BIGSdb Documentation

Release 1.14.0

Keith Jolley

Contents

1	1.1 1.2 1.3 1.4 1.5 1.6 1.7	BIGSdb Loci Alleles Schemes Profiles Classification groups Sequence tags Sets	3 3 3 4 4 4 4 5
2	BIGS 2.1	Sdb dependencies Required packages	7 7 7
		2.1.2 Perl modules	7 8
3	3.1 3.2 3.3 3.4 3.5 3.6 3.7 3.8 3.9 3.10 3.11	Site-specific configuration Setting up the offline job manager Setting up the submission system Periodically delete temporary files Prevent preference database getting too large Purging old jobs from the jobs database Log file rotation Upgrading BIGSdb Running the BIGSdb RESTful interface	9 10 11 12 13 14 14 14 15
4	Data 4.1 4.2 4.3	Creating databases Database-specific configuration XML configuration attributes used in config.xml 4.3.1 Isolate database XML attributes Special values 4.3.2 Sequence definition database XML attributes	19 20 20 20 25 26 28

		4.4.1 Apache authentication	28
		4.4.2 Built-in authentication	28
	4.5	Setting up the admin user	29
	4.6	Updating PubMed citations	29
5	Admi	inistrator's guide	31
	5.1		31
	5.2		31
	5.3		31
	5.4	Locus and scheme permissions (sequence definition database)	34
	5.5		36
			36
	5.6		36
	5.7		37
	5.8		38
	5.9		38
			38
			39
	5.10		39
	5.11	11 6	39
	5.11		39
		– 1	39
			41
	5.12		41
	3.12		41
			42
		e e e e e e e e e e e e e e e e e e e	43
			45
			45 45
	5.13		45 45
	3.13	E	45 45
		1	
			45
		E	48
			49
		e	49
	~ 1.4		55 55
	5.14	\mathcal{E}	57
	5.15		60
		1	60
	5.16		64
	5.16	6 6	67 - 2
	5.17		70
			72
			74
			75
	5.18	<u>i</u> <u>i</u>	76
		1	76
			78
			78
			78
	5.19	C 1	78
	5.20		79
	5.21		80
		5.21.1 Defining classification scheme in sequence definition database	80

		5.21.2 Defining classification scheme in isolate database
	5 22	Defining new loci based on annotated reference genome
		Genome filtering
	3.23	5.23.1 Filtering by <i>in silico</i> PCR
		5.23.2 Filtering by <i>in silico</i> hybridization
	5.24	Setting locus genome positions
	5.25	Defining composite fields
	5.26	Extended provenance attributes (lookup tables)
	5.27	Sequence bin attributes
	5.28	Checking external database configuration settings
	5.29	Exporting table configurations
	5.30	Authorizing third-party client software to access authenticated resources
	5.50	Tradionizing time party effect software to access authenticated resources
6	Cura	tor's guide
	6.1	Adding new sender details
	6.2	Adding new allele sequence definitions
		6.2.1 Single allele
		6.2.2 Batch adding multiple alleles
		Upload using a spreadsheet
		Upload using a FASTA file
	6.3	Updating and deleting allele sequence definitions
	6.4	Retiring allele identifiers
	6.5	Updating locus descriptions
	6.6	Adding new scheme profile definitions
	6.7	Updating and deleting scheme profile definitions
	6.8	Adding isolate records
	6.9	Updating and deleting single isolate records
	6.10	Batch updating multiple isolate records
	6.11	Deleting multiple isolate records
	6.12	Linking isolate records to publications
	6.13	Uploading sequence contigs linked to isolate records
		6.13.1 Select isolate from drop-down list
		6.13.2 Select from isolate query
		6.13.3 Upload options
	6.14	Automated web-based sequence tagging
	6.15	Projects
		6.15.1 Creating the project
		6.15.2 Explicitly adding isolates to a project
	6.16	Isolate record versioning
		6
7	Cura	ting submitted data 163
	7.1	Alleles
		7.1.1 Individual allele curation
		7.1.2 Batch allele curation
		7.1.3 Rejecting sequences
		7.1.4 Requesting additional information
		7.1.5 Closing the submission
	7.2	Profiles
		7.2.1 Individual profile curation
		7.2.2 Batch profile curation
		7.2.3 Rejecting profiles
		7.2.4 Requesting additional information
		7.2.5 Closing the submission

	7.3	Isolates177.3.1 Requesting additional information187.3.2 Closing the submission18	30
8	Offlir	e curation tools	33
_	8.1	Automated offline sequence tagging	
	8.2	Defining exemplar alleles	
	8.3	Automated offline allele definition	
	8.4	Cleanly interrupting offline curation	
	8.5	Uploading contigs from the command line	
9		tion downloads 19	
	9.1	Allele sequence definitions	
		9.1.1 Scheme tree	
		9.1.2 Alphabetical list	
		9.1.3 All loci by scheme	
	0.0	9.1.4 Download locus table	
	9.2	Scheme profile definitions	15
10	Data	ecords 19)7
		Solate records	
		10.1.1 Projects	
		10.1.2 Provenance metadata	
		10.1.3 Publications)()
		10.1.4 Sequence bin summary)1
		10.1.5 Scheme and locus data)1
	10.2	Allele definition records)1
	10.3	Sequence tag records)2
	10.4	Profile records)3
	10.5	Sequence bin records)4
11	0	*** J.4**	
11	11.1	ing data Querying sequences to determine allele identity	
	11.1	11.1.1 Querying whole genome data	
	11.2	Searching for specific allele definitions	
	11,2	11.2.1 General (all loci) sequence attribute search	
		11.2.2 Locus-specific sequence attribute search	
	11 3	Browsing scheme profile definitions	
	11.4	Querying scheme profile definitions	
	11.5	Investigating allele differences	
		11.5.1 Sequence similarity	
		11.5.2 Sequence comparison	22
	11.6	Browsing isolate data	24
	11.7	Querying isolate data	26
		11.7.1 Query by allele designation/scheme field	29
		11.7.2 Query by allele designation count	30
		11.7.3 Query by allele designation status	32
		11.7.4 Query by sequence tag count	33
		11.7.5 Query by tagged sequence status	34
		11.7.6 Query by list of attributes	36
		11.7.7 Query filters	
	11.8	Querying by allelic profile	38
	110	Retrieving isolates by linked publication	1
	11.9		
		User-configurable options	

		1.10.2 Main results table	16
		1.10.3 Isolate record display	
		1.10.4 Provenance field display	19
		1.10.5 Query filters	50
		1.10.6 Modifying locus and scheme display options	52
10	D (
12		nalysis plugins 25	
	12.1	cocus explorer	
		2.1.1 Polymorphic site analysis	
		2.1.2 Codon usage	
	12.2	2.1.3 Aligned translations	
		wo field breakdown	
		cheme and allele breakdown	
		equence bin breakdown	
		Genome comparator	
	12.0	2.6.1 Analysis using defined loci	
		2.6.2 Analysis using annotated reference genome	
		2.6.3 Include in identifiers fieldset	
		2.6.4 Reference genome fieldset	
		2.6.5 Parameters/options fieldset	
		2.6.6 Distance matrix calculation fieldset	
		2.6.7 Alignments fieldset	
		2.6.8 Core genome analysis fieldset	
		2.6.9 Filter fieldset	
		2.6.10 Understanding the output	
		Distance matrix	
		Unique strains	
	12.7	BLAST	
		2.7.1 Include in results table fieldset	37
		2.7.2 Parameters fieldset	38
		2.7.3 No matches	
		2.7.4 Filter fieldset	
		BURST	
		Codon usage	
		Unique combinations	
		olymorphisms	
	12.12	resence/absence	
		2.12.1 Options	
	12.13	ag status)3
13	Data	aport plugins 30)7
		solate record export	
		3.1.1 Advanced options)9
		3.1.2 Molecular weight calculation	
	13.2	equence export	10
		3.2.1 Aligning sequences	12
	13.3	Contig export	12
		3.3.1 Filtering by tagged status of contigs	15
14	Subm	ting data using the submission system 31	17
17		Registering a user account	
		Allele submission	
		4.2.1 Start 31	

		14.2.2 14.2.3	Select the submission locus	
		14.2.4	Paste in sequence(s)	
		14.2.5	Add message to curator	
		14.2.6	Add supporting files	
	4.4.0	14.2.7	Finalize submission	
	14.3		submission	
		14.3.1	Start	
		14.3.2	Paste in profile(s)	
		14.3.3	Add message to curator	
		14.3.4	Add supporting files	
		14.3.5	Finalize submission	
	14.4		submission	
		14.4.1	Start	
		14.4.2	Paste in isolate data	335
		14.4.3	Add message to curator	337
		14.4.4	Add supporting files	338
		14.4.5	Finalize submission	340
	14.5	Genome	e submission	342
	14.6	Removi	ng submissions from your notification list	342
15	REST		lication Programming Interface (API)	345
	15.1		additional/optional parameters	
	15.2		es	345
		15.2.1	GET / or /db - List site resources	
		15.2.2	GET /db/{database} - List database resources	347
		15.2.3	GET /db/{database}/loci - List loci	347
		15.2.4	GET /db/{database}/loci/{locus} - Retrieve locus record	348
		15.2.5	GET /db/{database}/loci/{locus}/alleles - Retrieve list of alleles defined for a locus	349
		15.2.6	GET /db/{database}/loci/{locus}/alleles_fasta - Download alleles in FASTA format	349
		15.2.7	GET /db/{database}/loci/{locus}/alleles/{allele_id} - Retrieve full allele information	350
		15.2.8	GET /db/{database}/schemes - List schemes	350
		15.2.9	GET /db/{database}/schemes/{scheme_id} - Retrieve scheme information	350
		15.2.10	GET /db/{database}/schemes/{scheme_id}/fields/{field} - Retrieve information about	
			scheme field	351
		15.2.11	GET /db/{database}/schemes/{scheme_id}/profiles - List allelic profiles defined for scheme	351
			GET /db/{database}/schemes/{scheme id}/profiles csv - Download allelic profiles in CSV	
			(tab-delimited) format	352
		15.2.13	GET /db/{database}/schemes/{scheme_id}/profiles/{profile_id} - Retrieve allelic profile	
				352
		15.2.14		353
			GET /db/{database}/isolates/{isolate_id} - Retrieve isolate record	353
			GET /db/{database}/isolates/{isolate_id}/allele_designations - Retrieve list of allele desig-	
			nation records	354
		15.2.17	GET /db/{database}/isolates/{isolate_id}/allele_designations/{locus} - Retrieve full allele	
			designation record	355
		15 2 18	GET /db/{database}/isolates/{isolate_id}/allele_ids - Retrieve allele identifiers	355
			GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_designations - Re-	555
		10.2.17	trieve scheme allele designation records	356
		15 2 20	GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_ids - Retrieve list of	550
		13.2.20	scheme allele identifiers	356
		15 2 21	GET /db/{database}/isolate_id}/contigs - Retrieve list of contigs	357
			GET /db/{database}/isolates/{isolate_id}/contigs_fasta - Download contigs in FASTA format	
			GET /db/{database}/isolates/{isolate_id}/contigs_fasta - Download contigs in FASTA format GET /db/{database}/contigs/{contig_id} - Retrieve contig record	
		13.4.43	OLI Tau (autouse frenings) (coning in freniere coning record	550

18	Datal	base sch	ema	379
1,	17.1 17.2	Query of Sequence	perators	375 375
17	Appe	ndix		375
	16.3	Admini	stration	372
			ion	
16			ked questions (FAQs)	371
		15.3.6	Accessing protected resources	
		15.3.5	Getting a session token	
		15.3.4	Getting an access token	
		15.3.3	Getting user authorization	
		15.3.2	Getting a request token	
	10.0	15.3.1	Developer sign up to get a consumer key	
	15.3	Authent	ication	
		15.2.38	DELETE /db/{database}/submissions/{submission_id}/files/{filename} - Delete submission supporting file	365
		15.2.37	GET /db/{database}/submissions/{submission_id}/files/{filename} - Download submission supporting file	364
			POST /db/{database}/submissions/{submission_id}/files - Upload submission supporting file	364
		15000	loaded for submission	364
		15.2.35	GET /db/{database}/submissions/{submission_id}/files - Retrieve list of supporting files up-	264
			dence	364
		15.2.34	POST /db/{database}/submissions/{submission_id}/messages - Add submission correspon-	
			spondence	363
			GET /db/{database}/submissions/{submission_id}/messages - Retrieve submission corre-	200
			DELETE /db/{database}/submissions/{submission_id} - Delete submission record	363
			GET /db/{database}/submissions/{submission_id} - Retrieve submission record	362
			POST /db/{database}/submissions - create new submission	361
		15 2 20	GET /db/{database}/submissions - retrieve list of submissions	360
		13.2.20	project	360
			GET /db/{database}/projects/{project_id} - Retrieve project information	359
			GET /db/{database}/projects - Retrieve list of projects	
			GET /db/{database}/users/{user_id} - Retrieve user information	
		15.2.24	GET /db/{database}/fields - Retrieve list of isolate provenance field descriptions	

Gene-by-gene population annotation and analysis

BIGSdb is software designed to store and analyse sequence data for bacterial isolates. Any number of sequences can be linked to isolate records - these can be small contigs assembled from dideoxy sequencing through to whole genomes (complete or multiple contigs generated from parallel sequencing technologies such as 454 or Illumina).

BIGSdb extends the principle of MLST to genomic data, where large numbers of loci can be defined, with alleles assigned by reference to sequence definition databases (which can also be set up with BIGSdb). Loci can also be grouped into schemes so that types can be defined by combinations of allelic profiles, a concept analogous to MLST.

The software has been released under the GNU General Public Licence version 3. The latest version of this documentation can be found at http://bigsdb.readthedocs.org/.

Contents 1

2 Contents

Concepts and terms

1.1 BIGSdb

BIGSdb is the software platform - not a specific database. There are many instances of BIGSdb databases, so referring to 'the BIGSdb' is meaningless.

1.2 Loci

Loci are regions of the genome that are identified by similarity to a known sequence. They can be defined by DNA or peptide sequence. They are often complete coding sequences (genes), but may represent gene fragments (such as used in MLST), antigenic peptide loops, or indeed any sequence feature.

In versions of BIGSdb prior to 1.8.0, an isolate record could only have one live *allele* designation for a locus (inactive/pending designations could be stored within the database but were unavailable for querying or analysis purposes). Since biology is rarely so clean, and some genomes may contain more than one copy of a gene, later versions of the software allow multiple allele designations for a locus, all of which can be queried and analysed.

Paralogous loci can be difficult to differentiate by sequence similarity alone. Because of this, loci can be further defined by context, where in silico PCR or hybridization reactions can be performed to *filter the genome* to specific regions based on sequence external to the locus.

1.3 Alleles

Alleles are instances of loci. Every unique sequence, either DNA or peptide depending on the locus, is defined as a new allele and these are defined in a sequence definition database, where they are given an allele identifier. These identifiers are usually integers, but can be text strings. Allele identifiers in text format can be constrained by length and formatting.

When a specific allele of a locus is identified within the sequence data of an isolate record, the allele designation, i.e. identifier, is associated with the isolate record. This efficiently stores the sequence variation found within an isolate. Two isolates with the same allele designation for a locus have identical sequences at that locus. Once the sequence variation within a genome has been reduced to a series of allele designations, genomes can be efficiently compared by identifying which loci vary between them.

It is important to note that allele identifiers are usually arbitrary and are allocated sequentially in the order of discovery. Alleles with adjacent identifiers may vary by a single nucleotide or by many.

1.4 Schemes

Schemes are collections of loci that may be associated with additional field values. At their simplest they just group loci together. Example uses of simple schemes include:

- Antibiotic resistance genes
- · Genes involved in specific biochemical pathways
- Antigens
- Vaccine components
- Whole genome MLST (wgMLST)

When schemes are associated with additional fields, one of these fields must be the primary key, i.e. its value uniquely defines a particular combination of alleles at its member loci. The pre-eminent example of this is MLST - where a sequence type (ST) is the primary key that uniquely defines combinations of alleles that make up the MLST profiles. Additional fields can also then be included. The values for these need not be unique. In the MLST example, a field for clonal complex can be included, and the same value for this can be set for multiple STs.

1.5 Profiles

Profiles are instances of *schemes*. A profile consists of a set of *allele identifiers* for the *loci* that comprise the scheme. If the scheme has a primary key field, e.g. sequence type (ST) in MLST schemes, then the unique combination of alleles in a complete profile can be defined by the value of this field.

1.6 Classification groups

Classification groups are a way to cluster scheme profiles using a specified threshold of pairwise allelic mismatches. Currently, single-linkage clustering is supported whereby each member of a group must have no more than the specified number of allelic differences with at least one other member of the group.

1.7 Sequence tags

Sequence tags record locus position within an isolate record's sequence bin. The process of creating these tags, is known as *tag-scanning*. A sequence tag consists of:

- sequence bin id this identifies a particular contig
- · locus name
- · start position
- · end position
- flag to indicate if sequence is reversed
- flag to indicate if sequence is complete and does not continue off the end of the contig

1.8 Sets

Sets provide a means to take a large database with multiple loci and/or schemes and present a subset of these as though it was a complete database. The loci and schemes chosen to belong to a set can be renamed when used with this set. The rationale for this is that in a database with disparate isolates and a large number of loci, the naming of these loci may have to be long to specify a species name. For example, you may have a database that contains multiple MLST schemes for different species, but since these schemes may use different fragments of the same genes they may have to be named something like 'Streptococcus_pneumoniae_MLST_aroE' to uniquely specify them. If we define a set for 'Streptococcus pneumoniae' we can then choose to only include S. pneumoniae loci and therefore shorten their names, e.g. to 'aroE'.

Additional metadata fields can also be associated with each set so it is possible to have a database containing genomes from multiple species and a generic set of metadata, then have additional specific metadata fields for particular species or genera. These additional fields only become visible and searchable when the specific set that they belong to has been selected.

1.8. Sets 5

BIGSdb dependencies

2.1 Required packages

BIGSdb requires a number of software components to be installed:

2.1.1 Linux packages

- Apache2 web server with mod_perl2
- PostgreSQL database
- Perl 5.10+
- BioPerl
- BLAST+
- EMBOSS
 - infoalign use to extract alignment stats in Genome Comparator.
 - sixpack used to translate sequences in multiple reading frames.
 - stretcher used for sequence alignment in allele query.
- Ipcress part of exonerate package used to simulate PCR reactions which can be used to filter the genome to predicted amplification products.
- Xvfb X virtual framebuffer needed to support SplitsTree in command line mode as used in Genome Comparator.

2.1.2 Perl modules

These are included with most Linux distributions.

- DBI Database independent interface module used to interact with databases.
- DBD-Pg PostgreSQL database driver for DBI.
- XML::Parser::perlSAX part of libxml-perl Used to parse XML configuration files.
- Log::Log4perl Configurable status and error logging.
- Log::Dispatch::File Object for logging to file.
- Error Exception handling.

- Config::Tiny Configuration file handling.
- Bio::Biblio This used to be part of BioPerl but will need to be installed separately if using BioPerl 1.6.920 or later.
- IO::String
- Data::UUID Globally unique identifer handling for preference storage.
- List::MoreUtils (version 0.28+).
- Time::Duration [optional] Used by Job Viewer to display elapsed time in rounded units.
- Excel::Writer::XLSX Used to export data in Excel format.
- Parallel::ForkManager Required for multi-threading autotagger and autodefiner scripts.
- Net::Oauth Required for REST authentication (this needs to be installed even if you are not using REST).
- Crypt::Eksblowfish::Bcrypt Used for password hashing.
- Mail::Sender [optional] Used to send E-mail messages by submission system.
- Email::Valid [optional] Used to validate E-mails sent by job manager.

2.1.3 Optional packages

Installing these packages will enable extra functionality, but they are not required by the core BIGSdb package.

- ChartDirector library used for generating charts. Used by some plugins.
- ImageMagick mogrify used by some plugins.
- MAFFT 6.8+ sequence alignment used by some plugins.
- Muscle sequence alignment used by some plugins.
- Splitstree4 used by GenomeComparator plugin.

Installation and configuration of BIGSdb

3.1 Software installation

BIGSdb consists of two main Perl scripts, bigsdb.pl and bigscurate.pl, that run the query and curator's interfaces respectively. These need to be located somewhere within the web cgi-bin directories. In addition, there are a large number of library files, used by both these scripts, that are installed by default in /usr/local/lib/BIGSdb. Plugin scripts are stored within a 'Plugins' sub-directory of this library directory.

All databases on a system can use the same instance of the scripts, or alternatively any database can specify a particular path for each script, enabling these script directories to be protected by apache htaccess directives.

- Software requirements
- Download from SourceForge.net or GitHub.
- 1. Unpack the distribution package in a temporary directory:

```
gunzip bigsdb_1.x.x.tar.gz
tar xvf bigsdb_1.x.x.tar
```

- 2. Copy the bigsdb.pl and bigscurate.pl scripts to a subdirectory of your web server's cgi-bin directory. Make sure these are readable and executable by the web server daemon.
- 3. Copy the contents of the lib directory to /usr/local/lib/BIGSdb/. Make sure you include the Plugins and Offline directories which are subdirectories of the main lib directory.
- 4. Copy the contents of the javascript directory to a javascript directory within the web root tree, i.e. accessible from http://your_website/javascript/.
- 5. Copy the contents of the css directory to a css directory within the web root tree, i.e. accessible from http://your_website/css/.
- 6. Copy the contents of the fonts directory to a fonts directory within the web root tree, i.e. accessible from http://your_website/fonts/.
- 7. Copy the images directory to the root directory of your website, i.e. accessible from http://your_website/images/.
- 8. Copy the contents of the conf directory to /etc/bigsdb/. Check the paths of helper applications and database names in the bigsdb.conf file and modify for your system.
- 9. Create a PostgreSQL database user called apache this should not have any special priveleges. First you will need to log in as the postgres user:

```
sudo su postgres
```

Then use the createuser command to do this, e.g.

```
createuser apache
```

From the psql command line, set the apache user password:

```
psql
ALTER ROLE apache WITH PASSWORD 'remote';
```

10. Create PostgreSQL databases called bigsdb_auth, bigsdb_prefs and bigsdb_refs using the scripts in the sql directory. Create the database using the createdb command and set up the tables using the psql command.

```
createdb bigsdb_auth
psql -f auth.sql bigsdb_auth
createdb bigsdb_prefs
psql -f prefs.sql bigsdb_prefs
createdb bigsdb_refs
psql -f refs.sql bigsdb_refs
```

- 11. Create a writable temporary directory in the root of the web site called tmp, i.e. accessible from http://your_website/tmp.
- 12. Create a log file, bigsdb.log, in /var/log owned by the web server daemon, e.g.

```
touch /var/log/bigsdb.log
chown www-data /var/log/bigsdb.log
```

(substitute www-data for the web daemon user).

3.2 Configuring PostgreSQL

PostgreSQL can be configured in many ways and how you do this will depend on your site requirements.

The following security settings will allow the appropriate users 'apache' and 'bigsdb' to access databases without allowing all logged in users full access. Only the UNIX users 'postgres' and 'webmaster' can log in to the databases as the Postgres user 'postgres'.

You will need to edit the pg_hba.conf and pg_ident.conf files. These are found somewhere like /etc/postgresql/9.1/main/

pg_hba.conf

```
# Database administrative login by UNIX sockets
       all
                  postgres
                                                    ident map=mymap
# TYPE DATABASE USER
                              CIDR-ADDRESS
                                                    METHOD
# "local" is for Unix domain socket connections only
local all all
                                                    ident map=mymap
# IPv4 local connections:
                              127.0.0.1/32
                                                    md5
       all
# IPv6 local connections:
                   all
                                                   md5
host
       all
                              ::1/128
```

pg_ident.conf

# MAPNAME	SYSTEM-USERNAME	PG-USERNAME
mymap	postgres	postgres
mymap	webmaster	postgres
mymap	www-data	apache

mymap	bigsdb	bigsdb	
mymap	bigsdb	apache	

You may also need to change some settings in the postgresql.conf file. As an example, a configuration for a machine with 16GB RAM, allowing connections from a separate web server may have the following configuration changes made:

```
listen_addresses = '*'
max_connections = 200
shared_buffers = 1024Mb
work_mem = 8Mb
effective_cache_size = 8192Mb
stats_temp_directory = '/dev/shm'
```

Setting stats_temp_directory to /dev/shm makes use of a ramdisk usually available on Debian or Ubuntu systems for frequently updated working files. This reduces a lot of unnecessary disk access.

See Tuning Your PostgreSQL Server for more details.

Restart PostgreSQL after any changes, e.g.

```
/etc/init.d/postgresql restart
```

3.3 Site-specific configuration

Site-specific configuration files are located in /etc/bigsdb by default.

- bigsdb.conf main configuration file
- logging.conf error logging settings. See log4perl project website for advanced configuration details.

3.4 Setting up the offline job manager

To run plugins that require a long time to complete their analyses, an offline job manager has been developed. The plugin will save the parameters of a job to a job database and then provide a link to the job status page. An offline script, run frequently from CRON, will then process the job queue and update status and outputs via the job status page.

1. Create a 'bigsdb' UNIX user, e.g.:

```
sudo useradd -s /bin/sh bigsdb
```

2. As the postgres user, create a 'bigsdb' user and create a bigsdb_jobs database using the jobs.sql SQL file, e.g.:

```
createuser bigsdb [no need for special priveleges]
createdb bigsdb_jobs
psql -f jobs.sql bigsdb_jobs
```

From the psql command line, set the bigsdb user password::

```
psql
ALTER ROLE bigsdb WITH PASSWORD 'bigsdb';
```

3. Set up the jobs parameters in the /etc/bigsdb/bigsdb.conf file, e.g.:

```
jobs_db=bigsdb_jobs
max_load=8
```

The jobs script will not process a job if the server's load average (over the last minute) is higher than the max_load parameter. This should be set higher than the number of processor cores or you may find that jobs never run on a busy server. Setting it to double the number of cores is probably a good starting point.

- 4. Copy the job_logging.conf file to the /etc/bigsdb directory.
- 5. Set the script to run frequently (preferably every minute) from CRON. Note that CRON does not like '.' in executable filenames, so either rename the script to 'bigsjobs' or create a symlink and call that from CRON, e.g.:

```
copy bigsjobs.pl to /usr/local/bin
sudo ln -s /usr/local/bin/bigsjobs.pl /usr/local/bin/bigsjobs
```

You should install xvfb, which is a virtual X server that may be required for third party applications called from plugins. This is required, for example, for calling splitstree4 from the Genome Comparator plugin.

Add the following to /etc/crontab::

```
* * * * bigsdb xvfb-run -a /usr/local/bin/bigsjobs
```

(set to run every minute from the 'bigsdb' user account).

If you'd like to run this more frequently, e.g. every 30 seconds, multiple entries can be added to CRON with an appropriate sleep prior to running, e.g.:

```
* * * * bigsdb xvfb-run -a /usr/local/bin/bigsjobs
* * * * * bigsdb sleep 30;xvfb-run -a /usr/local/bin/bigsjobs
```

6. Create a log file, bigsdb_jobs.log, in /var/log owned by 'bigsdb', e.g.:

```
sudo touch /var/log/bigsdb_jobs.log
sudo chown bigsdb /var/log/bigsdb_jobs.log
```

3.5 Setting up the submission system

The submission system allows users to submit new data to the database for curation. Submissions are placed in a queue for a curator to upload. All communication between submitters and curators can occur via the submission system.

1. Create a writable submissions directory in the root of the web site called submissions, i.e. accessible from http://your_website/submissions. This is used for file uploads. The directory should be writable by the Apache web daemon (user 'www-data' on Debian/Ubuntu systems). If you are running the RESTful interface the directory should also be writable by the bigsdb user. To ensure this, make the directory group-writable and add the bigsdb user to the apache group ('www-data' on Debian/Ubuntu systems). If you will be allowing submissions via the RESTful interface, you should also add the apache user ('www-data' on Debian/Ubuntu systems) to the bigsdb group, e.g.

```
sudo usermod -a -G www-data bigsdb
sudo usermod -a -G bigsdb www-data
```

The actual directory can be outside of the web root and made accessible using a symlink provided your Apache configuration allows this, e.g. the default location is /var/submissions symlinked to /var/www/submissions (assuming your web site is located in /var/www), e.g.

```
sudo touch /var/submissions
sudo chown www-data:www-data /var/submissions
sudo chmod 775 /var/submissions
sudo ln -s /var/submissions /var/www
```

- 2. Set the submission_dir location in bigsdb.conf.
- 3. Set the smtp_server in bigsdb.conf to the IP or DNS name of your organisation's SMTP relay. Depending on how your E-mail system is configured, you may be able to use the localhost address (127.0.0.1).
- 4. Make sure the curate_script and query_script values are set in bigsdb.conf. These point to the web-accessible location of the web scripts and are required to allow curators to be directed between the web interfaces as needed.
- 5. Set submissions="yes" in the system tag of the *database config.xml file* of each database for which submissions should be enabled.

3.6 Periodically delete temporary files

There are two temporary directories (one public, one private) which may accumulate temporary files over time. Some of these are deleted automatically when no longer required but some cannot be cleaned automatically since they are used to display results after clicking a link or to pass the database query between pages of results.

The easiest way to clean the temp directories is to run a cleaning script periodically, e.g. create a root-executable script in /etc/cron.hourly containing the following::

```
#!/bin/sh
#Remove temp BIGSdb files from secure tmp folder older than 1 week.
find /var/tmp/ -name '*BIGSdb_*' -type f -mmin +10080 -exec rm -f {} \; 2>/dev/null

#Remove .jnlp files from web tree older than 1 day
find /var/www/tmp/ -name '*.jnlp' -type f -mmin +1440 -exec rm -f {} \; 2>/dev/null

#Remove other tmp files from web tree older than 1 week
find /var/www/tmp/ -type f -mmin +10080 -exec rm -f {} \; 2>/dev/null
```

3.7 Prevent preference database getting too large

The preferences database stores user preferences for BIGSdb databases running on the site. Every user will have a globally unique identifier (guid) stored in this database along with a datestamp indicating the last access time. On public databases that do not require logging in, this guid is stored as a cookie on the user's computer. Databases that require logging in use a combination of database and username as the identifier. Over time, the preferences database can get quite large since every unique user will result in an entry in the database. Since many of these entries represent casual users, or even web indexing bots, they can be periodically cleaned out based on their last access time. A weekly CRON job can be set up to remove any entries older than a defined period. For example, the following line entered in /etc/crontab will remove the preferences for any user that has not accessed any database in the past 6 months (the script will run at 6pm every Sunday).

```
#Prevent prefs database getting too large
00 18 * * 0 postgres psql -c "DELETE FROM guid WHERE last_accessed < NOW() - INTERVAL '6 month
```

3.8 Purging old jobs from the jobs database

If you are running the offline job manager, the jobs database (default bigsdb_jobs) contains the parameters and output messages of these jobs. Job output files are only *usually kept on the server for 7 days* so there is no point keeping the database entries for longer than this. These can be purged with a daily cron job, e.g. set the following in /etc/crontab (the script will run at 5am every day).

```
#Purge jobs older than 7 days from the jobs database.
00 5 * * * postgres psql -c "DELETE FROM jobs where (stop_time IS NOT NULL AND stop_time < now
```

3.9 Log file rotation

Set the log file to auto rotate by adding a file called 'bigsdb' with the following contents to /etc/logrotate.d:

```
/var/log/bigsdb.log {
 weekly
 rotate 4
 compress
 copytruncate
 missingok
 notifempty
 create 640 root adm
/var/log/bigsdb_jobs.log {
 weekly
 rotate 4
 compress
 copytruncate
 missingok
 notifempty
 create 640 root adm
```

3.10 Upgrading BIGSdb

Major version changes, e.g. 1.7 -> 1.8, indicate that there has been a change to the underlying database structure for one or more of the database types. Scripts to upgrade the database are provided in sql/upgrade and are named by the database type and version number. For example, to upgrade an isolate database (bigsdb_isolates) from version 1.7 to 1.8, log in as the postgres user and type:

```
psql -f isolatedb_v1.8.sql bigsdb_isolates
```

Upgrades are sequential, so to upgrade from a version earlier than the last major version you would need to upgrade to the intermediate version first, e.g. to go from 1.6 -> 1.8, requires upgrading to 1.7 first.

Minor version changes, e.g. 1.8.0 -> 1.8.1, have no modifications to the database structures. There will be changes to the Perl library modules and possibly to the contents of the Javascript directory, images directory and CSS files. The version number is stored with the bigsdb.pl script, so this should also be updated so that BIGSdb correctly reports its version.

3.11 Running the BIGSdb RESTful interface

BIGSdb has an Application Programming Interface (API) that allows third-party applications to access the data within the databases. The script that runs this is called bigsrest.pl. This is a Dancer2 application that can be run using a wide range of options, e.g. as a stand-alone script, using Perl webservers with plackup, or from apache. Full documentation for deploying Dancer2 applications can be found online.

The script requires a new database that describes the resources to make available. This is specified in the bigsdb.conf file as the value of the 'rest_db' attribute. By default, the database is named bigsdb_rest.

A SQL file to create this database can be found in the sql directory of the download archive. It is called rest.sql. To create the database, as the postgres user, navigate to the sql directory and type

```
createdb bigsdb_rest
psql -f rest.sql bigsdb_rest
```

This database will need to be populated using psql or any tool that can be used to edit PostgreSQL databases. The database contains three tables that together describe and group the databases resources that will be made available through the API. The tables are:

- resources
 - this contains two fields (both compulsory):
 - * **dbase_config** the name of the database configuration used with the database. This is the same as the name of the directory that contains the config.xml file in the /etc/bigsdb/dbases directory.
 - * **description** short description of the database.
- groups (used to group related resources together)
 - this contains two fields (compulsory fields shown in bold):
 - * **name** short name of group. This is usually a single word and is also the key that links resources to groups.
 - * **description** short description of group.
 - * long_description fuller description of group.
- group_resources (used to add resources to groups)
 - this contains two fields (both compulsory)
 - * group_name name of group. This must already exist in the groups table.
 - * dbase_config the name of database resource. This must already exist in the resources table.

For example, to describe the PubMLST resources for Neisseria, connect to the bigsdb_rest database using psql,

```
psql bigsdb_rest
```

Then enter the following SQL commands. First add the database resources:

```
INSERT INTO resources (dbase_config,description) VALUES
('pubmlst_neisseria_seqdef','Neisseria sequence/profile definitions');
INSERT INTO resources (dbase_config,description) VALUES
('pubmlst_neisseria_isolates','Neisseria isolates');
```

Then create a 'neisseria' group that will contain these resources:

```
INSERT INTO groups (name, description) VALUES
('neisseria','Neisseria spp.');
```

Finally, add the database resources to the group:

```
INSERT INTO group_resources (group_name,dbase_config) VALUES
('neisseria','pubmlst_neisseria_seqdef');
INSERT INTO group_resources (group_name,dbase_config) VALUES
('neisseria','pubmlst_neisseria_isolates');
```

The REST API will need to run on its own network port. By default this is port 3000. To run as a stand-alone script, from the script directory, as the bigsdb user, simply type:

```
./bigsrest.pl
```

This will start the API on port 3000. You will be able to check that this is running using a web browser by navigating to http://localhost:3000 on the local machine, or using the server IP address from a remote machine. You may need to modify your server firewall rules to allow connection to this port.

Running as a stand-alone script is useful for testing, but you can achieve much better performance using a Perl webserver with plackup. There are various options to choose. PubMLST uses Starman.

To run the API using Starman, type the following as the bigsdb user:

description "Start BIGSdb REST interface"

```
plackup -a /var/rest/bigsrest.pl -s Starman -E deployment
```

where the value of -a refers to the location of the bigsrest.pl script. Starman defaults to using port 5000.

Different Linux distributions use different means to control services/daemons. To start the REST interface on system boot on systems using upstart, create a file called bigsdb-rest.conf in /etc/init. The contents of this file should be something like (modify file paths as appropriate):

```
version "1.0"
author "Keith Jolley"
start on runlevel [12345]
## tell upstart we're creating a daemon
expect fork
script
exec su -s /bin/sh -c 'exec "$0" "$@"' bigsdb -- /usr/local/bin/plackup -a /var/rest/bigsrest.pl -s send script
```

3.11.1 Proxying the API to use a standard web port

Usually you will want your API to be available on the standard web port 80. To do this you will need to set up a virtual host using a different domain name from your web site to proxy the API port. For example, PubMLST has a separate domain 'http://rest.pubmlst.org' for its API. This is set up as a virtual host directive in apache with the following configuration file:

```
<VirtualHost *>
   ServerName rest.pubmlst.org
   DocumentRoot /var/rest
   ServerAdmin keith.jolley@zoo.ox.ac.uk
   <Directory /var/rest>
    AllowOverride None
   Require all granted
   </Directory>
```

```
ProxyPass / http://rest.pubmlst.org:5000/
ProxyPassReverse / http://rest.pubmlst.org:5000/

<Proxy *>
    Order allow,deny
    Allow from all
    </Proxy>

ErrorLog /var/log/apache2/rest.pubmlst.org-error.log
CustomLog /var/log/apache2/rest.pubmlst.org-access.log common

</VirtualHost>
```

Database setup

There are two types of BIGSdb database:

- sequence definition databases, containing
 - allele sequences and their identifiers
 - scheme data, e.g. MLST profile definitions
- · isolate databases, containing
 - isolate provenance metadata
 - genome sequences
 - allele designations for loci defined in sequence definition databases.

These two databases are independent but linked. A single isolate database can communicate with multiple sequence definition databases and vice versa. Different access restrictions can be placed on different databases.

Databases are described in XML files telling BIGSdb everything it needs to know about them. Isolate databases can have any fields defined for the isolate table, allowing customisation of metadata - these fields are described in the XML file (config.xml) and must match the fields defined in the database itself.

4.1 Creating databases

There are templates available for the sequence definition and isolate databases. These are SQL scripts found in the sql directory.

To create a database, you will need to log in as the postgres user and use these templates. For example to create a new sequence definition database called bigsdb_test_seqdef, navigate to the sql directory and log in as the postgres user, e.g.

```
sudo su postgres
```

then

```
createdb bigsdb_test_seqdef
psql -f seqdef.sql bigsdb_test_seqdef
```

Create an isolate database the same way:

```
createdb bigsdb_test_isolates
psql -f isolatedb.sql bigsdb_test_isolates
```

The standard fields in the isolate table are limited to essential fields required by the system. To add new fields, you need to log in to the database and alter this table. For example, to add fields for country and year, first log in to the newly created isolate database as the postgres user:

```
psql bigsdb_test_isolates
```

and alter the isolate table:

```
ALTER TABLE isolates ADD country text;
ALTER TABLE isolates ADD year int;
```

Remember that any fields added to the table need to be described in the config.xml file for this database.

4.2 Database-specific configuration

Each BIGSdb database on a system has its own configuration directory, by default in /etc/bigsdb/dbases. The database has a short configuration name used to specify it in a web query and this matches the name of the configuration sub-directory, e.g. http://pubmlst.org/cgi-bin/bigsdb/bigsdb.pl?db=pubmlst_neisseria_isolates is the URL of the front page of the PubMLST Neisseria isolate database whose configuration settings are stored in /etc/bigsdb/dbases/pubmlst_neisseria_isolates. This database sub-directory contains a number of files (hyperlinks lead to the files used on the Neisseria database):

- config.xml the database configuration file. Fields defined here correspond to fields in the isolate table of the database.
- banner.html optional file containing text that will appear as a banner within the database index pages. HTML markup can be used within this text.
- header.html HTML markup that is inserted at the top of all pages. This can be used to set up site-specific menubars and logos.
- footer.html HTML markup that is inserted at the bottom of all pages.
- curate_header.html HTML markup that is inserted at the top of all curator's interface pages.
- curate_footer.html HTML markup that is inserted at the bottom of all curator's interface pages.

4.3 XML configuration attributes used in config.xml

The following lists describes the attributes used in the config.xml file that is used to describe databases.

4.3.1 Isolate database XML attributes

Please note that database structure described by the field and sample elements must match the physical structure of the database isolate and sample tables respectively. Required attributes are in **bold**:

```
<db>
```

Top level element. Contains child elements: system, field and sample.:

<system>

· authentication

- Method of authentication: either 'builtin' or 'apache'. See *user authentication*.

• db

- Name of database on system.

dbtype

- Type of database: either 'isolates' or 'sequences'.

· description

- Description of database used throughout interface.
- align_limit
 - Overrides the sequence export record alignment limit in the Sequence Export plugin. Default: '200'.
- all_plugins
 - Enable all appropriate plugins for database: either 'yes' or 'no', default 'no'.
- · annotation
 - Semi-colon separated list of accession numbers with descriptions (separated by a l), eg. 'AL157959|Z2491;AM421808|FAM18;NC_002946|FA 1090;NC_011035|NCCP11945;NC_014752|020-06'. Currently used only by Genome Comparator plugin.
- · cache_schemes
 - Enable automatic refreshing of scheme field caches when batch adding new isolates: either 'yes' or 'no', default 'no'.
 - See scheme caching.
- · codon_usage_limit
 - Overrides the record limit for the Codon Usage plugin. Default: '500'.
- contig_analysis_limit
 - Overrides the isolate number limit for the Contig Export plugin. Default: '1000'.
- curate_path_includes
 - Partial path of the bigscurate.pl script used to curate the database. See user authentication.
- curate_script
 - Relative web path to curation script. Default 'bigscurate.pl' (version 1.11+).
 - This is only needed if automated submissions are enabled. If bigscurate.pl is in a different directory from bigsdb.pl, you need to include the whole web path, e.g. /cgi-bin/private/bigsdb/bigscurate.pl.
- · curators only
 - Set to 'yes' to prevent ordinary authenticated users having access to database configuration. This is only
 effective if read_access is set to 'authenticated_users'. This may be useful if you have different configurations for curation and querying with some data hidden in the configuration used by standard users. Default
 'no'.
- daily_rest_submissions_limit
 - Overrides the limit on number of submissions that can be made to the database via the RESTful interface. This is useful to prevent flooding of the submission system by aberrant scripts. Default: '100'.
- · default_access

- The default access to the database configuration, either 'allow' or 'deny'. If 'allow', then specific users can be denied access by creating a file called 'users.deny' containing usernames (one per line) in the configuration directory. If 'deny' then specific users can be allowed by creating a file called 'users.allow' containing usernames (one per line) in the configuration directory. See *default access*.
- default_seqdef_config
 - Isolate databases only: Name of the default sequel database configuration used with this database. Used
 to automatically fill in details when adding new loci.
- default_seqdef_dbase
 - Isolate databases only: Name of the default seqdef database used with this database. Used to automatically fill in details when adding new loci.
- default_seqdef_script
 - Isolate databases only: URL of BIGSdb script running the seqdef database (default: '/cgi-bin/bigsdb/bigsdb.pl').
- export_limit
 - Overrides the default allowed number of data points (isolates x columns) to export. Default: '25000000'.
- fieldgroup1 fieldgroup10
 - Allows multiple fields to be queried as a group. Value should be the name of the group followed by a colon
 (:) followed by a comma-separated list of fields to group, e.g. identifiers:id,strain,other_name.
- genome_comparator_limit
 - Overrides the isolate number limit for the Genome Comparator plugin. Default: '1000'.
- genome_comparator_max_ref_loci
 - Overrides the limit on number of loci allowed in a reference genome. Default: '10000'.
- · hide unused schemes
 - Sets whether a scheme is shown in a main results table if none of the isolates on that page have any data for the specific scheme: either 'yes' or 'no', default 'no'.
- host
 - Host name/IP address of machine hosting isolate database, default 'localhost'.
- job_priority
 - Integer with default job priority for offline jobs (default:5).
- job quota
 - Integer with number of offline jobs that can be queued or currently running for this database.
- labelfield
 - Field that is used to describe record in isolate info page, default 'isolate'.
- · locus_aliases
 - Display locus aliases and use them in dropdown lists by default: must be either 'yes' or 'no', default 'no'.
 This option can be overridden by a user preference.
- locus_superscript_prefix
 - Superscript the first letter of a locus name if it is immediately following by an underscore, e.g. f_abcZ would be displayed as fabcZ within the interface: must be either 'yes' or 'no', default 'no'. This can be used to designate gene fragments (or any other meaning you like).

- · maindisplay_aliases
 - Default setting for whether isolates aliases are displayed in main results tables: either 'yes' or 'no', default 'no'. This setting can be overridden by individual user preferences.
- · noshow
 - Comma-separated list of fields not to use in breakdown statistic plugins.
- no_publication_filter
 - Isolate databases only: Switches off display of publication filter in isolate query form by default: either 'yes' or 'no', default 'no'.
- only_sets
 - Don't allow option to view the 'whole database' only list sets that have been defined: either 'yes' or 'no', default 'no'.
- · password
 - Password for access to isolates database, default 'remote'.
- port
 - Port number that the isolate host is listening on, default '5432'.
- privacy
 - Displays E-mail address for sender in isolate information page if set to 'no'. Default 'yes'.
- · query_script
 - Relative web path to bigsdb script. Default 'bigsdb.pl' (version 1.11+).
 - This is only needed if automated submissions are enabled. If bigsdb.pl is in a different directory from bigscurate.pl, you need to include the whole web path, e.g. /cgi-bin/bigsdb/bigsdb.pl.
- read_access
 - Describes who can view data: either 'public' for everybody or 'authenticated_users' for anybody who has been able to log in. Default 'public'.
- script_path_includes
 - Partial path of the bigsdb.pl script used to access the database. See *user authentication*.
- seqbin_size_threshold
 - Sets the size values in Mbp to enable for the *seqbin filter*.
 - Example: segbin size threshold="0.5,1,2,4".
- seq_export_limit
 - Overrides the sequence export limit (records x loci) in the Sequence Export plugin. Default: '1000000'.
- sets
 - Use sets: either 'yes' or 'no', default 'no'.
- set id
 - Force the use of a specific set when accessing database via this XML configuration: Value is the name of the set.
- start_id

- Defines the minimum record id to be used when uploading new isolate records. This can be useful when it is anticipated that two databases may be merged and it would be easier to do so if the id numbers in the two databases were different. Default: '1'.

· submissions

- Enable automated submission system: either 'yes' or 'no', default 'no' (version 1.11+).
- The curate_script and query_script paths should also be set, either in the bigsdb.conf file (for site-wide configuration) or within the system attribute of config.xml.
- submissions_deleted_days
 - Overrides the default number of days before closed submissions are deleted from the system. Default: '90'.
- · tblastx_tagging
 - Sets whether tagging can be performed using TBLASTX: either 'yes' or 'no', default 'no'.
- user
 - Username for access to isolates database, default 'apache'.
- user_job_quota
 - Integer with number of offline jobs that can be queued or currently running for this database by any specific user this parameter is only effective if users have to log in.
- · view
 - Database view containing isolate data, default 'isolates'.
- · views
 - Comma-separated list of views of the isolate table defined in the database. This is used to set a view for a set.
- webroot
 - URL of web root, which can be relative or absolute. The bigsdb.css stylesheet file should be located in this
 directory. Default '/'.

<field>

Element content: Field name + optional list < optlist> of allowed values, e.g.:

```
<field type="text" required="no" length="40" maindisplay="no"
  web="http://somewebsite.com/cgi-bin/script.pl?id=[?]" optlist="yes">epidemiology
  <optlist>
        <option>carrier</option>
        <option>healthy contact</option>
        <option>sporadic case</option>
        <option>endemic</option>
        <option>epidemic</option>
        <option>pandemic</option>
        <option>pandemic</option>
        <option>pandemic</option>
        <option>pandemic</option>
        </optlist>
        </field>
```

- type
 - Data type: int, text, float, bool, or date.
- comments * optional

- Comments about the field. These will be displayed in the field description plugin and as tooltips within the curation interface.
- · curate_only
 - Set to 'yes' to hide field on an isolate information page in the standard interface. The field will be visible if the page is accessed via the curator's interface (version 1.10.0+).
- · default
 - Default value. This will be entered automatically in the web form but can be overridden.
- dropdown
 - Select if you want this field to have its own dropdown filter box on the query page. If the field has an option list it will use the values in it, otherwise all values defined in the database will be included: 'yes' or 'no', default 'no'. This setting can be overridden by individual user preferences.
- · length
 - Length of field, default 12.
- maindisplay
 - Sets if field is displayed in the main table after a database search, 'yes' or 'no', default 'yes'. This setting can be overridden by individual user preferences.
- max
 - Maximum value for integer types. Special values such as CURRENT_YEAR can be used.
- min
 - Minimum value for integer types.
- · optlist
 - Sets if this field has a list of allowed values, default 'no'. Surround each option with an <option> tag.
- regex
 - Regular expression used to constrain field values, e.g. regex="^[A-Z].*\$" forces the first letter of the value to be capitalized.
- required
 - Sets if data is required for this field, 'yes' or 'no', default 'yes'.
- · userfield
 - Select if you want this field to have its own dropdown filter box of users (populated from the users table): 'yes' or 'no', default 'no'.
- web
 - URL that will be used to hyperlink field values. If [?] is included in the URL, this will be substituted for the actual field value.

Special values

The following special variables can be used in place of an actual value:

• CURRENT_YEAR: the 4 digit value of the current year

<sample>

Element content: Sample field name + optional list <optlist> of allowed values. Attributes are essentially the same as isolate field attributes, but refer to the samples table rather than the isolates table.

The sample table, if defined, must include isolate_id and sample_id fields, which must also be described in the XML file. These must be set as integer fields.

4.3.2 Sequence definition database XML attributes

Required attributes are in **bold**.

<db>

Top level element. Contains child elements: system, field and sample.

<system>

authentication

- Method of authentication: either 'builtin' or 'apache'. See user authentication.

db

- Name of database on system.

dbtype

- Type of database: either 'isolates' or 'sequences'.

description

- Description of database used throughout interface.
- align_limit
 - Overrides the sequence export record alignment limit in the Sequence Export plugin. Default: '200'.
- allele_comments
 - Enable comments on allele sequences: either 'yes' or 'no', default 'no'.
 - This is not enabled by default to discourage the practice of adding isolate information to allele definitions (this sort of information belongs in an isolate database).
- allele_flags
 - Enable flags to be set for alleles: either 'yes' or 'no', default 'no'.
- curate_path_includes
 - Partial path of the bigscurate.pl script used to curate the database. See *user authentication*.
- curate script
 - Relative web path to curation script. Default 'bigscurate.pl' (version 1.11+).
 - This is only needed if automated submissions are enabled. If bigscurate.pl is in a different directory from bigsdb.pl, you need to include the whole web path, e.g. /cgi-bin/private/bigsdb/bigscurate.pl.
- curators only
 - Set to 'yes' to prevent ordinary authenticated users having access to database configuration. This is only effective if read_access is set to 'authenticated_users'. This may be useful if you have different configurations for curation and querying with some data hidden in the configuration used by standard users. Default 'no'.
- daily_rest_submissions_limit

Overrides the limit on number of submissions that can be made to the database via the RESTful interface.
 This is useful to prevent flooding of the submission system by aberrant scripts. Default: '100'.

· diploid

- Allow IUPAC 2-nuclotide ambiguity codes in allele definitions for use with diploid typing schemes: either 'yes' or 'no', default 'no'.
- disable_seq_downloads
 - Prevent users or curators from downloading all alleles for a locus (admins always can). 'yes' or 'no', default 'no'.
- job_priority
 - Integer with default job priority for offline jobs (default:5).
- job_quota
 - Integer with number of offline jobs that can be queued or currently running for this database.
- profile_submissions
 - Enable profile submissions (automated submission system): either 'yes' or 'no', default 'no' (version 1.11+).
 - To enable, you will also need to set submissions="yes". By default, profile submissions are disabled since generally new profiles should be accompanied by representative isolate data, and the profile can be extracted from that.
- · query_script
 - Relative web path to bigsdb script. Default 'bigsdb.pl' (version 1.11+).
 - This is only needed if automated submissions are enabled. If bigsdb.pl is in a different directory from bigscurate.pl, you need to include the whole web path, e.g. /cgi-bin/bigsdb/bigsdb.pl.
- · read access
 - Describes who can view data: either 'public' for everybody, or 'authenticated_users' for anybody who has been able to log in. Default 'public'.
- script_path_includes
 - Partial path of the bigsdb.pl script used to access the database. See *user authentication*.
- seq_export_limit
 - Overrides the sequence export limit (records x loci) in the Sequence Export plugin. Default: '1000000'.
- · sets
 - Use sets: either 'yes' or 'no', default 'no'.
- set_id
 - Force the use of a specific set when accessing database via this XML configuration: Value is the name of the set.
- · submissions
 - Enable automated submission system: either 'yes' or 'no', default 'no' (version 1.11+).
 - The curate_script and query_script paths should also be set, either in the bigsdb.conf file (for site-wide configuration) or within the system attribute of config.xml.
- submissions deleted days

- Overrides the default number of days before closed submissions are deleted from the system. Default: '90'.
- · user_job_quota
 - Integer with number of offline jobs that can be queued or currently running for this database by any specific user - this parameter is only effective if users have to log in.
- · webroot
 - URL of web root, which can be relative or absolute. The bigsdb.css stylesheet file should be located in this
 directory. Default '/'.

4.4 User authentication

You can choose whether to allow Apache to handle your authentication or use built-in authentication.

4.4.1 Apache authentication

Using apache to provide your authentication allows a flexible range of methods and back-ends (see the Apache authentication HowTo for a start, or any number of tutorials on the web).

At its simplest, use a .htaccess file in the directory containing the bigscurate.pl (and bigsdb.pl for restriction of read-access) script or by equivalent protection of the directory in the main Apache server configuration. It is important to note however that, by default, any BIGSdb database can be accessed by any instance of the BIGSdb script (including one which may not be protected by a .htaccess file, allowing public access). To ensure that only a particular instance (protected by a specific htaccess directive) can access the database, the following attributes can be set in the system tag of the database XML description file:

- script_path_includes: the BIGSdb script path must contain the value set.
- curate_path_includes: the BIGSdb curation script path must contain the value set.

For public databases, the 'script_path_includes' attribute need not be set.

To use apache authentication you need to set the authentication attribute in the system tag of the database XML configuration to 'apache'.

4.4.2 Built-in authentication

BIGSdb has its own built-in authentication, using a separate database to store password and session hashes. The advantages of using this over many forms of apache authentication are:

- Users are able to update their own passwords.
- Passwords are not transmitted over the Internet in plain text.

When a user logs in, the server provides a random one-time session variable and the user is prompted to enter their username and password. The password is encrypted within the browser using a Javscript one-way hash algorithm, and this is combined with the session variable and hashed again. This hash is passed to the server. The server compares this hash with its own calculated hash of the stored encrypted password and session variable that it originally sent to the browser. Implementation is based on perl-md5-login.

To use built-in authentication you need to set the authentication attribute in the system tag of the database XML configuration to 'builtin'.

4.5 Setting up the admin user

The first admin user needs to be manually added to the users table of the database. Connect to the database using psql and add the following (changing details to suit the user).:

```
INSERT INTO users (id, user_name, surname, first_name, email, affiliation, status, date_entered, datestamp, curator) VALUES (1, 'keith', 'Jolley', 'Keith', 'keith.jolley@zoo.ox.ac.uk', 'University of Oxford, UK', 'admin', 'now', 'now', 1);
```

If you are using built-in authentication, set the password for this user using the *add_user.pl* script. This hashes the password and stores this within the authentication database. Other users can be added by the admin user from the curation interface accessible from http://your_website/cgi-bin/private/bigscurate.pl?db=test_db (or wherever you have located your bigscurate.pl script).

4.6 Updating PubMed citations

Publications listed in PubMed can be associated with individual isolate records, profiles, loci and sequences. Full citations for these are stored within a local reference database, enabling these to be displayed within isolate records and searching by publication and author. This local database is populated by a script that looks in BIGSdb databases for PubMed records not locally stored and then requests the full citation record from the PubMed database.

The script is called getrefs.pl and can be found in the scripts/maintenance directory. This script needs to know which BIGSdb databases and tables it needs to search for PubMed ids. These are listed in a configuration file (usually called getrefs.conf) which contains two columns - the first is the name of the database, the second is a comma-separated list of tables to search, e.g.

```
pubmlst_bigsdb_neisseria_isolatesrefspubmlst_bigsdb_neisseria_seqdefprofile_refs, sequence_refs, locus_refs
```

The script can be called as follows:

```
getrefs.pl getrefs.conf
```

Run either as the 'postgres' user or an account that is allowed to connect as the postgres user.

This should be run periodically from a CRON job, e.g. every hour.

Administrator's guide

Please note that links displayed within the curation interface will vary depending on database contents and the permissions of the curator.

5.1 Types of user

There are four types of user in BIGSdb:

- User can view data but never modify it. Users should be created for every submitter of data so that records can be tracked, even if they do not actually use the database.
- Submitter (isolate databases only) can add and modify their own isolate data and data submitted by anybody else that is in the same *user group* as them but not anyone elses. A limited range of *Individual permissions* can be set for each submitter, so their roles can be controlled. A submitter with no specific permissions set has no more power than a standard user.
- Curator can modify data but does not have full control of the database. *Individual permissions* can be set for
 each curator, so their roles can be controlled. A curator with no specific permissions set has no more power than
 a standard user.
- Admin has full control of the database, including setting permissions for curators and setting user passwords
 if built-in authentication is in use.

5.2 User groups

User groups allow submitter accounts to be grouped such that the submitter can edit isolates where the sender is either themselves or any member of a user group to which they belong.

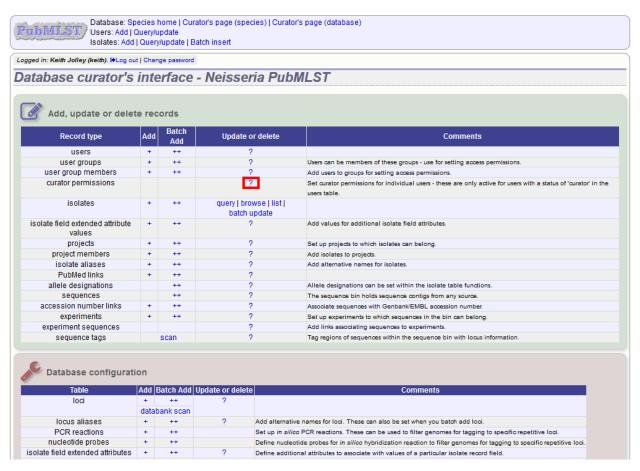
5.3 Curator permissions

Individual permissions can be set for each curator:

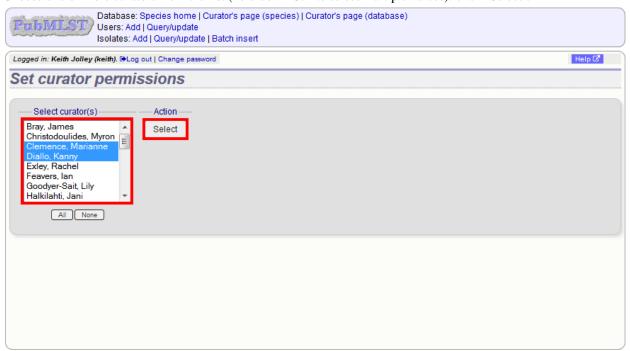
- disable_access if set to true, this user is completely barred from access.
- modify_users allowed to add or modify user records. They can change the status of users, but can not revoke admin privileges from an account. They can also not raise the status of a user to admin level.
- modify_usergroups allowed to add or modify user groups and add users to these groups.

- set_user_passwords allowed to modify other users' passwords (if built-in authentication is in use).
- modify_loci allowed to add or modify loci.
- modify_schemes allowed to add or modify schemes.
- modify_sequences allowed to add sequences to the sequence bin (for isolate databases) or new allele definitions (for sequence definition databases).
- modify_experiments define new experiments that can be used to group contigs uploaded to the sequence bin.
- modify_isolates allowed to add or modify isolate records.
- modify_projects allowed to create projects, modify their descriptions and add or remove isolate records to these.
- modify_composites allowed to add or modify composite fields (fields made up of other fields, including scheme fields defined in external databases). Composite fields involve defining regular expressions that are evaluated by Perl this can be dangerous so this permission should be granted with discretion.
- modify_field_attributes allow user to create or modify secondary field attributes (lookup tables) for isolate record fields.
- modify_value_attributes allow user to add or modify secondary field values for isolate record fields.
- modify_probes allow user to define PCR or hybridization reactions to filter tag scanning.
- tag_sequences allowed to tag sequences with locus information.
- designate_alleles allowed to manually designate allele numbers for isolate records.
- modify_profiles allowed to add or modify scheme profiles (only used in a sequence definitions database).

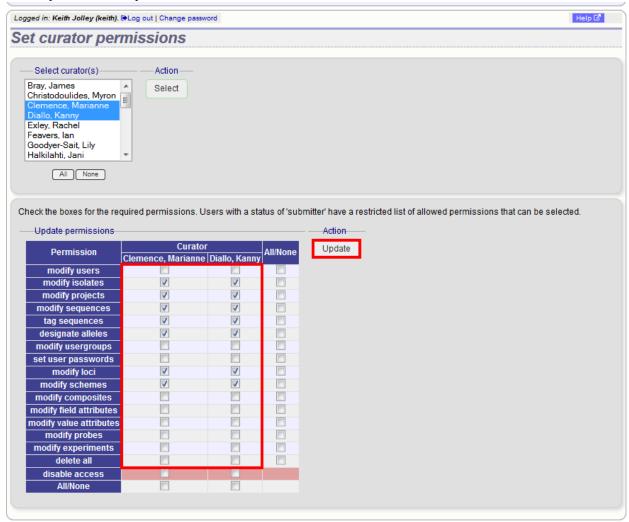
Permissions can be set by clicking the '?' button next to 'curator permissions' on the curator's interface:



Choose one or more curators from the list (hold down Ctrl to select multiple values). click 'Select'.



Click the appropriate checkboxes to modify permissions. There are also 'All/None' buttons to facilitate quicker selection of options. Click 'Update'.

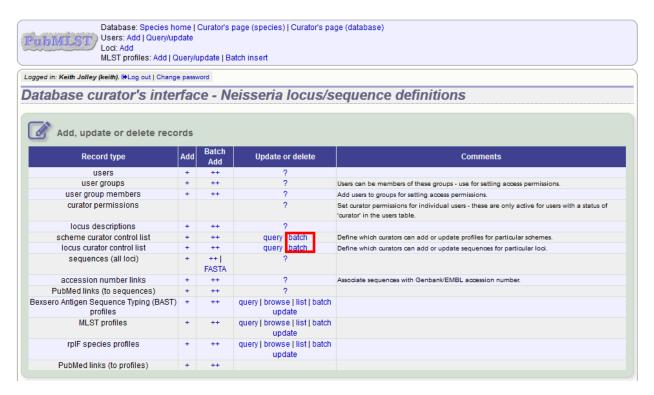


The 'disable access' option provides a quick way to disable access to a curator. This will not be selected by the 'All/None' buttons.

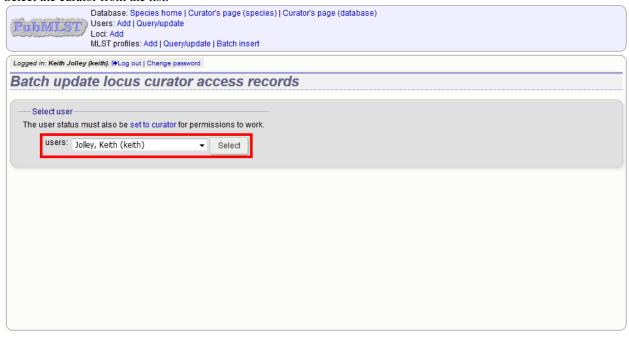
5.4 Locus and scheme permissions (sequence definition database)

To be allowed to define alleles or scheme profiles, curators must be granted specific permission for each locus and scheme by adding their user id number to the 'locus curator' and 'scheme curator' lists.

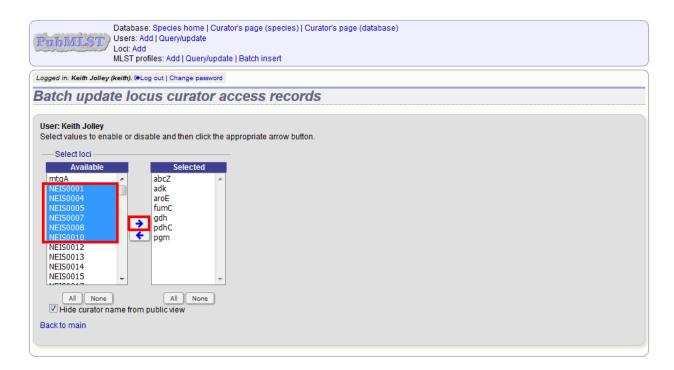
The easiest way to modify these lists is to use the batch update link next to 'locus curator control list' and 'scheme curator control list':



Select the curator from the list:



Then select loci/schemes that the user is allowed to curate in the left hand 'Available' list, and click the right button to move these to the 'Selected' list:



If you uncheck the 'Hide curator name from public view' checkbox, the curator name and E-mail address will appear alongside loci in the download table on the website.

5.5 Controlling access

5.5.1 Restricting particular configurations to specific user accounts

Suppose you only wanted specific users to access a database configuration.

In the config.xml, add the following directive:

```
default_access="deny"
```

This tells BIGSdb to deny access to anybody unless their account name appears within a file called users.allow within the config directory. The users.allow file should contain one username per line.

Alternatively, you can deny access to specific users, while allowing every other authenticated user. In config.xml, add the following directive:

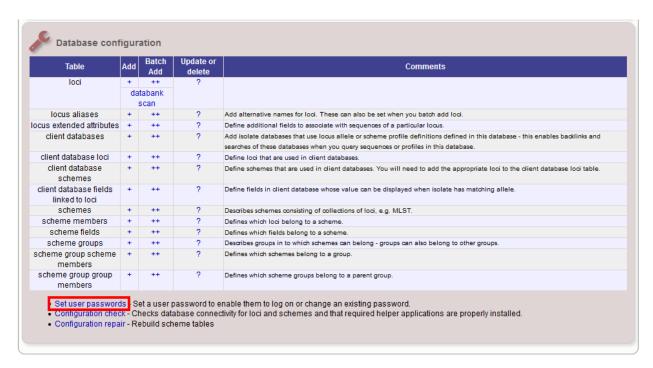
```
default_access="allow"
```

This tells BIGSdb to allow access to anybody unless their account name appears within a file called users.deny within the config directory. The users.deny file should contain one username per line.

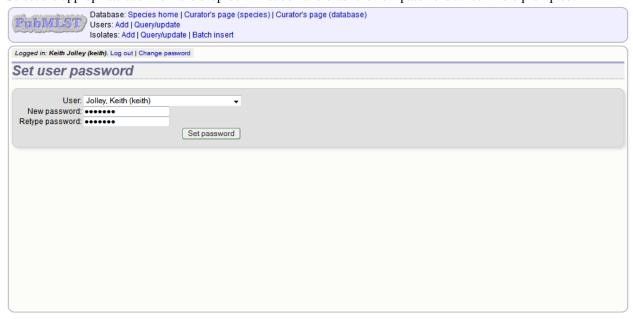
5.6 Setting user passwords

Please note that these instructions only apply if using the built-in BIGSdb authentication system.

If you are an administrator or a curator with specific permission to change other users' passwords, you should see a link to 'set user passwords' at the bottom of the curator's index page. Click this.



Select the appropriate user from the drop-down list box and enter the new password twice where prompted.



Click 'Set password' and the password will be updated.

5.7 Setting the first user password

To set the first administrator's password for a new database, use the add_user.pl script found in the scripts/maintenance directory:

```
add_user.pl [-a] -d <dbase> -n <username> -p <password>
```

The first user account needs to be added to the database *manually*.

5.8 Enabling plugins

Some plugins can be enabled/disabled for specific databases. If you look in the get_attributes function of the specific plugin file and see a value for system_flag, this value can be used in the system tag of the database configuration XML file to enable the plugin.

For example, the get_attributes function of the BURST plugin looks like:

```
sub get_attributes {
      my %att = (
                        => 'BURST',
              name
              author => 'Keith Jolley',
              affiliation => 'University of Oxford, UK',
              email => 'keith.jolley@zoo.ox.ac.uk',
              description => 'Perform BURST cluster analysis on query results query results',
              category => 'Cluster',
              buttontext => 'BURST',
              menutext => 'BURST',
              module
                         => 'BURST',
              version
                        => '1.0.0',
              dbtype => 'isolates, sequences',
              section => 'postquery',
order => 10,
              system_flag => 'BURST',
                       => 'query',
              input
              requires => 'mogrify',
                        => 2,
              min
                         => 1000
              max
      );
      return \%att;
```

The 'system_flag' attribute is set to 'BURST', so this plugin can be enabled for a database by adding:

```
BURST="yes"
```

to the system tag of the database XML file. If the system_flag value is not defined then the plugin is always enabled if it is installed on the system.

5.9 Temporarily disabling database updates

There may be instances where it is necessary to temporarily disable database updates. This may be during periods of server or database maintenance, for instance when running on a backup database server.

Updates can be disabled on a global or database-specific level.

5.9.1 Global

In the /etc/bigsdb/bigsdb.conf file, add the following line:

```
disable_updates=yes
```

An optional message can also be displayed by adding a disable_update_message value, e.g.

```
disable_update_message=The server is currently undergoing maintenance.
```

5.9.2 Database-specific

The same attributes described above for use in the bigsdb.conf file can also be used within the system tag of the database config.xml file, e.g.

```
<system
  db="bigsdb_neisseria"
  dbtype="isolates"
  ...
  disable_updates="yes"
  disable_update_message="The server is currently undergoing maintenance."</pre>
```

5.10 Host mapping

During periods of server maintenance, it may be necessary to map a database host to an alternative server. This would allow a backup database server to be used while the primary database server is unavailable. In this scenario, you would probably also want to *disable updates*.

Host mapping can be achieved by editing the /etc/bigsdb/host_mapping.conf file. Each host mapping is placed on a single line, with the current server followed by any amount of whitespace and then the new mapped host, e.g.

```
#Existing_host Mapped_host
server1 server2
localhost server2
```

[Lines beginning with a hash are comments and are ignored.]

This configuration would use server2 instead of server 1 or localhost wherever they are defined in the database configuration (either host attribute in the database config.xml file, or within the loci or schemes tables).

5.11 Improving performance

5.11.1 Use mod_perl

The single biggest improvement to speed can be obtained by running BIGSdb under mod_perl. There's very little point trying anything else until you have mod_perl set up and running - this can improve start-up performance a hundred-fold since the script isn't compiled on each page access but persists in memory.

5.11.2 Cache scheme definitions within an isolate database

If you have a large number of allelic profiles defined for a scheme, you can cache these definitions within an isolate database to speed up querying of isolates by scheme criteria (e.g. by ST for a MLST scheme).

To do this use the update_scheme_caches.pl script found in the scripts/maintenance directory, e.g. to cache all schemes in the pubmlst_bigsdb_neisseria_isolates database

```
update_scheme_caches.pl --database pubmlst_bigsdb_neisseria_isolates
```

This script creates indexed tables within the isolate database called temp_scheme_X and temp_isolates_scheme_fields_1 (where X is the scheme_id). If these table aren't present, they are created as temporary tables every time a query is performed that requires a join against scheme definition data. This requires importing all profile definitions from the definitions database and determining scheme field values for all isolates. This may sound like it would be slow but caching only has a noticeable effect once you have >5000 profiles.

You are able to update the cache for a single scheme, or a list of schemes, and choose the method of update. For large schemes, such as cgMLST, a full refresh may take a long time, so you may wish to only perform this infrequently (perhaps once a week) with more regular 'daily' or 'daily_replace' updates. A full list of options available are shown by typing

```
update_scheme_caches.pl --help
NAME
   update_scheme_caches.pl - Update scheme field caches
SYNOPSIS
    update_scheme_caches.pl --database NAME [options]
OPTIONS
--database NAME
   Database configuration name.
--help
   This help page.
--method METHOD
   Update method - the following values are allowed:
   full: Completely recreate caches
   incremental: Only add values for records not in cache.
   daily: Only add values for records not in cache updated today.
   daily_replace: Refresh values only for records updated today.
--quiet
   Don't output progress messages.
--schemes SCHEMES
   Comma-separated list of scheme ids to use.
    If left empty, all schemes will be updated.
```

Note that you will need to run this script periodically as a CRON job to refresh the cache. Admins can also refresh the caches manually from a link on the curators' page. This link is only present if the caches have been previously generated.



You can also set cache_schemes="yes" in the system tag of config.xml to enable automatic refreshing of the caches (using the 'daily' method) when batch adding new isolates (you should still periodically run the update_scheme_caches.pl script via CRON to ensure any changes in the sequence definition database are picked up).

If queries are taking longer than 5 seconds to perform and a cache is not in place, you will see a warning message in bigsdb.log suggesting that the caches be set up. Unless you see this warning regularly, you probably don't need to do this.

5.11.3 Use a ramdisk for the secure temporary directory

If you are running BIGSdb on a large server with lots of RAM, you could use some of this as a ramdisk for temporary files. Debian/Ubuntu systems make available up to half the system RAM as a ramdisk mounted under /run/shm (or /dev/shm) by default. Set the secure_tmp_dir to this RAM disk and you should see significant improvement in operations requiring the writing of lots of temporary files, e.g. tag scanning and the Genome Comparator plugin. This is only likely to be appropriate if you have very large amounts of RAM available. As an example, the server hosting the PubMLST databases is a dedicated machine with 1TB RAM with temporary files rarely using more than 50GB space.

5.12 Dataset partitioning

5.12.1 Sets

Sets provide a means to partition the database in to manageable units that can appear as smaller databases to an end user. Sets can include constrained groups of isolates, loci, and schemes from the complete database and also include additional metadata fields only applicable to that set.

See also:

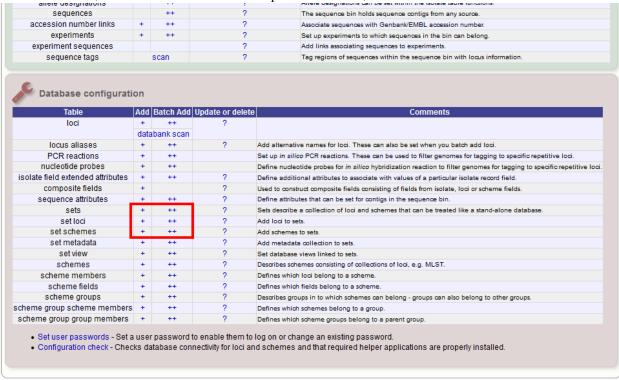
Sets (concept)

5.12.2 Configuration of sets

First sets need to be enabled in the XML configuration file (config.xml) of the database. Add the following attribute to the system tag:

sets="yes"

With this attribute, the curation interface now has options to add sets, and then add loci or schemes to these sets.



The name of a locus or scheme to use within a set can be defined in the set_name field when adding loci or schemes to a set. Common names can also be set for loci - equivalent to the common name used within the loci table.

Now when a user goes to the contents page of the database they will be presented with a dropdown menu of datasets and can choose either the 'whole database' or a specific set. This selection is remembered between sessions.



Alternatively, a specific set can be selected within the XML config file so that only a specific set is available when accessed via that configuration. In that case, the user would be unaware that the database contains anything other than the loci and schemes available within the set.

To specify this, add the following attributute to the system tag:

```
set_id="1"
```

where the value is the name of the set.

Note: If the set_id attribute is set, database configuration settings in the curator's interface are disabled. This is because when the configuration is constrained to a set, only loci and schemes already added to the set are visible, so functionality to edit schemes/loci would become very confusing. To modify these settings, you either need to access the interface from a different configuration, i.e. an alternative config.xml with the set_id attribute not set, or temporarily remove the set_id directive from the current config.xml while making configuration changes.

5.12.3 Set metadata

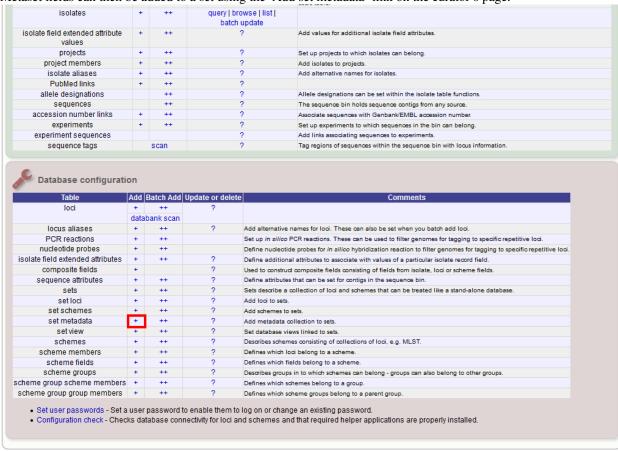
Additional metadata fields can be set within the XML configuration file. They are specified as belonging to a metaset by prefixing the field name with 'meta_NAME:' where NAME is the name of the metaset, e.g.

```
<field type="text" required="no" length="30" maindisplay="no"
    optlist="yes">meta_1:clinical_outcome
  <optlist>
```

```
<option>no sequeleae</option>
  <option>hearing loss</option>
  <option>amputation</option>
  <option>death</option>
  </optlist>
</field>
```

Metaset fields can be defined just like any other *provenance field* and their position in the output is determined by their position in the XML file.

Metaset fields can then be added to a set using the 'Add set metadata' link on the curator's page.



A new database table needs to be added for each metaset. This should contain all the fields defined for a metaset. The table should also contain an isolate_id field to act as the foreign key linking to the isolate table, e.g. the SQL would look something like the following:

```
CREATE TABLE meta_1 (
isolate_id integer NOT NULL,
town text,
clinical_outcome text,
PRIMARY KEY (isolate_id),
CONSTRAINT m1_isolate_id FOREIGN KEY (isolate_id) REFERENCES isolates
ON DELETE CASCADE
ON UPDATE CASCADE
);
GRANT SELECT, UPDATE, INSERT, DELETE ON meta_1 TO apache;
```

The above creates the database table for a metaset called '1', defining new text fields for 'town' and 'clinical_outcome'.

5.12.4 Set views

Finally the isolate record table can be partitioned using database views and these views associated with a set. Create views using something like the following:

```
CREATE VIEW spneumoniae AS SELECT * FROM isolates WHERE species = 'Streptococcus pneumoniae'; GRANT SELECT ON spneumoniae TO apache;
```

Add the available views to the XML file as a comma separated list in the system tag 'views' attribute:

```
<system
....
sets="yes"
views="spneumoniae, saureus"
>
</system>
```

Set the view to the set by using the 'Add set view' link on the curator's page.

5.12.5 Using only defined sets

The only_sets attribute can be set to 'yes' to disable the option for 'Whole database' so that only sets can be viewed, e.g.

```
<system
....
sets="yes"
only_sets="yes"
>
</system>
```

5.13 Adding new loci

See also:

Loci (concept)

5.13.1 Sequence definition databases

Single locus

Click the add (+) loci link on the curator's interface contents page.



Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:

- id The name of the locus.
 - Allowed: any value starting with a letter or underscore.
- data_type Describes whether the locus is defined by nucleotide or peptide sequence.
 - Allowed: DNA/peptide.
- allele_id_format The format for allele identifiers.
 - Allowed: integer/text.
- length varies Sets whether alleles can vary in length.
 - Allowed: true/false.
- coding_sequence Sets whether the locus codes for a protein.
 - Allowed: true/false.
- formatted_name Name with HTML formatting (optional).
 - This allows you to add formatting such as bold, italic, underline and superscripting to locus names as they
 appear in the web interface.
 - Allowed: valid HTML.
- common_name The common name for the locus (optional).
 - Allowed: any value.
- formatted_common_name Common name with HTML formatting (optional).
 - Allowed: valid HTML.
- allele_id_regex Regular expression to enforce allele id naming (optional).
 - ^: the beginning of the string
 - \$:the end of the string
 - d: digit

- D: non-digit
- s: white space character
- S: non white space character
- w: alpha-numeric plus '_'
- .: any character
- *: 0 or more of previous character
- +: 1 or more of previous character
- e.g. ^Fd-d+\$ states that an allele name must begin with a F followed by a single digit, then a dash, then
 one or more digits, e.g. F1-12
- length Standard length of locus (required if length_varies is set to false.
 - Allowed: any integer.
- min_length Minimum length of locus (optional).
 - Allowed: any integer.
- max length Maximum length of locus (optional).
 - Allowed: any integer (larger than the minimum length).
- orf Open reading frame of locus (optional).
 - 1-3 are the forward reading frame, 4-6 are the reverse reading frames.
 - Allowed: 1-6.
- genome_position The start position of the locus on a reference genome (optional).
 - Allowed: any integer.
- match_longest Specifies whether in a sequence query to only return the longest match (optional).
 - This is useful for some loci that can have some sequences shorter than others, e.g. peptide loci defining antigenic loops. This can lead to instances of one sequence being longer than another but otherwise being identical. In these cases, usually the longer sequence is the one that should be matched.
 - Allowed: true/false.
- full name Full name of the locus (optional).
 - Allowed: any value.
- product Name of gene product (optional).
 - Allowed: Any value.
- description Description of the locus (optional).
 - Allowed: any value.
- aliases Alternative names for the locus (optional).
 - Enter each alias on a separate line in the text box.
 - Allowed: any value.
- pubmed_ids PubMed ids of publications describing the locus (optional).
 - Enter each PubMed id on a separate line in the text box.
 - Allowed: any integer.

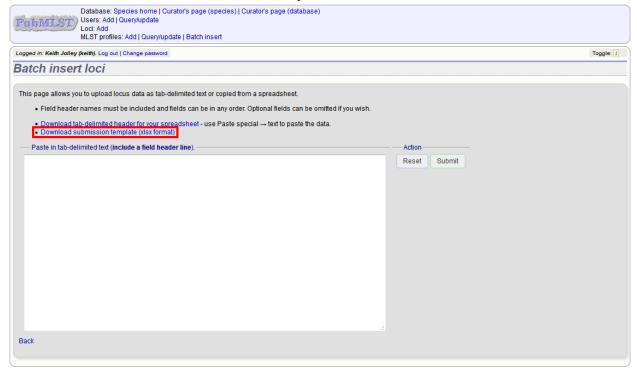
- links Hyperlinks pointing to additional resources to display in the locus description (optional).
 - Enter each link on a separate line in the format with the URL first, followed by a | then the description (URLIdescription).

Batch adding

Click the batch add (++) loci link on the curator's interface contents page.



Click the link to download a header line for an Excel spreadsheet:



Fill in the spreadsheet using the fields described for adding single loci.

Fill in the spreadsheet fields using the table above as a guide, then paste the completed table into the web form and press 'Submit query'.

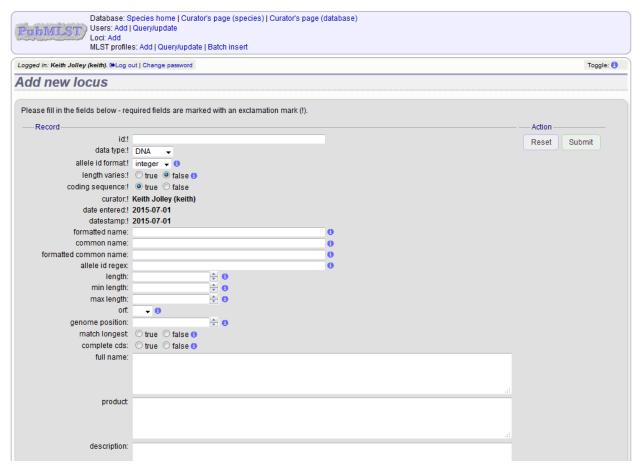
5.13.2 Isolate databases

Single locus

Click the add (+) loci link on the curator's interface contents page.



Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:



- id The name of the locus
 - Allowed: any value starting with a letter or underscore.
- data_type Describes whether the locus is defined by nucleotide or peptide sequence.
 - Allowed: DNA/peptide.
- allele_id_format The format for allele identifiers.
 - Allowed: integer/text.
- length_varies Sets whether alleles can vary in length.
 - Allowed: true/false.
- coding sequence Sets whether the locus codes for a protein.
 - Allowed: true/false.
- flag_table Set to true to specify that the sequence definition database contains an allele flag table (which is the case for BIGSdb version 1.4 onwards).
 - Allowed: true/false.
- isolate_display Sets how alleles for this locus are displayed in a detailed isolate record this can be overridden by user preference.
 - Allowed: allele only/sequence/hide.
- main_display Sets whether or not alleles for this locus are displayed in a main results table by default this can be overridden by user preference.

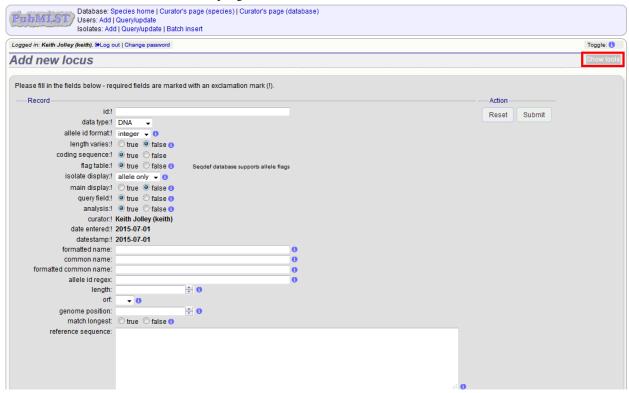
- Allowed: true/false.
- query_field Sets whether or not alleles for this locus can be used in queries by default this can be overridden by user preference.
 - Allowed: true/false.
- analysis Sets whether or not alleles for this locus can be used in analysis functions by default this can be overridden by user preference.
 - Allowed: true/false.
- formatted_name Name with HTML formatting (optional).
 - This allows you to add formatting such as bold, italic, underline and superscripting to locus names as they
 appear in the web interface.
 - Allowed: valid HTML.
- common_name The common name for the locus (optional).
 - Allowed: any value.
- formatted_common_name Common name with HTML formatting (optional).
 - Allowed: valid HTML.
- allele_id_regex Regular expression to enforce allele id naming.
 - ^: the beginning of the string
 - \$:the end of the string
 - d: digit
 - D: non-digit
 - s: white space character
 - S: non white space character
 - w: alpha-numeric plus '_'
 - .: any character
 - *: 0 or more of previous character
 - +: 1 or more of previous character
 - e.g. ^Fd-d+\$ states that an allele name must begin with a F followed by a single digit, then a dash, then one or more digits, e.g. F1-12
- length Standard length of locus (required if length varies is set to false).
 - Allowed: any integer.
- orf Open reading frame of locus (optional). 1-3 are the forward reading frame, 4-6 are the reverse reading frames.
 - Allowed: 1-6.
- genome_position The start position of the locus on a reference genome.
 - Allowed: any integer.
- match_longest Only select the longest exact match when tagging/querying.

- This is useful for some loci that can have some sequences shorter than others, e.g. peptide loci defining antigenic loops. This can lead to instances of one sequence being longer than another but otherwise being identical. In these cases, usually the longer sequence is the one that should be matched.
- Allowed: true/false.
- reference_sequence Sequence used by the automated sequence comparison algorithms to identify sequences matching this locus. This is only used if a sequence definition database has not been set up for this locus.
- pcr_filter Set to true if this locus is further defined by genome filtering using in silico PCR.
 - Allowed: true/false.
- probe_filter Set to true if this locus is further defined by genome filtering using in silico hybdridization.
 - Allowed: true/false.
- dbase_name Name of database (system name).
 - Allowed: any text.
- dbase_host Resolved name of IP address of database host leave blank if running on the same machine as the isolate database.
 - Allowed: network address, e.g. 129.67.26.52 or zoo-oban.zoo.ox.ac.uk
- dbase_port Network port on which the sequence definition database server is listening leave blank unless using a non-standard port (5432).
 - Allowed: integer.
- dbase_user Name of user with permission to access the sequence definition database depending on the database configuration you may be able to leave this blank.
 - Allowed: any text (no spaces).
- dbase_password Password of database user again depending on the database configuration you may be able to leave this blank.
 - Allowed: any text (no spaces).
- dbase_table Table in the sequence definition database that contains allele sequences for this locus. If the definition database uses BIGSdb this will be 'sequences'.
 - Allowed: any text (no spaces).
- dbase_id_field Primary field in sequence database that defines allele. If the definition database uses BIGSdb this will be 'allele_id'.
 - Allowed: any text (no spaces).
- dbase_id2_field Secondary field in sequence database that defines alleles. If dbase_id_field uniquely defines
 alleles for this locus then this should be left blank. If the definition database uses BIGSdb this will be 'locus'.
 - Allowed: any text (no spaces).
- dbase_id2_value Secondary field value in sequence database that defines alleles. If dbase_id_field uniquely
 defines alleles for this locus then this should be left blank. If the definition database uses BIGSdb this will be
 the name of the locus used in the id field
 - Allowed: any text (no spaces).
- dbase_seq_field Field in sequence database containing allele sequence. If the definition database uses BIGSdb this will be 'sequence'.
 - Allowed: any text (no spaces).

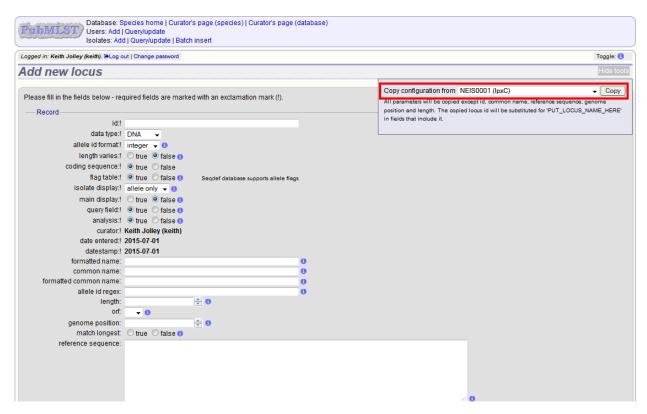
- description_url The URL used to hyperlink to locus information in the isolate information page. This can either be a relative (e.g. /cgi-bin/...) or an absolute (containing http://) URL.
 - Allowed: any valid URL.
- url The URL used to hyperlink to information about the allele. This can either be a relative or absolute URL. If [?] (including the square brackets) is included then this will be substituted for the allele value in the resultant URL. To link to the appropriate allele info page on a corresponding seqdef database you would need something like /cgi-bin/bigsdb/bigsdb.pl?db=pubmlst_neisseria_seqdef&page=alleleInfo&locus=abcZ&allele_id=[?].
 - Allowed: any valid URL.
- submission_template Sets whether or not a column for this locus is included in the Excel submission template.
 - Allowed: true/false (default: false)

Using existing locus definition as a template

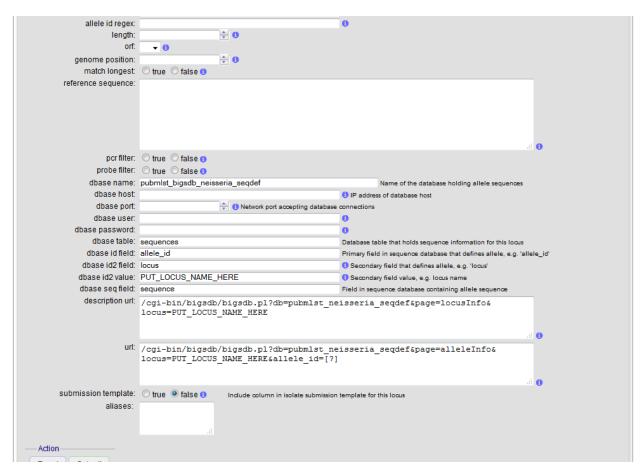
When defining a new locus in the isolate database, it is possible to use an existing locus record as a template. To do this, click the 'Show tools' link in the top-right of the screen:



This displays a drop-down box containing existing loci. Select the locus that you wish to use as a template, and click 'Copy'.



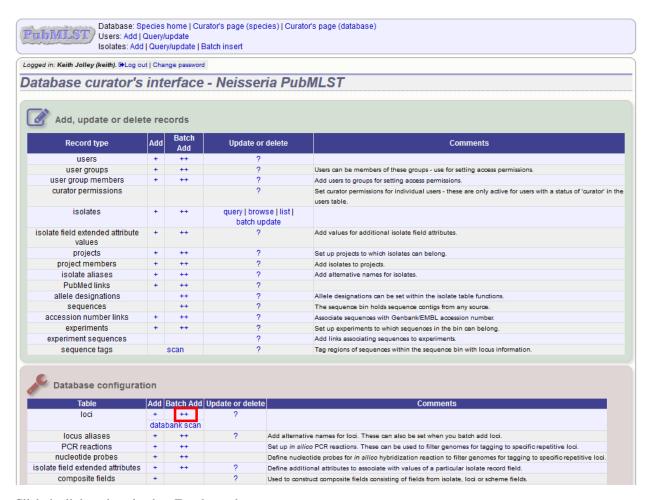
The configuration will be copied over to the web form, with the exception of name fields. Some fields will require you to change the value 'PUT_LOCUS_NAME_HERE' with the value you enter in the id field. These are usually the dbase_id2_value, description_url and url fields:



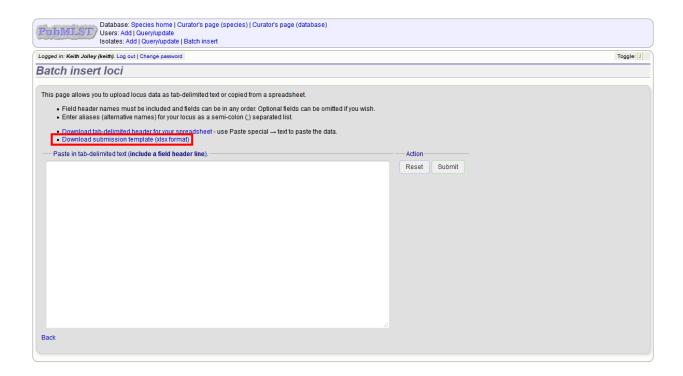
Complete the form and click 'Submit'.

Batch adding

Click the batch add (++) loci link on the curator's interface contents page.



Click the link to download an Excel template:

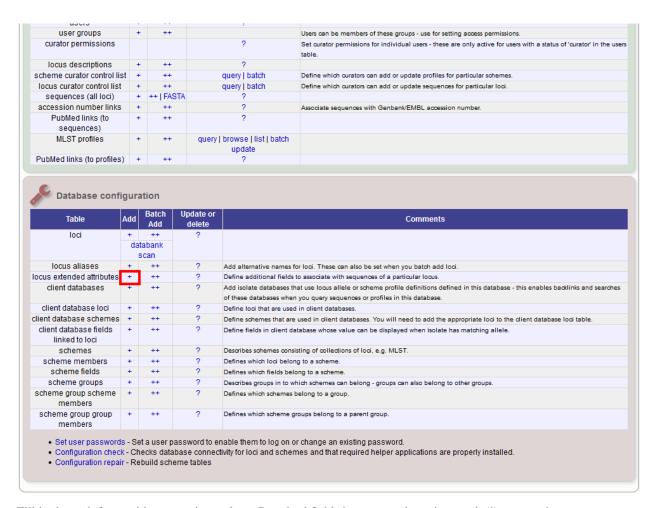


Fill in the spreadsheet fields using the *table above as a guide*, then paste the completed table into the web form and press 'Submit query'.

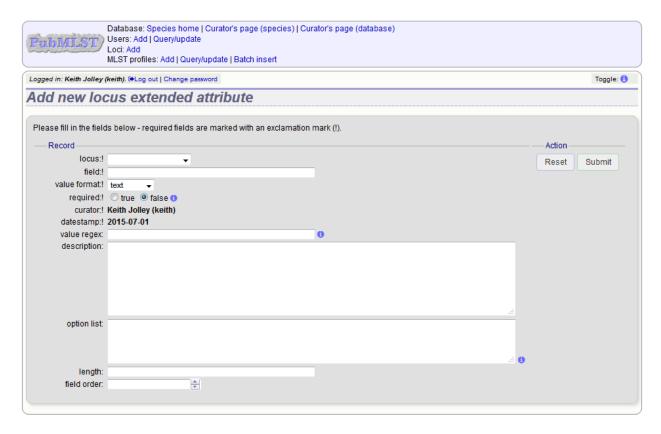
5.14 Defining locus extended attributes

You may want to add additional metadata for the allele definitions of some loci. Since these are likely to be specific to each locus, they cannot be defined generically within the standard locus definition. We can, instead, define extended attributes. Examples of these include higher order grouping of antigen sequences, antibody reactivities, identification of important mutations, or cross-referencing of alternative nomenclatures.

To add extended attributes for a locus, click add (+) locus extended attributes in the sequence definition database curator's interface contents page.



Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:



- locus Select locus from dropdown box.
 - Allowed: existing locus name.
- field Name of extended attributes.
 - Allowed: any value.
- value_format Data type of attribute.
 - Allowed: integer/text/boolean.
- required Specifies whether the attribute value but be defined when adding a new sequence.
 - Allowed: true/false.
- value_regex Regular expression to enforce allele id naming (optional).
 - ^: the beginning of the string
 - \$:the end of the string
 - d: digit
 - D: non-digit
 - s: white space character
 - S: non white space character
 - w: alpha-numeric plus '_'
 - .: any character
 - *: 0 or more of previous character
 - +: 1 or more of previous character

- description Description that will appear within the web form when adding new sequences (optional).
 - Allowed: any value.
- option_list Pipe (l) separated list of allowed values (optional).
- length Maximum length of value (optional).
 - Allowed: any integer.
- field order Integer that sets the position of the field within scheme values in any results (optional).
 - Allowed: any integer.

Once extended attributes have been defined, they will appear in the web form when adding new sequences for that locus. The values are searchable when using a *locus-specific sequence query*, and they will appear within query results and allele information pages.

5.15 Defining schemes

Schemes are collections of loci that may be associated with particular fields - one of these fields can be a primary key, i.e. a field that uniquely defines a particular combination of alleles at the associated loci.

A specific example of a scheme is MLST - see workflow for setting up a MLST scheme.

To set up a new scheme, you need to:

- 1. Add a new scheme description.
- 2. Define loci as 'scheme members'.
- 3. Add 'scheme fields' associated with the scheme.

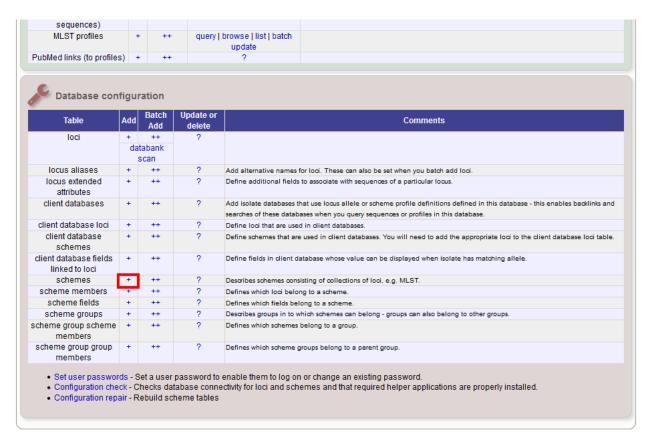
See also:

Schemes (concept)

5.15.1 Sequence definition databases

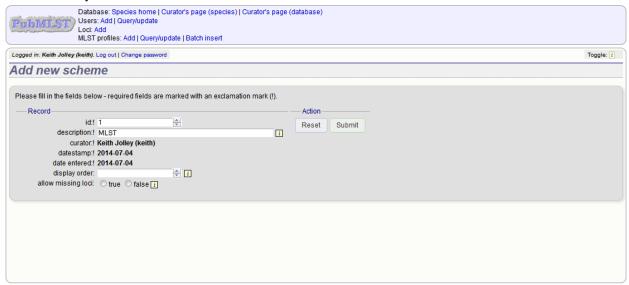
As with all configuration, tables can be populated using the batch interface (++) or one at a time (+). Details for the latter are described below:

Click the add (+) scheme link on the curator's interface contents page.

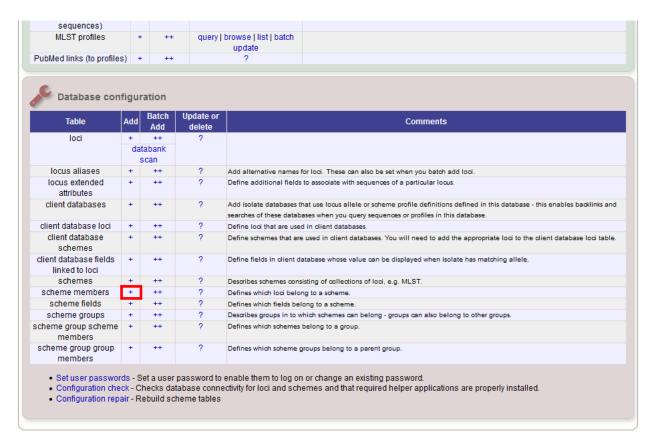


Fill in the scheme description in the web form. The next available scheme id number will be filled in already.

The display_order field accepts an integer that can be used to order the display of schemes in the interface - this can be left blank if you wish.



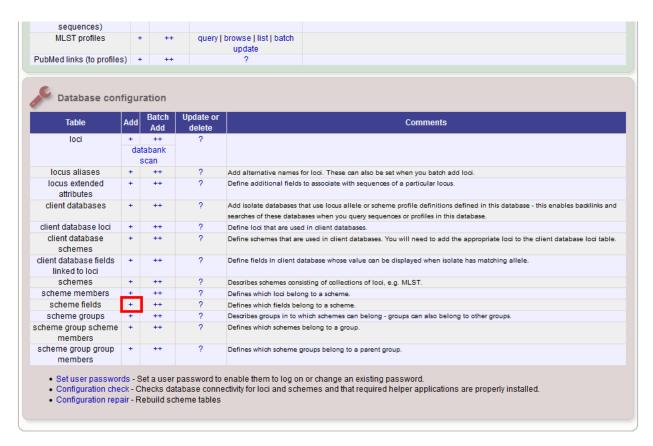
To add loci to the scheme, click the add (+) scheme members link on the curator's interface contents page.



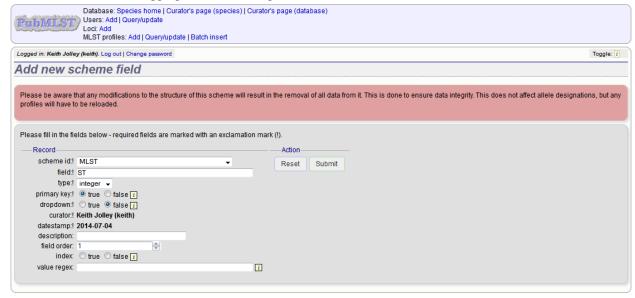
Select the scheme name and a locus that you wish to add to the scheme from the appropriate drop-down boxes. *Loci need to have already been defined*. The field_order field allows you to set the display order of the locus within a profile - if these are left blank that alphabetical ordering is used.



To add scheme fields, click the add (+) scheme fields link on the curator's interface contents page.



Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:



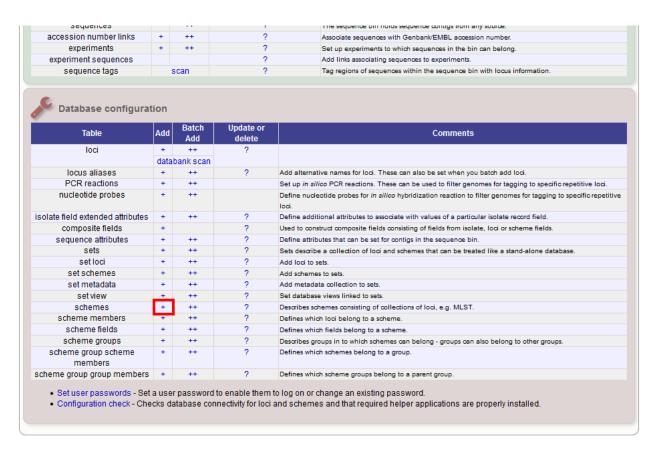
- scheme_id Dropdown box of scheme names.
 - Allowed: selection from list.
- field Field name.
 - Allowed: any value.
- type Format for values.

- Allowed: text/integer/date.
- primary_key Set to true if field is the primary key. There can only be one primary key for a scheme.
 - Allowed: true/false.
- dropdown Set to true if a dropdown box is displayed in the query interface, by default, for values of this field to be quickly selected. This option can be overridden by user preferences.
 - Allowed: true/false.
- description This field isn't currently used.
- field_order Integer that sets the position of the field within scheme values in any results.
 - Allowed: any integer.
- value_regex Regular expression to enforce field values.
 - ^: the beginning of the string
 - \$:the end of the string
 - d: digit
 - D: non-digit
 - s: white space character
 - S: non white space character
 - w: alpha-numeric plus '_'
 - .: any character
 - *: 0 or more of previous character
 - +: 1 or more of previous character

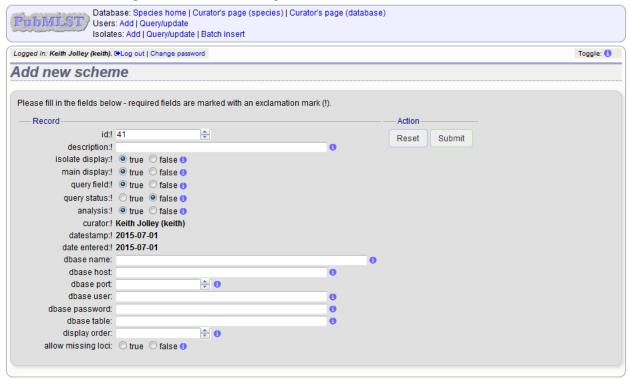
5.15.2 Isolate databases

As with all configuration, tables can be populated using the batch interface (++) or one at a time (+). Details for the latter are described below:

Click the add (+) scheme link on the curator's interface contents page.



Fill in the scheme description in the web form. Required fields have an exclamation mark (!) next to them:



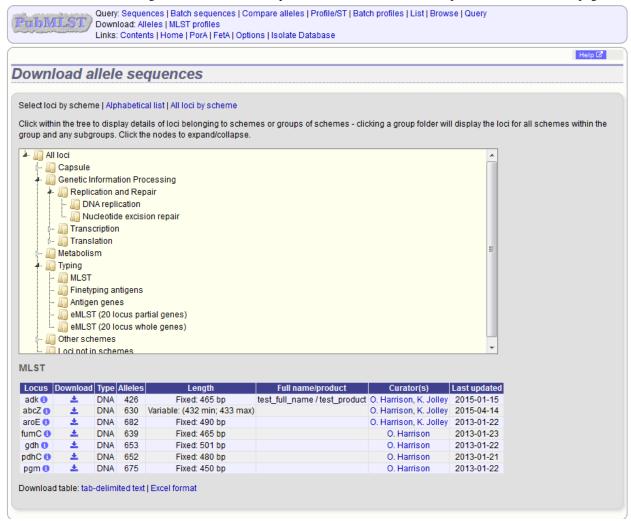
id - Index number of scheme - the next available number will be entered automatically.

- Allowed: any positive integer.
- description Short description this is used in tables so make sure it's not too long.
 - Allowed: any text.
- isolate_display Sets whether or not fields for this scheme are displayed in a detailed isolate record this can be overridden by user preference.
 - Allowed: allele only/sequence/hide.
- main_display Sets whether or not fields for this scheme are displayed in a main results table by default this can be overridden by user preference.
 - Allowed: true/false.
- query_field Sets whether or not fields for this scheme can be used in queries by default this can be overridden
 by user preference.
 - Allowed: true/false.
- query_status Sets whether a dropdown list box should be displayed in the query interface to filter results based on profile completion for this scheme this can be overridden by user preference.
 - Allowed: true/false.
- analysis Sets whether or not alleles for this locus can be used in analysis functions by default this can be overridden by user preference.
 - Allowed: true/false.
- dbase name Name of segdef database (system name) containing scheme profiles (optional).
 - Allowed: any text.
- dbase_host Resolved name of IP address of database host leave blank if running on the same machine as the isolate database (optional).
 - Allowed: network address, e.g. 129.67.26.52 or zoo-oban.zoo.ox.ac.uk
- dbase_port Network port on which the sequence definition database server is listening leave blank unless using a non-standard port, 5432 (optional).
 - Allowed: integer.
- dbase_user Name of user with permission to access the sequence definition database depending on the database configuration you may be able to leave this blank (optional).
 - Allowed: any text (no spaces).
- dbase_password Password of database user again depending on the database configuration you may be able to leave this blank (optional).
 - Allowed: any text (no spaces).
- dbase_table Table in the sequence definition database that contains profiles for this scheme. If the definition database uses BIGSdb this will be 'scheme_X' where X is the scheme id number in the sequef database.
 - Allowed: any text (no spaces).
- display_order Integer reflecting the display position for this scheme within the interface (optional).
 - Allowed: any integer.
- allow_missing_loci Allow profile definitions to contain '0' (locus missing) or 'N' (any allele).

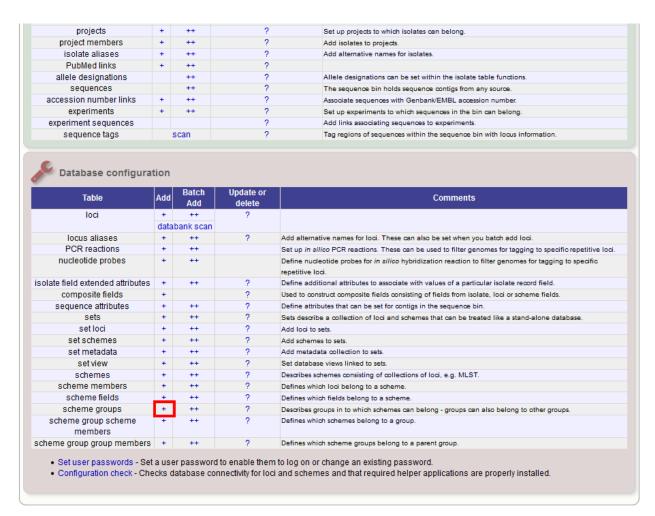
5.16 Organizing schemes into hierarchical groups

Schemes can be organized in to groups, and these groups can in turn be members of other groups. This faciliates hierarchical ordering of loci, but with the flexibility to allow loci and schemes to belong to multiple groups.

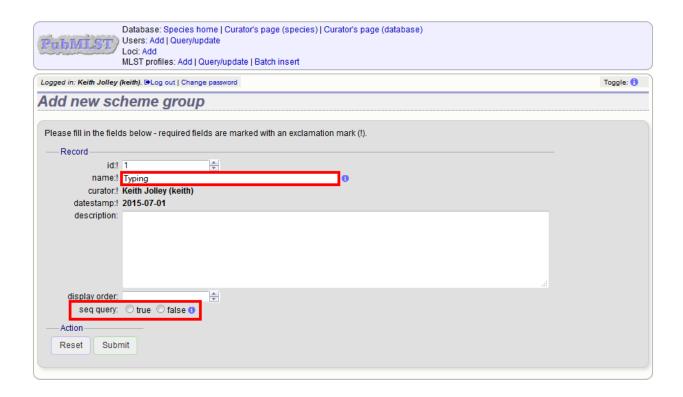
This hierarchical structuring can be seen in various places within BIGSdb, for example the *allele download* page.



Scheme groups can be added in both the sequence definition and isolate databases. To add a new group, click the add (+) scheme group link on the curator's contents page.



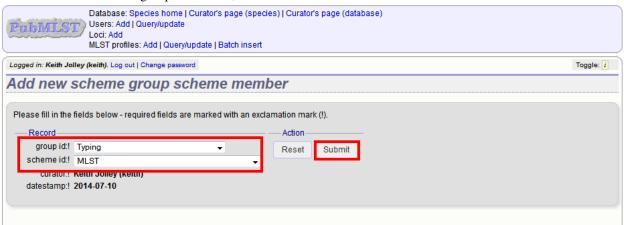
Enter a short name for the group - this will appear within drop-down list boxes and the hierarchical tree, so it needs to be fairly short.



If you are creating a scheme group in the sequence definition database, there is an additional field called 'seq_query'. Set this to true to add the scheme group to the dropdown lists in the *sequence query* page. This enables all loci belonging to schemes within the group to be queried together.

Schemes can be added to groups by clicking the add (+) scheme group scheme members link. Database configuration Update or Batch Comments Add delete loci databank scan locus aliases Add alternative names for loci. These can also be set when you batch add loci. locus extended Define additional fields to associate with sequences of a particular locus. attributes client databases 2 Add isolate databases that use locus allele or scheme profile definitions defined in this database - this enables backlinks and searches of these databases when you query sequences or profiles in this database. client database loci Define loci that are used in client databases. client database Define schemes that are used in client databases. You will need to add the appropriate loci to the client database loci schemes client database fields Define fields in client database whose value can be displayed when isolate has matching allele linked to loci schemes Describes schemes consisting of collections of loci, e.g. MLST. scheme members Defines which loci belong to a scheme ? scheme fields Defines which fields belong to a scheme. Describes groups in to which schemes can belong - groups can also belong to other groups. scheme groups scheme group scheme + Defines which schemes belong to a group. members scheme group group Defines which scheme groups belong to a parent group. members . Set user passwords - Set a user password to enable them to log on or change an existing password. . Configuration check - Checks database connectivity for loci and schemes and that required helper applications are properly installed. Configuration repair - Rebuild scheme tables

Select the scheme and the group to add it to, then click 'Submit'.



Scheme groups can be added to other scheme groups in the same way by clicking the add (+) scheme group group members link.

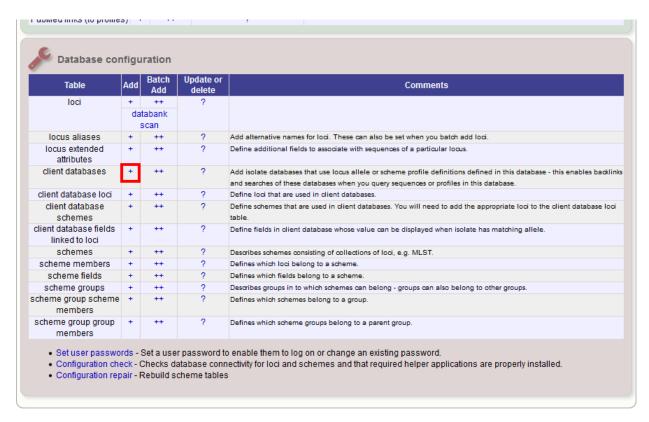
5.17 Setting up client databases

Sequence definition databases can have any number of isolate databases that connect as clients. Registering these databases allows the software to perform isolate data searches relevant to results returned by the sequence definition database, for example:

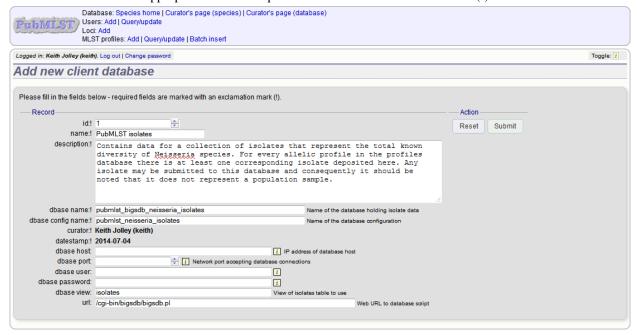
- Determine the number of isolates that a given allele is found in and link to these.
- Determine the number of isolates that a given scheme field, e.g. a sequence type, is found in and link to these.
- Retrieve specific attributes of isolates that have a given allele, e.g. species that have a particular 16S allele, or penicillin resistance given a particular penA allele.

Multiple client databases can be queried simultaneously.

To register a client isolate database for a sequence definition database, click the add (+) client database link on the curator's interface contents page.



Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:

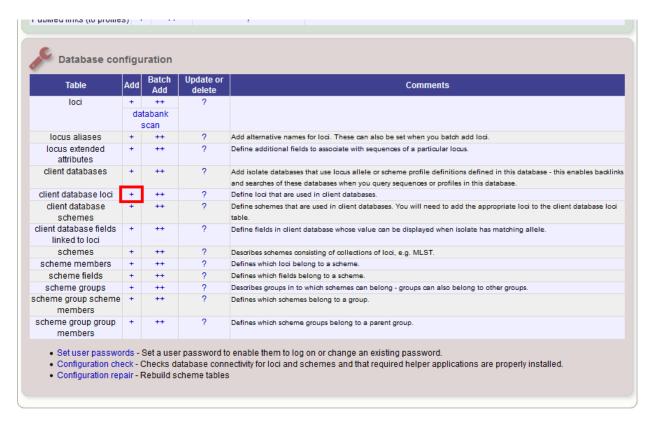


- id Index number of client database. The next available number is entered automatically but can be overridden.
 - Allowed: any positive integer.
- name Short description of database. This is used within the interface result tables so it is better to make it as short as possible.

- Allowed: any text.
- description Longer description of database.
 - Allowed: any text.
- dbase_name Name of database (system name).
 - Allowed: any text.
- dbase_config_name Name of database configuration this is the text string that appears after the db= part of script URLs.
 - Allowed: any text (no spaces)
- dbase_host Resolved name of IP address of database host (optional).
 - Allowed: Network address, e.g. 129.67.26.52 or zoo-oban.zoo.ox.ac.uk
 - Leave blank if running on the same machine as the seqdef database.
- dbase_port Network port on which the client database server is listening (optional).
 - Allowed: integer.
 - Leave blank unless using a non-standard port (5432).
- dbase_user Name of user with permission to access the client database.
 - Allowed: any text (no spaces).
 - Depending on the database configuration you may be able to leave this blank.
- dbase_password Password of database user
 - Allowed: any text (no spaces).
 - Depending on the database configuration you may be able to leave this blank.
- url URL of client database bigsdb.pl script
 - Allowed: valid script path.
 - This can be relative (e.g. /cgi-bin/bigsdb/bigsdb.pl) if running on the same machine as the seqdef database or absolute (including http://) if on a different machine.

5.17.1 Look up isolates with given allele

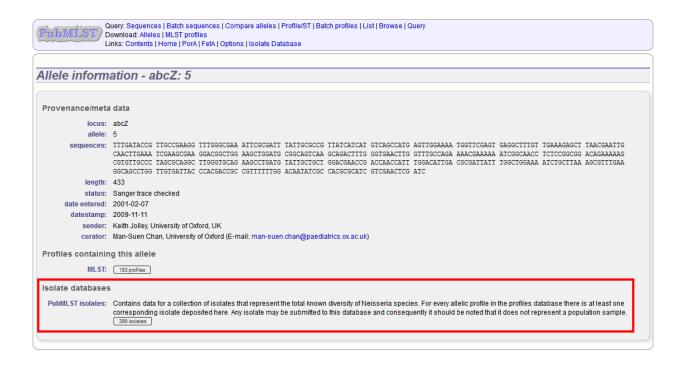
To link a locus, click the add (+) client database loci link on the curator's interface contents page.



Link the locus to the appropriate client database using the dropdown list boxes. If the locus is named differently in the client database, fill this name in the locus_alias.



Now when information on a given allele is shown following a query, the software will list the number of isolates with that allele and link to a search on the database to retrieve these.



5.17.2 Look up isolates with a given scheme primary key

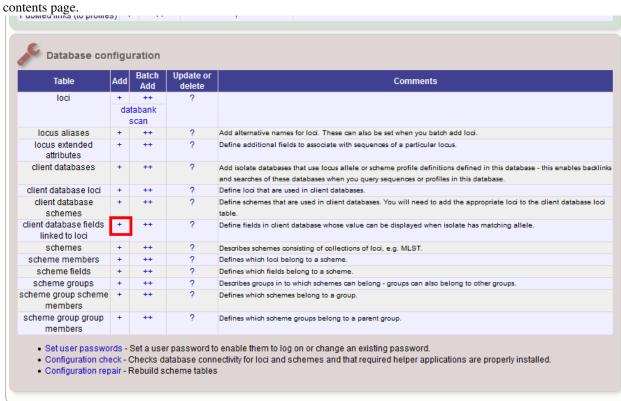
Setting this up is identical to setting up for alleles (see above) except you click on the add (+) client database schemes link and choose the scheme and client databases in the dropdown list boxes.

Now when information on a given scheme profile (e.g. MLST sequence type) is shown following a query, the software will list the number of isolates with that profile and link to a search on the database to retrieve these.

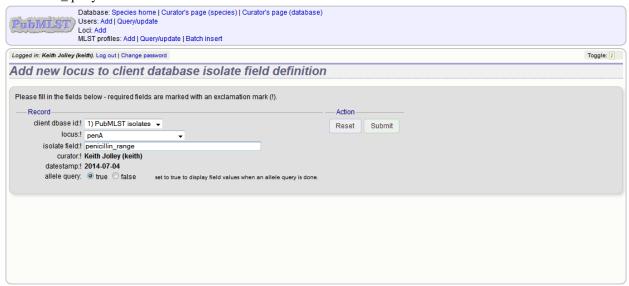


5.17.3 Look up specific isolate database fields linked to a given allele

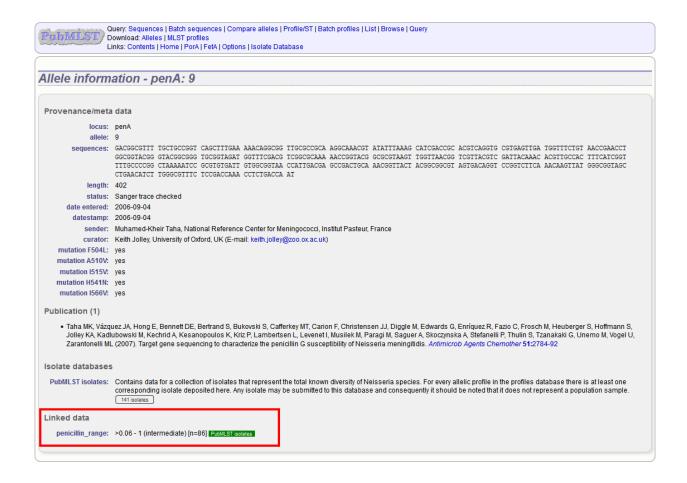
To link an allele to an isolate field, click the add (+) 'client database fields linked to loci' link on the curator's interface contents page



Select the client database and locus from the dropdown lists and enter the isolate database field that you'd like to link. The 'allele_query' field should be set to true.



Now, in the allele record or following a sequence query that identifies an allele, all values for the chosen field from isolates with the corresponding allele are shown.



5.18 Rule-based sequence queries

The RuleQuery plugin has been designed to extract information from a pasted-in genome sequence, look up scheme fields and client database fields, and then format the output in a specified manner.

Rules are written in Perl, allowing the full power of this scripting language to be utilised. Helper functions that perform specific actions are available to the script (see example).

Please note that direct access to the database is prevented as are system calls.

5.18.1 Example rule code

An example can be found on the Neisseria sequence database that takes a genome sequence and determines a fine type and antibiotic resistance.

The code for this rule is as follows:

```
#Clinical identification rule

#Update job viewer status
update_status({stage=>'Scanning MLST loci'});

#Scan genome against all scheme 1 (MLST) loci
scan_scheme(1);
```

```
#Update job viewer status
update_status({percent_complete=>30, stage=>'Scanning PorA and FetA VRs'});
#Scan genome against the PorA VR and FetA VR loci
scan_locus($_) foreach qw(PorA_VR1 PorA_VR2 FetA_VR);
Add text to main output
append_html("<h1>Strain type</h1>");
#Set variables for the scanned results. These can be found in the
#$results->{'locus'} hashref
my %alleles;
$alleles($_) = $results->{'locus'}->{$_} // 'ND' foreach gw(PorA_VR1 PorA_VR2);
$alleles{'FetA_VR'} = $results->{'locus'}->{'FetA_VR'} // 'F-ND';
#Scheme field values are automatically determined if a complete
#profile is available. These are stored in the $results->{'scheme'} hashref
my $st = $results->{'scheme'}->{1}->{'ST'} // 'ND';
append_html("P1.$alleles{'PorA_VR1'}, $alleles{'PorA_VR2'}; $alleles{'FetA_VR'}; ST-$st ");
#Reformat clonal complex using a regular expression, e.g.
#'ST-11 clonal complex/ET-37 complex' gets rewritten to 'cc11'
my $cc = $results->{'scheme'}->{1}->{'clonal_complex'} // '-';
c = s/ST-(S+) complex.*/cc$1/;
append_html("($cc)");
if ($st eq 'ND') {
 append html("ST not defined. If individual MLST loci have been found "
  . "they will be displayed below:");
 #The get_scheme_html function automatically formats output for a scheme.
 #Select whether to display in a table rather than a list, list all loci, and/or list fields.
 append_html(get_scheme_html(1, {table=>1, loci=>1, fields=>0}));
#Antibiotic resistance
update_status({percent_complete=>80, stage=>'Scanning penA and rpoB'});
scan_locus($_) foreach qw(penA rpoB);
if (defined $results->{'locus'}->{'penA'} || defined $results->{'locus'}->{'rpoB'} ){
 append html("<h1>Antibiotic resistance</h1>");
 if (defined $results->{'locus'}->{'penA'}) {
   append_html("<i>penA</i> allele: $results->{'locus'}->{'penA'}");
    #If a client isolate database has been defined and values have been defined in
    #the client_dbase_loci_fields table, the values for a field in the isolate database can be
    #retrieved based on isolates that have a particular allele designated.
    #The min_percentage attribute states that only values that are represented by at least that
    #proportion of all isolates that had a value set are returned (null values are ignored).
   my $range = get_client_field(1,'penA','penicillin_range',{min_percentage => 75});
   append_html(" (penicillin MIC: $range->[0]->{'penicillin_range'})") if @$range;
   append_html("");
 if (defined $results->{'locus'}->{'rpoB'}){
   append_html("<i>rpoB</i> allele: $results->{'locus'}->{'rpoB'}");
   my $range = get_client_field(1,'rpoB','rifampicin_range',{min_percentage => 75});
   append_html(" (rifampicin MIC: $range->[0]->{'rifampicin_range'})") if @$range;
    append_html("");
```

```
append_html("");
}
```

Rule files

The rule file is placed in a rules directory within the database configuration directory, e.g. /etc/bigsdb/dbase/pubmlst_neisseri_seqdef/rules. Rule files are suffixed with '.rule' and their name should be descriptive since it is used within the interface, i.e. the above rule file is named Clinical_identification.rule (underscores are converted to spaces in the web interface).

Linking to the rule query

Links to the rule query are not automatically placed within the web interface. The above rule query can be called using the following URL:

http://pubmlst.org/perl/bigsdb/bigsdb.pl?db=pubmlst_neisseria_seqdef&page=plugin&name=RuleQuery&ruleset=Clinical_identification

To place a link to this within the database contents page an HTML file called job_query.html can be placed in a contents directory within the database configuration directory, e.g. in /etc/bigsdb/dbases/pubmlst_neisseria_seqdef/contents/job_query.html. This file should contain a list entry (i.e. surrounded with and

Adding descriptive text

Descriptive text for the rule, which will appear on the rule query page, can be placed in a file called description.html in a directory with the same name as the rule within the rule directory, e.g. in /etc/bigsdb/dbases/pubmlst neisseria segdef/rules/Clinical identification/description.html.

5.19 Workflow for setting up a MLST scheme

The workflow for setting up a MLST scheme is as follows (the example seqdef database is called seqdef_db):

Seqdef database

- 1. Create appropriate loci
- 2. Create new scheme 'MLST'
- 3. Add scheme_field 'ST' with primary_key=TRUE (add clonal_complex if you want; set this with primary_key=FALSE)
- 4. Add each locus as a scheme member
- 5. You'll then be able to add profiles

Isolate database

- 1. Create the same loci with the following additional parameters (example locus 'atpD')
- dbase name: segdef db
- dbase table: sequences
- · dbase_id_field: allele_id
- dbase_id2_field: locus

- dbase_id_value: atpD
- dbase_seq_field: sequence
- url: something like /cgi-bin/bigsdb/bigsdb.pl?db=seqdef_db&page=alleleInfo&locus=atpD&allele_id=[?]
- 2. Create scheme 'MLST' with:
- · dbase name: segdef db
- dbase_table: mv_scheme_1 (or whatever the id of your seqdef scheme is)
- 3. Add scheme_field ST as before
- 4. Add loci as scheme_members

5.20 Automated assignment of scheme profiles

It is not practical to define cgMLST profiles via the web interface. A script is provided in the scripts/automation directory of the BIGSdb package called define_profiles.pl that can be used to scan an isolate database and automatically define cgMLST profiles in the corresponding sequence definition database.

The script is run as follows:

```
define_profiles.pl --database <name> --scheme <scheme_id>
```

A full list of options can be found by typing:

```
define_profiles.pl --help
NAME
    define_profiles.pl - Define scheme profiles found in isolate database
   define_profiles.pl --database NAME --scheme SCHEME_ID [options]
OPTIONS
--cache
   Update scheme field cache in isolate database.
--database NAME
   Database configuration name.
--help
   This help page.
--exclude_isolates LIST
   Comma-separated list of isolate ids to ignore.
--exclude_projects LIST
   Comma-separated list of projects whose isolates will be excluded.
--ignore_multiple_hits
   Set allele designation to 'N' if there are multiple designations set for
   a locus. The default is to use the lowest allele value in the profile
   definition.
--isolates LIST
    Comma-separated list of isolate ids to scan (ignored if -p used).
```

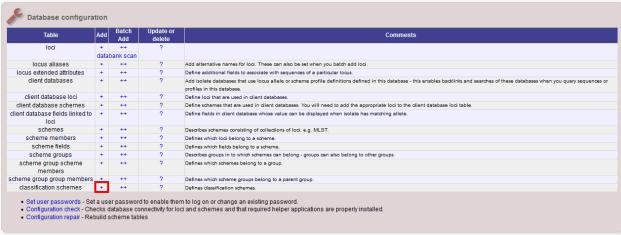
```
--isolate_list_file FILE
   File containing list of isolate ids (ignored if -i or -p used).
--max ID
   Maximum isolate id.
--min TD
   Minimum isolate id.
--min_size SIZE
   Minimum size of seqbin (bp) - limit search to isolates with at least this
--missing NUMBER
   Set the number of loci that are allowed to be missing in the profile. If
   the remote scheme does not allow missing loci then this number will be set
   to 0. Default=0.
--projects LIST
   Comma-separated list of project isolates to scan.
--scheme SCHEME ID
   Scheme id number.
```

5.21 Scheme profile clustering - setting up classification schemes

Classification groups are a way to cluster scheme profiles using a specified threshold of pairwise allelic mismatches. Any number of different classification schemes can sit on top of a standard scheme (such as cgMLST), allowing different similarity thresholds to be pre-determined. Currently, single-linkage clustering is supported whereby each member of a group must have no more than the specified number of allelic differences with at least one other member of the group.

5.21.1 Defining classification scheme in sequence definition database

Once a scheme has been defined, add a classification scheme by clicking the add classification schemes (+) link on the curator's interface contents page.



Select the underlying scheme and enter a name for the classification scheme, the number of mismatches allowed

in order to include a scheme profile in a group, and a description. An example name for such a scheme could be 'Nm_cgc_25' indicating that this is a classification scheme for *Neisseria meningitidis* core genome cluster with a threshold of 25 mismatches.

You can additionally choose whether a relative threshold is used to calculate the number of mismatches to account for missing loci in pairwise comparisons. In this case, in order to be grouped, the number of matching alleles must exceed:

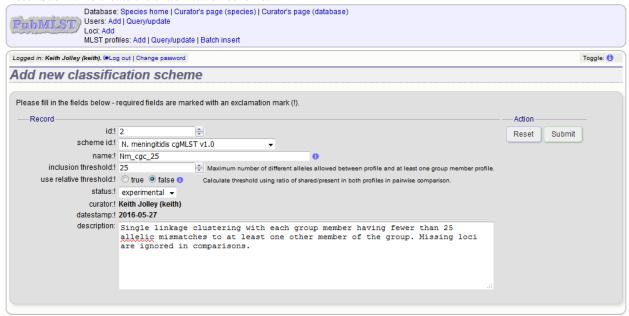
```
(number of common loci x (total loci - defined threshold)) / total loci
rather than
total loci - defined threshold
```

when an absolute threshold is used.

As this threshold has to be calculated for each pairwise comparison, clustering using relative thresholds is slower than using an absolute value, and probably makes little real world difference.

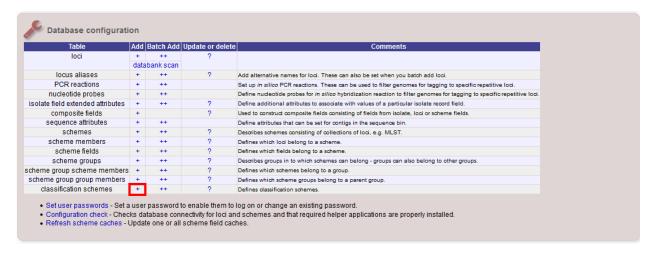
The status can be 'experimental' or 'stable'. The status of a scheme will be shown in the web interface to indicate that any groupings are subject to change and do not form part of the stable nomenclature.

Press 'Submit' to create the classification scheme.



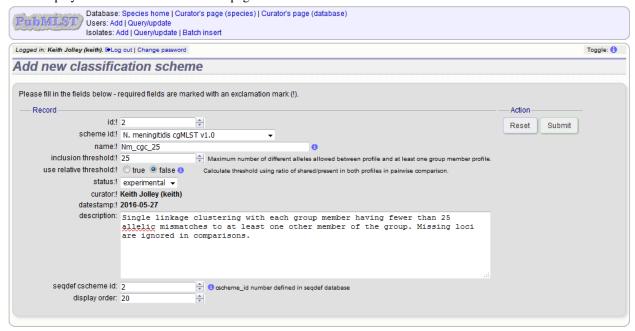
5.21.2 Defining classification scheme in isolate database

Duplicate the scheme definition from the sequence definition database. Click the add classification schemes (+) link on the curator's interface contents page.



Enter the same details used in the sequence definition database. If a different id number is used in the isolate and sequence definition databases, you can set the seqdef id in the seqdef_cscheme_id field (the default is to use the same id).

You can also define a display order - this is an integer field on which the ordering of classification schemes is sorted when displayed in the isolate information page.



It is a good idea to check the configuration.

5.21.3 Clustering

Clustering is performed using the cluster.pl script found in the scripts/automation directory of the BIGSdb package. It should be run by the bigsdb user account (or any account with access to the databases).

Currently only single-linkage clustering is supported.

The script is run as follows from the command line:

cluster.pl --database <database configuration> --cscheme <classification scheme id>

A full list of options can be found by typing:

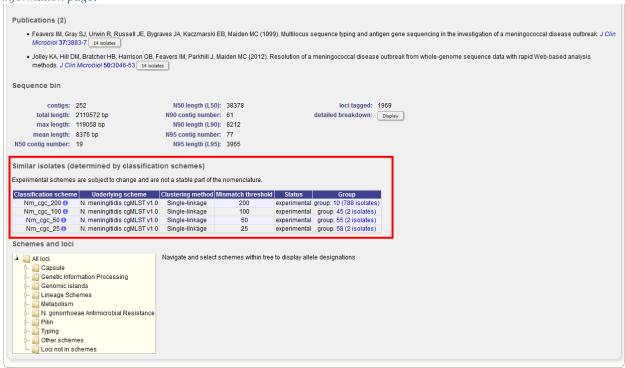
```
cluster.pl --help
NAME
    cluster.pl - Cluster cgMLST profiles using classification groups.

SYNOPSIS
    cluster.pl --database NAME --cscheme_id SCHEME_ID [options]

OPTIONS
--cscheme CLASSIFICATION_SCHEME_ID
    Classification scheme id number.
--database NAME
    Database configuration name.
--help
    This help page.
--reset
    Remove all groups and profiles currently defined for classification group.
```

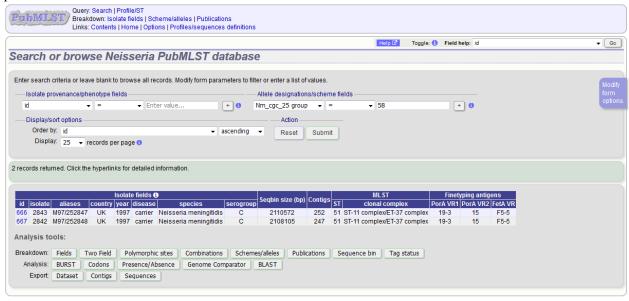
Note: Note that for classification schemes to be accessible within the isolate database, *scheme cache tables* must be generated and kept up-to-date.

Where an isolate has been clustered in to a group with other isolates, this information is available in the *isolate* information page.



Clicking the hyperlinks will take you to a table containing matching isolates, from where standard analyses can be

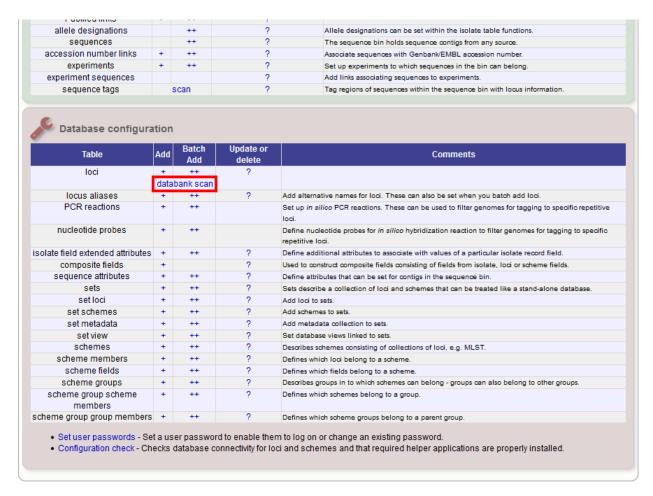
performed.



5.22 Defining new loci based on annotated reference genome

An annotated reference genome can be used as the basis of defining loci. The 'Databank scan' function will create an upload table suitable for pasting directly in to the batch locus add form of the *sequence definition* or *isolate* databases.

Click 'Database scan' on the curator's contents pag.



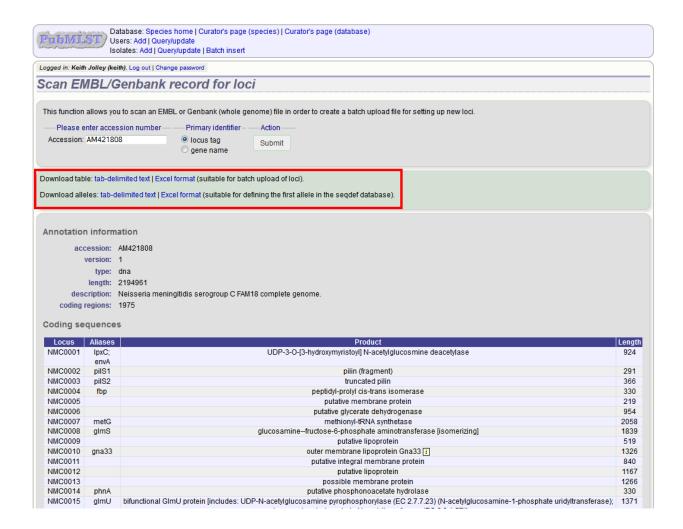
Enter an EMBL or Genbank accession number for a complete annotated genome and press 'Submit'.



A table of loci will be generated provided a valid accession number is provided.



Tab-delimited text and Excel format files will be created to be used as the basis for upload files for the sequence definition and isolate databases. Batch sequence files, in text and Excel formats, are also created for defining the first allele once the locus has been set up in the sequence definition database.



5.23 Genome filtering

Within a genome there may be multiple loci that share allele pools. If an allele sequence is tagged from a genome using only BLAST then there is no way to determine which locus has been identified. It is, however, possible to further define loci by their context, i.e. surrounding sequence.

5.23.1 Filtering by in silico PCR

Provided a locus can be predicted to be specifically amplified by a PCR reaction, the genome can be filtered to only look at regions prediced to fall within amplification products of one or more PCR reactions. Since this is *in silico* we don't need to worry about problems such as sequence secondary structure and primers can be any length.

To define a PCR reaction that can be linked to a locus definition, click the add (+) PCR reaction link on the curator's main page.

Locus 1 and locus 2 share allele pool

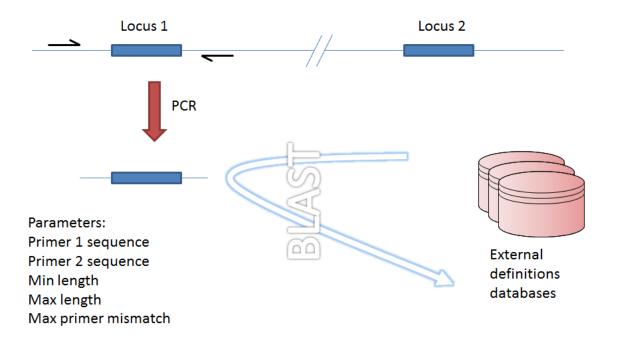
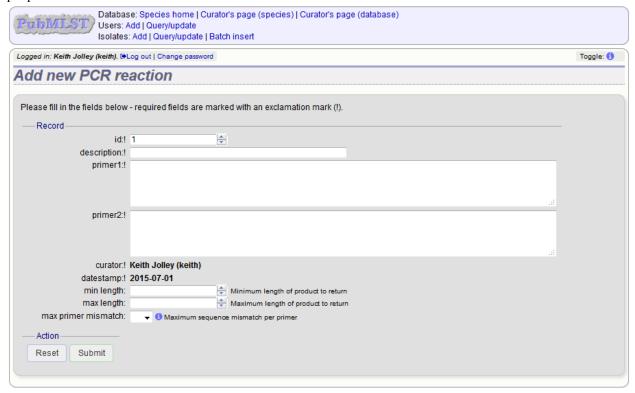


Fig. 5.1: Genome filtering by in silico PCR.

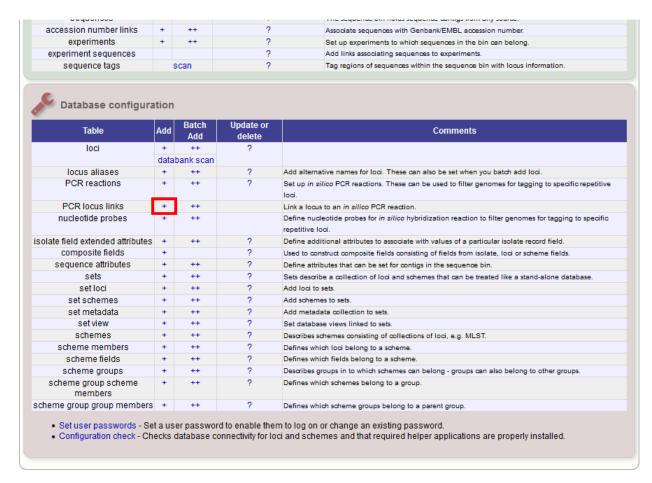
allele designations		++	?	Allele designations can be set within the isolate table functions.
sequences		++	?	The sequence bin holds sequence contigs from any source.
accession number links	+	++	?	Associate sequences with Genbank/EMBL accession number.
experiments	+	++	?	Set up experiments to which sequences in the bin can belong.
experiment sequences			?	Add links associating sequences to experiments.
sequence tags	8	scan	?	Tag regions of sequences within the sequence bin with locus information.
Database configura	ition Add	Batch	Update or	Comments
		Add	delete	
loci	+	++	?	
		bank scan		
locus aliases	+	++	?	Add alternative names for loci. These can also be set when you batch add loci.
PCR reactions	+	++		Set up in silico PCR reactions. These can be used to filter genomes for tagging to specific repetitive loci.
nucleotide probes	+	++		Define nucleotide probes for <i>in silico</i> hybridization reaction to filter genomes for tagging to specific repetitive loci.
solate field extended attributes	+	++	?	Define additional attributes to associate with values of a particular isolate record field.
composite fields	+		?	Used to construct composite fields consisting of fields from isolate, loci or scheme fields.
sequence attributes	+	++	?	Define attributes that can be set for contigs in the sequence bin.
sets	+	++	?	Sets describe a collection of loci and schemes that can be treated like a stand-alone database.
set loci	+	++	?	Add loci to sets.
set schemes	+	++	?	Add schemes to sets.
set metadata	+	++	?	Add metadata collection to sets.
set view	+	++	?	Set database views linked to sets.
schemes	+	++	?	Describes schemes consisting of collections of loci, e.g. MLST.
scheme members	+	++	?	Defines which loci belong to a scheme.
scheme fields	+	++	?	Defines which fields belong to a scheme.
scheme groups	+	++	?	Describes groups in to which schemes can belong - groups can also belong to other groups.
scheme group scheme members	+	++	?	Defines which schemes belong to a group.
cheme group group members	+	++	?	Defines which scheme groups belong to a parent group.
				m to log on or change an existing password. ci and schemes and that required helper applications are properly installed.

In the resulting web form you can enter values for your two primer sequences (which can be any length), the minimum and maximum lengths of reaction products you wish to consider and a value for the allowed number of mismatches per primer.



- id PCR reaction identifier number.
 - Allowed: integer.
- description Description of PCR reaction product.
 - Allowed: any text.
- primer1 Primer 1 sequences
 - Allowed: nucleotide sequence (IUPAC ambiguous characters allowed).
- primer2 Primer 2 sequence.
 - Allowed: nucleotide sequence (IUPAC ambiguous characters allowed).
- min_length Minimum length of predicted PCR product.
 - Allowed: integer.
- max_length Maximum length of predicted PCR product.
- max_primer_mismatch Number of mismatches allowed in primer sequence.
 - Allowed: integer.
 - Do not set this too high or the simulation will run slowly.

Associating this with a particular locus is a two step process. First, create a locus link by clicking the add (+) PCR locus link on the curator's main page. This link will only appear once a PCR reaction has been defined.



Select the locus and PCR reaction name from the dropdown lists to create the link. You also need to edit the locus table and set the pcr filter field to 'true'.

Now when you next perform tag scanning there will be an option to use PCR filtering.

5.23.2 Filtering by in silico hybridization

An alternative is to define a locus by proximity to a single probe sequence. This is especially useful if you have multiple contigs and the locus in question may be at the end of a contig so that it doesn't have upstream or downstream sequence available for PCR filtering.

The process is very similar to setting up PCR filtering, but this time click the nucleotide probe link on the curator's content page.

Enter the nucleotide sequence and a name for the probe. Next you need to link this to the locus in question. Click the add (+) probe locus links link on the curator's main page. This link will only appear once a probe has been defined.

Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:

- probe_id Dropdown list of probe names.
 - Allowed: selection from list.
- locus Dropdown list of loci.
 - Allowed: selection from list.
- max_distance Minimum distance of probe from end of locus.

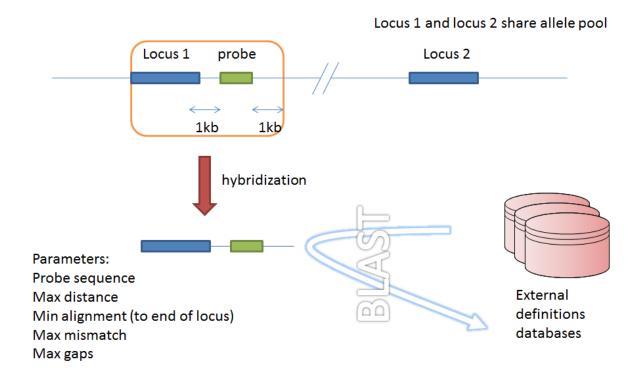


Fig. 5.2: Filtering by in silico hybridization

- Allowed: any positive integer.
- min_alignment Minimum length of alignment allowed.
 - Allowed: any positive integer.
- max_mismatch Maximum number of mismatches allowed in alignment.
 - Allowed: any positive integer.
- max_gaps Maximum number of gaps allowed in alignment.
 - Allowed: any positive integer.

Finally edit the locus table and set the probe_filter field for the specified locus to 'true'.

Now when you next perform tag scanning there will be an option to use probe hybridization filtering.

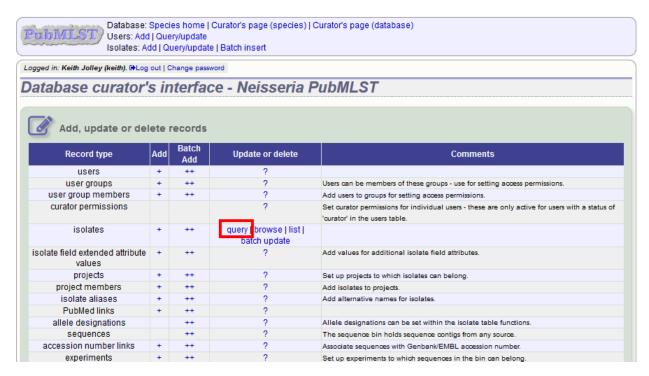
5.24 Setting locus genome positions

The genome position for a locus can be set directly by editing the locus record. To batch update multiple loci based on a tagged genome, however, a much easier way is possible. For this method to work, the reference genome must be represented by a single contig.

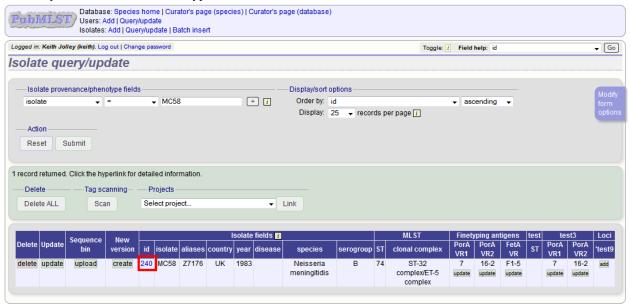
From the curator's main page, you need to do a query to find the isolate that you will base your numbering on. Click 'isolate query' to take you to a standard query form.

+	++	?	Associate sequences with Genbank/EMBL accession number.
+	++	?	Set up experiments to which sequences in the bin can belong.
		?	Add links associating sequences to experiments.
sequence tags scan ?		?	Tag regions of sequences within the sequence bin with locus information.
tion			
Add	Batch Add	Update or delete	Comments
+ ++		?	
databank scan			
+	++	?	Add alternative names for loci. These can also be set when you batch add loci.
+	++	?	Set up in silico PCR reactions. These can be used to filter genomes for tagging to specific repetitive loci.
+	++		Link a locus to an in silico PCR reaction.
+	++		Define nucleotide probes for in silico hybridization reaction to filter genomes for tagging to specific repetitive loci.
+	++	?	Define additional attributes to associate with values of a particular isolate record field.
+		?	Used to construct composite fields consisting of fields from isolate, loci or scheme fields.
+	++	?	Define attributes that can be set for contigs in the sequence bin.
+	++	?	Sets describe a collection of loci and schemes that can be treated like a stand-alone database.
+	++	?	Add loci to sets.
+	++	?	Add schemes to sets.
+	++	?	Add metadata collection to sets.
+	++	?	Set database views linked to sets.
+	++	?	Describes schemes consisting of collections of loci, e.g. MLST.
+	++	?	Defines which loci belong to a scheme.
+	++	?	Defines which fields belong to a scheme.
+	++	?	Describes groups in to which schemes can belong - groups can also belong to other groups.
+	++	?	Defines which schemes belong to a group.
+	++	2	Defines which scheme groups belong to a parent group.
	+ s s s s s s s s s s s s s s s s s s s	ion Scan ion Add Batch Add + ++ databank scan + ++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++	+ ++ ++ ? scan ? scan ? ion Add Batch Add H A

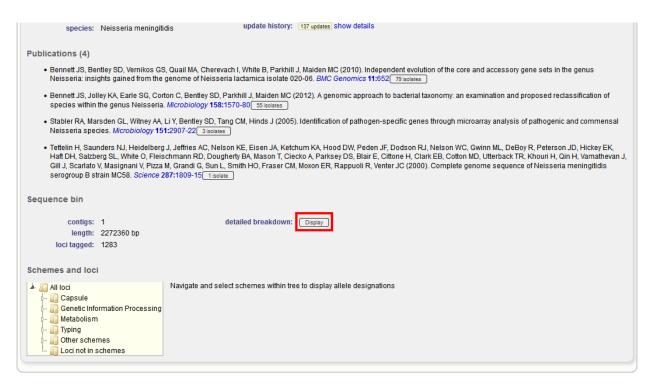
experiment sequences			?	Add links associating sequences to experiments.
sequence tags	scan		?	Tag regions of sequences within the sequence bin with locus information.
Database configura	ation			
		Deteb	Hedete es	
Table	Add	Batch Add	Update or delete	Comments
loci	+	++	?	
	databank scan			
locus aliases	+	++	?	Add alternative names for loci. These can also be set when you batch add loci.
PCR reactions	+	++	?	Set up in silico PCR reactions. These can be used to filter genomes for tagging to specific repetitive loci.
PCR locus links	+	++		Link a locus to an in silico PCR reaction.
nucleotide probes	+	++	?	Define nucleotide probes for <i>in silico</i> hybridization reaction to filter genomes for tagging to specific repetitive loci.
probe locus links	+	++		Link a locus to an in silico hybridization reaction.
olate field extended attributes	+	++	?	Define additional attributes to associate with values of a particular isolate record field.
composite fields	+		?	Used to construct composite fields consisting of fields from isolate, loci or scheme fields.
sequence attributes	+	++	?	Define attributes that can be set for contigs in the sequence bin.
sets	+	++	?	Sets describe a collection of loci and schemes that can be treated like a stand-alone database.
set loci	+	++	?	Add loci to sets.
set schemes	+	++	?	Add schemes to sets.
set metadata	+	++	?	Add metadata collection to sets.
setview	+	++	?	Set database views linked to sets.
schemes	+	++	?	Describes schemes consisting of collections of loci, e.g. MLST.
scheme members	+	++	?	Defines which loci belong to a scheme.
scheme fields	+	++	?	Defines which fields belong to a scheme.
scheme groups	+	++	?	Describes groups in to which schemes can belong - groups can also belong to other groups.
scheme group scheme members	+	++	?	Defines which schemes belong to a group.
heme group group members	+	++	?	Defines which scheme groups belong to a parent group.
•		•		m to log on or change an existing password. ci and schemes and that required helper applications are properly installed.



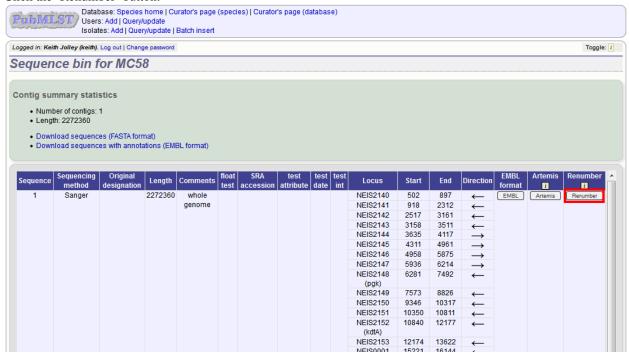
Perform your search and click the hyperlinked id number of the record.



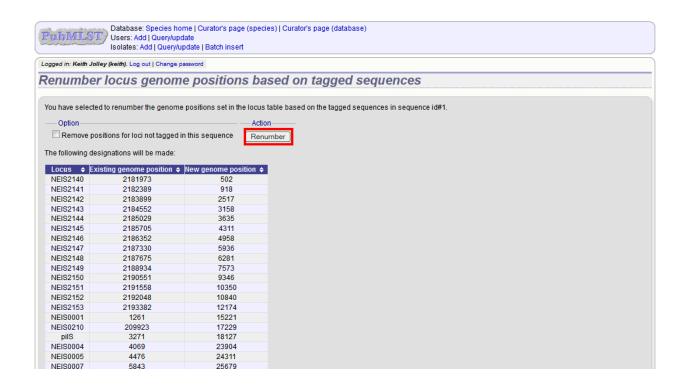
In the isolate record, click the sequence bin 'Display' button to bring up details of the isolate contigs.



Click the 'Renumber' button:



A final confirmation screen is displayed with the option to remove existing numbering that doesn't appear within the reference genome. Click 'Renumber'.



5.25 Defining composite fields

Composite fields are virtual fields that don't themselves exist within the database but are made up of values retrieved from other fields or schemes and formatted in a particular way. They are used for display and analysis purposes only and can not be searched against.

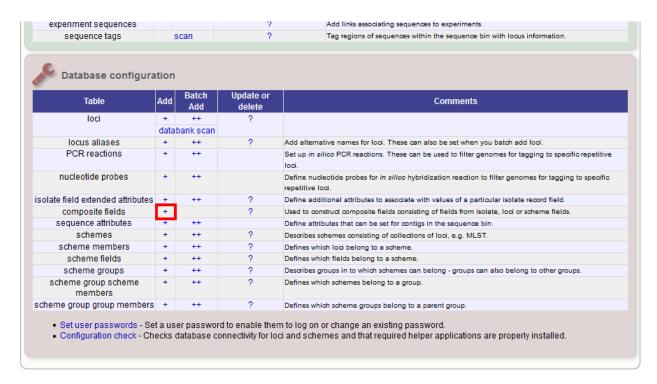
One example of a composite field is used in the Neisseria PubMLST database which has a strain designation composite field made up of serogroup, PorA VR1 and VR2, FetA VR, ST and clonal complex designations in the format:

[serogroup]: P1.[PorA_VR1],[PorA_VR2]: [FetA_VR]: ST-[ST] ([clonal_complex])

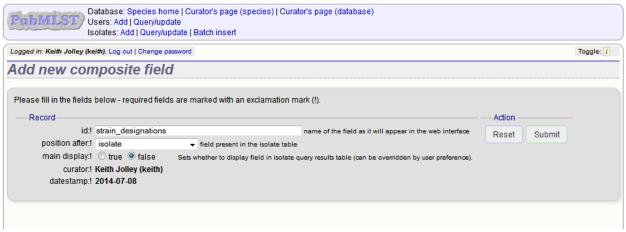
e.g. A: P1.5-2,10: F1-5: ST-4 (cc4)

Additionally, the clonal complex field in the above example is converted using a regular expression from 'ST-4 complex/subgroup IV' to 'cc4'.

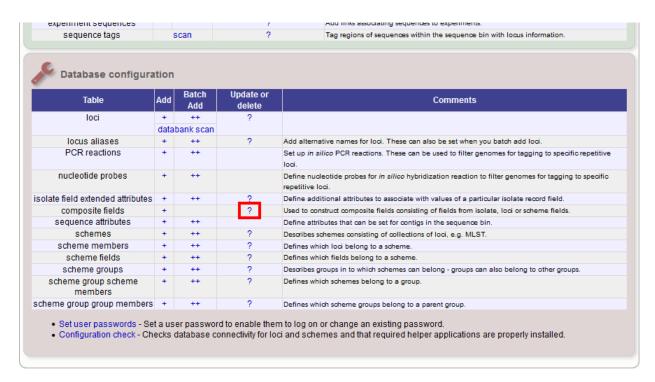
Composite fields can be added to the database by clicking the add (+) composite fields link on the curator's main page.



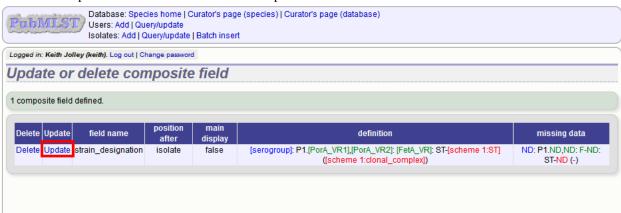
Initially you just enter a name for the composite field and after which field it should be positioned. You can also set whether or not it should be displayed by default in main results tables following a query - this is overrideable by user preferences.



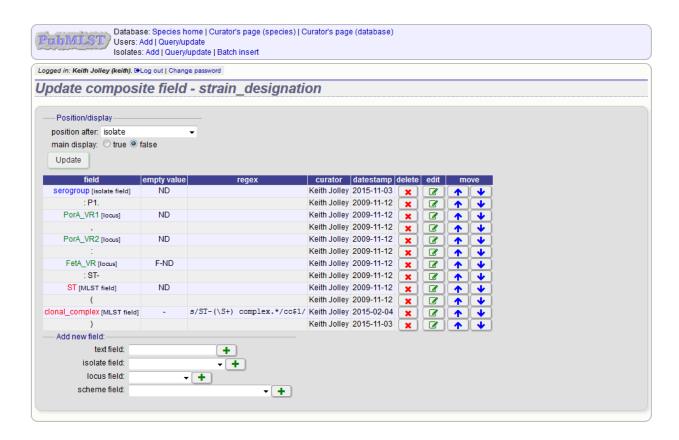
Once the field has been created it needs to be defined. This can be done from query composite field link on the main curator's page.



Select the composite field from the list and click 'Update'.



From this page you can build up your composite field from snippets of text, isolate field, locus and scheme field values. Enter new values in the boxes at the bottom of the page.



Once a field has been added to the composite field, it can be edited by clicking the 'edit' button next to it to add a regular expression to modify its value by specific rules, e.g. in the clonal complex field above, the regular expression is set as:

```
s/ST-(\S+) complex.*/cc$1/
```

which extracts one or more non-space characters following the 'ST-' in a string that then contains the work 'complex', and appends this to 'cc' to produce the final string.

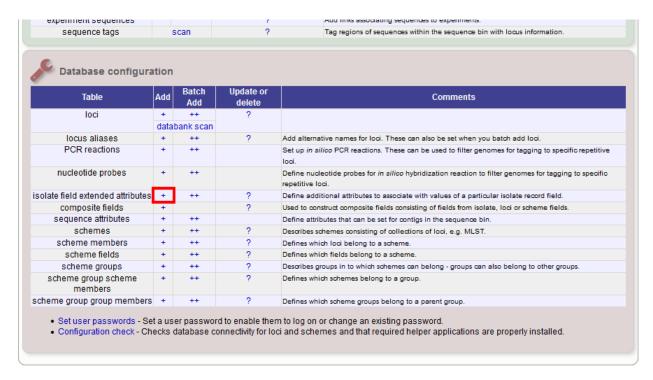
This will convert 'ST-4 complex/subgroup IV' to 'cc4'.

You can also define text to be used for when the field value is missing, e.g. 'ND'.

5.26 Extended provenance attributes (lookup tables)

Lookup tables can be associated with an isolate database field such that the database can be queried by extended attributes. An example of this is the relationship between continent and country - every country belongs to a continent but you wouldn't want to store the continent with each isolate record (not only could data be entered inconsistently but it's redundant). Instead, each record may have a country field and the continent is then determined from the lookup table, allowing, for example, a search of isolates limited to those from Europe.

To set up such an extended attribute, click the add (+) isolate field extended attributes link on the curator's main page.

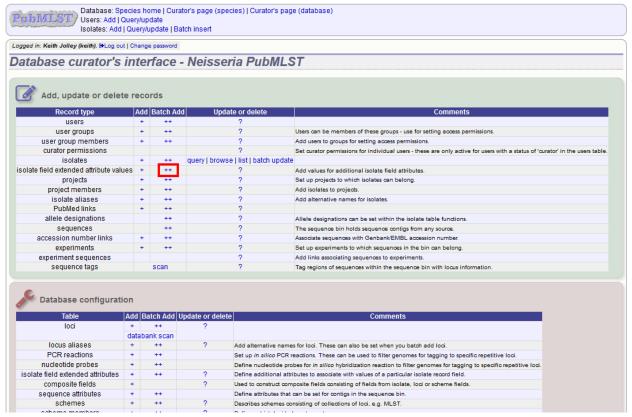


Fill in the web form with appropriate values. Required fields have an exclamation mark (!) next to them:

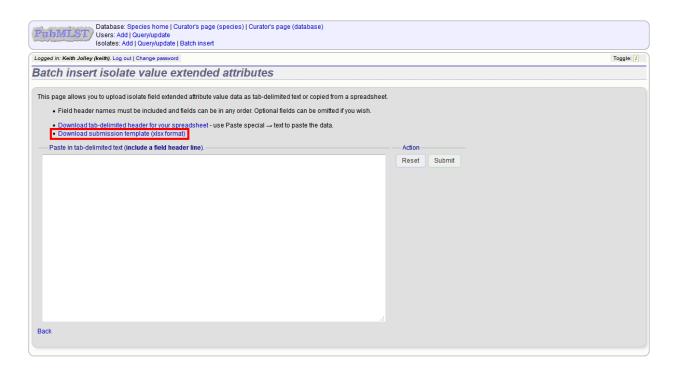
- isolate_field Dropdown list of isolate fields.
 - Allowed: selection from list.
- attribute Name of extended attribute, e.g. continent.
 - Allowed: any text (no spaces).
- value_format Format for values.
 - Allowed: integer/float/text/date.
- value_regex Regular expression to enforce allele id naming.
 - ^: the beginning of the string
 - \$:the end of the string
 - d: digit
 - D: non-digit
 - s: white space character
 - S: non white space character
 - w: alpha-numeric plus '_'
 - .: any character
 - *: 0 or more of previous character
 - +: 1 or more of previous character
 - e.g. ^Fd-d+\$ states that a value must begin with a F followed by a single digit, then a dash, then one or more digits, e.g. F1-12
- description Long description this isn't currently used but may be in the future.
 - Allowed: any text.

- url URL used to hyperlink values in the isolate information page. Instances of [?] within the URL will be substituted with the value.
 - Allowed: any valid URL (either relative or absolute).
- length Maximum length of extended attribute value.
 - Allowed: any positive integer.
- field_order Integer that sets the order of the field following it's parent isolate field.
 - Allowed: any integer.

The easiest way to populate the lookup table is to do a batch update copied from a spreadsheet. Click the batch add (++) isolate field extended attribute values link on the curator's main page (this link will only appear once an extended attribute has been defined).



Download the Excel template:



Fill in the columns with your values, e.g.

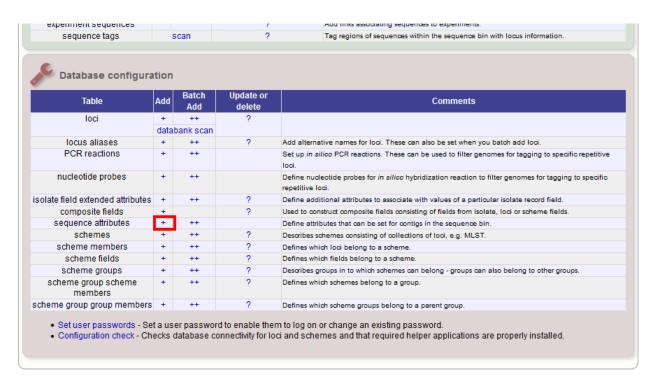
isolate_field	attribute	field_value	value
country	continent	Afghanistan	Asia
country	continent	Albania	Europe
country	continent	Algeria	Africa
country	continent	Andorra	Europe
country	continent	Angola	Africa

Paste from the spreadsheet in to the upload form and click 'Submit'.

5.27 Sequence bin attributes

It is possible that you will want to store extended attributes for sequence bin contigs when you upload them. Examples may be read length, assembler version, etc. Since there are almost infinite possibilities for these fields, and they are likely to change over time, they are not hard-coded within the database. An administrator can, however, create their own attributes for a specific database and these will then be available in the web form when uploading new contig data. The attributes are also searchable.

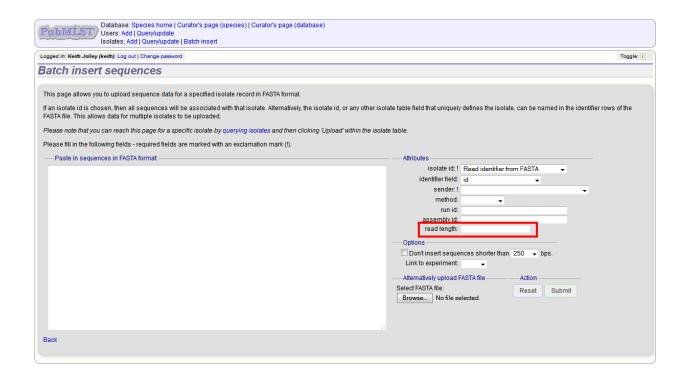
To set up new attributes, click the add (+) 'sequence attributes' link on the isolate database curator's index page.



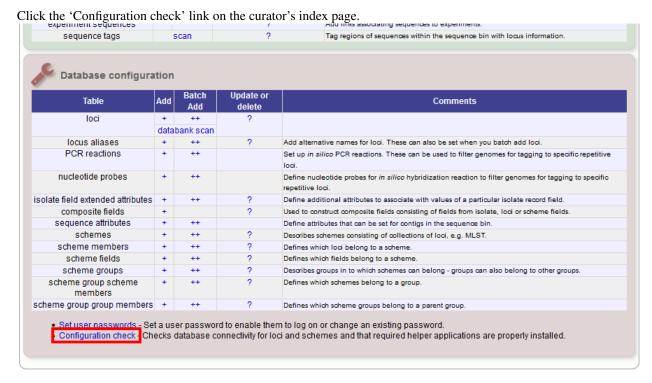
Enter the name of the attribute as the 'key', select the type of data (text, integer, float, date) and an optional short description. Click 'Submit'.



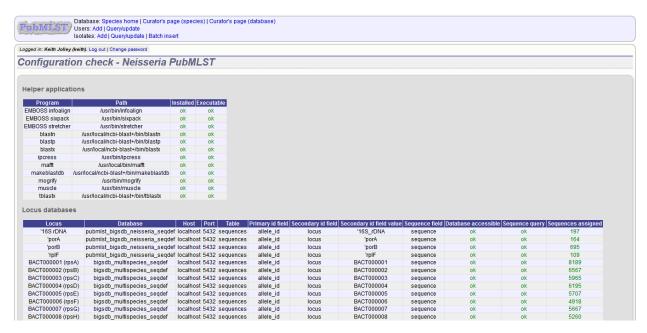
This new attribute will then be available when *uploading contig data*.



5.28 Checking external database configuration settings



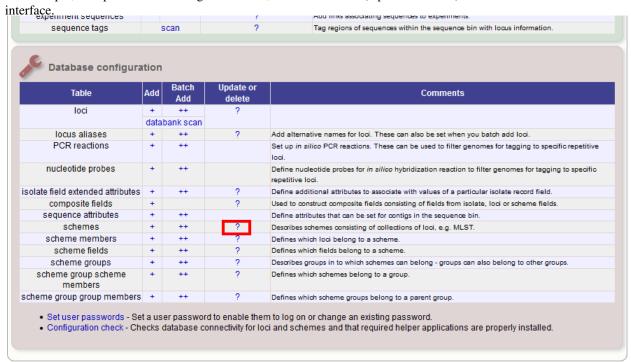
The software will check that required helper applications are installed and executable and, in isolate databases, test every locus and scheme external database to check for connectivity and that data can be retrieved.



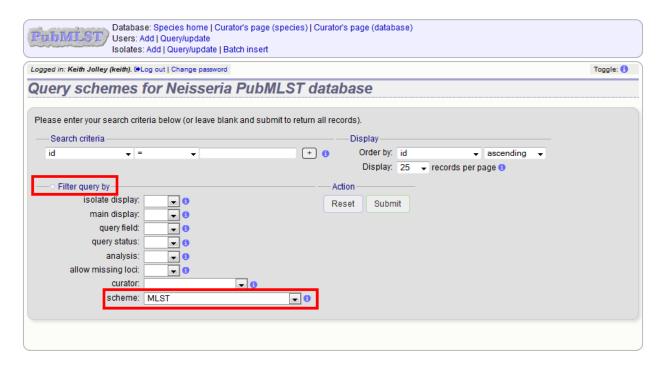
Any problems will be highlighted with a red X.

5.29 Exporting table configurations

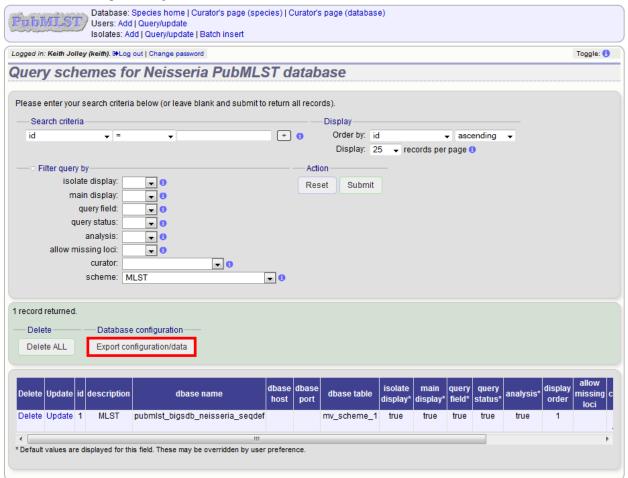
Sometimes it is useful to transfer configurations between different databases or to export a configuration for troubleshooting. Data from most of the tables can be exported in tab-delimited text format suitable for batch uploading. For example, to export scheme configuration data, click the '?' link (Update or delete) next to schemes in the curator's



Expand the filters and select the required scheme in the dropdown box, then press submit.



Click the button 'Export configuration/data'.



The three tables that are used to define a scheme (schemes, scheme_members and scheme_fields) are displayed in a format suitable for copy and pasting.

```
schemes
id description dbase_name dbase_host dbase_port dbase_user dbase_password dbase_table isolate_di:
1 MLST pubmlst_bigsdb_neisseria_seqdef
                                                     mv_scheme_1 1 1 1 1 1 2 2012-03-22
scheme_members
scheme_id locus profile_name field_order curator datestamp
1 abcZ 1 2 2009-11-12
1 adk
          2 2 2009-11-12
1 aroE
          3 2 2009-11-12
          4 2 2009-11-12
1 fumC
          5 2 2009-11-12
1 adh
1 pdhC 6 2 2009-11-12
1 pgm 7 2 2009-11-12
scheme_fields
scheme_id field type primary_key description field_order url isolate_display main_display
1 ST integer 1 1 /cgi-bin/bigsdb/bigsdb.pl?page=profileInfo&db=pubmlst_neisseria_seqdef&scheme
1 \quad {\tt clonal\_complex \ text} \quad 0 \qquad \quad 2 \qquad \quad 1 \quad 1 \quad 1 \quad 2 \quad 2009 \hbox{--} 11 \hbox{--} 16
```

5.30 Authorizing third-party client software to access authenticated resources

If you are running the *RESTful API*, you will need to specifically authorize client software to connect to authenticated resources. This involves creating a client key and a client secret that is used to sign requests coming from the application. The client key and secret should be provided to the application developer.

There is a script to do this in the scripts/maintenace directory of the download archive. The script is called create_client_credentials and should be run by the postgres user. A full list of options can be found by typing:

-i, --insert

Add credentials to authentication database. This will fail if a matching application version already exists (use --update in this case to overwrite existing credentials).

-u, --update

Update exisitng credentials in the authentication database.

-v, --version VERSION

Version of application (optional).

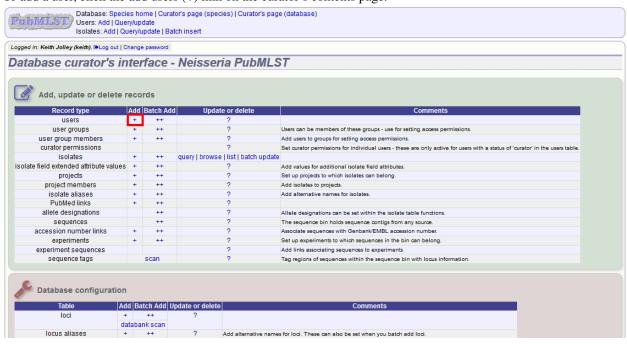
Curator's guide

Please note that links displayed within the curation interface will vary depending on database contents and the permissions of the curator.

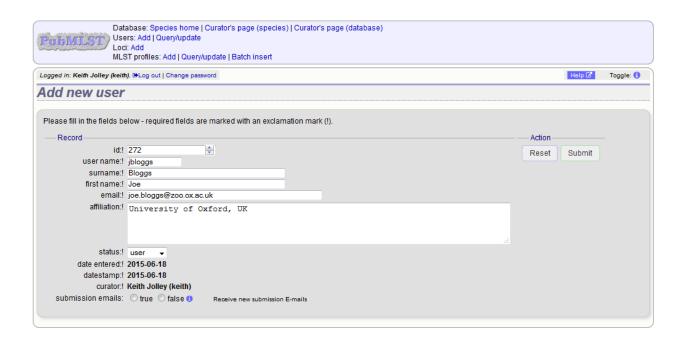
6.1 Adding new sender details

All records within the databases are associated with a sender. Whenever somebody new submits data, they should be added to the users table so that their name appears in the dropdown lists on the data upload forms.

To add a user, click the add users (+) link on the curator's contents page.



Enter the user's details in to the form.



Normally the status should be set as 'user'. Only admins and curators with special permissions can create users with a status of curator or admin.

If the submission system is in operation there will be an option at the bottom called 'submission_emails'. This is to enable users with a status of 'curator' or 'admin' to receive E-mails on receipt of new submissions. It is not relevant for users with a status of 'user' or 'submitter'.

6.2 Adding new allele sequence definitions

6.2.1 Single allele

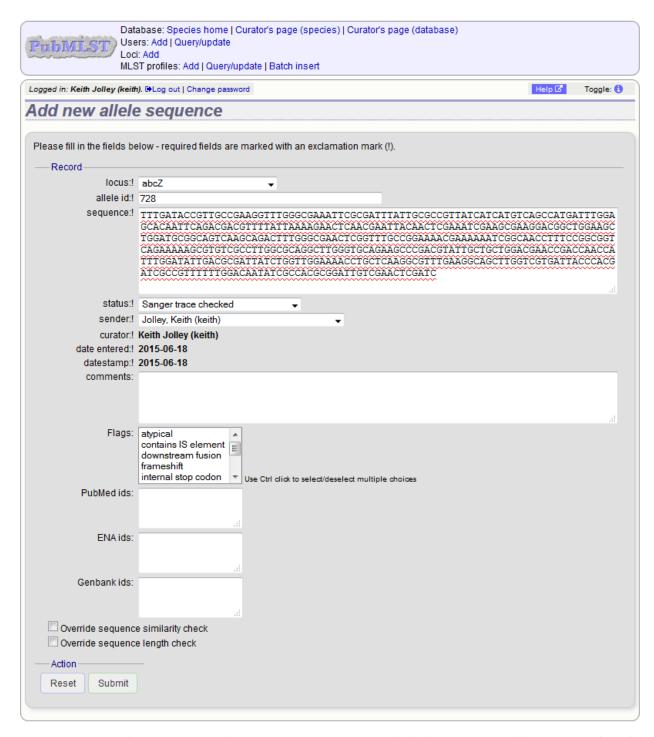
To add a single new allele, click the sequences (all loci) add (+) link on the curator's main page - if only a few loci are defined with permission for the current user to curate then they will be listed individually and the specific locus allele addition links can also be used.



Select the locus from the dropdown list box. The next available allele id will be entered automatically (if the allele id format is set to integer). Paste the sequence in to form, set the status and select the sender name from the dropdown box. If the sender does not appear in the box, you will need to add them to the registered users.

The status reflects the level of curation that the curator has done personally - the curator should not rely on assurances from the submitter. The status can either be:

- Sanger trace checked
 - Sequence trace files have been assembled and inspected by the curator.
- WGS: manual extract (BIGSdb)
 - The sequence has been extracted manually from a BIGSdb database *by the curator* . There may be some manual intervention to identify the start and stop sites of the sequence.
- WGS: automated extract (BIGSdb)
 - The sequences have been generated by a BIGSdb tag scanning run and have had no manual inspection or intervention.
- WGS: visually checked
 - Short read data has been inspected visually using an alignment program by the curator.
- WGS: automatically checked
 - The sequences have been checked by an automated algorithm that assesses the quality of the data to ensure it meets specified criteria.
- · unchecked
 - If none of the above match, then the sequence should be entered as unchecked.



Press submit. By default, the system will test whether your sequence is similar enough to existing alleles defined for that locus. The sequence will be rejected if it isn't considered similar enough. This test can be overridden by checking the 'Override sequence similarity check' checkbox at the bottom. It will also check that the sequence length is within the allowed range for that locus. These checks can also be overridden by checking the 'Override sequence length check' checkbox, allowing the addition of unusual length alleles.

See also:

allele sequence flags

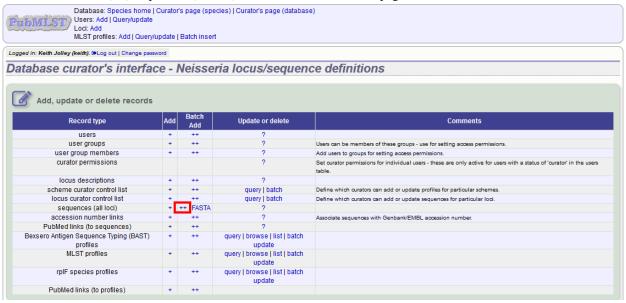
Sequences can also be associated with PubMed, ENA or Genbank id numbers by entering these as lists (one value per line) in the appropriate form box.

6.2.2 Batch adding multiple alleles

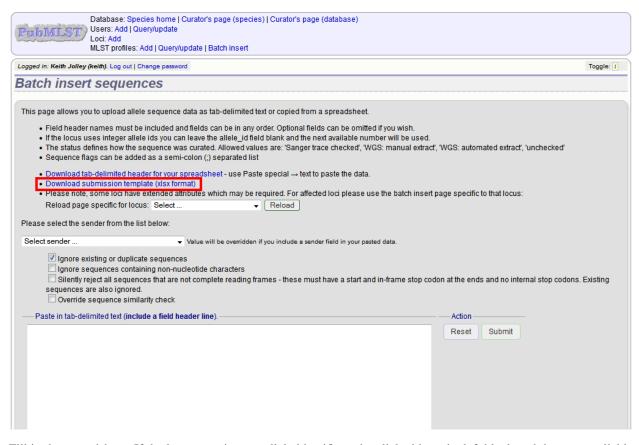
There are two methods of batch adding alleles. You can either upload a spreadsheet with all fields in tabular format, or you can upload a FASTA file provided all sequences are for the same locus and have the same status.

Upload using a spreadsheet

Click the batch add (++) sequences (all loci) link on the curator's main page.



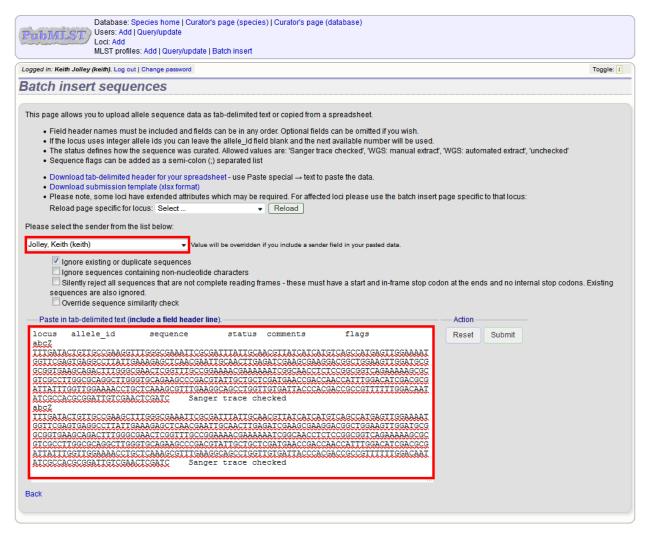
Download a template Excel file from the following page.



Fill in the spreadsheet. If the locus uses integer allele identifiers, the allele_id can be left blank and the next available number will be used automatically. Paste the entire sheet in to the web form and select the sender from the dropdown box.

Additionally, there are a number of options available. Some of these will ignore sequences if they don't match certain criteria - this is useful when sequence data has been extracted from genomes automatically. Available options are:

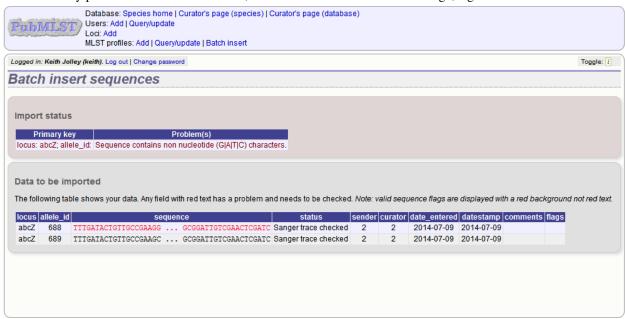
- Ignore existing or duplicate sequences.
- Ignore sequences containing non-nucleotide characters.
- Silently reject all sequences that are not complete reading frames these must have a start and in-frame stop codon at the ends and no internal stop codons. Existing sequences are also ignored.
- Override sequence similarity check.



Press submit. You will be presented with a page indicating what data will be uploaded. This gives you a chance to back out of the upload. Click 'Import data'.



If there are any problems with the submission, these should be indicated at this stage, e.g.:

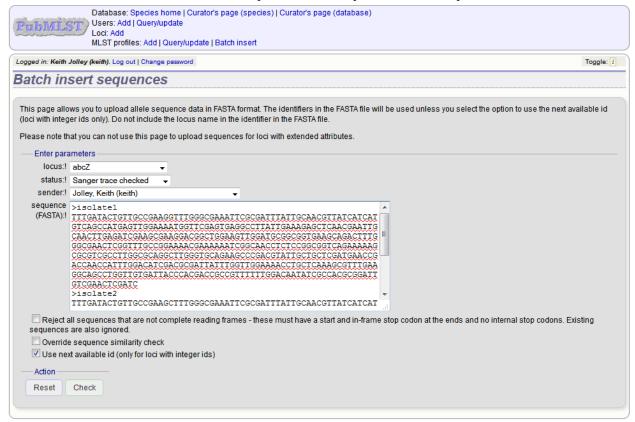


Upload using a FASTA file

Uploading new alleles from a FASTA file is usually more straightforward than generating an Excel sheet. Click 'FASTA' upload on the curator's contents page.



Select the locus, status and sender from the dropdown boxes and paste in the new sequences in FASTA format.



For loci with integer ids, the next available id number will be used by default (and the identifier in the FASTA file will be ignored). Alternatively, you can indicate the allele identifier within the FASTA file (do not include the locus name in this identifier).

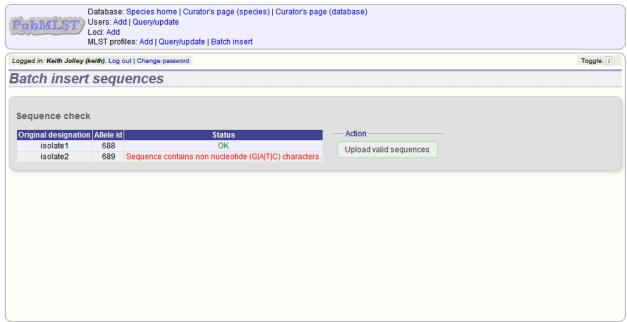
As with the spreadsheet upload, you can select options to ignore selected sequences if they don't match specific criteria.

Click 'Check'.

The sequences will be checked. You will be presented with a page indicating what data will be uploaded. This gives you a chance to back out of the upload. Click 'Upload valid sequences'.



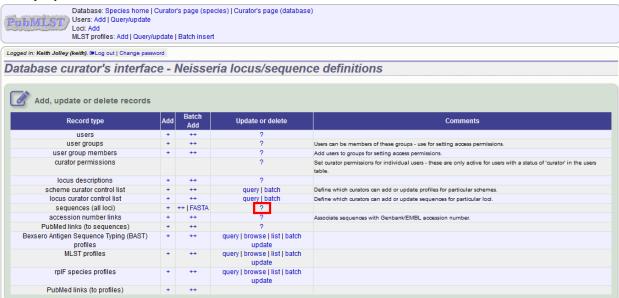
Any invalid sequences will be indicated in this confirmation page and these will not be uploaded (you can still upload the others), e.g.



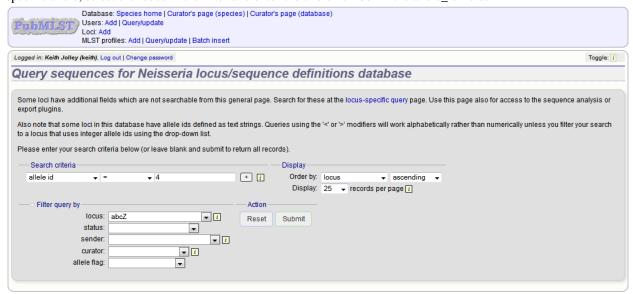
6.3 Updating and deleting allele sequence definitions

Note: You cannot update the sequence of an allele definition. This is for reasons of data integrity since an allele may form part of a scheme profile and be referred to in multiple databases. If you really need to change a sequence, you will have to remove the allele definition and then re-add it. If the allele is a member of a scheme profile, you will also have to remove that profile first, then re-create it after deleting and re-adding the allele.

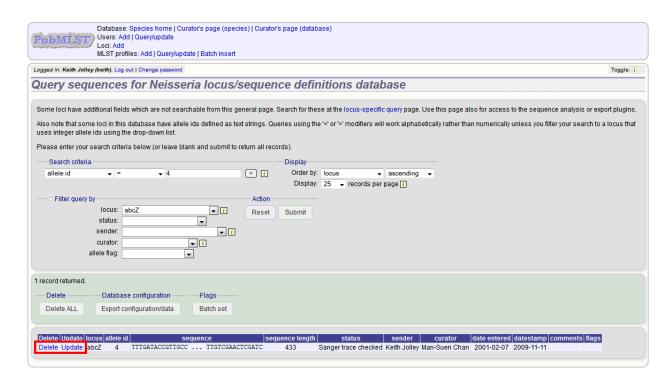
In order to update or delete an allele, first you must select it. Click the query (?) sequences (all loci) link - if only a few loci are defined with permission for the current user to curate then they will be listed individually and the specific locus query links can also be used.



Either search for specific attributes in the search form, or leave it blank and click 'Submit' to return all alleles. For a specific allele, select the locus in the filter and enter the allele number in the allele id field.



Click the appropriate link to either update the allele attributes or to delete it. If you have appropriate permissions, there may also be a link to 'Delete ALL'. This allows you to quickly delete all alleles returned from a search.

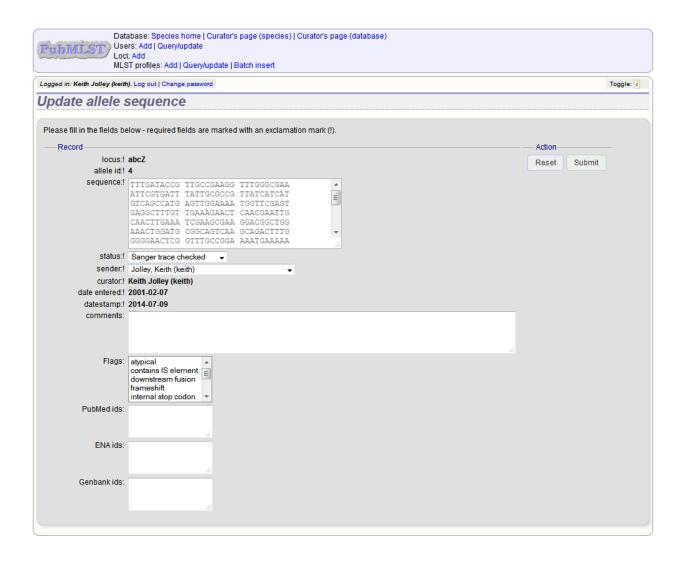


If you choose to delete, you will be presented with a final confirmation screen. To go ahead, click 'Delete!'. Deletion will not be possible if the allele is part of a scheme profile - if it is you will need to delete any profiles that it is a member of first.



If instead you clicked 'Update', you will be able to modify attributes of the sequence, or link PubMed, ENA or Genbank records to it. You will not be able to modify the sequence itself.

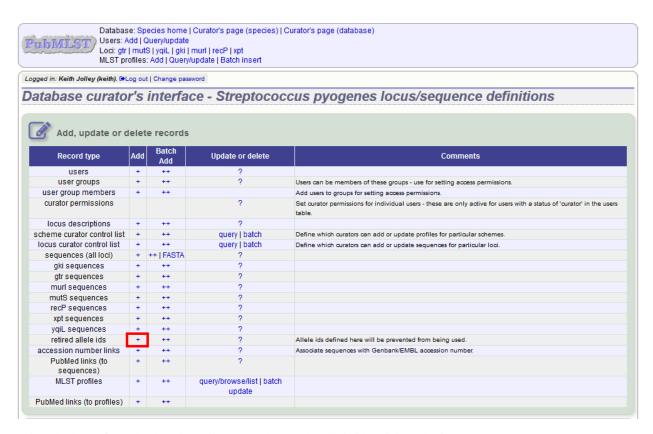
Note: Adding flags and comments to an allele record requires that this feature is enabled in the *database configuration*.



6.4 Retiring allele identifiers

Sometimes there is a requirement to prevent the automated assignment of a particular allele identifier - an allele with that identifier may have been commonly used and has since been removed. Reassignment of the identifier to a new sequence may lead to confusion, so in this instance, it would be better to prevent this.

You can retire an allele identifier by clicking the 'Add' retired allele ids link on the sequence database curators' page.



Select the locus from the dropdown list box and enter the allele id. Click 'Submit'.



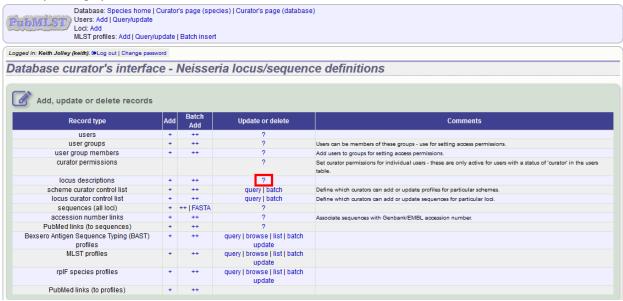
You cannot retire an allele that already exists, so you must delete it before retiring it. Once an identifier is retired, you will not be able to create a new allele with that name.

6.5 Updating locus descriptions

Loci in the sequence definitions database can have a description associated with them. This may contain information about the gene product, the biochemical reaction it catalyzes, or publications providing more detailed information etc. This description is accessible from various pages within the interface such as an *allele information page* or from the *allele download page*.

Note: In recent versions of BIGSdb, a blank description record is created when a new locus is defined. The following instructions assume that this is the case. It is possible for this record to be deleted or it may never have existed if the locus was created using an old version of BIGSdb. If the record does not exist, it can be added by clicking the Add (+) button next to 'locus descriptions'. Fill in the fields in the same way as described below.

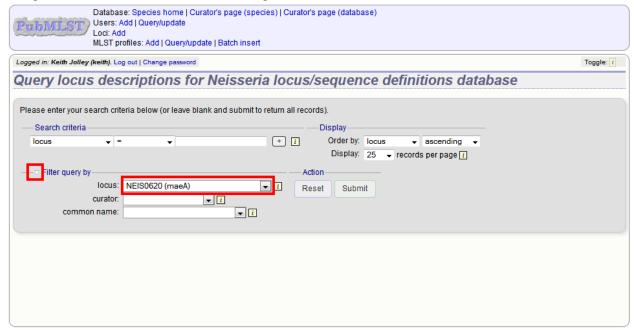
To edit a locus description, first you need to find it. Click the update/delete (?) button next to 'locus descriptions' on the sequence database curator's page (depending on the permissions set for your user account not all the links shown here may be displayed).



Either enter the name of the locus in the query box:

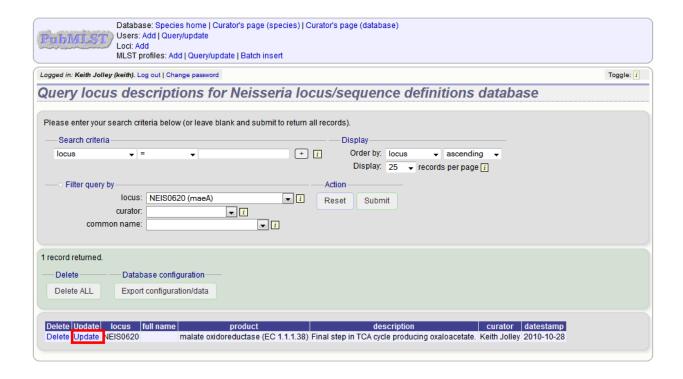


or expand the filter list and select it from the dropdown box:



Click 'Submit'.

If the locus description exists, click the 'Update' link (if it doesn't, see the note above).



Fill in the form as needed:



• full_name

The full name of the locus - often this can be left blank as it may be the same as the locus name. An example of where it is appropriately used is where the locus name is an abbreviation, e.g. PorA_VR1 - here we could enter 'PorA variable region 1'. This should not be used for the 'common name' of the locus (which is defined within the locus record itself) or the gene product.

• product

The name of the protein product of a coding sequence locus.

· description

This can be as full a description as possible. It can include the specific part of the biochemical pathway the gene product catalyses or may provide background information, as appropriate.

aliases

These are alternative names for the locus as perhaps found in different genome annotations. Don't duplicate the locus name or common name defined in the locus record. Enter each alias on a separate line.

· Pubmed_ids

Enter the PubMed id of any paper that specifically describes the locus. Enter each id on a separate line. The software will retrieve the full citation from PubMed (this happens periodically so it may not be available for display immediately).

Links

Enter links to additional web-based resources. Enter the URL first followed by a pipe symbol (I) and then the description.

Click 'Submit' when finished.

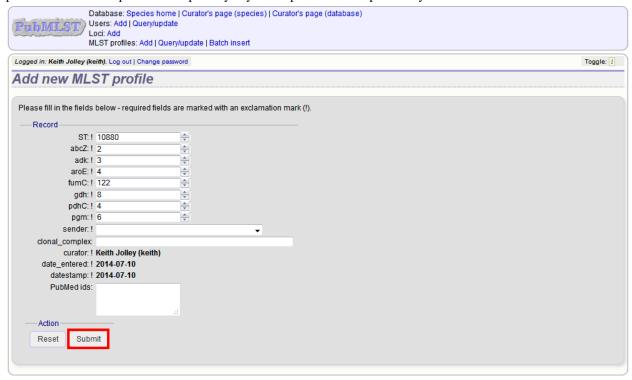
6.6 Adding new scheme profile definitions

Provided a scheme has been set up with at least one locus and a scheme field set as a primary key, there will be links on the curator's main page to add profiles for that scheme.

To add a single profile you can click the add (+) profiles link next to the scheme name (e.g. MLST):

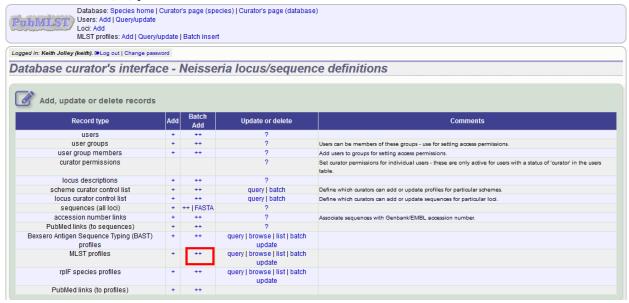


A form will be displayed with the next available primary key number already entered (provided integers are used for the primary key format). Enter the new profile, associated scheme fields, and the sender, then click 'Submit'. The new profile will be added provided the primary key or the profile has not previously been entered.

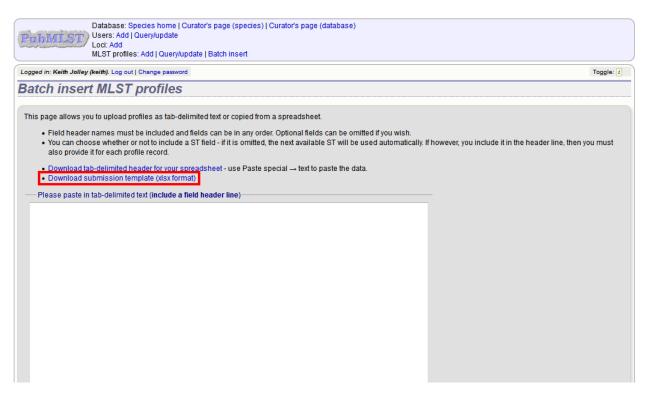


More usually, profiles are added in a batch mode. It is often easier to do this even for a single profile since it allows copying and pasting data from a spreadsheet.

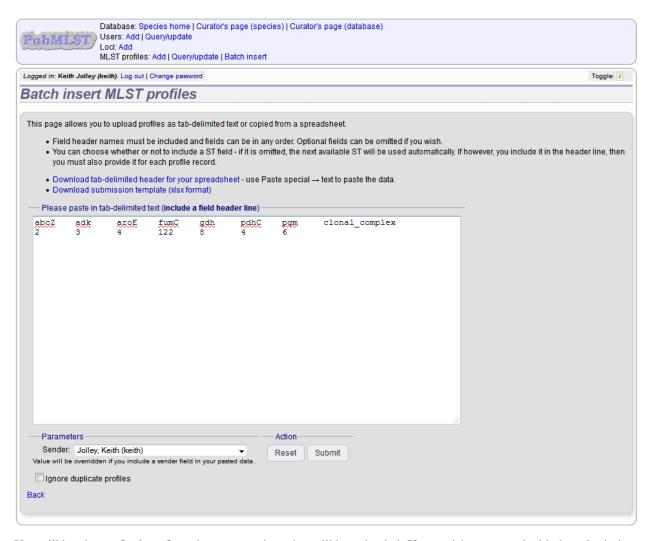
Click the batch add (++) profiles link next to the scheme name:



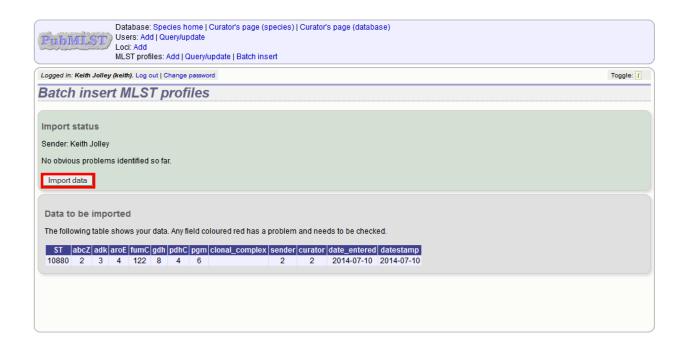
Click the 'Download submission template (xlsx format)' link to download an Excel submission template.



Fill in the spreadsheet using the copied template, then copy and paste the whole spreadsheet in to the large form on the upload page. If the primary key has an integer format, you can exclude this column and the next available number will be used automatically. If the column is included, however, a value must be set. Select the sender from the dropdown list box and then click 'Submit'.

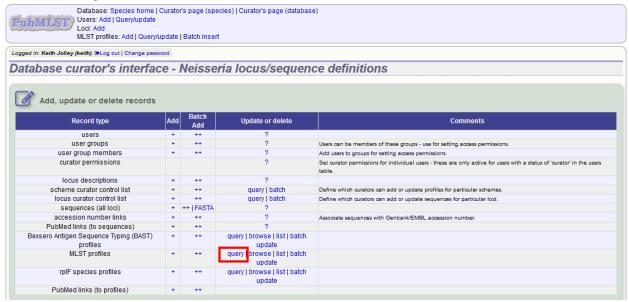


You will be given a final confirmation page stating what will be uploaded. If you wish to proceed with the submission, click 'Import data'.

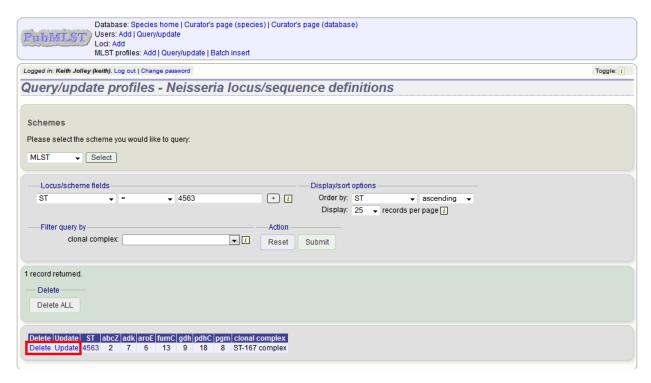


6.7 Updating and deleting scheme profile definitions

In order to update or delete a scheme profile, first you must select it. Click the query (?) profiles link next to the scheme name (e.g. MLST):



Search for your profile by entering search criteria (alternatively you can use the browse or list query functions).

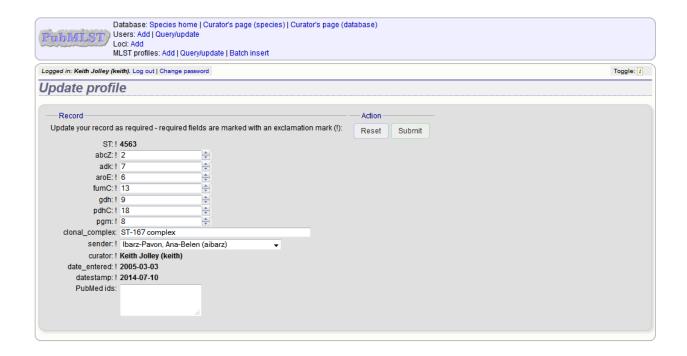


To delete the profile, click the 'Delete' link next to the profile. Alternatively, if your account has permission, you may be able to 'Delete ALL' records retrieved from the search.

For deletion of a single record, the full record will be displayed. Confirm deletion by clicking 'Delete!'.

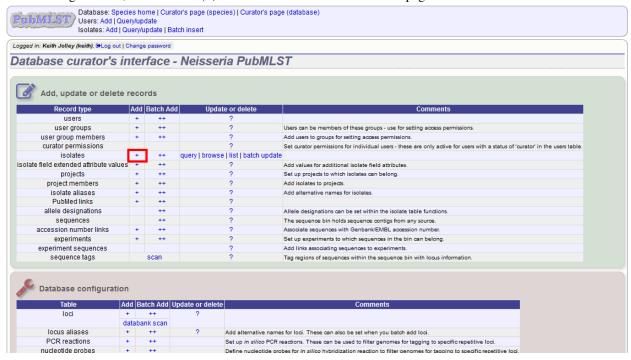


To modify the profile, click the 'Update' link next to the profile following the query. A form will be displayed - make any changes and then click 'Update'.

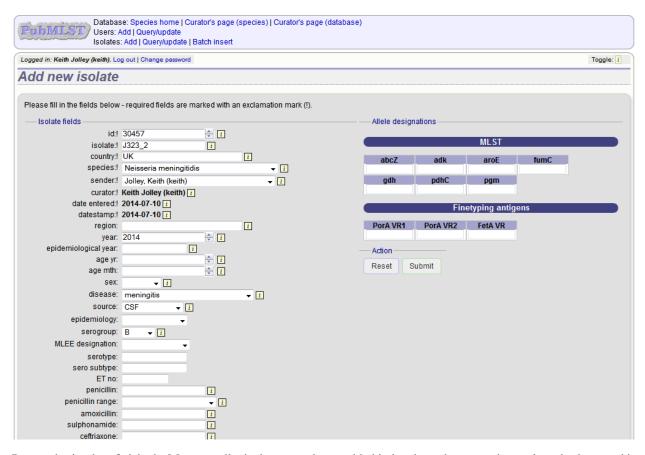


6.8 Adding isolate records

To add a single record, click the add (+) isolates link on the curator's index page.



The next available id will be filled in automatically but you are free to change this. Fill in the individual fields. Required fields are listed first and are marked with an exclamation mark (!). Some fields may have drop-down list boxes of allowed values. You can also enter allele designations for any loci that have been defined.

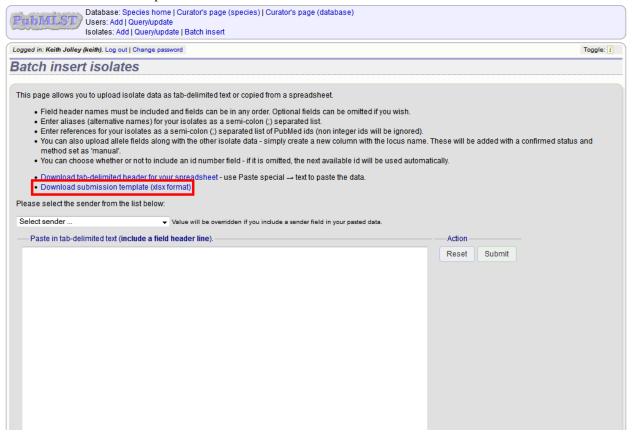


Press submit when finished. More usually, isolate records are added in batch mode, even when only a single record is added, since the submission can be prepared in a spreadsheet and copied and pasted.

Select batch add (++) isolates link on the curator's index page.

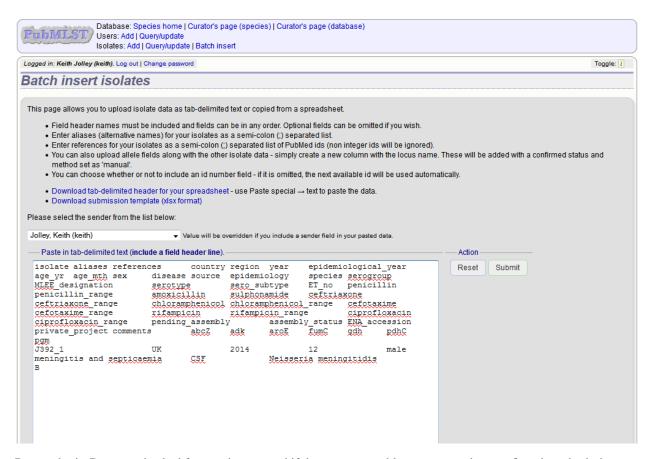


Download a submission template in Excel format from the link.

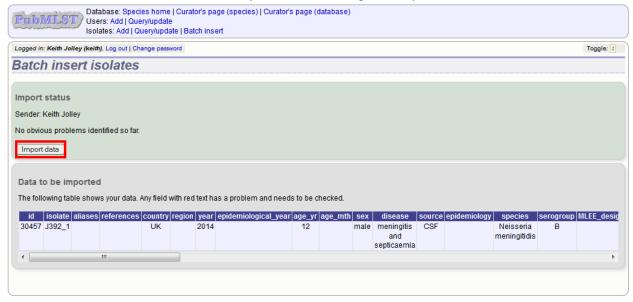


Prepare your data in the spreadsheet - the column headings must match the database fields. In databases with large numbers of loci, there won't be columns for each of these. You can, however, manually add locus columns.

Pick a sender from the drop-down list box and paste the data from your spreadsheet in to the web form. The next available isolate id number will be used automatically (this can be overridden if you manually add an id column).



Press submit. Data are checked for consistency and if there are no problems you can then confirm the submission.



Any problems with the data will be listed and highlighted within the table. Fix the data and resubmit if this happens.

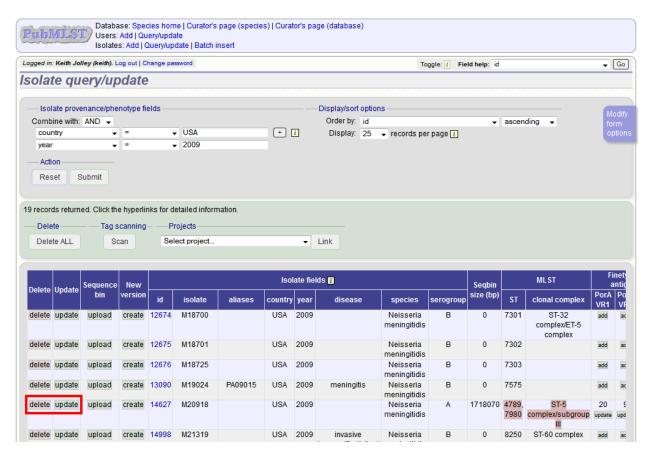


6.9 Updating and deleting single isolate records

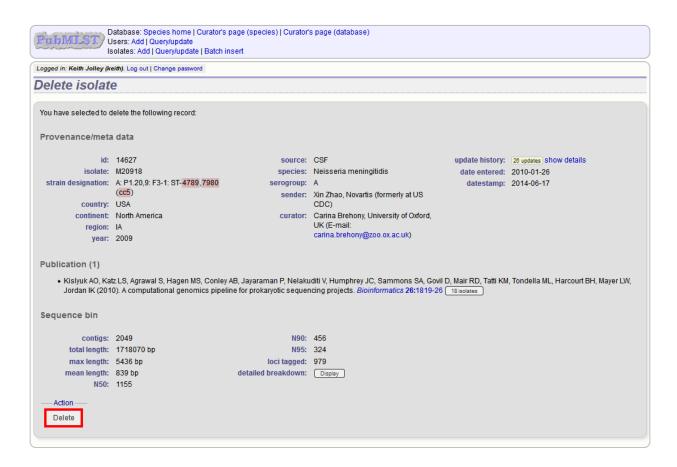
First you need to locate the isolate record. You can either browse or use a search or list query.



The query interface is the same as the *public query interface*. Following a query, a results table of isolates will be displayed. There will be delete and update links for each record.

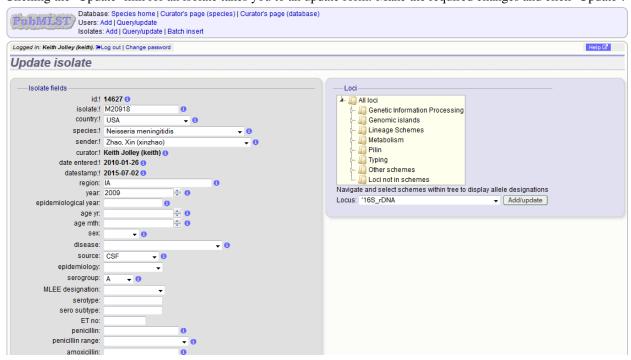


Clicking the 'Delete' link takes you to a page displaying the full isolate record.



Pressing 'Delete' from this record page confirms the deletion.

Clicking the 'Update' link for an isolate takes you to an update form. Make the required changes and click 'Update'.

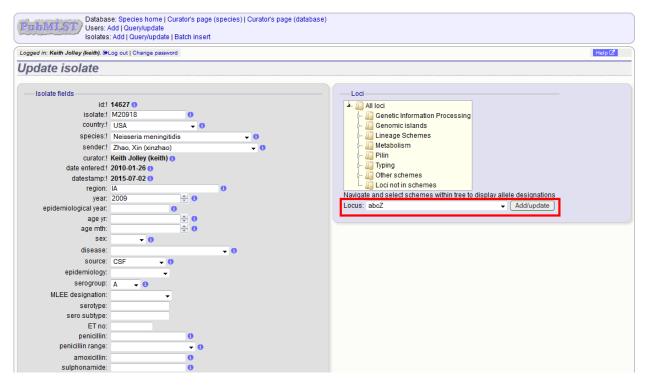


Allele designations can also be updated by clicking within the scheme tree and selecting the 'Add' or 'Update' link

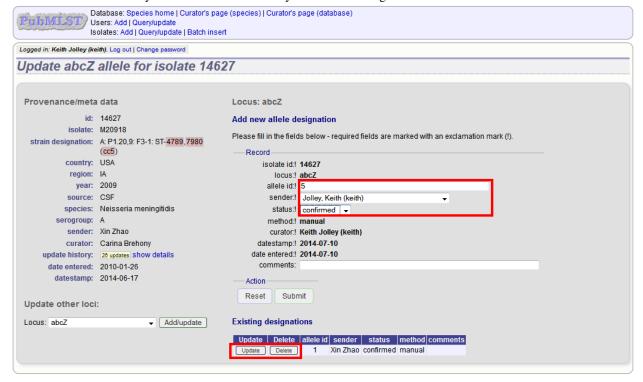
next to a displayed locus.



Schemes will only appear in the tree if data for at least one of the loci within the scheme has been added. You can additionally add or update allelic designations for a locus by choosing a locus in the drop-down list box and clicking 'Add/update'.



The allele designation update page allows you to modify an existing designation, or alternatively add additional designations. The sender, status (confirmed/provisional) and method (manual/automatic) needs to be set for each designation (all pending designations have a provisional status). The method is used to differentiate designations that have been determined manually from those determined by an automated algorithm.



6.10 Batch updating multiple isolate records





Prepare your update data in 3 columns in a spreadsheet:

- 1. Unique identifier field
- 2. Field to be updated
- 3. New value for field

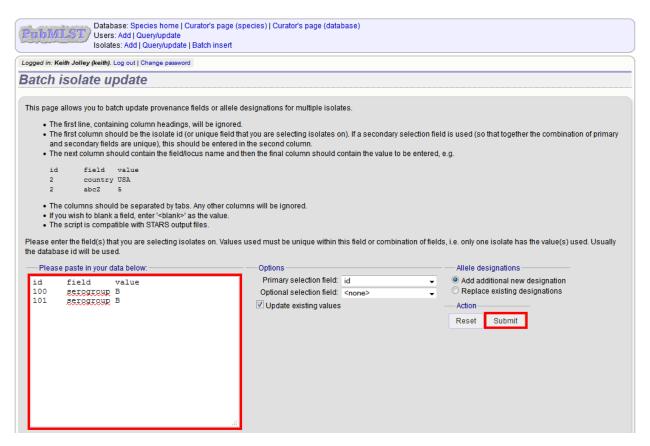
You should also include a header line at the top - this isn't used so can contain anything but it should be present.

Columns must be tab-delimited which they will be if you copy and paste directly from the spreadsheet.

So, to update isolate id-100 and id-101 to serogroup B you would prepare the following:

id	field	value
100	serogroup	В
101	serogroup	В

Select the field you are using as a unique identifier, in this case id, from the drop-down list box, and paste in the data. If the fields already have values set, you should also check the 'Update existing values' checkbox. Press 'submit'.



A confirmation page will be displayed if there are no problems. If there are problems, these will be listed. Press 'Upload' to upload the changes.



You can also use a secondary selection field such that a combination of two fields uniquely defines the isolate, for example using country and isolate name.

So, for example, to update the serogroups of isolates CN100 and CN103, both from the UK, select the appropriate primary and secondary fields and prepare the data as follows:

isolate	country	field	value
CN100	UK	serogroup	В
CN103	UK	serogroup	В

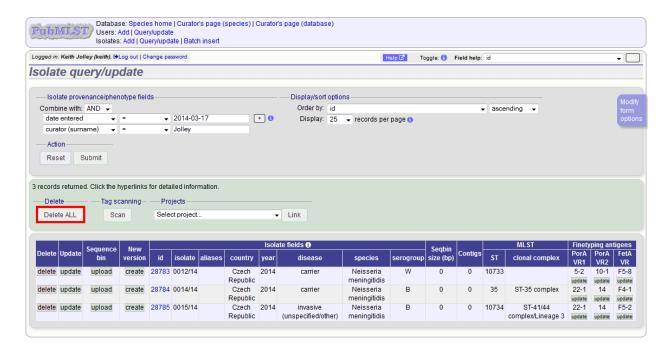
6.11 Deleting multiple isolate records

Note: Please note that standard curator accounts may not have permission to delete multiple isolates. Administrator accounts are always able to do this.

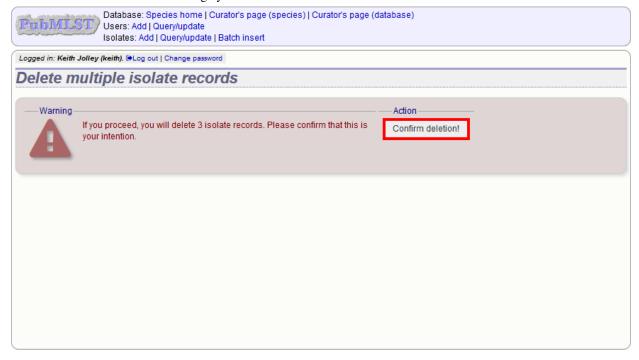
Before you can delete multiple records, you need to search for them. From the curator's main page, click the Query isolates link:



Enter search criteria that specifically return the isolates you wish to delete. Click 'Delete ALL'.



You will have a final chance to change your mind:

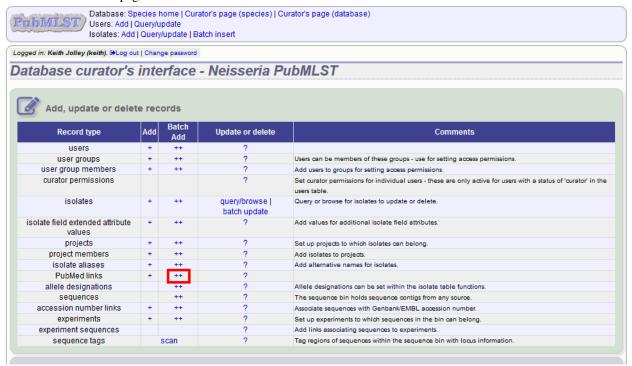


Click 'Confirm deletion!'.

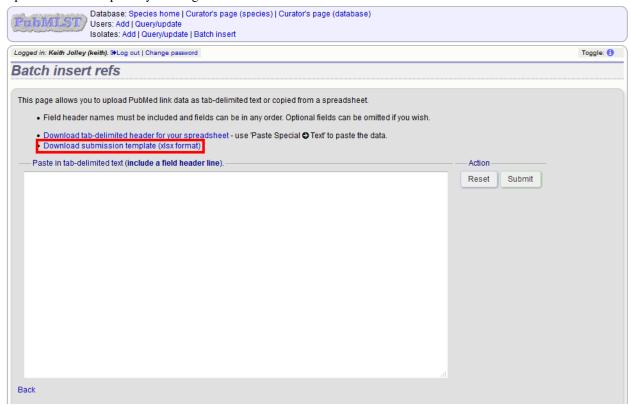
6.12 Linking isolate records to publications

Isolates can be associated with publications by adding PubMed id(s) to the record. This can be done when *adding the isolate*, where lists of PubMed ids can be entered in to the web form.

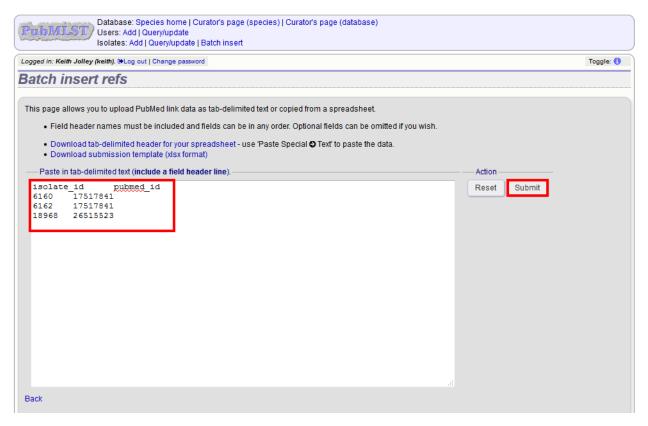
They can also be associated in batch after the upload of isolate records. Click the PubMed links batch add (++) link on the curator's main page.



Open the Excel template by clicking the link.



The Excel template has two columns, isolate_id and pubmed_id. Simply fill this in with a line for each record and then paste the entire spreadsheet in to the web form and press submit.

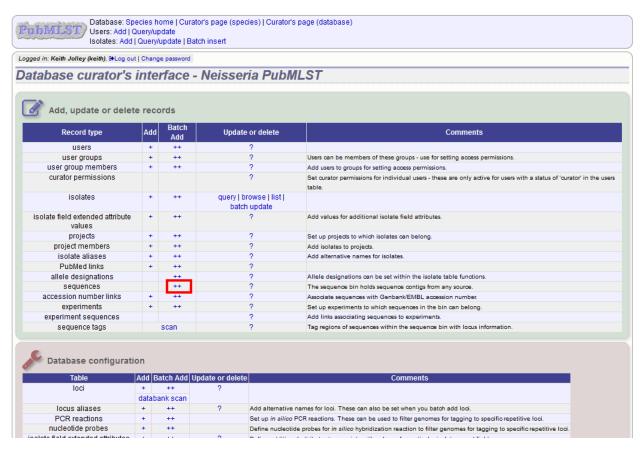


To ensure that publication information is stored locally and available for searching, the references database needs to be *updated regularly*.

6.13 Uploading sequence contigs linked to isolate records

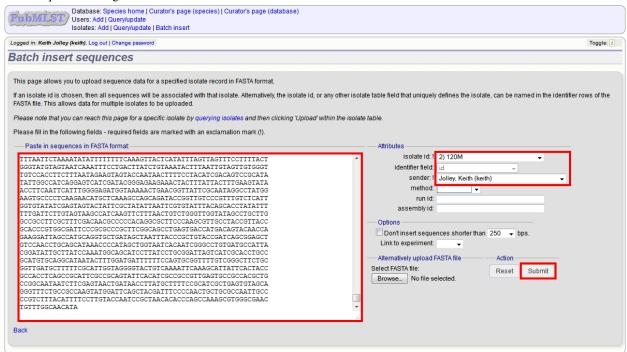
6.13.1 Select isolate from drop-down list

To upload sequence data, click the sequences batch add (++) link on the curator's main page.

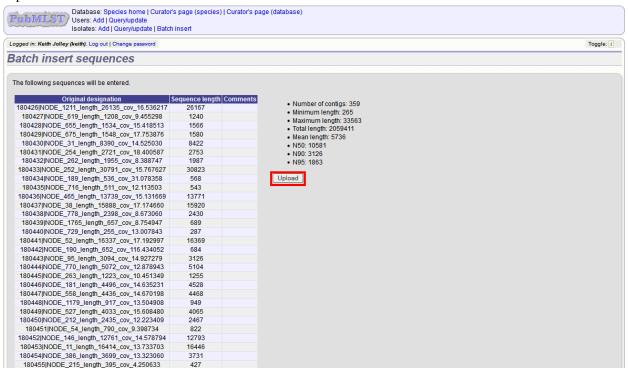


Select the isolate that you wish to link the sequence to from the dropdown list box. You also need to enter the person who sent the data. Optionally, you can add the sequencing method used.

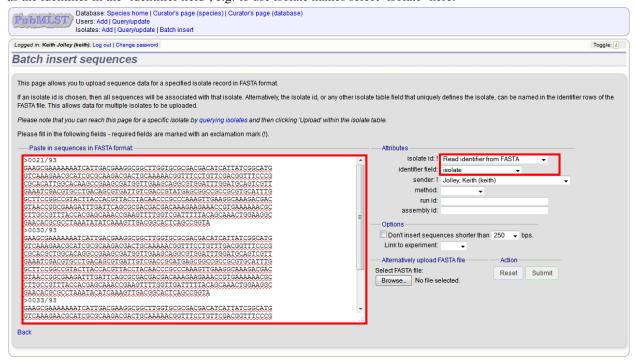
Paste sequence contigs in FASTA format in to the form.



Click 'Submit'. A summary of the number of isolates and their lengths will be displayed. To confirm upload, click 'Upload'.



It is also possible to upload data for multiple isolates at the same time, but these must exist as single contigs for each isolate. To do this, select 'Read identifier from FASTA' in the isolate id field and select the field that you wish to use as the identifier in the 'identifier field', e.g. to use isolate names select 'isolate' here.



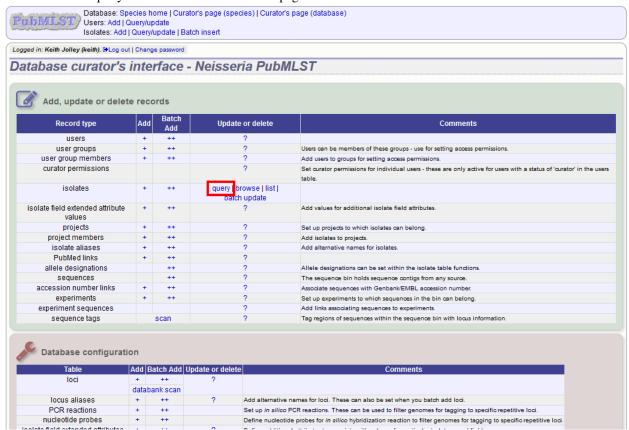
Provided the identifier used uniquely identifies the isolate you will get a confirmation screen. If the isolate name does not do this you'll probably have to use the database id number instead. Click 'Upload' to confirm.



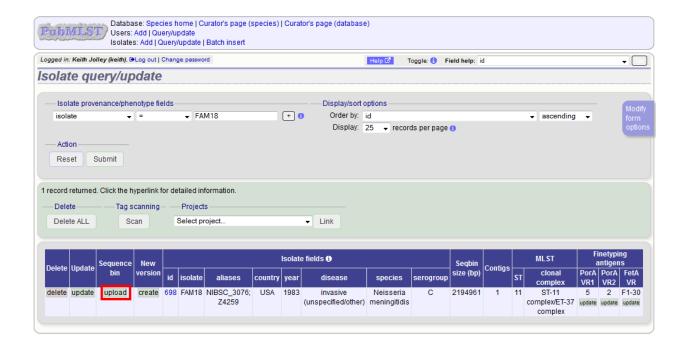
6.13.2 Select from isolate query

As an alternative to selecting the isolate from a dropdown list (which can become unwieldy for large databases), it is also possible to upload sequence data following an isolate query.

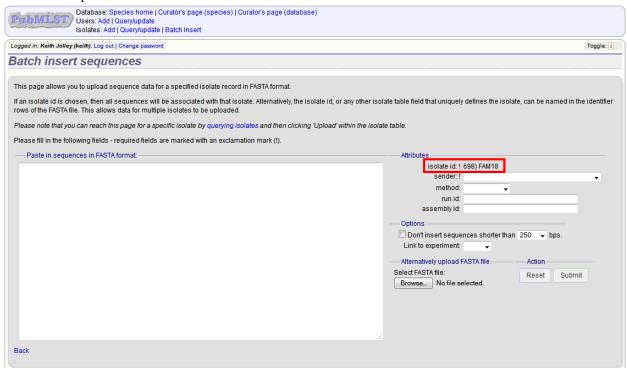
Click the isolate query link from the curator's main page.



Enter your search criteria. From the list of isolates displayed, click the 'Upload' link in the sequence bin column of the appropriate isolate record.



The same upload form as detailed above is shown. Instead of a dropdown list for isolate selection, however, the chosen isolate will be pre-selected.



6.13.3 Upload options

On the upload form, you can select to filter out short sequences from your contig list.

If your database has experiments defined (experiments are used for grouping sequences and can be used to filter the sequences used in *tag scanning*), you can also choose to upload your contigs as part of an experiment. To do this,

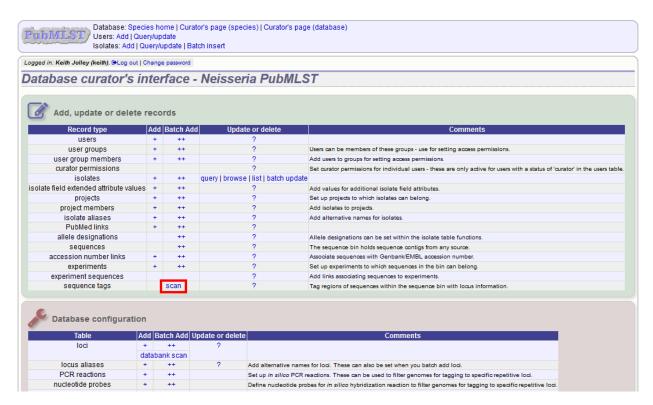
Database: Species home | Curator's page (species) | Curator's page (database) Users: Add | Query/update Isolates: Add | Query/update | Batch insert Logged in: Keith Jolley (keith). Log out | Change password Toggle: i Batch insert sequences If an isolate id is chosen, then all sequences will be associated with that isolate. Alternatively, the isolate id, or any other isolate table field that uniquely defines the isolate, can be named in the identifier rows of the FASTA file. This allows data for multiple isolates to be uploaded. Please note that you can reach this page for a specific isolate by querying isolates and then clicking 'Upload' within the isolate table Please fill in the following fields - required fields are marked with an exclamation mark (!). Paste in sequences in FASTA format: - Attributes isolate id: ! Read identifier from FASTA identifier field: id sender:! method: run id: assembly id: Don't insert sequences shorter than 250 Link to experiment: Alternatively upload FASTA file Action Select FASTA file: Reset Submit Browse... No file selected.

select the experiment from the dropdown list box.

6.14 Automated web-based sequence tagging

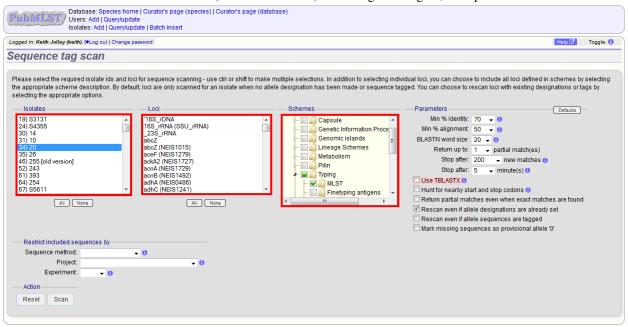
Sequence tagging, or tag-scanning, is the process of identifying alleles by scanning the sequence bin linked to an isolate record. Defined loci can either have a single reference sequence, that is defined in the locus table, or they can be linked to an external database that contains the sequences for known alleles. The tagging function uses BLAST to identify sequences and will tag the specific sequence region with locus information and an allele designation if a matching allele is identified by reference to an external database.

Select 'scan' sequence tags on the curator's index page.



Next, select the isolates whose sequences you wish to scan against. Multiple isolates can be selected by holding down the Ctrl key. All isolates can be selected by clicking the 'All' button under the isolate selection list.

Select either individual loci or schemes (collections of loci) to scan against. Again, multiple selections can be made.



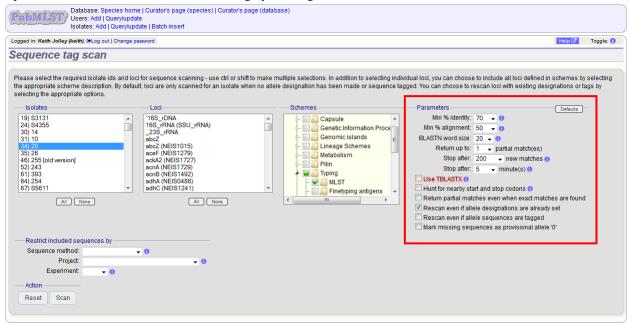
Choose your scan parameters. Lowering the value for BLASTN word size will increase the sensitivity of the search at the expense of time. Using TBLASTX is more sensitive but also much slower. TBLASTX can only be used to identify the sequence region rather than a specific allele (since it will only match the translated sequence and there may be multiple alleles that encode a particular peptide sequence).

By default, for each isolate only loci that have not had either an allele designation made or a sequence region scanned

will be scanned again. To rescan in these cases, select either or both the following:

- Rescan even if allele designations are already set
- · Rescan even if allele sequences are tagged

Options can be returned to their default setting by clicking the 'Defaults' button.



Press 'Scan'. The system takes approximately 1-2 seconds to identify each sequence (depending on machine speed and size of definitions databases). Any identified sequences will be listed in a table, with checkboxes indicating whether allele sequences or sequence regions are to be tagged.



Individual sequences can be extracted for inspection by clicking the 'extract \rightarrow ' link. The sequence (along with flanking regions) will be opened in another browser window or tab.

Checkboxes are enabled against any new sequence region or allele designation. You can also set a flag for a particular sequence to mark an attribute. These will be set automatically if these have been defined within the sequence definition database for an identified allele.

See also:

Sequence tag flags

Ensure any sequences you want to tag are selected, then press 'Tag alleles/sequences'.

If any new alleles are found, a link at the bottom will display these in a format suitable for automatic allele assignment by *batch uploading to sequence definition* database.

See also:

Offline curation tools

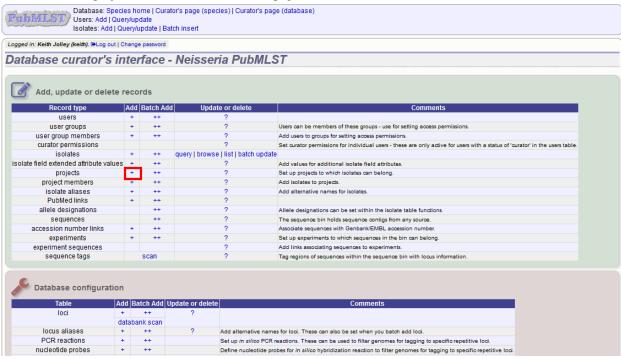
Automated offline sequence tagging

6.15 Projects

6.15.1 Creating the project

The first step in grouping by project is to set up a project.

Click the add (+) project link on the curator's main page.



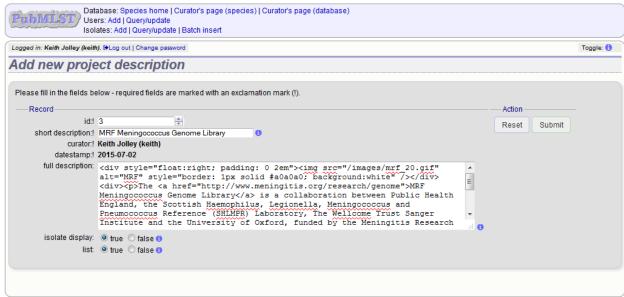
Enter a short description for the project. This is used in drop-down list boxes within the query interfaces, so make sure it is not too long.

You can also enter a full description. If this is added, the project description can displayed at the top of an isolate information page (but see 'isolate_display' flag below). The full description can include HTML formatting, including image links.

There are additionally two flags that affect how projects are listed:

- isolate_display Setting this is required for the project and its description to be listed at the top of an isolate record (default: false).
- list Setting this is required for the project to be listed in a page of projects linked from the main contents page.

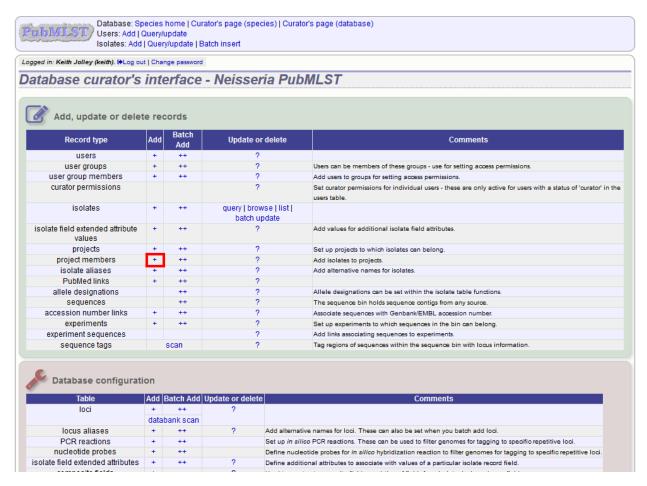
Click 'Submit'.



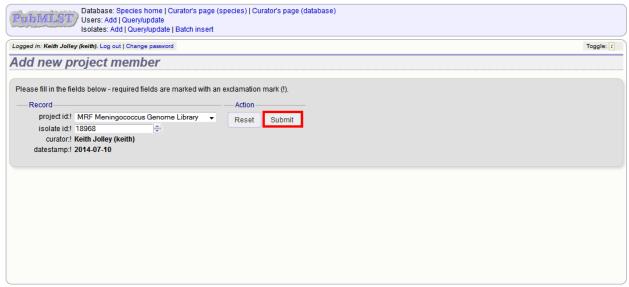
6.15.2 Explicitly adding isolates to a project

Explicitly adding isolates to the project can be done individually or in batch mode. To add individually, click the add (+) project member link on the curator's main page.

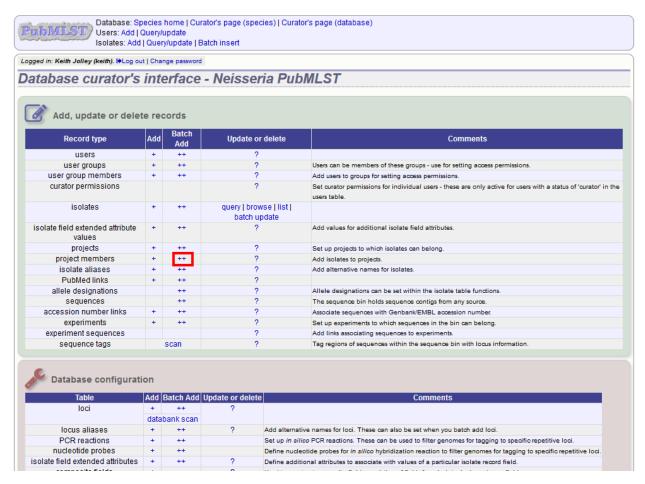
6.15. Projects 155



Select the project from the dropdown list box and enter the id of the isolate that you wish to add to the project. Click 'Submit'.

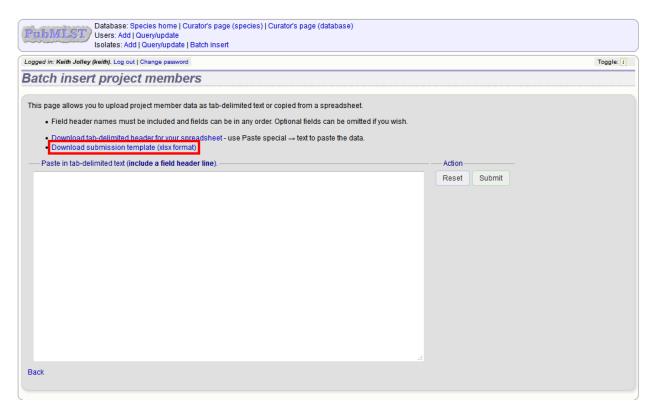


To add isolates in batch mode. Click the batch add (++) project members link on the curator's main page.



Download an Excel submission template:

6.15. Projects 157



You will need to know the id number of the project - this is the id that was used when you created the project. Fill in the spreadsheet, listing the project and isolate ids. Copy and paste this to the web upload form. Press 'Submit'.



6.16 Isolate record versioning

Versioning enables multiple versions of genomes to be uploaded to the database and be analysed separately. When a new version is created, a copy of the provenance metadata, and publication links are created in a new isolate record. The sequence bin and allele designations are not copied.

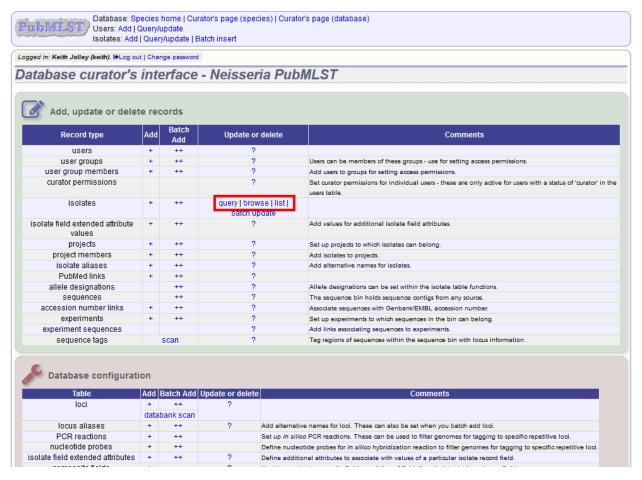
By default, old versions of the record are not returned from queries. Most query pages have a checkbox to 'Include old record versions' to override this.

Links to different versions are displayed within an isolate record:

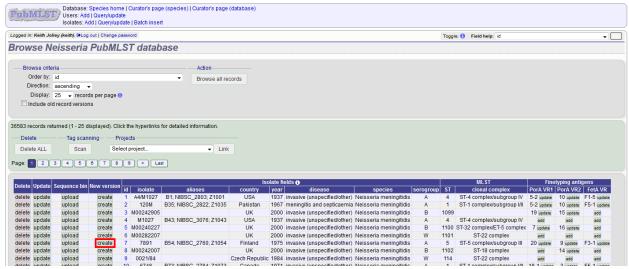


The different versions will also be listed in analysis plugins, with old versions identified with an [old version] designation after their name.

To create a new version of an isolate record, query or browse for the isolate:

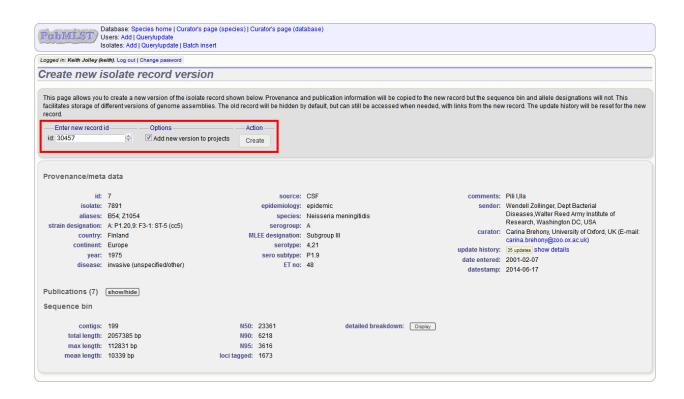


Click the 'create' new version link next to the isolate record:



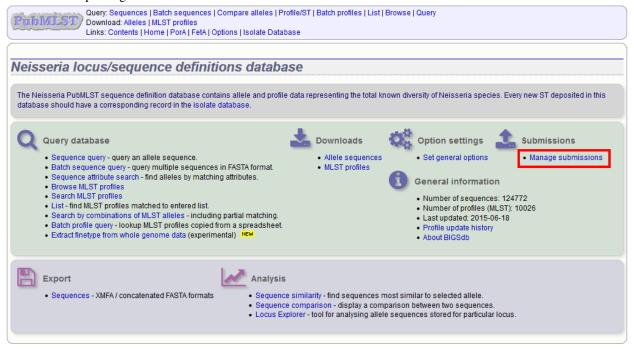
The isolate record will be displayed. The suggested id number for the new record will be displayed - you can change this. By default, the new record will also be added to any projects that the old record is a member of. Uncheck the 'Add new version to projects' checkbox to prevent this.

Click the 'Create' button.



Curating submitted data

Data may be submitted by users using the automated submission system if it has been enabled for a specific database. As a curator, you will be notified of pending submissions when you log in to the curator's interface or if you access the 'Manage submissions' links from the standard contents page. Additionally, if your user account has the 'submission_emails' flag set in the users' table you will also receive E-mail notification of new submissions for which you have sufficient privileges to curate.



Any submissions for which you have sufficient privileges to curate will be shown.



7.1 Alleles

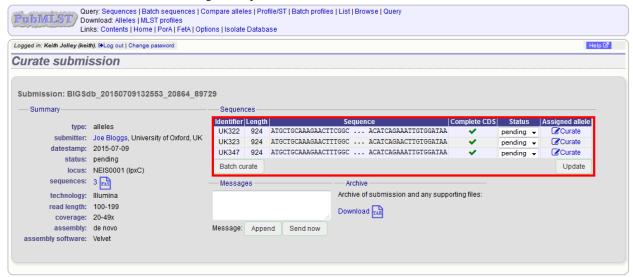
Click the link to the appropriate submission on the 'Manage submissions' page.



You will see a summary section that describes details about how the sequences were obtained. There should also be link here to download all the sequences in FASTA format.



There will also be a table summarizing the sequences in the submission and their current submission status.



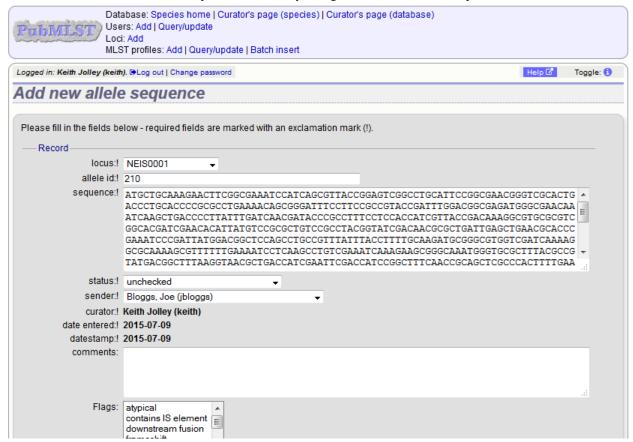
7.1.1 Individual allele curation

Individual sequences can be curated singly by clicking the 'Curate' links next to the sequence in the table. If you have supporting data attached to the submission, e.g. Sanger trace files then you may need to assess the submission based on the policy of the database.

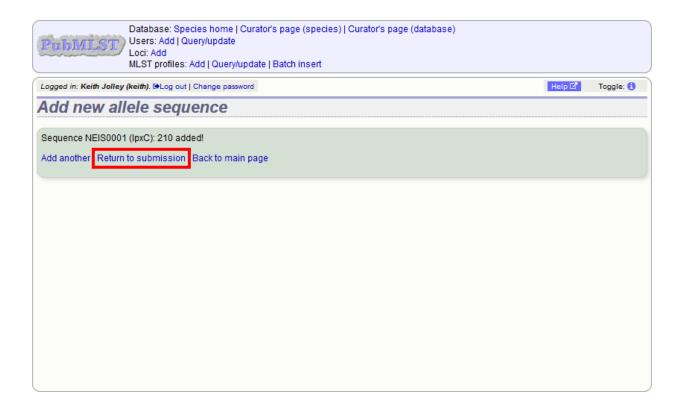
7.1. Alleles 165



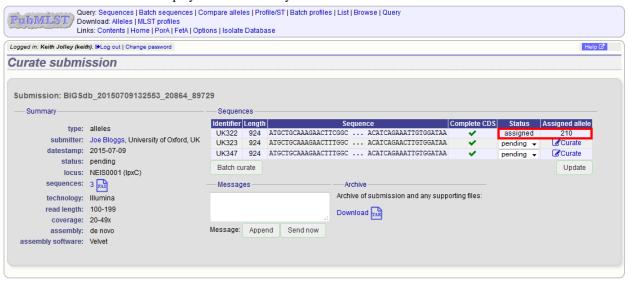
Clicking this link takes you to the curation interface *single sequence upload page*. The upload form will be filled with details from the submission. You may wish to manually change the status from the dropdown list of values.



Clicking 'Submit' from this form will define the new allele and add it to the database. A link on the confirmation page will take you back to the submission management page.



You will find that the status of the newly assigned sequence has changed in the summary table. The assigned value and status are determined on display and should always reflect the live database values.



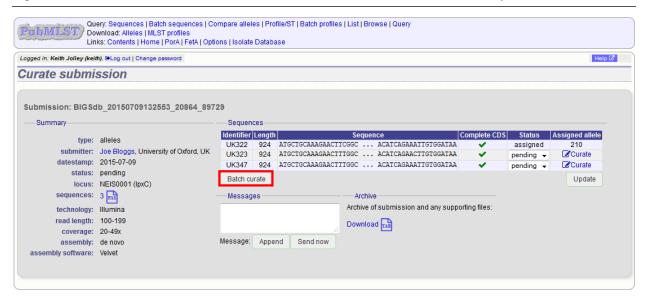
7.1.2 Batch allele curation

Often, you will want to batch upload submitted sequences. This can be done by clicking the 'Batch curate' button.

Note: Batch curation is only available for loci that do not have extended attributes defined. Entries for these loci

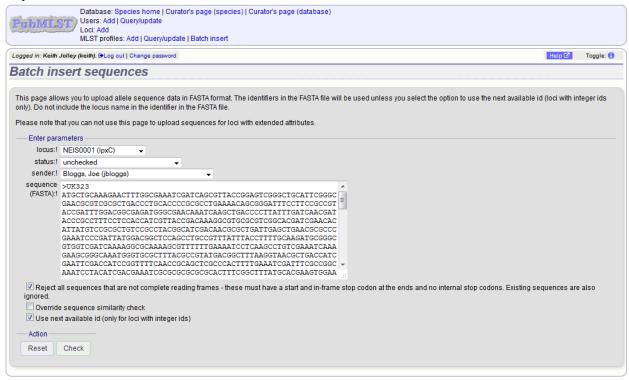
7.1. Alleles 167

require additional values set for these additional fields and so need to be handled individually.

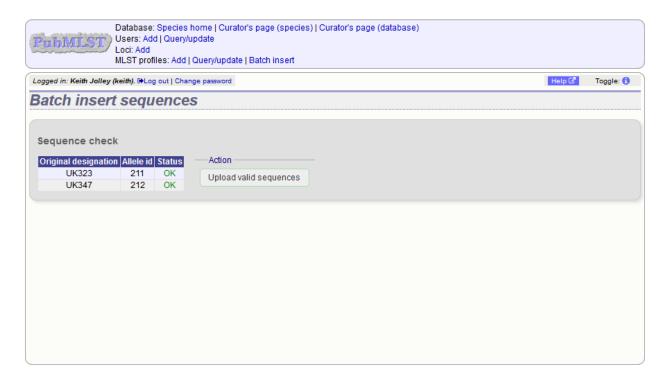


This takes you to the batch FASTA upload page in the curators' interface.

The upload form will be filled with details from the submission. You may wish to manually change the status from the dropdown list of values.



Click 'Check' on this form will perform some standard checks before allowing you to upload the sequences.

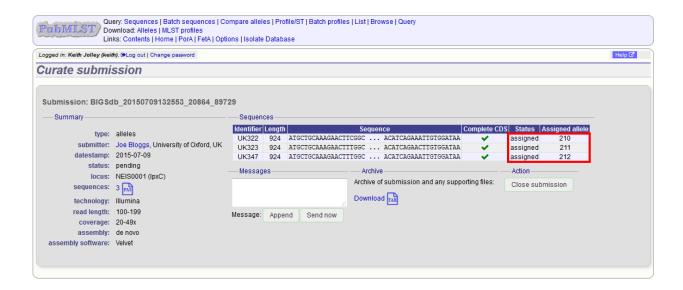


A link on the confirmation page will take you back to the submission management page.



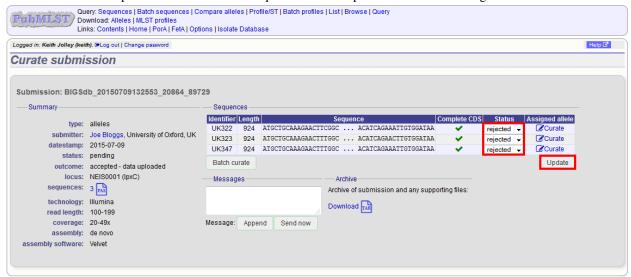
The status of the sequences should reflect their newly assigned status.

7.1. Alleles 169



7.1.3 Rejecting sequences

Sometimes you may need to reject all, or some of, the sequences in a submission. You can do this by changing the value in the status dropdown box next to each sequence. Click 'Update' to make the change.

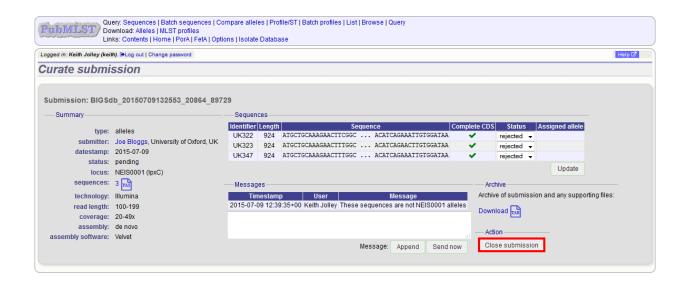


7.1.4 Requesting additional information

You can send a message to the submitter by entering it in the Messages box and clicking 'Send now'. This will append a message to the submission and send an update to the submitter so that they can respond.

7.1.5 Closing the submission

You can add a message to the submitter by entering it in the message box and clicking 'Append'. Once sequences have all been either assigned or rejected, the 'Close submission' button will be displayed. Click this to close the submission. The submitter will be notified of their submission status.



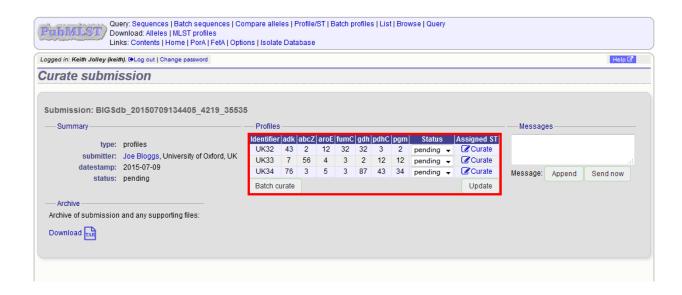
7.2 Profiles

Click the appropriate submission on the 'Manage submissions' page.



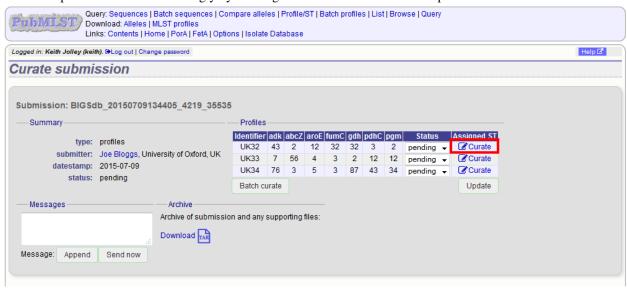
You will see a table summarizing the profiles in the submission and their current status.

7.2. Profiles 171

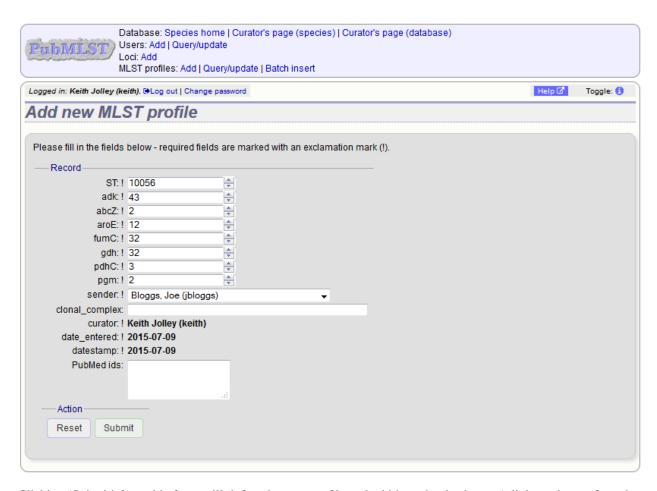


7.2.1 Individual profile curation

Individual profiles can be curated singly by clicking the 'Curate' links next to the profile in the table.



Clicking this link takes you to the curation interface *single profile upload page*. The upload form will be filled with details from the submission.

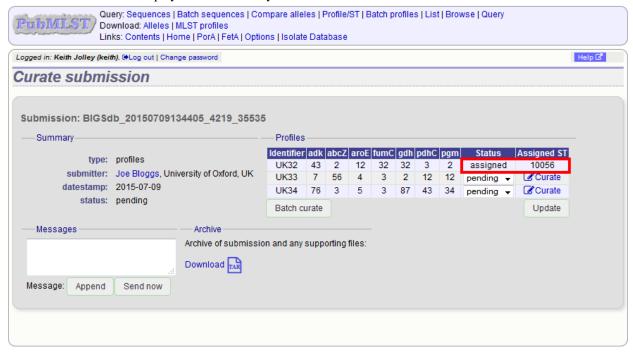


Clicking 'Submit' from this form will define the new profile and add it to the database. A link on the confirmation page will take you back to the submission management page.

7.2. Profiles 173

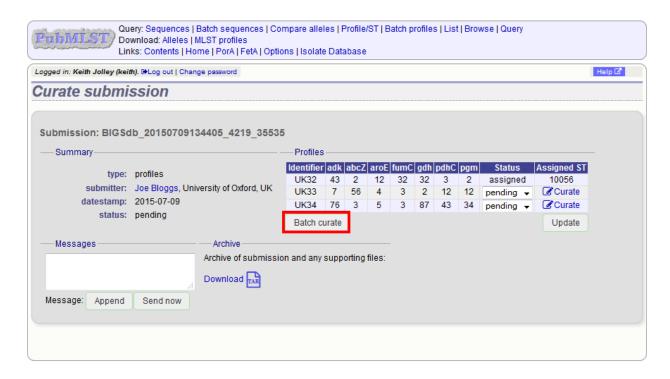


You will find that the status of the newly assigned profile has changed in the summary table. The assigned value and status are determined on display and should always reflect the live database values.



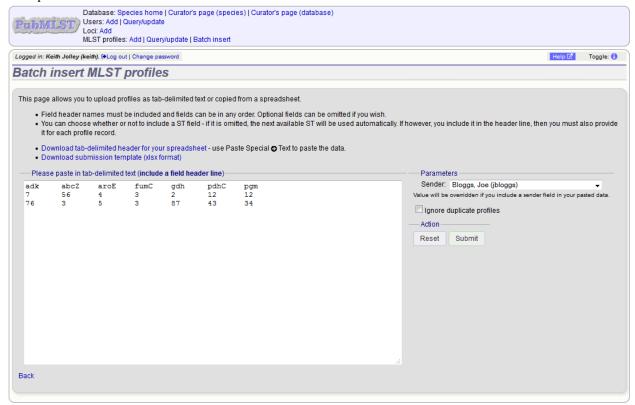
7.2.2 Batch profile curation

Often, you will want to batch upload submitted profiles. This can be done by clicking the 'Batch curate' button.



This takes you to the batch profile upload page in the curators' interface.

The upload form will be filled with details from the submission.

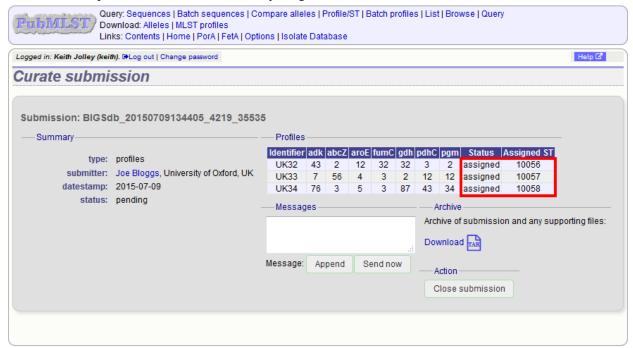


After upload, a link on the confirmation page leads back to the submission management page.

7.2. Profiles 175

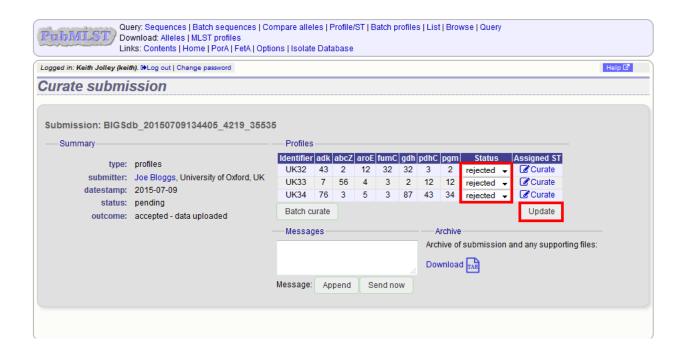


The status of the profiles should reflect their newly assigned status.



7.2.3 Rejecting profiles

Sometimes you may need to reject all, or some of, the profiles in the submission. This may be because isolate data had not been made available, against the policy of the database. You can do this by changing the value in the status dropdown box next to each profile. Click 'Update' to make the change.

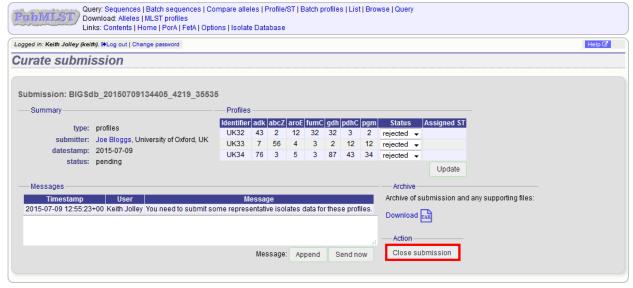


7.2.4 Requesting additional information

You can send a message to the submitter by entering it in the Messages box and clicking 'Send now'. This will append a message to the submission and send an update to the submitter so that they can respond.

7.2.5 Closing the submission

You can add a message to the submitter by entering it in the message box and clicking 'Append'. Once profiles have all been either assigned or rejected, the 'Close submission' button will be displayed. Click this to close the submission. The submitter will be notified of their submission status.



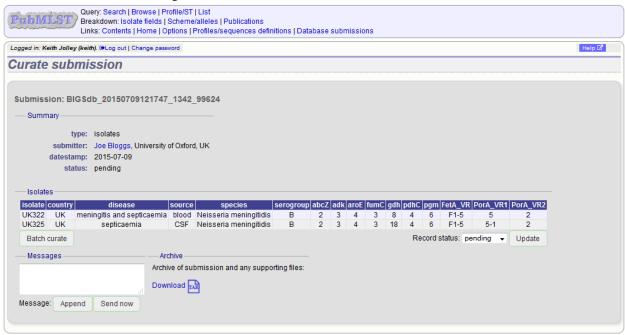
7.2. Profiles 177

7.3 Isolates

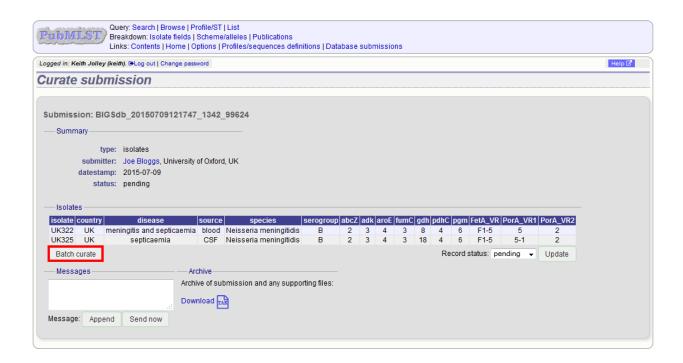
Clicking the appropriate submission on the 'Manage submissions' page.



You will see a table summarizing the submission.

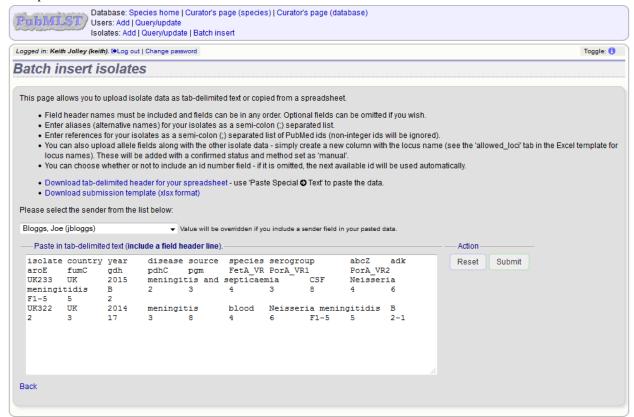


Click the 'Batch curate' button.



This will take you to the batch isolate upload page in the curators' interface.

The upload form will be filled with details from the submission.



Click submit to check and then import if there are no errors.

7.3. Isolates 179

After upload, a link on the confirmation page leads back to the submission management page.



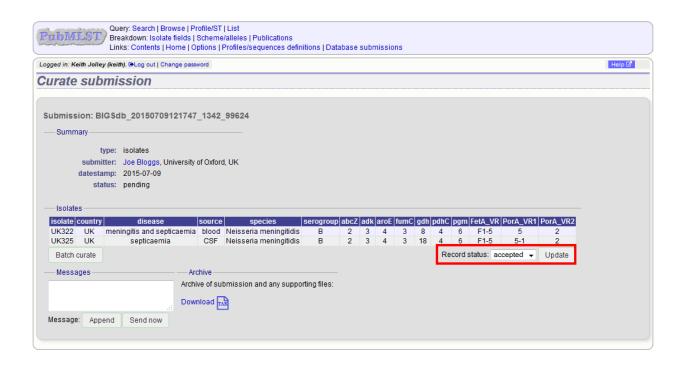
Note: Depending on the database policy, definitions of new scheme profiles, e.g. for MLST, may require submission of representative isolate records. Where this is the case, the curator will need to extract the new profile from the submitted record. The tab-delimited isolate text file can be downloaded from the archive of supporting files linked to the submission and used directly for *batch adding new profiles*. Alternatively, the curator could use the *Export functionality* of the database to generate the file required for batch profile definition after upload of the isolate data.

7.3.1 Requesting additional information

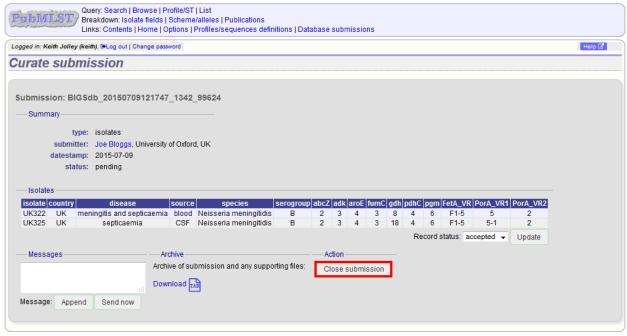
You can send a message to the submitter by entering it in the Messages box and clicking 'Send now'. This will append a message to the submission and send an update to the submitter so that they can respond.

7.3.2 Closing the submission

You can add a message to the submitter by entering it in the message box and clicking 'Append'. Change the record status to either 'accepted' or 'rejected' depending on whether you have accepted the submission. Click 'Update'.



The 'Close submission' button will now appear. Click this to close the submission. The submitter will be notified of their submission status.



7.3. Isolates 181

Offline curation tools

8.1 Automated offline sequence tagging

Sequence tagging is the process of identifying alleles by scanning the sequence bin linked to an isolate record. Loci need to be defined in an external sequence definition database that contains the sequences for known alleles. The tagging function uses BLAST to identify sequences and will tag the specific sequence region with locus information and an allele designation if a matching allele is identified by reference to an external database.

There is a script called 'autotag.pl' in the BIGSdb package. This can be used to tag genome sequences from the command line.

Before autotag.pl can be run for the first time, a log file needs to be created. This can be created if it doesn't already exist with the following:

```
sudo touch /var/log/bigsdb_scripts.log
sudo chown bigsdb /var/log/bigsdb_scripts.log
```

The autotag.pl script should be installed in /usr/local/bin. It is run as follows:

```
autotag.pl --database <database configuration>
```

where <database configuration> is the name used for the argument 'db' when using the BIGSdb application.

If you have multiple processor cores available, use the –threads option to set the number of jobs to run in parallel. Isolates for scanning will be split among the threads.

The script must be run by a user that can both write to the log file and access the databases, e.g. the 'bigsdb' user (see 'Setting up the offline job manager').

A full list of options can be found by typing:

```
autotag.pl --help

NAME
    autotag.pl - BIGSdb automated allele tagger

SYNOPSIS
    autotag.pl --database NAME [options]

OPTIONS
-0, --missing
    Marks missing loci as provisional allele 0. Sets default word size to 15.

-d, --database NAME
```

Database configuration name.

-e, --exemplar

Only use alleles with the 'exemplar' flag set in BLAST searches to identify locus within genome. Specific allele is then identified using a database lookup. This may be quicker than using all alleles for the BLAST search, but will be at the expense of sensitivity. If no exemplar alleles are set for a locus then all alleles will be used. Sets default word size to 15.

-f --fast

Perform single BLAST query against all selected loci together. This will take longer to return any results but the overall scan should finish quicker. This method will also use more memory — this can be used with ——exemplar to mitigate against this.

-h, --help

This help page.

-i, --isolates LIST

Comma-separated list of isolate ids to scan (ignored if -p used).

--isolate_list_file FILE

File containing list of isolate ids (ignored if -i or -p used).

-I, --exclude_isolates LIST

Comma-separated list of isolate ids to ignore.

-l, --loci LIST

Comma-separated list of loci to scan (ignored if -s used).

-L, --exclude_loci LIST

Comma-separated list of loci to exclude

-m, --min_size SIZE

Minimum size of seqbin (bp) - limit search to isolates with at least this much sequence.

-n, --new_only

New (previously untagged) isolates only. Combine with --new_max_alleles if required.

--new_max_alleles ALLELES

Set the maximum number of alleles that can be designated or sequences tagged before an isolate is not considered new when using the --new_only option.

-o, --order

Order so that isolates last tagged the longest time ago get scanned first (ignored if \neg r used).

--only_already_tagged

Only check loci that already have a tag present (but no allele designation). This must be combined with the --already_tagged option or no loci will match. This option is used to perform a catch-up scan where a curator has previously tagged sequence regions prior to alleles being defined, without the need to scan all missing loci.

-p, --projects LIST

```
Comma-separated list of project isolates to scan.
-P, --exclude_projects LIST
    Comma-separated list of projects whose isolates will be excluded.
-q, --quiet
    Only error messages displayed.
-r, --random
    Shuffle order of isolate ids to scan.
-R, --locus_regex REGEX
   Regex for locus names.
-s, --schemes LIST
    Comma-separated list of scheme loci to scan.
-t, --time MINS
    Stop after t minutes.
--threads THREADS
   Maximum number of threads to use.
-T, --already_tagged
    Scan even when sequence tagged (no designation).
-v, --view VIEW
    Isolate database view (overrides value set in config.xml).
-w, --word_size SIZE
   BLASTN word size.
-x, --min ID
   Minimum isolate id.
-y, --max ID
   Maximum isolate id.
```

8.2 Defining exemplar alleles

Exemplar alleles are a subset of the total number of alleles defined for a locus that encompass the known diversity within a specified identity threshold. They can be used to speed up *autotagging* as the BLAST queries are performed against exemplars to identify the locus region in the genome followed by a direct database lookup of the sequence found to identify the exact allele found. This is usually combined with the autotagger –fast option.

There is a script called 'find_exemplars.pl' in the BIGSdb scripts/maintenance directory.

A full list of options can be found by typing:

```
find_exemplars.pl --help

NAME
    find_exemplars.pl - Identify and mark exemplar alleles for use
    by tagging functions

SYNOPSIS
    find_exemplars.pl --database NAME [options]
```

```
OPTIONS
--database NAME
   Database configuration name.
--datatype DNA|peptide
   Only define exemplars for specified data type (DNA or peptide)
--exclude_loci LIST
   Comma-separated list of loci to exclude
--help
   This help page.
--loci LIST
   Comma-separated list of loci to scan (ignored if -s used).
--locus_regex REGEX
   Regex for locus names.
--schemes LIST
   Comma-separated list of scheme loci to scan.
--update
   Update exemplar flags in database.
--variation IDENTITY
   Value for percentage identity variation that exemplar alleles
   cover (smaller value will result in more exemplars). Default: 10.
```

8.3 Automated offline allele definition

There is a script called 'scannew.pl' in the BIGSdb scripts/automation directory. This can be used to identify new alleles from the command line. This can (optionally) upload these to a sequence definition database.

Before scannew.pl can be run for the first time, a log file needs to be created. This can be created if it doesn't already exist with the following:

```
sudo touch /var/log/bigsdb_scripts.log
sudo chown bigsdb /var/log/bigsdb_scripts.log
```

The autotag.pl script should be installed in /usr/local/bin. It is run as follows:

```
scannew.pl --database <database configuration>
```

where <database configuration> is the name used for the argument 'db' when using the BIGSdb application.

If you have multiple processor cores available, use the –threads option to set the number of jobs to run in parallel. Loci for scanning will be split among the threads.

The script must be run by a user that can both write to the log file and access the databases, e.g. the 'bigsdb' user (see 'Setting up the offline job manager').

A full list of options can be found by typing:

```
scannew.pl --help
```

```
NAME
 scannew.pl - BIGSdb automated allele definer
SYNOPSIS
 scannew.pl --database NAME [options]
OPTIONS
-a, --assign
   Assign new alleles in definitions database.
--allow_frameshift
   Allow sequences to contain a frameshift so that the length is not a
   multiple of 3, or an internal stop codon. To be used with
   --coding_sequences option to allow automated curation of pseudogenes.
   New alleles assigned will be flagged either 'frameshift' or 'internal stop
   codon' if appropriate. Essentially, combining these two options only
   checks that the sequence starts with a start codon and ends with a stop
   codon.
-A, --alignment INT
    Percentage alignment (default: 100).
-B, --identity INT
   Percentage identity (default: 99).
-c, --coding_sequences
   Only return complete coding sequences.
-d, --database NAME
   Database configuration name.
-h, --help
   This help page.
-i, --isolates LIST
   Comma-separated list of isolate ids to scan (ignored if -p used).
--isolate_list_file FILE
   File containing list of isolate ids (ignored if -i or -p used).
-I, --exclude_isolates LIST
   Comma-separated list of isolate ids to ignore.
-1, --loci LIST
   Comma-separated list of loci to scan (ignored if -s used).
-L, --exclude_loci LIST
    Comma-separated list of loci to exclude.
-m, --min_size SIZE
   Minimum size of seqbin (bp) - limit search to isolates with at least this
   much sequence.
-n, --new_only
   New (previously untagged) isolates only.
-o, --order
   Order so that isolates last tagged the longest time ago get scanned first
```

```
(ignored if -r used).
-p, --projects LIST
   Comma-separated list of project isolates to scan.
-P, --exclude_projects LIST
   Comma-separated list of projects whose isolates will be excluded.
-r, --random
   Shuffle order of isolate ids to scan.
-R, --locus_regex REGEX
   Regex for locus names.
-s, --schemes LIST
   Comma-separated list of scheme loci to scan.
-t, --time MINS
    Stop after t minutes.
--threads THREADS
   Maximum number of threads to use.
-T, --already_tagged
    Scan even when sequence tagged (no designation).
-v, --view VIEW
   Isolate database view (overrides value set in config.xml).
-w, --word_size SIZE
   BLASTN word size.
-x, --min ID
   Minimum isolate id.
-y, --max ID
   Maximum isolate id.
```

8.4 Cleanly interrupting offline curation

Sometimes you may wish to stop running autotagger or allele autodefiner jobs as they can be run for a long time and as CRON jobs. If these are running in single threaded mode, the easiest way is to simply send a kill signal to the process, i.e. identify the process id using 'top', e.g. 23232 and then

```
kill 23232
```

The scripts should respond to this signal within a couple of seconds, clean up all their temporary files and write the history log (where appropriate). Do not use 'kill -9' as this will terminate the processes immediately and not allow them to clean up.

If these scripts are running using multiple threads, then you need to cleanly kill each of these. The simplest way to terminate all autotagger jobs is to, type

```
pkill autotag
```

The parent process will wait for all forked processes to cleanly terminate and then exit itself.

Similarly, to terminate all allele autodefiner jobs, type

```
pkill scannew
```

8.5 Uploading contigs from the command line

There is a script called upload_contigs.pl in the BIGSdb scripts/maintenance directory. This can be used to upload contigs from a local FASTA file for a specified isolate record.

The upload contigs.pl script should be installed in /usr/local/bin. It is run as follows:

The script must be run by a user who has the appropriate database permissions and the local configuration settings should be modified to match the database user account to be used. The default setting uses the 'apache' user which is used by the BIGSdb web interface.

A full list of options can be found by typing:

```
upload_contigs.pl --help
NAME
    upload_contigs.pl - Upload contigs to BIGSdb isolate database
SYNOPSIS
    upload_contigs.pl --database NAME --isolate ID --file FILE
         --curator ID --sender ID [options]
OPTIONS
-a, --append
   Upload contigs even if isolate already has sequences in the bin.
-c, --curator ID
   Curator id number.
-d, --database NAME
   Database configuration name.
-f, --file FILE
   Full path and filename of contig file.
-h, --help
   This help page.
-i, --isolate ID
    Isolate id of record to upload to.
-m, --method METHOD
   Method, e.g. 'Illumina', default 'unknown'.
--min_length LENGTH
   Exclude contigs with length less than value.
-s, --sender ID
    Sender id number.
```

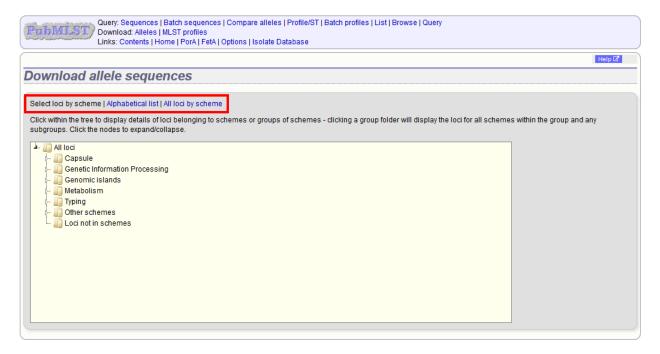
Definition downloads

The sequence definition database defines alleles, i.e. links an allele identifier to a sequence. It also defines scheme, e.g. MLST, profiles.

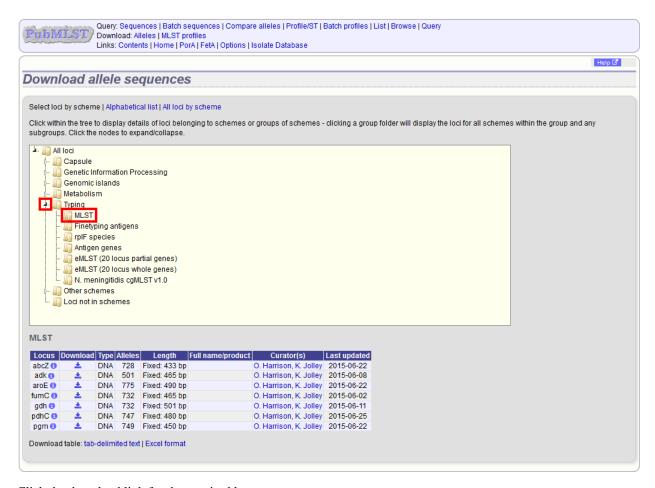
9.1 Allele sequence definitions

Click the 'Allele sequences' link in the 'Downloads' section. Depending on the database, you may see either a hierarchical scheme tree or a table of loci. You can choose to display links either by scheme using the scheme tree, as an alphabetical list or a page of all schemes, by selecting the approrpiate link at the top of the page.

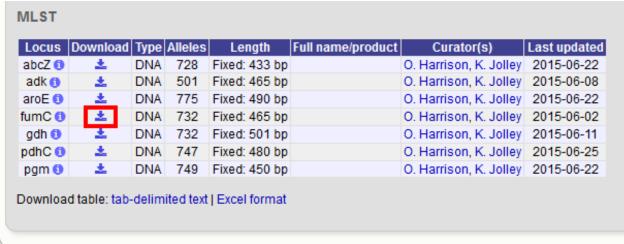
9.1.1 Scheme tree



You can drill down through the tree by clicking branch nodes. Clicking the labels of internal nodes will display tables of all schemes belonging to that scheme group. Clicking the labels of terminal nodes will display that single scheme table.



Click the download link for the required locus

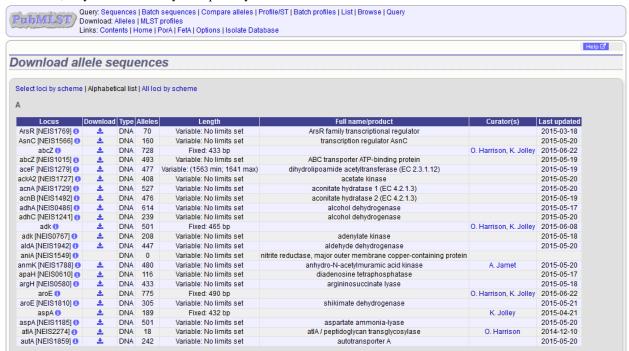


Alleles will be downloaded in FASTA format, e.g.

GCCTACAACCTCTTGCAATCCATCCGCCTGTTGGGCGACGCGTGCAACAGCTTCAACGAA CACTGCGCCGTCGGCATTGAACCCGTACCGGAAAAAATCGACTATTTCCTGCACCATTCC CTGATGCTCGTTACCGCGTTAAACCGCAAAATCGGTTACGAAAAC >fumC_2 GAAGCCTTGGGCGGACGCGATGCCGCCGTTGCCGCTTCGGGCGCATTGAAAACGCTGGCG TTGGGCGAAATCAAAATCCCCGAAAACGAGCCGGGTTCGTCCATCATGCCGGGCAAAGTC AACCCGACCCAATGCGAAGCGATGACCATGGTGTGCCCAAGTGTTCGGCAACGACGTT ACCATCGGCATGGCGGCGCGCGTCGGGCAATTTCGAGCTGAACGTCTATATGCCCGTTATC GCCTACAACCTCTTGCAATCCATCCGCCTCTTGGGCGACGCGTGCAACAGCTTCAACGAA CACTGCGCCATCGGCATCGAACCCGTACCGGAAAAAATCGACTATTTCCTGCACCATTCC CTGATGCTCGTTACCGCGTTAAACCGCAAAATCGGTTACGAAAAC >fumC 3 GAAGCCTTGGGCGGACGCGATGCCGCCGTTGCCGCTTCGGGCGCATTGAAAACGCTGGCG TTGGGCGAAATCAAAATCCCCGAAAACGAGCCGGGTTCGTCCATCATGCCGGGCAAAGTC AACCCGACCCAATGCGAAGCGATGACCATGGTGTGCCCAAGTGTTCGGCAACGACGTT ACCATCGGCATGGCGGCGCTCGGGCAATTTCGAGCTGAACGTCTATATGCCCGTTATC GCCTACAACCTCTTGCAATCCATCCGCCTGTTGGGCGACGCGTGCAACAGCTTCAACGAA CACTGCGCCGTCGGCATCGAACCCGTACCGGAAAAAATCGACTATTTCCTGCACCATTCC CTGATGCTGGTTACTGCGTTAAACCGTAAAATCGGCTACGAAAAC

9.1.2 Alphabetical list

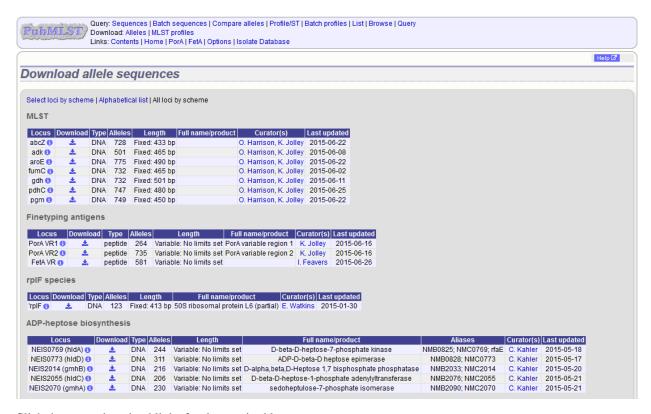
Loci can be displayed in an alphabetical list. Loci will be grouped in to tables by initial letter. If common names are set for loci, they will be listed by both primary and common names.



Click the download links for the required locus.

9.1.3 All loci by scheme

Loci can also be displayed by scheme with all schemes displayed.



Click the green download links for the required locus.

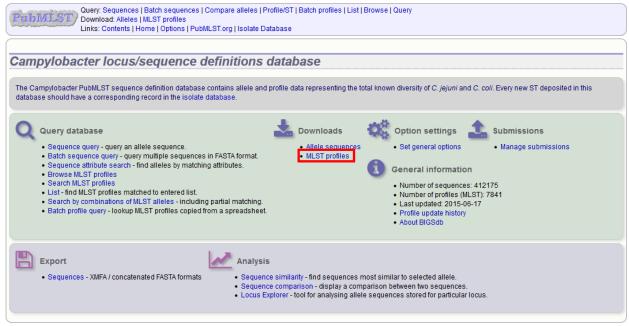
9.1.4 Download locus table

The locus table can be downloaded in tab-delimited text or Excel formats by clicking the links following table display.

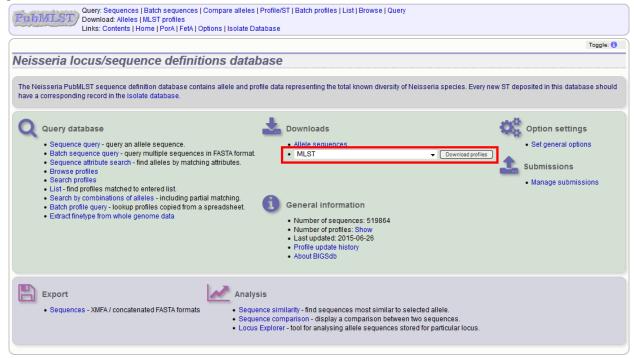


9.2 Scheme profile definitions

Scheme profiles, e.g. those for MLST, can be downloaded by clicking the appropriate link on the contents page.



If multiple schemes are available, you will need to select the scheme in the dropdown box and click 'Download profiles'



Profiles will be downloaded in tab-delimited format, e.g.

ST	abcZ	adk	aroE	fumC	gdh	pdhC	pgm	clonal_complex
1	1	3	1	1	1	1	3	ST-1 complex/subgroup I/II

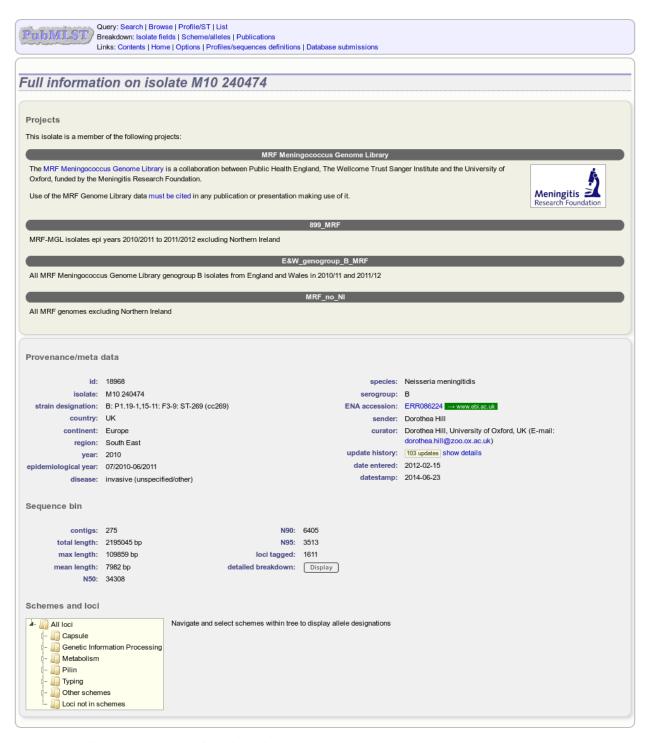
2	1	3	4	7	1	1	3	ST-1 complex/subgroup I/II
3	1	3	1	1	1	23	13	ST-1 complex/subgroup I/II
4	1	3	3	1	4	2	3	ST-4 complex/subgroup IV
5	1	1	2	1	3	2	3	ST-5 complex/subgroup III
6	1	1	2	1	3	2	11	ST-5 complex/subgroup III
7	1	1	2	1	3	2	19	ST-5 complex/subgroup III
8	2	3	7	2	8	5	2	ST-8 complex/Cluster A4
9	2	3	8	10	8	5	2	ST-8 complex/Cluster A4
10	2	3	4	2	8	15	2	ST-8 complex/Cluster A4
11	2	3	4	3	8	4	6	ST-11 complex/ET-37 complex
12	4	3	2	16	8	11	20	
13	4	10	15	7	8	11	1	ST-269 complex
14	4	1	15	7	8	11	1	ST-269 complex

Data records

Record pages for different types of data can be accessed following a query by clicking appropriate hyperlinks.

10.1 Isolate records

An Isolate record page displays everything known about an isolate.



Each record will have some or all of the following sections:

10.1.1 Projects



This displays a list of projects that the isolate is a member of. Only projects that have a full description will be displayed.

10.1.2 Provenance metadata

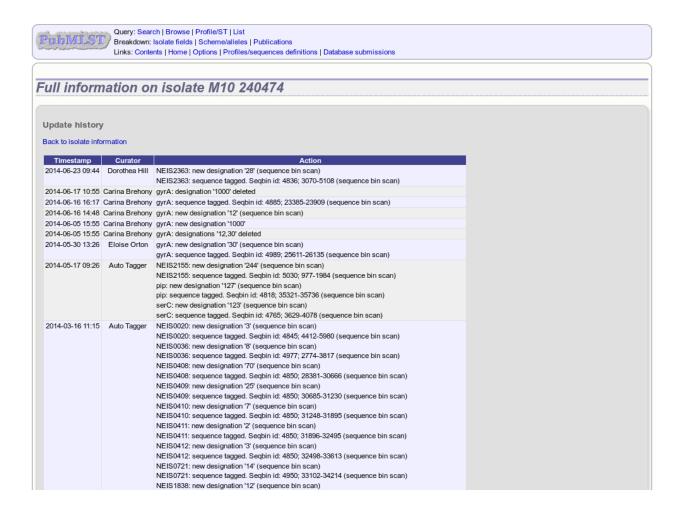


This section includes:

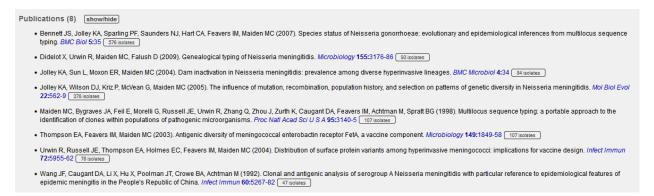
- · provenance fields
- · housekeeping data
 - who sent the isolate
 - who last curated
 - record creation times
 - last update times
 - links to update history

The update link displays page with exact times of who and when updated the record.

10.1. Isolate records



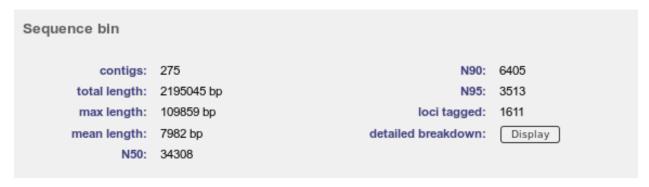
10.1.3 Publications



This section includes full citation for papers linked to the isolate record. Each citation has a button that will return a dataset of all isolates linked to the paper.

If there are five or more references they will be hidden by default to avoid cluttering the page too much. Click the 'Show/hide' button to display them in this case.

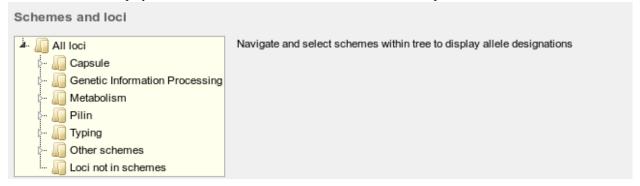
10.1.4 Sequence bin summary



This section contains basic statistics describing the sequence bin. Clicking the 'Display' button navigates to the *sequence bin record*.

10.1.5 Scheme and locus data

A hierarchical tree displays available schemes. Click within internal nodes to expand them.



Clicking any terminal node will display data available for a scheme or group of schemes.



Click an allele number within the scheme profile, will display the appropriate *allele definition record*. Clicking the green 'S' link will display the appropriate *sequence tag record*.

10.2 Allele definition records

An allele definition record displays information about a defined allele in a sequence definition database.



If the allele is a member of a scheme profile, e.g. MLST, this will be listed. In this case, there will be a button to display all profiles of that scheme that contain the allele.

Similarly, if a *client database* has been setup for the database and the allele has been identified in an isolate, there will be a button to display all isolates that have that allele.

10.3 Sequence tag records



A sequence tag record displays information about the location within a contig of a region associated with a locus. The nucleotide sequence will be displayed along with upstream and downstream flanking sequence. The length of these flanking sequences can be modified within the *general options*.

If the tag is for a DNA locus and it is marked as a coding sequence, the three-frame translation will also be displayed.

10.4 Profile records

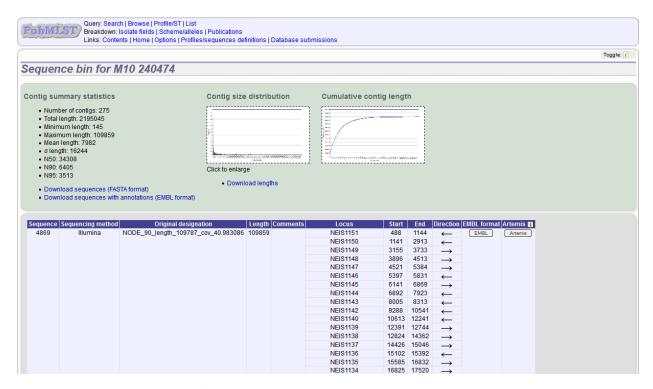


A profile record displays information about a scheme, e.g. MLST, profile. Each allele number within the profile will be hyperlinked. Clicking these will take you to the appropriate *allele definition record*.

If a *client database* has been setup for the database and an isolate has the profile, there will be a button to display all isolates that have the profile.

10.4. Profile records 203

10.5 Sequence bin records



A sequence bin record contains information about that contigs associated with an isolate record. This includes:

- · Number of contigs
- · Total length
- Minimum length
- Maximum length
- N50, N90 and N95 values
- Size distribution charts

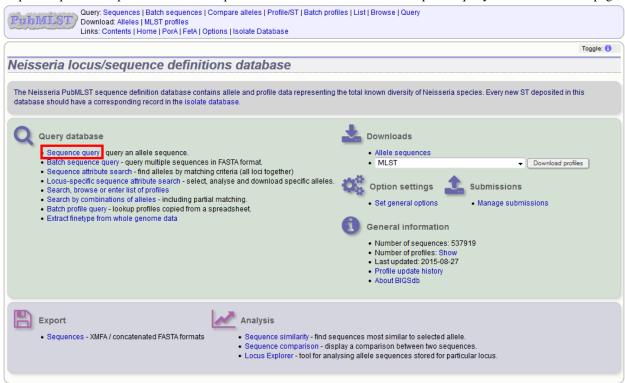
There are also links to download the contigs in FASTA or EMBL format.

Finally there is a table that shows the loci that are tagged on each contig. Individual contigs can also be downloaded in EMBL format.

Querying data

11.1 Querying sequences to determine allele identity

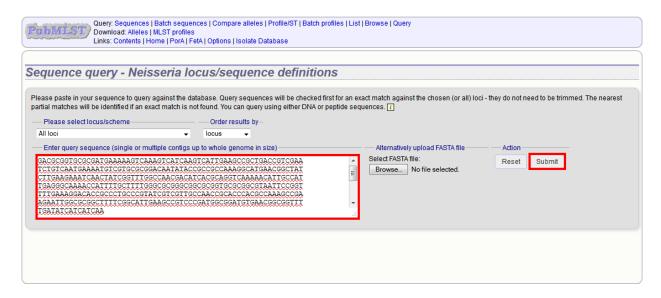
Sequence queries are performed in the sequence definition database. Click 'Sequence query' from the contents page.



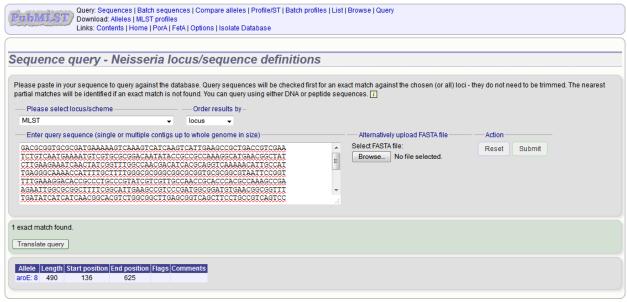
Paste your sequence in to the box - there is no need to trim. Normally, you can leave the locus setting on 'All loci' - the software should identify the correct locus based on your sequence. Sometimes, it may be quicker, however, to select the specific locus or scheme (e.g. MLST) that a locus belongs to.

Note: If the locus you are querying is a shorter version of another, e.g. an MLST fragment of a gene where the full length gene is also defined, you will need to select the specific locus or the scheme from the dropdown box. Leaving the selection on 'All loci' will return a match to the longer sequence in preference to the shorter one.

Click 'Submit'.



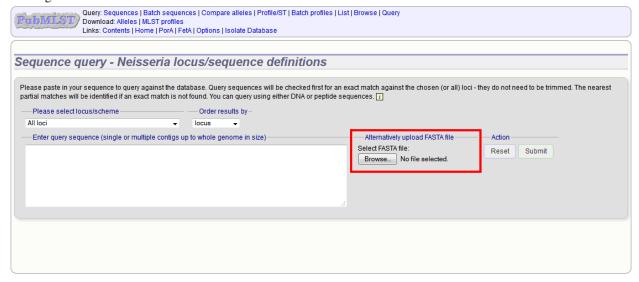
If an exact match is found, this will be indicated along with the start position of the locus within your sequence.



If only a partial match is found, the most similar allele is identified along with any nucleotide differences. The varying nucleotide positions are numbered both relative to the pasted in sequence and to the reference sequence. The start position of the locus within your sequence is also indicated.

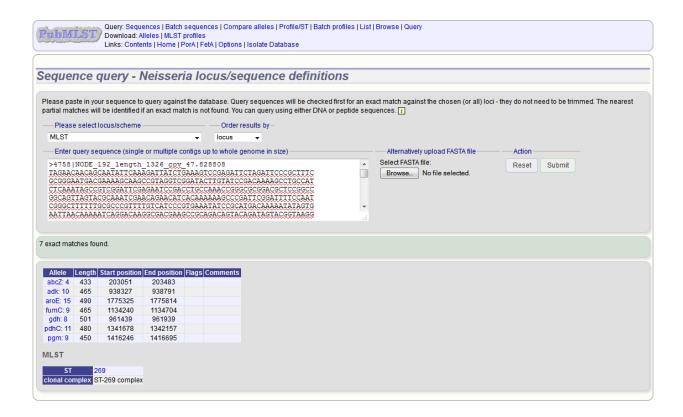


As an alternative to pasting a sequence in to the box, you can also choose to upload sequences in FASTA format by clicking the file browse button.



11.1.1 Querying whole genome data

The sequence query is not limited to single genes. You can also paste or upload whole genomes - these can be in multiple contigs. If you select a specific scheme from the dropdown box, all loci belonging to that scheme will be checked (although only exact matches are reported for a locus if one of the other loci has an exact match). If all loci are matched, scheme fields will also be returned if these are defined. This, for example, enables you to identify the MLST sequence type of a genome in one step.

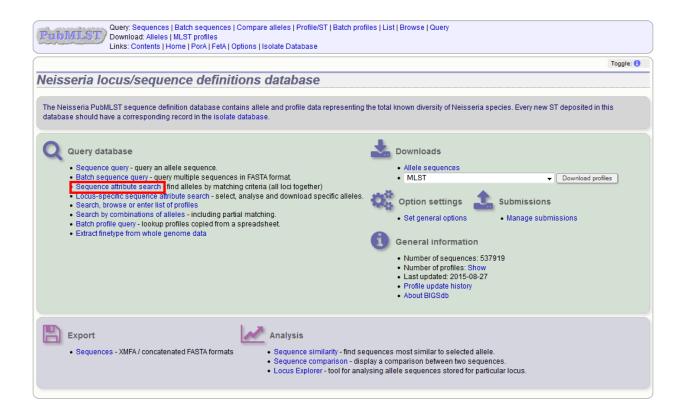


11.2 Searching for specific allele definitions

There are two query pages available that allow searching for specific allele definitions. The first allows querying of all loci together by criteria that are common to all. The second is a locus-specific attribute query that can search on any extended attributes that may be defined for a locus. This locus-specific query also allows you to paste in lists of alleles for download or analysis.

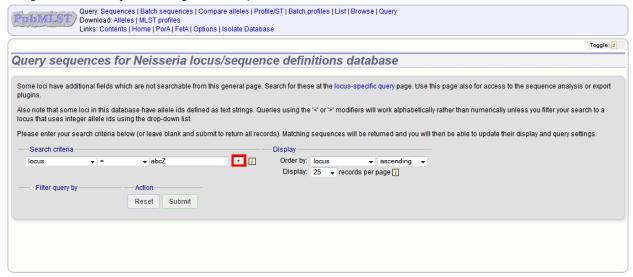
11.2.1 General (all loci) sequence attribute search

To retrieve specific allele designations, click 'Sequence attribute search' on a sequence definition database contents page.



Enter your query using the dropdown search box - additional terms can be added by clicking the '+' button.

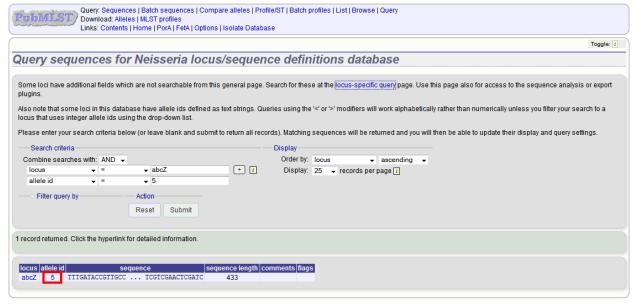
Designations can be queried using *standard operators*.



Click submit.

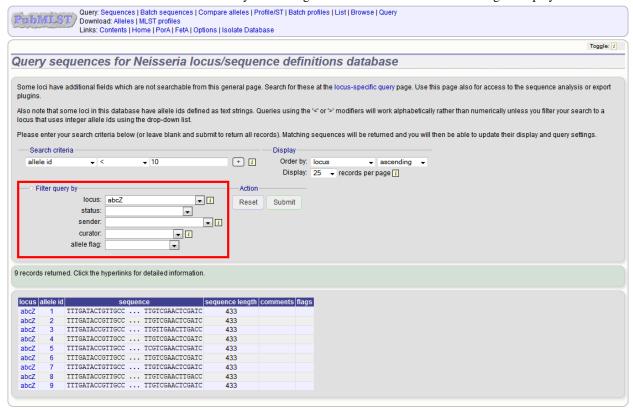


Click the hyperlinked results to display allele records.



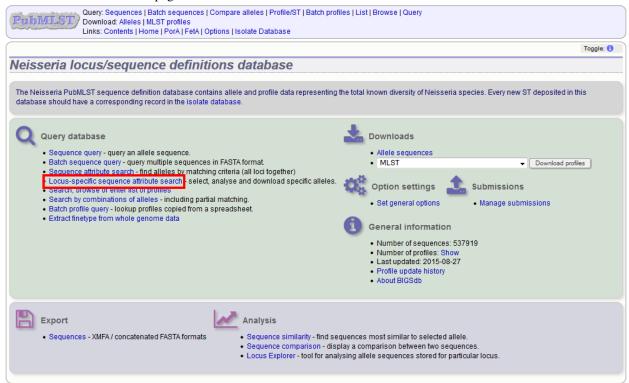


Various search criteria can also be selected by combining with filters. Click the filter heading to display these.

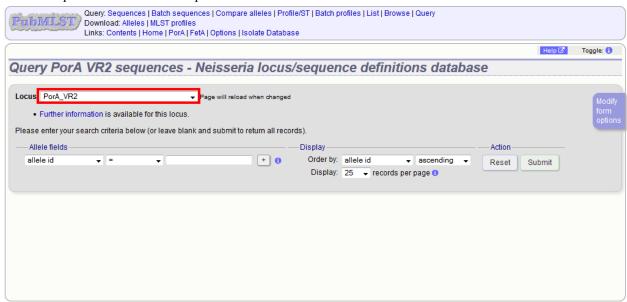


11.2.2 Locus-specific sequence attribute search

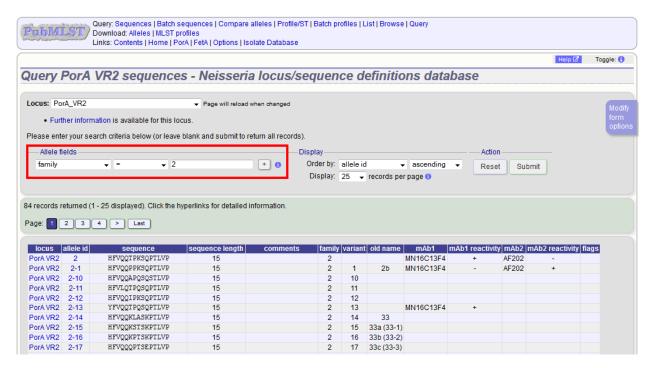
Some loci have *extended attribute fields*. To query these, click 'Locus-specific sequence attribute search' on a sequence definition database contents page.



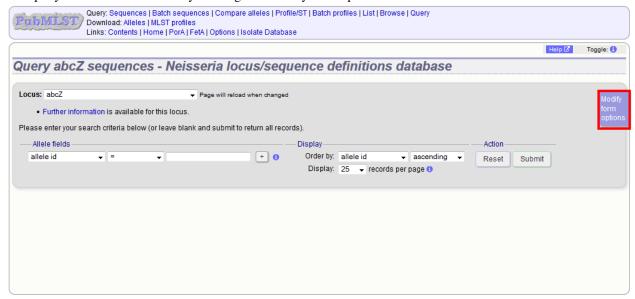
Pick the required locus from the dropdown box.



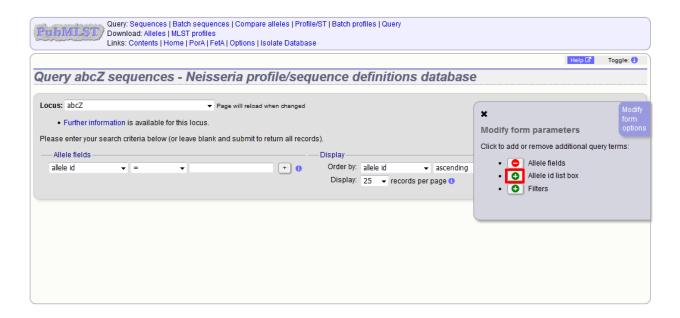
The fields specific for that locus will be added to the dropdown query boxes.



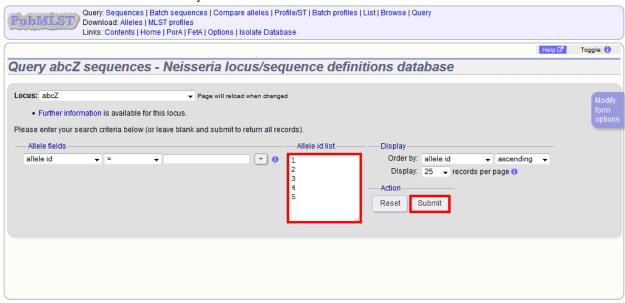
The query form can be modified by clicking the 'Modify form options' tab:



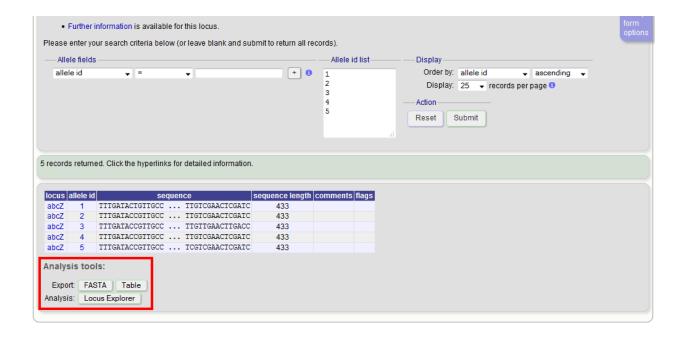
A list box can be added by clicking the 'Show' button for 'Allele id list box'.



Close the form modification tab and you can now enter a list of allele ids for retrieval.

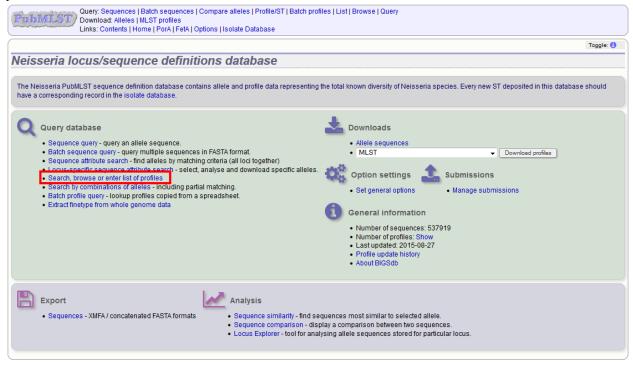


Various analysis and export options will be available for use on the retrieved sequences. These include FASTA output and *Locus Explorer* analysis.

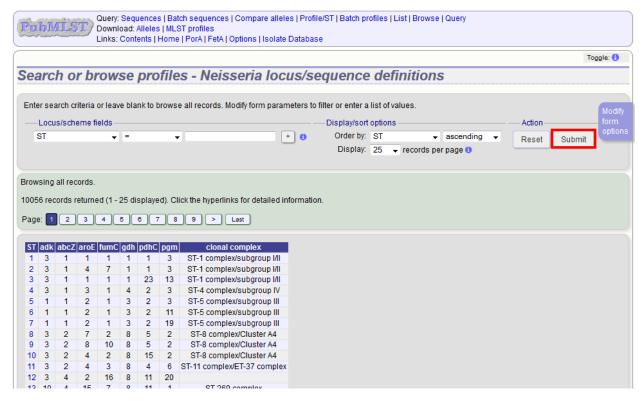


11.3 Browsing scheme profile definitions

If a sequence definition database has schemes defined that include a primary key field, i.e. collections of loci that together create profiles, e.g. for MLST, these can be browsed by clicking the link to 'Search, browse or enter list of profiles'.



Leave query form fields blank (the display of these may vary depending on modification options set by the user). Choose the field to order the results by, the number of results per page to display, and click 'Submit'.

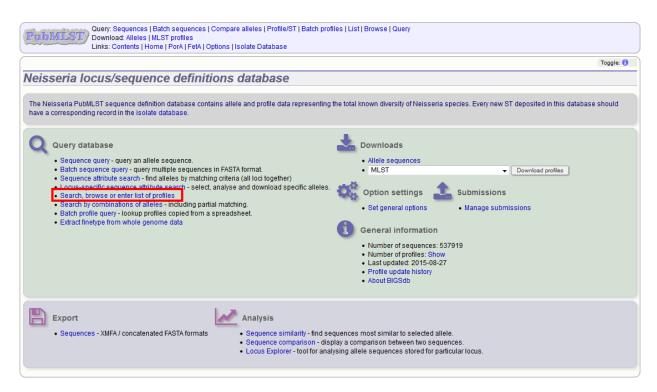


Clicking the hyperlink for any profile will display full information about the profile.

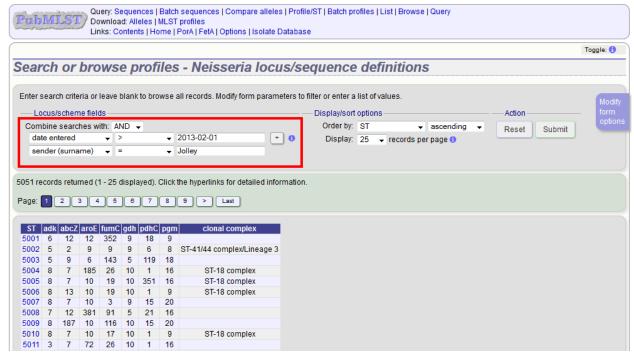


11.4 Querying scheme profile definitions

Click the link to 'Search, browse of enter list of profiles'.

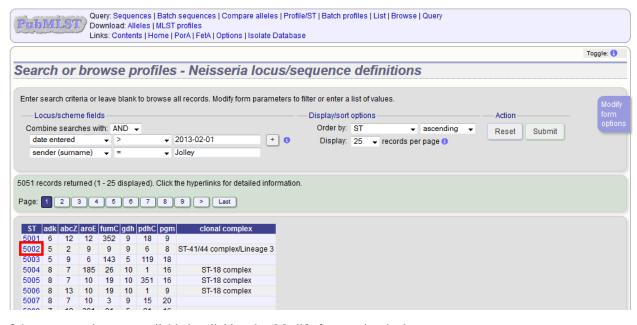


Enter the search criteria you wish to search on. You can add search criteria by clicking the '+' button in the 'Locus/scheme fields' section. These can be combined using 'AND' or 'OR'.



Each field can be queried using standard operators.

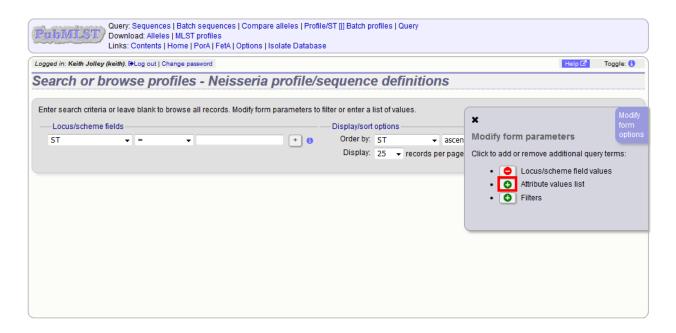
Clicking the hyperlink for any profile will display full information about the profile.



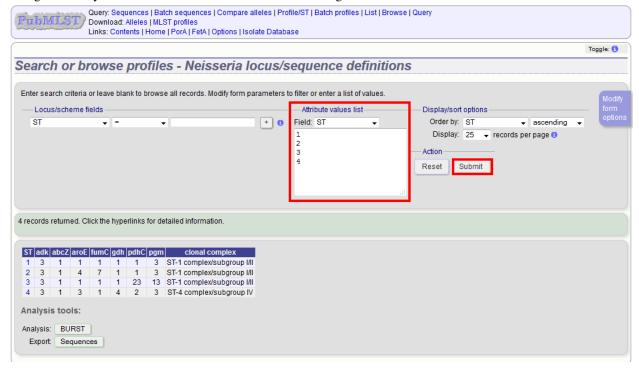
Other query options are available by clicking the 'Modify form options' tab.



For example, you can enter a list of attributes to query on by clicking the 'Show' button next to 'Attribute values list'.

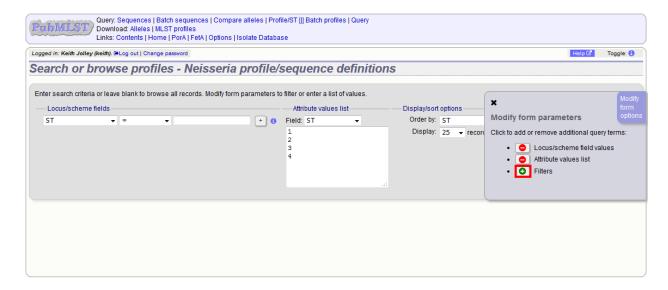


A list box will appear within the page. Hide the form modification tab by clicking the 'X' in the corner or the purple tab again. Now you can choose the attribute to search on along with a list of values.

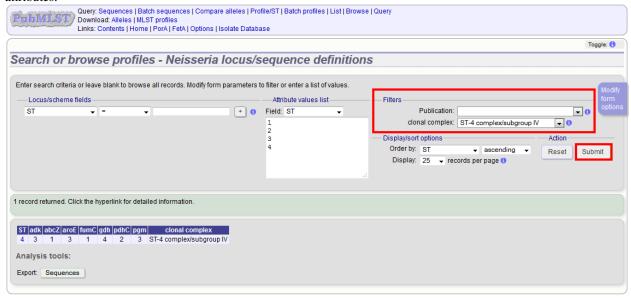


List values will be combined with any other attributes entered in the query form allowing complex queries can be constructed.

You can also add filters to the form by again clicking the 'Modify form options' tab and selecting 'Filters'.



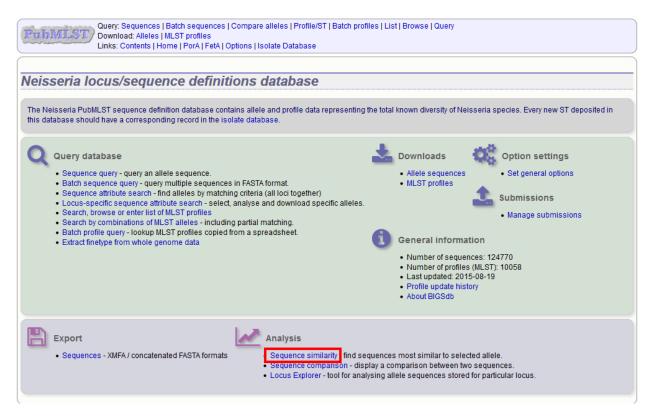
Available filters will vary depending on the database. These will be combined with other query criteria or lists of attributes.



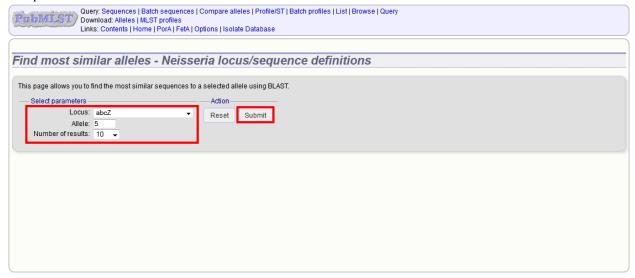
11.5 Investigating allele differences

11.5.1 Sequence similarity

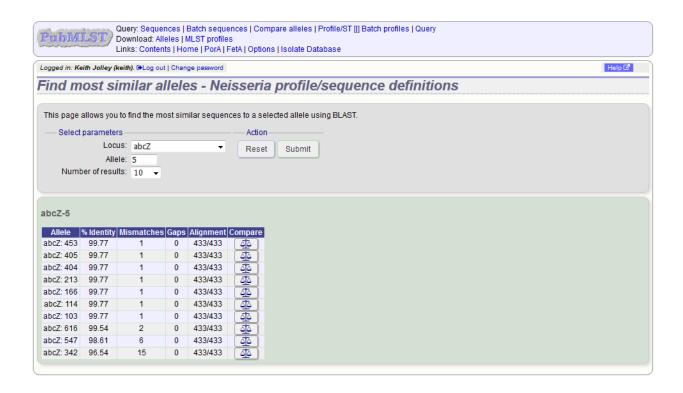
To find sequences most similar to a selected allele within a sequence definition database, click 'Sequence similarity' on the contents page.



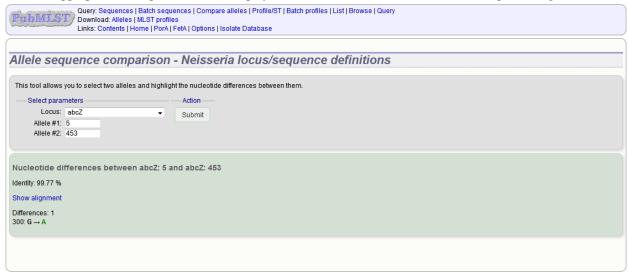
Enter the locus and allele identifer of the sequence to investigate and the number of nearest matches you'd like to see, then press submit.



A list of nearest alleles will be displayed, along with the percentage identity and number of gaps between the sequences.

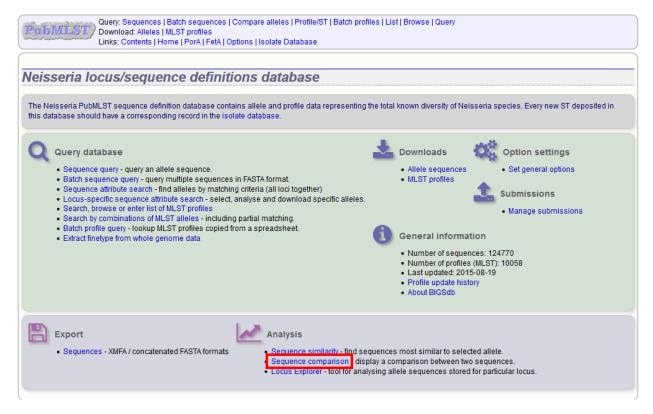


Click the appropriate 'Compare' button to display a list of nucleotide differences and/or a sequence alignment.

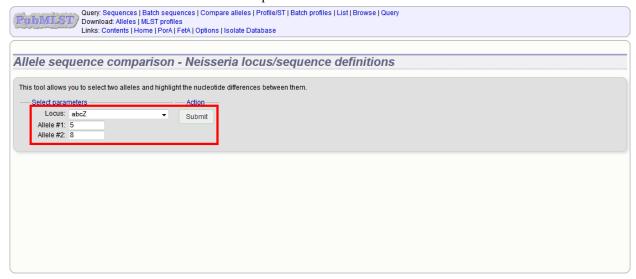


11.5.2 Sequence comparison

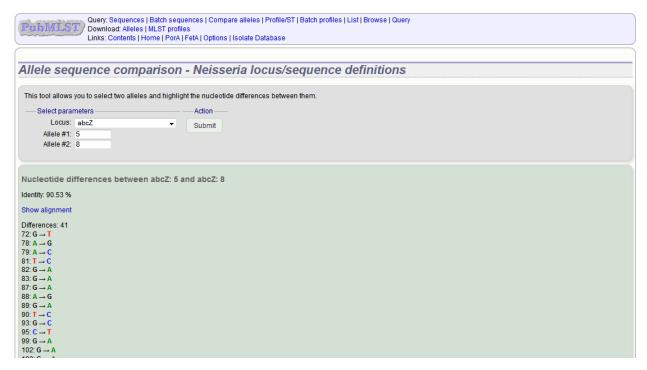
To directly compare two sequences click 'Sequence comparison' from the contents page of a sequence definition database.



Enter the locus and two allele identifiers to compare. Press submit.



A list of nucleotide differences and/or an alignment will be displayed.

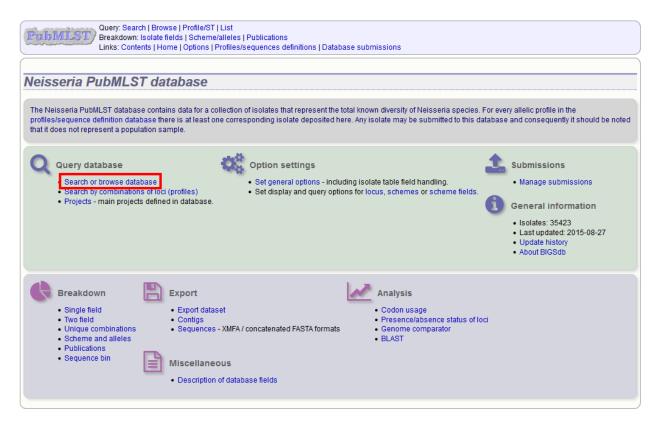


See also:

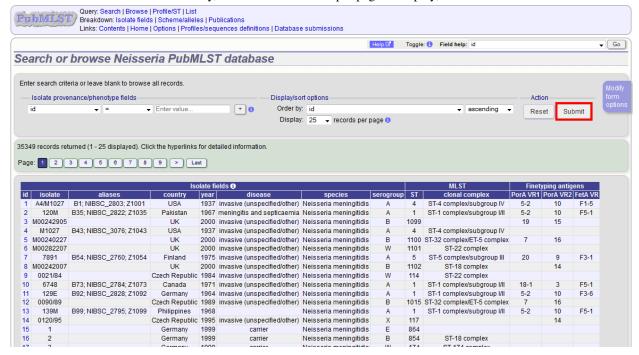
Locus explorer plugin.

11.6 Browsing isolate data

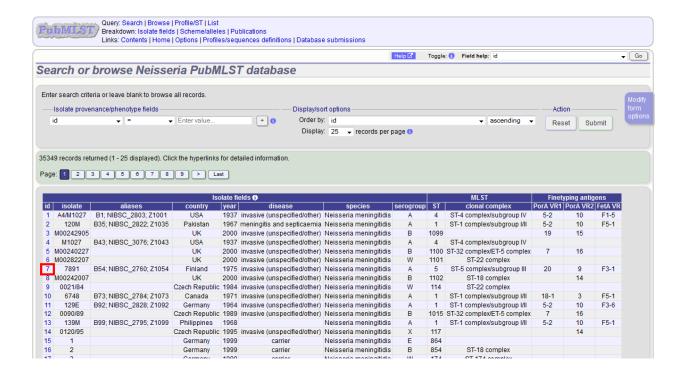
Isolate records can be browsed by clicking the link to 'Search or browse database'.



Leave query form fields blank (the display of these may vary depending on modification options set by the user). Choose the field to order the results by, the number of results per page to display, and click 'Submit'.

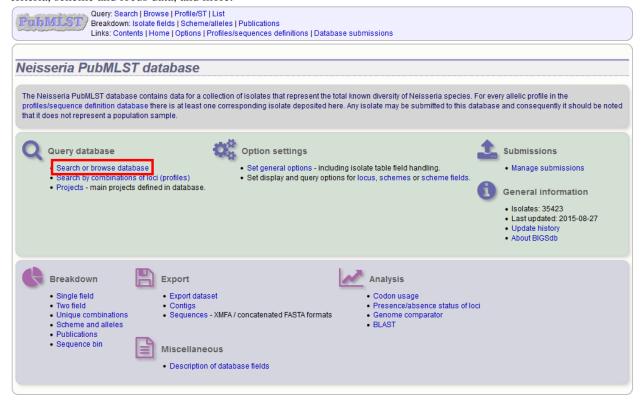


Clicking the hyperlink for any record will display full information about the profile.



11.7 Querying isolate data

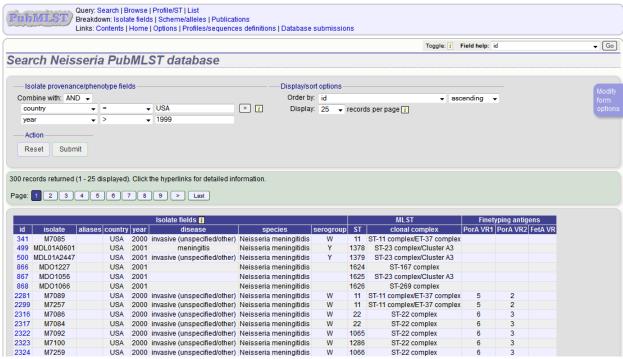
The 'Search or browse database' page of an isolate database allows you to also search by combinations of provenance criteria, scheme and locus data, and more.



To start with, only one provenance field search box is displayed but more can be added by clicking the '+' button (highlighted). These can be linked together by 'and' or 'or'.

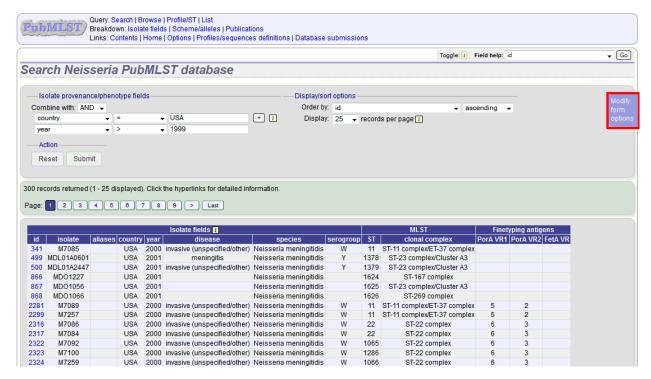


After the search has been submitted, the results will be displayed in a table.



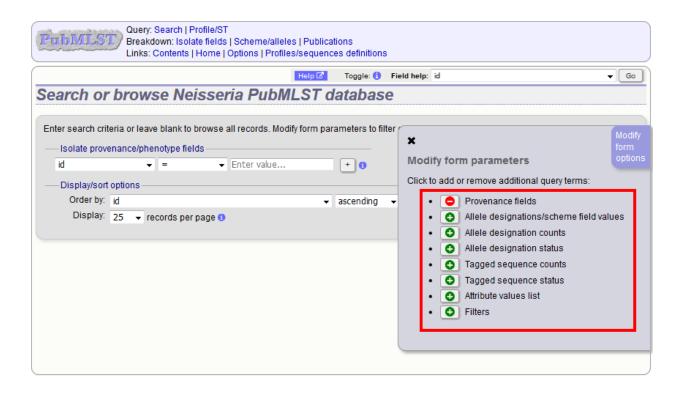
Each field can be queried using *standard operators*.

More search features are available by clicking the 'Modify form options' tab on the right-hand side of the screen.



A tab will be displayed. Different options will be available here depending on the database. Queries will be combined from the values entered in all form sections. Possible options are:

- · Provenance fields
 - Search by combination of provenance field values, e.g. country, year, sender.
- Allele designations/scheme field values
 - Search by combination of allele designations and/or scheme fields e.g. ST, clonal complex information.
- Allele designation status
 - Search by whether allele designation status is confirmed or provisional.
- Tagged sequence status
 - Search by whether tagged sequence data is available for a locus. You can also search by sequence flags.
- · Attribute values list
 - Enter a list of values for any provenance field, locus, or scheme field.
- Filters
 - Various filters may be available, including
 - * Publications
 - * Projects
 - * MLST profile completion status
 - * Clonal complex
 - * Sequence bin size
 - * Inclusion/exclusion of old versions

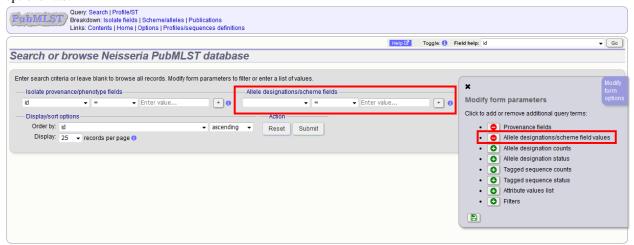


If the interface is modified, a button to save options becomes available within the tab. If this is clicked, the modified form will be displayed the next time you go to the query page.

11.7.1 Query by allele designation/scheme field

Queries can be combined with allele designation/scheme field values.

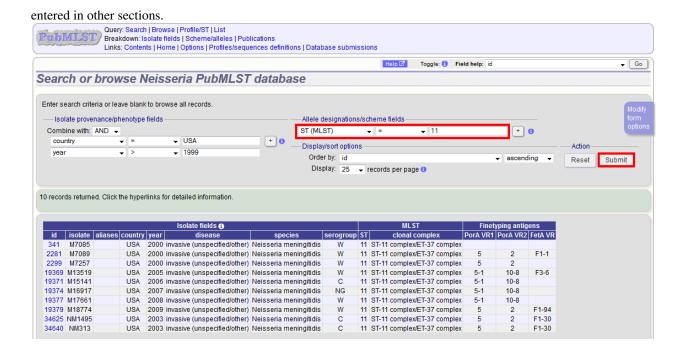
Make sure that the allele designation/scheme field values fieldset is displayed by selecting it in the 'Modify form options' tab.



Designations can be queried using standard operators.

Additional search terms can be combined using the '+' button.

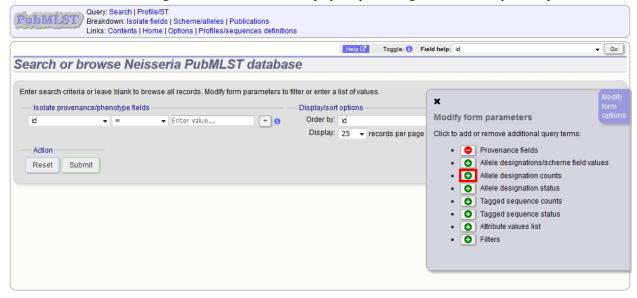
Add your search terms and click 'Submit'. Allele designation/scheme field queries will be combined with terms



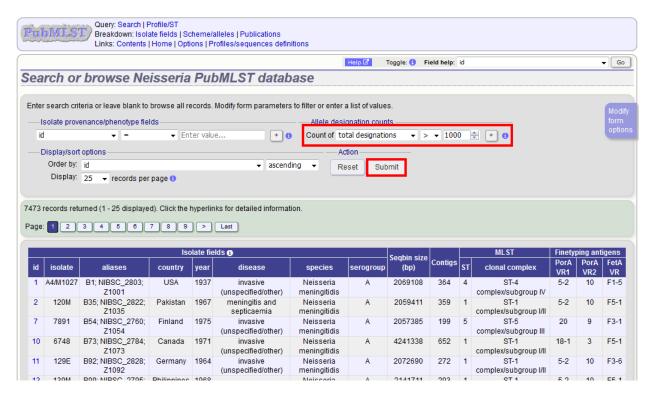
11.7.2 Query by allele designation count

Queries can be combined with counts of the total number of designations or for individual loci.

Make sure that the allele designation counts fieldset is displayed by selecting it in the 'Modify form options' tab.

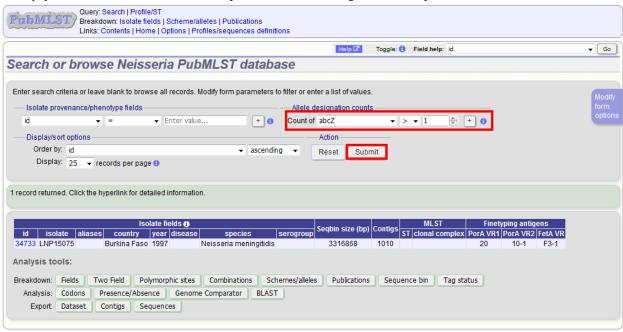


For example, to find all isolates that have designations at >1000 loci, select 'total designations > 1000', then click 'Submit'.



You can also search for isolates where any isolate has a particular number of designations. Use the term 'any locus' to do this.

Finally, you can search for isolates with a specific number of designations at a specific locus.



Additional search terms can be combined using the '+' button. Designation count queries will be combined with terms entered in other sections.

Note: Searches for 'all loci' with counts that include zero, e.g. 'count of any locus = 0' or with a '<' operator are not supported. This is because such searches have to identify every isolate for which one or more loci are missing. In databases with thousands of loci this can be a very expensive database query.

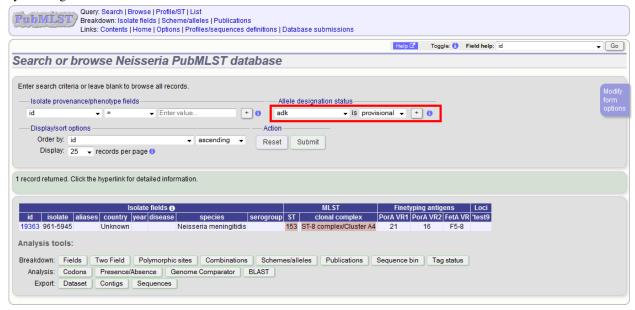
11.7.3 Query by allele designation status

Allele designations can be queried based on their status, i.e. whether they are confirmed or provisional. Queries will be combined from the values entered in all form sections.

Make sure that the allele designation staus fieldset is displayed by selecting it in the 'Modify form options' tab.



Select a locus from the dropdown box and either 'provisional' or 'confirmed'. Additional query fields can be displayed by clicking the '+' button. Click 'Submit'.

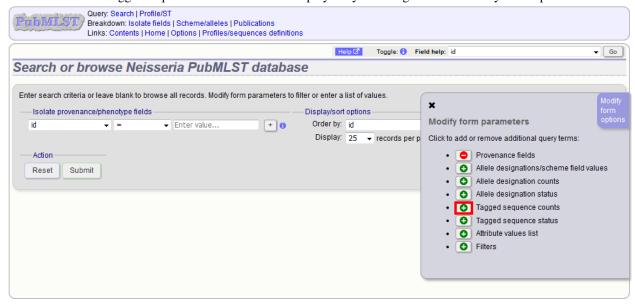


Provisional allele designations are marked within the results tables with a pink background. Any scheme field designations that depend on the allele in question, e.g. a MLST ST, will also be marked as provisional.

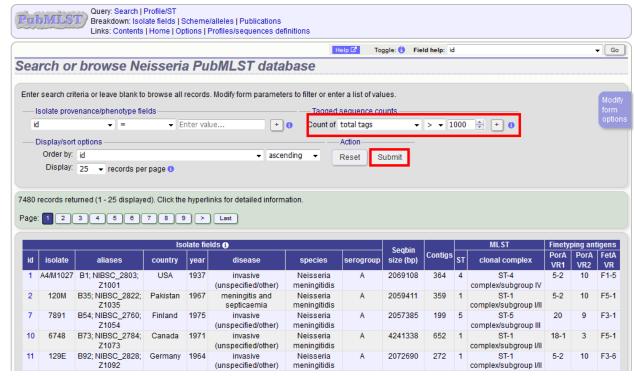
11.7.4 Query by sequence tag count

Queries can be combined with counts of the total number of tags or for individual loci.

Make sure that the tagged sequence counts fieldset is displayed by selecting it in the 'Modify form options' tab.



For example, to find all isolates that have sequence tags at >1000 loci, select 'total tags > 1000', then click 'Submit'.



You can also search for isolates where any isolate has a particular number of sequence tags. Use the term 'any locus' to do this.

Finally, you can search for isolates with a specific number of tags at a specific locus.



Additional search terms can be combined using the '+' button. Sequence tag count queries will be combined with terms entered in other sections.

Note: Searches for 'all loci' with counts that include zero, e.g. 'count of any locus = 0' or with a '<' operator are not supported. This is because such searches have to identify every isolate for which one or more loci are not tagged. In databases with thousands of loci this can be a very expensive database query.

11.7.5 Query by tagged sequence status

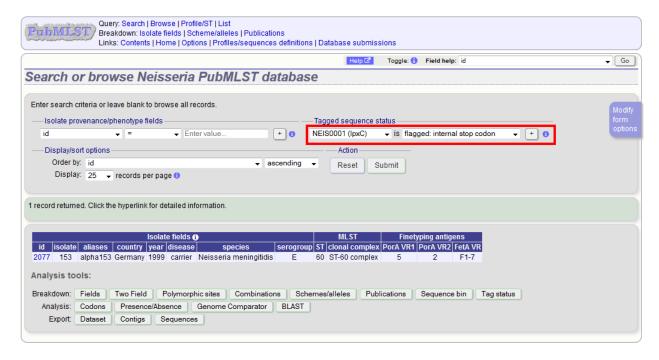
Sequence tags identify the region of a contig within an isolate's sequence bin entries that correspond to a particular locus. The presence or absence of these tags can be queried as can whether or not the sequence has an a flag associated with. These flags designate specific characteristics of the sequences. Queries will be combined from the values entered in all form sections.

Make sure that the tagged sequences status fieldset is displayed by selecting it in the 'Modify form options' tab.



Select a specific locus in the dropdown box (or alternatively 'any locus') and a status. Available status values are:

- · untagged
 - The locus has not been tagged within the sequence bin.
- · tagged
 - The locus has been tagged within the sequence bin.
- complete
 - The locus sequence is complete.
- incomplete
 - The locus sequence is incomplete normally because it continues beyond the end of a contig.
- · flagged: any
 - The sequence for the locus has a flag set.
- · flagged: none
 - The sequence for the locus does not have a flag set.
- flagged: <specific flag>
 - The sequence for the locus has the specific flag chosen.



See also:

Sequence tag flags

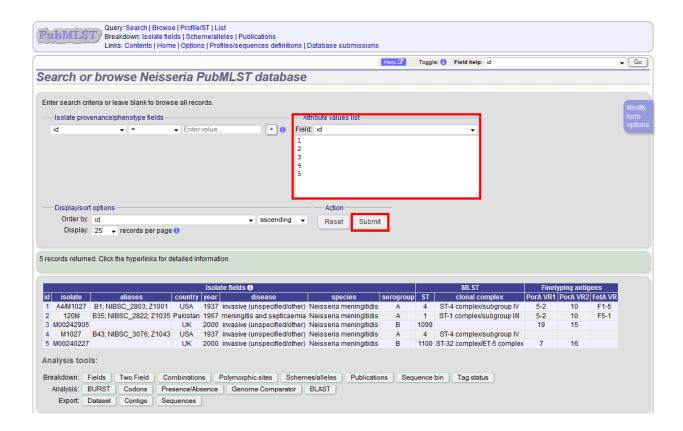
11.7.6 Query by list of attributes

The query form can be modified with a list box in to which a list of values for a chosen attribute can be entered - this could be a list of ids, isolate names, alleles or scheme fields. This list will be combined with any other criteria or filter used on the page.

If the list box is not shown, add it by selecting it in the 'Modify form options' tab.



Select the attribute to query and enter a list of values.



11.7.7 Query filters

There are various filters that can additionally be applied to queries, or the filters can be applied solely on their own so that they filter the entire database.

Make sure that the filters fieldset is displayed by selecting it in the 'Modify form options' tab.



The filters displayed will depend on the database and what has been defined within it. Common filters are:

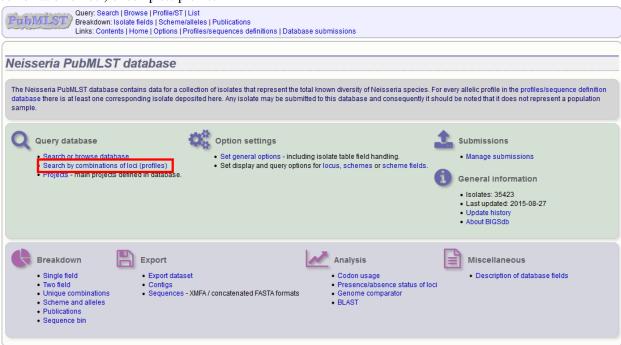
- Publication Select one or more publication that has been linked to isolate records.
- Project Select one or more project that isolates belong to.

- Profile completion This is commonly displayed for MLST schemes. Available options are:
 - complete All loci of the scheme have alleles designated.
 - incomplete One or more loci have not yet been designated.
 - partial The scheme is incomplete, but at least one locus has an allele designated.
 - started At least one locus has an allele designated. The scheme mat be complete or partial.
 - not started The scheme has no loci with alleles designated.
- Sequence bin Specify whether any sequence data has been associated with a record. Specific threshold values may be selected if these have been *set up for the database*.
- Provenance fields Dropdown list boxes of values for specific provenance fields may be present if set for the database. Users can choose to *add additional filters*.

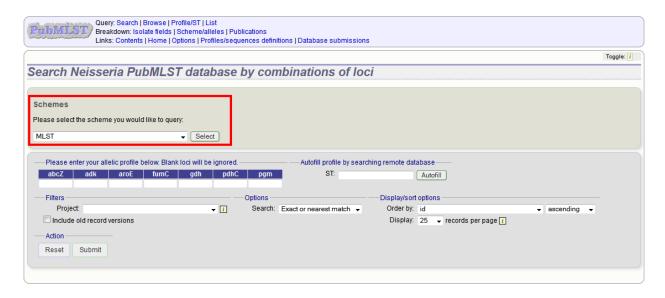
11.8 Querying by allelic profile

If a scheme, such as MLST, has been defined for an isolate database it is possible to query the database against complete or partial allelic profiles. Even if no scheme is defined, queries can be made against all loci. This can also be done in sequence definition databases if the scheme has a primary key field defined.

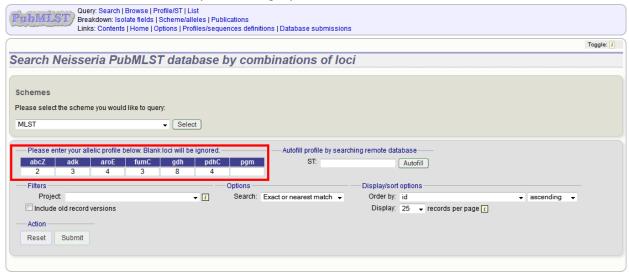
On the index page, click 'Search by combinations of loci (profiles)' for any defined scheme. Enter either a partial (any combination of loci) or complete profile.



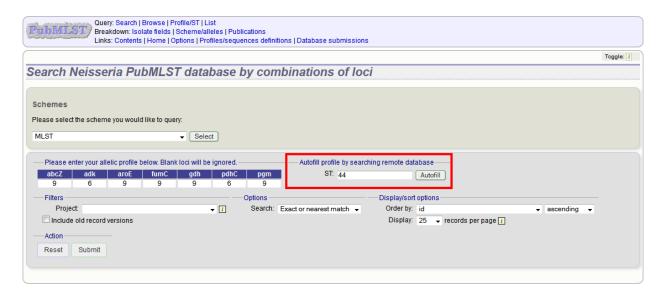
If multiple schemes are defined, you may have to select the scheme you wish to query in the 'Schemes' dropdown box and click 'Select'.



Enter the combination of alleles that you want to query for. Fields can be left blank.



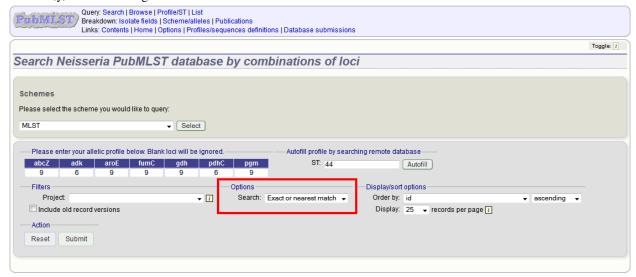
Alternatively, for scheme profiles, you can enter a primary key value (e.g. ST) and select 'Autofill' to automatically fill in the associated profile.



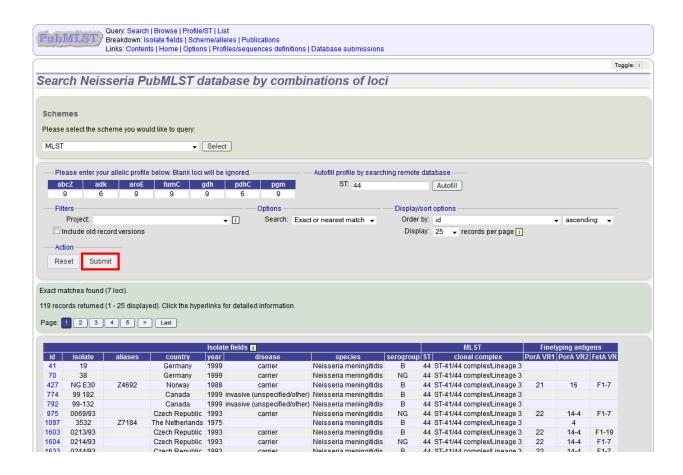
Select the number of loci that you'd like to match in the options dropdown box. Available options are:

- · Exact or nearest match
- · Exact match only
- x or more matches
- y or more matches
- · z or more matches

Where x,y, and z will range from n-1 to 1 where n is the number of loci in the scheme.

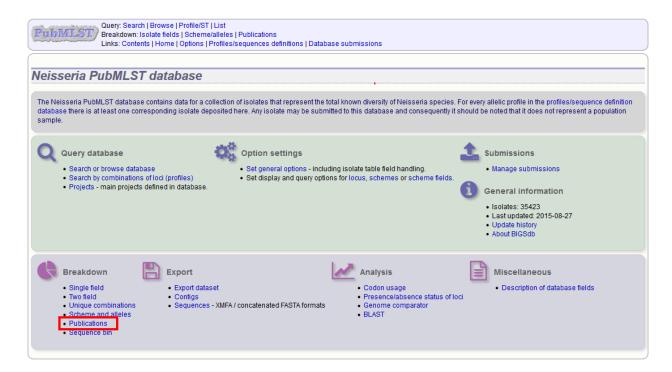


Click 'Submit'.

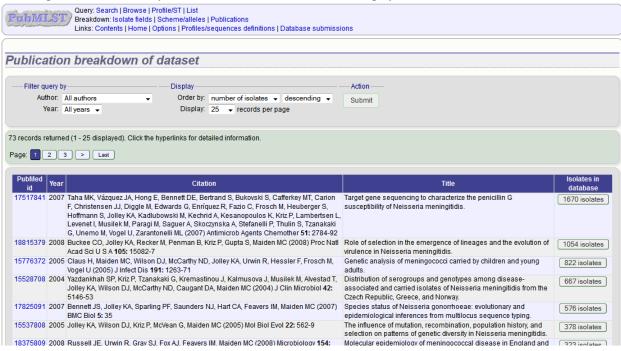


11.9 Retrieving isolates by linked publication

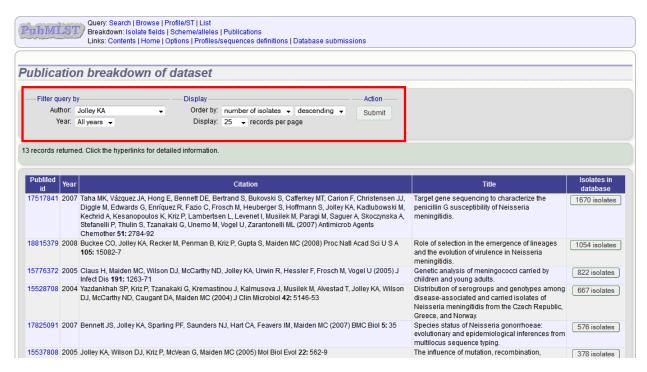
Click 'Publications' in the Breakdown section of the contents page.



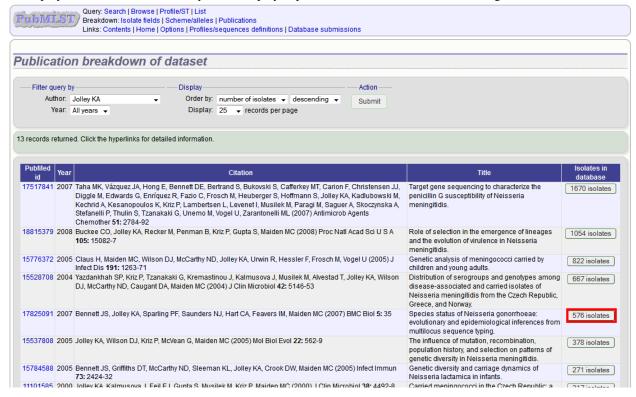
A list of publications linked by isolates within the database will be displayed.



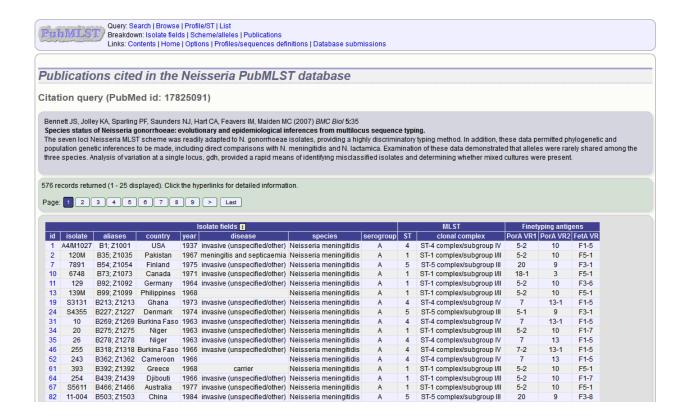
These can be filtered by author and/or year, and the sort order changed.



To display the isolate records for any of the displayed publications, click the button to the right of the citation.



The abstract of the paper will be displayed (if available), along with all isolates linked to it.



11.10 User-configurable options

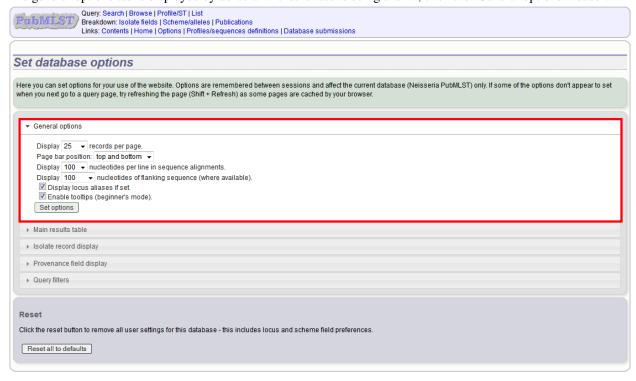
The BIGSdb user interface is configurable in a number of ways. Choices made are remembered between sessions. If the database requires you to log on, the options are associated with your user account, whereas if it is a public database, that you haven't logged in to, the options are associated with a browser cookie so they will be remembered if you connect from the same computer (using the same browser).

Most options are set by clicking the 'Set general options' link on the database contents page. Most of the available options are visible for isolate databases, whereas sequence definition databases have fewer available.



11.10.1 General options

The general options tab is displayed by default. If another tab is being shown, click the 'General options' header.



The general tab allows the following options to be modified:

• Records per page

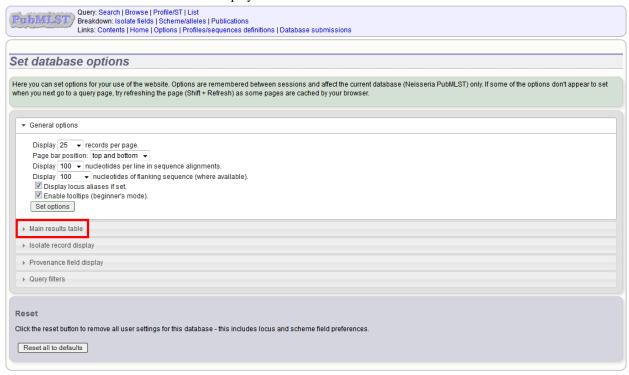
- Page bar position
- Nucleotides per line Some analyses display sequence alignments. This option allows you to set the width of these alignments so suit your display.
- Flanking sequence length This sets the length of flanking sequence upstream and downstream of a particular locus that is included whenever a sequence is displayed. Flanking sequences are displayed fainter that the locus sequence.
- Locus aliases Loci can have multiple names (aliases). Setting this option will display all alternative names in results tables.
- Tooltips (beginner's mode) Most query forms have help available in the form of information tooltips. These can be switched on/off here. They can also be toggled off by clicking the Toggle: 'i' button at the top-right of the display of some pages.

Click 'Set options' to remember any changes you make.

11.10.2 Main results table

The 'main results table' tab contains options for the display of paged results following a query.

Click the 'Main results table' header to display the tab.



The 'main results table' tab will scroll up.



This tab allows the following options to be modified:

- Hyperlink allele designations Hyperlinks point to an information page about the particular allele definition. Depending on the locus, these may exist on a different website.
- Differentiate provisional allele designations Allele designations can be set as confirmed or provisional, usually
 depending on the method of assignment. Selecting this option will display provisional designations in a different
 colour to confirmed designations.
- Information about sequence bin records Creates a tooltip that displays details about sequence tags corresponding to a locus.
- Sequence bin records Displays a tooltip linking to the sequence tag if available.
- Sequence bin size Displays the size of the sum of all contigs associated with each isolate record.
- Contig count Displays the number of contigs associated with each isolate record.
- Publications Displays citations with links to PubMed for each record.

11.10.3 Isolate record display

The 'isolate record display' tab contains options for the display of a full isolate record.

Click the 'Isolate record display' tab to display the tab.



The 'Isolate record display' tab will scroll up.



This tab allows the following options to be modified:

Differentiate provisional allele designations - Allele designations can be set as confirmed or provisional, usually
depending on the method of assignment. Selecting this option will display provisional designations in a different
colour to confirmed designations.

- Display sender, curator and last updated records Displays a tooltip containing sender information next to each allele designation.
- Sequence bin information Displays a tooltip with information about the position of the sequence if tagged within the sequence bin.
- Allele flags Displays information about whether alleles have flags defined in sequence definition databases.
- Display full information about sample records Used when the database is used as part of a basic laboratory
 information management system (LIMS). This option will display records of samples available for the displayed
 isolate.

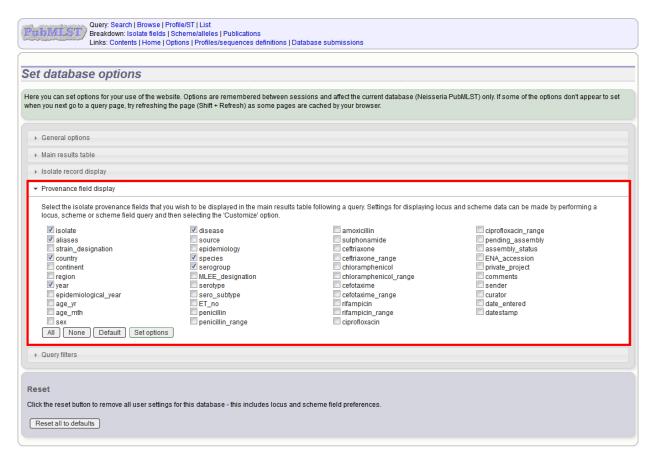
11.10.4 Provenance field display

The 'provenance field display' tab contains checkboxes for fields to display in the main results table.

Click the 'Provenance field display' tab to display the tab.



The 'Provenance field display' tab will scroll up.



Some fields will be checked by default - these are defined during database setup (maindisplay option).

Check any fields that you wish to be displayed and then click 'Set options'. You can return to the default selection by clicking 'Default' followed by 'Set options'.

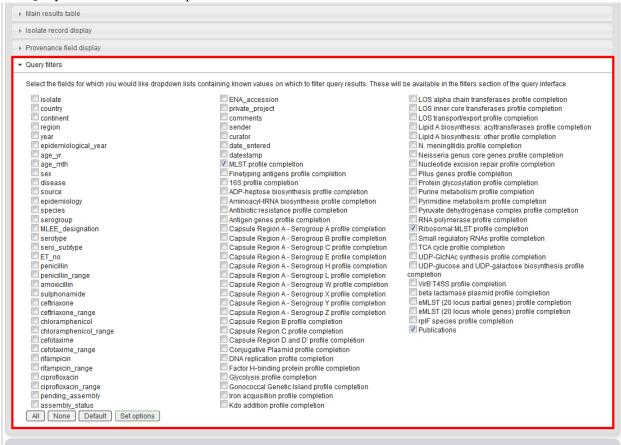
11.10.5 Query filters

The 'query filters' tab contains checkboxes for provenance fields and scheme completion status. Checking these results in drop-down list box filters appearing in the query page *filters fieldset*.

Click the 'Query filters' tab to display the tab.



The 'Query filters' tab will scroll up.



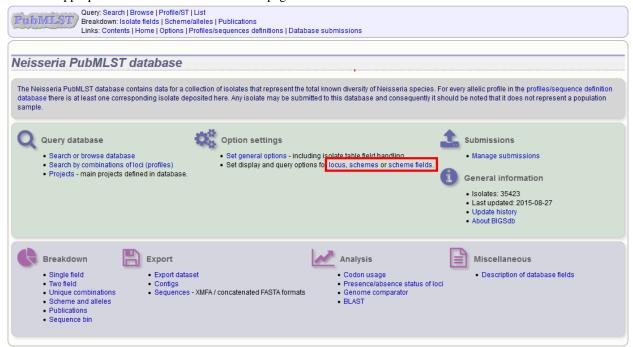
A list of possible filters appears. Click any checkbox for a filter you would like to make available. Click 'Set options' when done. You can return to the default selection by clicking 'Default' followed by 'Set options'.

11.10.6 Modifying locus and scheme display options

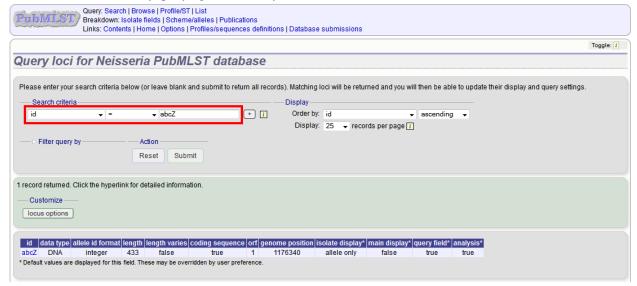
Whether or not loci, schemes or scheme fields are displayed in result tables, isolate records, or within query dropdown boxes can all be set with default options when first defined. These attributes can, however, be overridden by a user, and these selections will be remembered between sessions.

The procedure to modify these attributes is the same for locus, schemes or scheme fields, so the steps for loci will be demonstrated only.

Click the appropriate link on the isolate contents page.

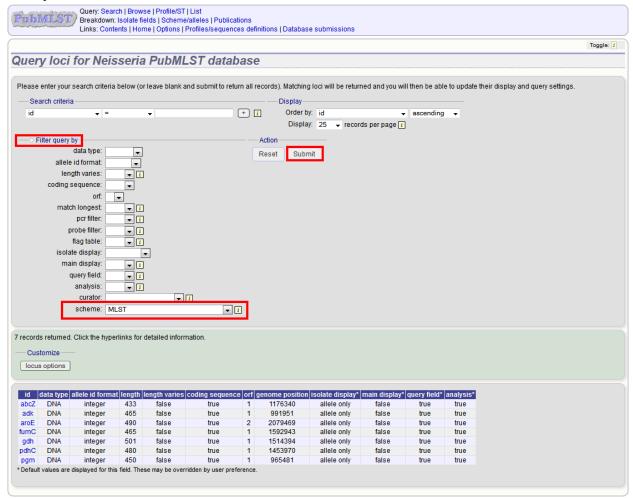


Either select the locus id by querying for it directly.

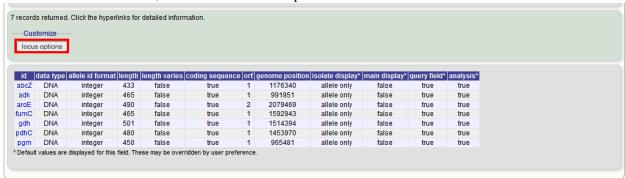


Designations can be queried using standard operators.

Alternatively, you can search by filtering loci by schemes. Click the 'Filter query by' header and select the scheme in the dropdown box.



Once loci have been selected, click Customize 'locus options'.



You can then choose to add or remove individual loci from the selection by clicking the appropriate checkboxes. At the bottom of the page are a number of attributes that you can change - clicking 'Change' will affect all selected loci.

Possible options for loci are:

• isolate_display - Sets how the locus is displayed within an isolate record:

- allele only display only identifier
- sequence display the full sequence
- hide don't show at all
- main_display Sets whether the locus is displayed in the main results table following a query.
- query_field Sets whether the locus appears in dropdown list boxes to be used within queries.
- analysis Sets whether the locus can be used in data analysis functions.

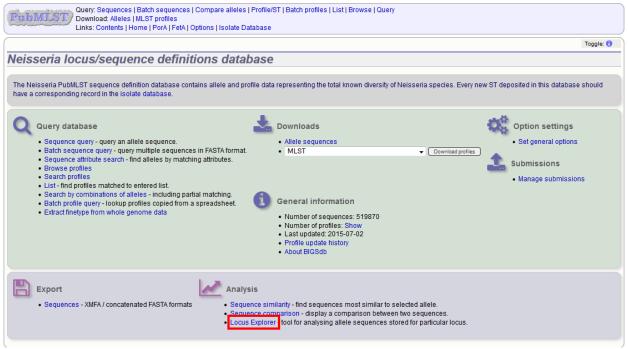
Note: Settings for loci can be overridden by those set for schemes that they are members of. For example, if you set a locus to be displayed within a main results table, but that locus is a member of a scheme and you set that scheme not to be displayed, then the locus will not be shown. Conversely, if you set a scheme to be displayed, but set its member locus not to be shown, then that locus will not be displayed (but other loci and scheme fields may be, depending on their independent settings).

Data analysis plugins

12.1 Locus explorer

The locus explorer is a sequence definition database plugin. It can create schematics showing the polymorphic sites within a locus, calculate the GC content and generate aligned translated sequences.

Click 'Locus Explorer' from the sequence definition database contents page.

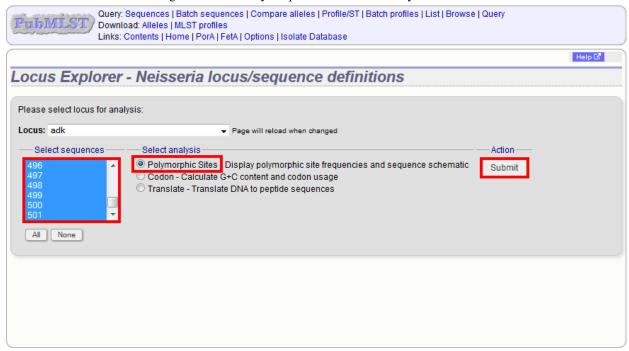


12.1.1 Polymorphic site analysis

Select the locus you would like to analyse in the Locus dropdown box. The page will reload.

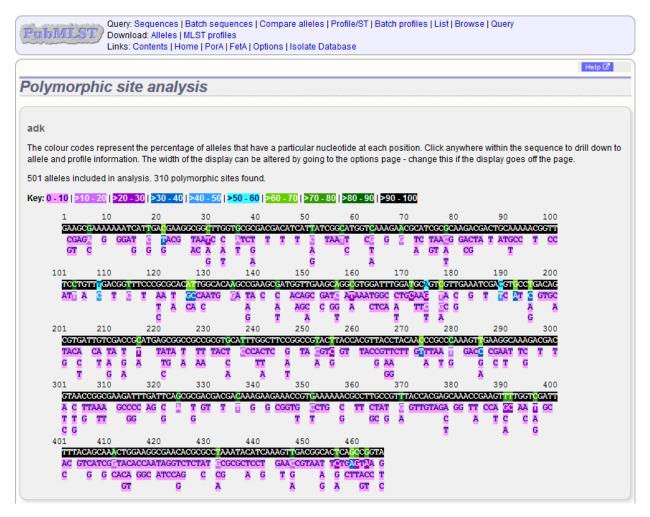


Select the alleles that you would like to include in the analysis. Variable length loci are limited to 2000 sequences or fewer since these need to be aligned. Select 'Polymorphic Sites' in the Analysis selection and click 'Submit'.

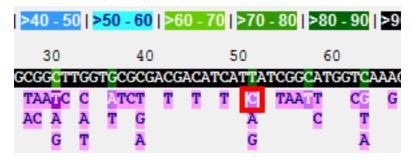


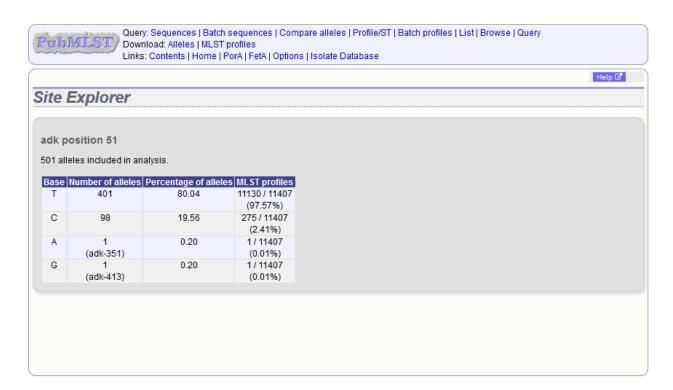
If an alignment is necessary, the job will be submitted to the job queue and the analysis performed. If no alignment is necessary, then the analysis is shown immediately.

The first part of the page shows the schematic.

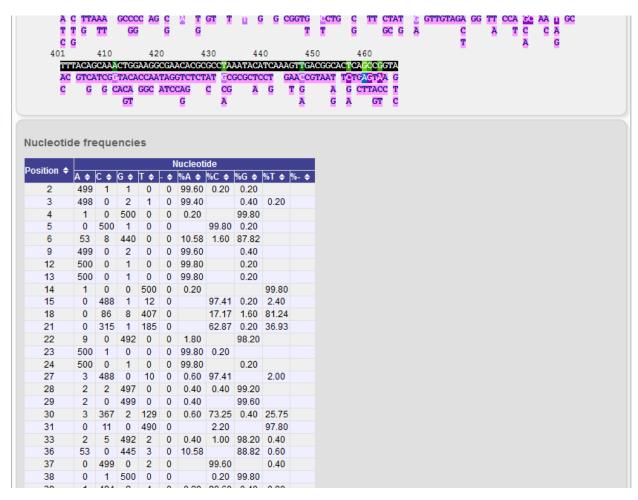


Clicking any of the sequence bases will calculate the exact frequencies of the different nucleotides at that position.





The second part of the page shows a table listing nucleotide frequencies at each of the variable positions.

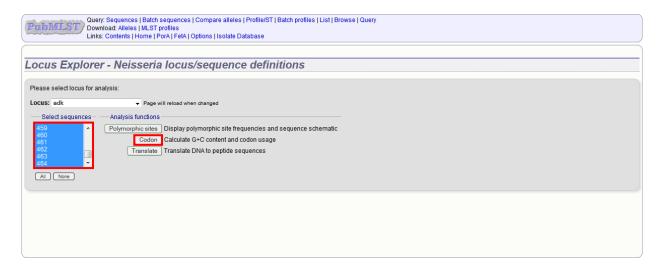


See also:

- Investigating allele differences.
- Polymorphism analysis following isolate query.

12.1.2 Codon usage

Select the alleles that you would like to include in the analysis. Again, variable length loci are limited to 200 sequences or fewer since these need to be aligned. Click 'Codon'.



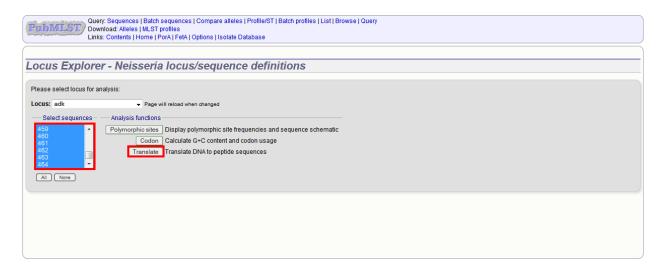
The GC content of the alleles will be determined and a table of the codon frequencies displayed.



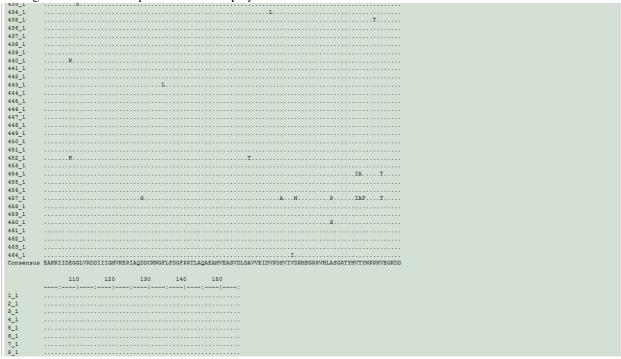
12.1.3 Aligned translations

If a DNA coding sequence locus is selected, an aligned translation can be produced.

Select the alleles that you would like to include in the analysis. Again, variable length loci are limited to 200 sequences or fewer since these need to be aligned. Click 'Translate'.



An aligned amino acid sequence will be displayed.



If there appear to be a lot of stop codons in the translation, it is possible that the orf value in the *locus definition* is not set correctly.

12.2 Field breakdown

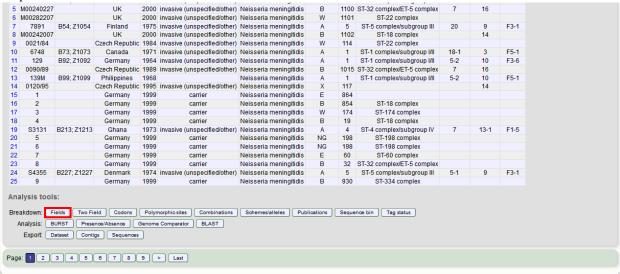
The field breakdown plugin for isolate databases displays the frequency of each value for fields stored in the isolates table. *Allele and scheme field breakdowns* are handled by a different plugin.

The breakdown function can be selected for the whole database by clicking the 'Single field' link in the Breakdown section of the main contents page.

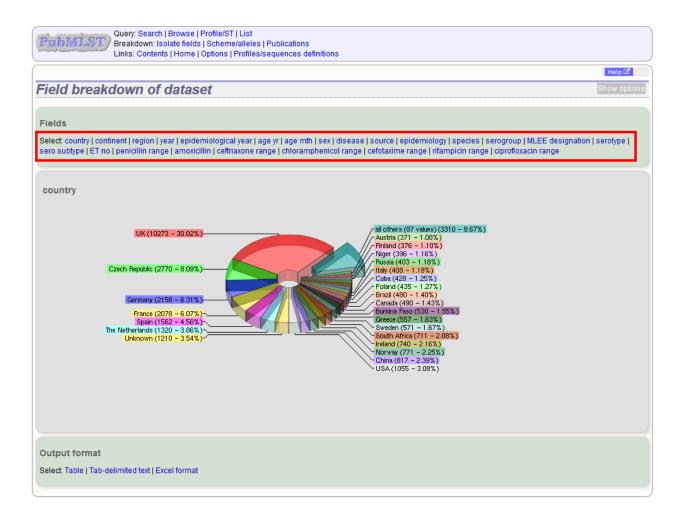
12.2. Field breakdown 261



Alternatively, a breakdown can be displayed of the dataset returned from a query by clicking the 'Fields' button in the Breakdown list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.

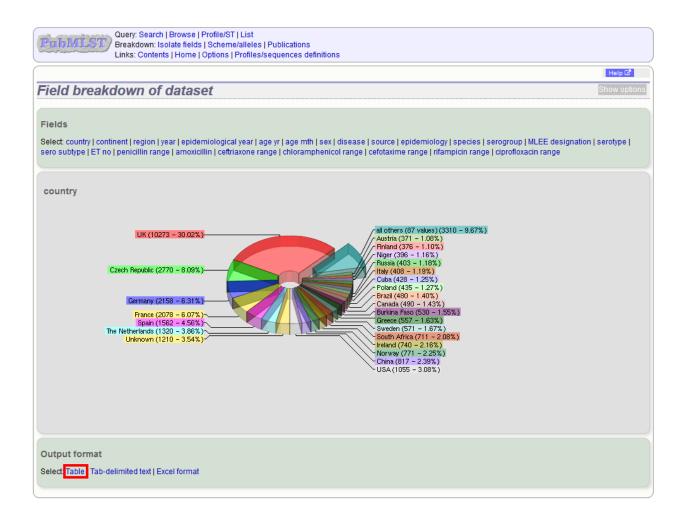


A series of charts will be displayed. Pick the field to display from the list at the top.

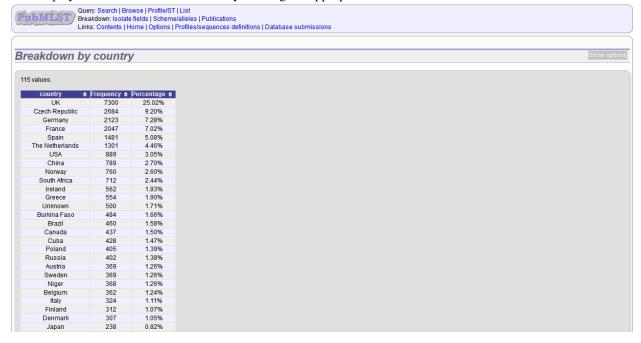


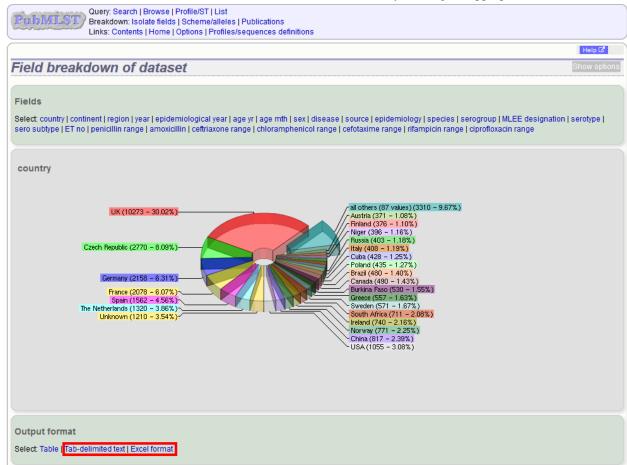
The values used to generate the chart can be displayed or extracted by clicking the 'Table' link at the bottom of the page.

12.2. Field breakdown 263



This displays a table that can be ordered by clicking the appropriate header.





The data can also be downloaded in tab-delimited text or Excel formats by clicking the appropriate links.

12.3 Two field breakdown

The two field breakdown plugin displays a table breaking down one field against another, e.g. breakdown of serogroup by year.

The analysis can be selected for the whole database by clicking the 'Two field breakdown' link on the main contents page.



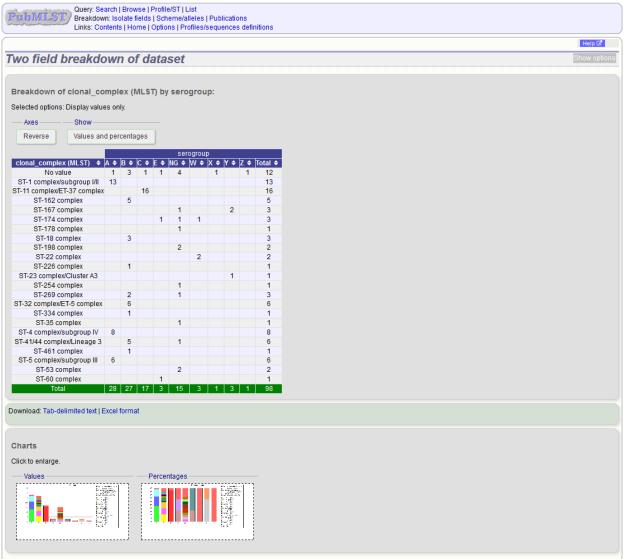
Alternatively, a two field breakdown can be displayed of the dataset returned from a query by clicking the 'Two field' button in the Breakdown list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



Select the two fields you wish to breakdown and how you would like the values displayed (percentage/absolute values and totaling options).



Click submit. The breakdown will be displayed as a table. Bar charts will also be displayed provided the number of returned values for both fields are fewer than 30.



The table values can be exported in a format suitable for copying in to a spreadsheet by clicking 'Download as tabdelimited text' underneath the table.

Note: The job will be submitted to the offline job queue if the query returns 10,000 or more isolates. In this case, the buttons to reverse the axes or to change whether values or percentages are shown will not be available.

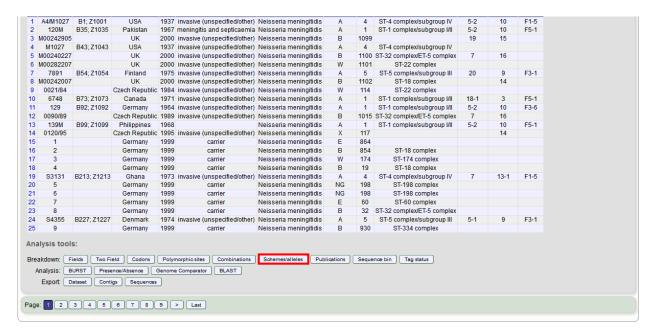
12.4 Scheme and allele breakdown

The scheme and allele breakdown plugin displays the frequency of each allele and scheme field (e.g. ST or clonal complex).

The function can be selected for the whole database by clicking the 'Scheme and allele breakdown' link on the main contents page.



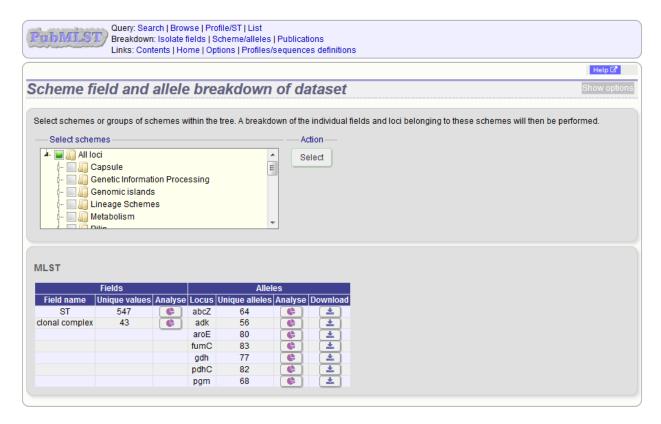
Alternatively, a breakdown can be displayed of the dataset returned from a query by clicking the 'Schemes/alleles' button in the Breakdown list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



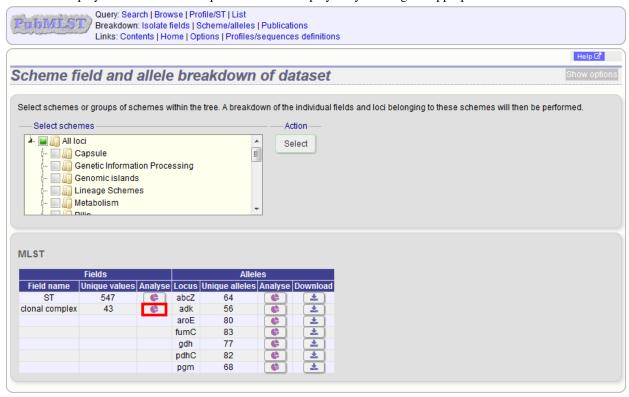
A scheme tree is shown. Select any combination of schemes to analyse. Click 'Select'.



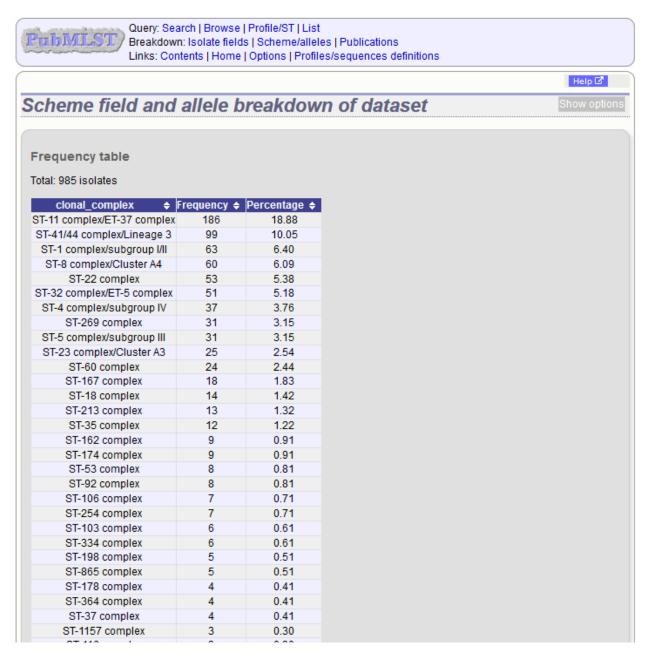
A table showing the number of unique values for each locus and scheme field will be displayed.



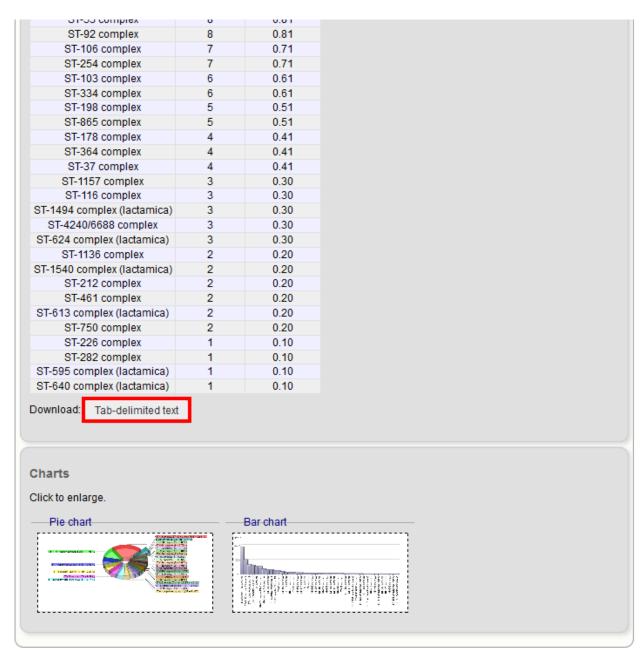
A detailed display of allele or field frequencies can be displayed by clicking the appropriate 'Breakdown' button.



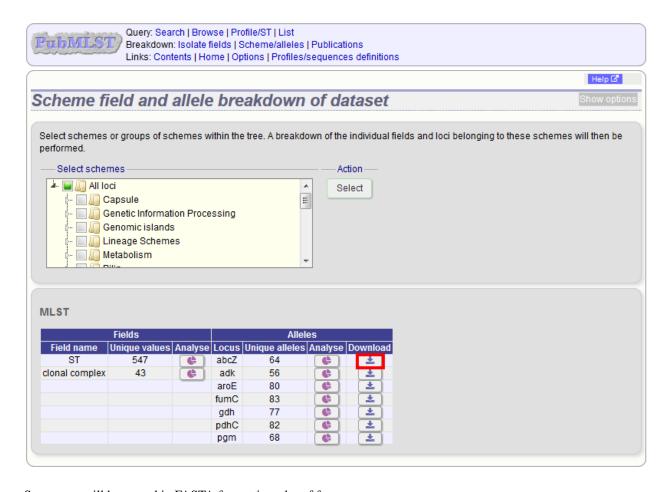
The sorting of the table can be changed by clicking the appropriate header - this toggles between ascending and descending order.



The table values can be exported in a format suitable for copying in to a spreadsheet by clicking the 'Tab-delimited text' button.



You can also download the sequences for alleles designated in the dataset for the loci belonging to the scheme by clicking the appropriate 'Download' button in the first results table.



Sequences will be served in FASTA format in order of frequency.

>2 TTTGATACCGTTGCCGAAGGTTTGGGTGAAATTCGCGATTTATTGCGCCGTTACCACCGC GTCGGCCATGAGTTGGAAAACGGTTCGGGTGAGGCTTTGTTGAAAGAACTCAACGAATTA CAACTTGAAATCGAAGCGAAGGACGGCTGGAAGCTGGATGCGGCAGTCAAGCAGACTTTG GGGGAACTCGGTTTGCCGGAAAACGAAAAAATCGGCAACCTTTCCGGCGGTCAGAAAAAG CGTGTCGCCTTGGCGCAGGCTTGGGTGCAGAAGCCCGACGTATTGCTGCTGGACGAACCG ACCAACCATTTGGATATCGACGCGATTATTTGGCTGGAAAATCTGCTCAAAGCGTTTGAA GGCAGCTTGGTTGTGATTACCCACGACCGCCGTTTTTTTGGACAATATCGCCACGCGGATT GTCGAACTCGATC >1 TTTGATACTGTTGCCGAAGGTTTGGGCGAAATTCGCGATTTATTGCGCCGTTATCATCAT GTCAGCCATGAGTTGGAAAATGGTTCGAGTGAGGCCTTATTGAAAGAGCTCAACGAATTG CAACTTGAGATCGAAGCGAAGGACGGCTGGAAGTTGGATGCGGCGGTGAAGCAGACTTTG GGCGAACTCGGTTTGCCGGAAAACGAAAAATCGGCAACCTCTCCGGCGGTCAGAAAAAG ACCAACCATTTGGACATCGACGCGATTATTTGGTTGGAAAACCTGCTCAAAGCGTTTGAA GGCAGCCTGGTTGTGATTACCCACGACCGCCGTTTTTTTGGACAATATCGCCACGCGGATT GTCGAACTCGATC TTTGATACCGTTGCCGAAGGTTTGGGCGAAATTCGTGATTTATTGCGCCGTTATCATCAT GTCAGCCATGAGTTGGAAAATGGTTCGAGTGAGGCTTTGTTGAAAGAACTCAACGAATTG CAACTTGAAATCGAAGCGAAGGACGGCTGGAAACTGGATGCGGCAGTCAAGCAGACTTTG GGGGAACTCGGTTTGCCGGAAAATGAAAAAATCGGCAACCTTTCCGGCGGTCAGAAAAAG CGCGTCGCCTTGGCTCAGGCTTGGGTGCAAAAGCCCGACGTATTGCTGCTGGACGAGCCG ACCAACCATTTGGATATCGACGCGATTATTTGGCTGGAAAATCTGCTCAAAGCGTTTGAA

GGCAGCTTGGTTGTGATTACCCACGACCGCCGTTTTTTGGACAATATCGCCACGCGGATTGTCGAACTCGATC

12.5 Sequence bin breakdown

The sequence bin breakdown plugin calculates statistics based on the number and length of contigs in the sequence bin as well as the number of loci tagged for an isolate record.

The function can be selected by clicking the 'Sequence bin' link on the Breakdown section of the main contents page.



Alternatively, it can be accessed following a query by clicking the 'Sequence bin' button in the Breakdown list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.

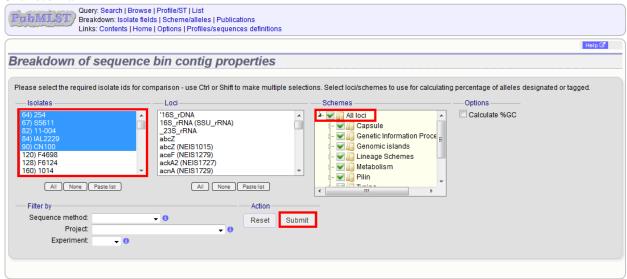


Select the isolate records to analyse - these will be pre-selected if you accessed the plugin following a query. You

can also select loci and/or schemes which will be used to calculate the totals and percentages of loci designated and tagged. This may be useful as a guide to assembly quality if you use a scheme of core loci where a good assembly would be expected to include all member loci. To determine the total of all loci designated or tagged, click 'All loci' in the scheme tree.

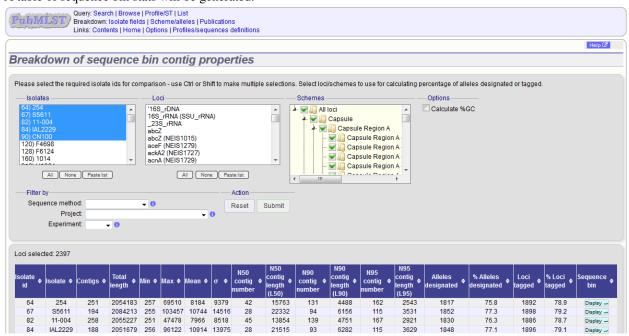
There is also an option to determine the mean G+C content of the sequence bin of each isolate.

Click submit.

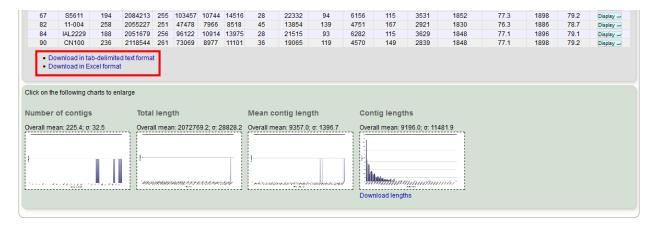


If there are fewer than 100 isolates selected, the table will be generated immediately. Otherwise it will be submitted to the job queue.

A table of sequence bin stats will be generated.



You can choose to export the data in tab-delimited text or Excel formats by clicking the appropriate link at the bottom of the table.



Sequence bin records can also be accessed by clicking the 'Display' button for each row of the table.



12.6 Genome comparator

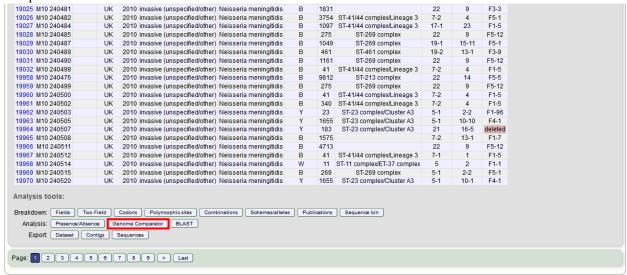
Genome Comparator is an optional plugin that can be enabled for specific databases. It is used to compare whole genome data of isolates within the database using either the database defined loci or the coding sequences of an annotated genome as the comparator.

Output is equivalent to a whole genome MLST profile, a distance matrix calculated based on allelic differences and a NeighborNet graph generated from this distance matrix.

Genome Comparator can be accessed on databases where it is enabled from the contents page by clicking the 'Genome Comparator' link.

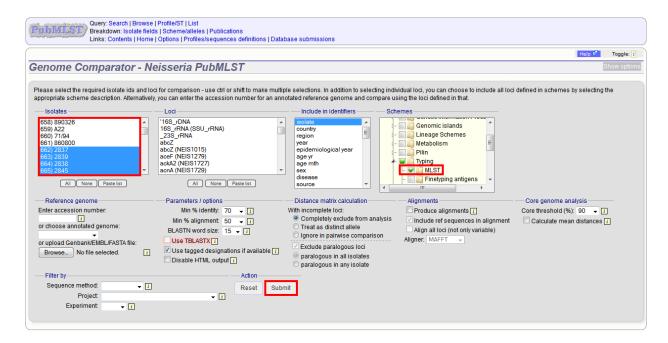


Alternatively, it can be accessed following a query by clicking the 'Genome Comparator' button at the bottom of the results table. Isolates with sequence data returned in the query will be automatically selected within the Genome Comparator interface.



12.6.1 Analysis using defined loci

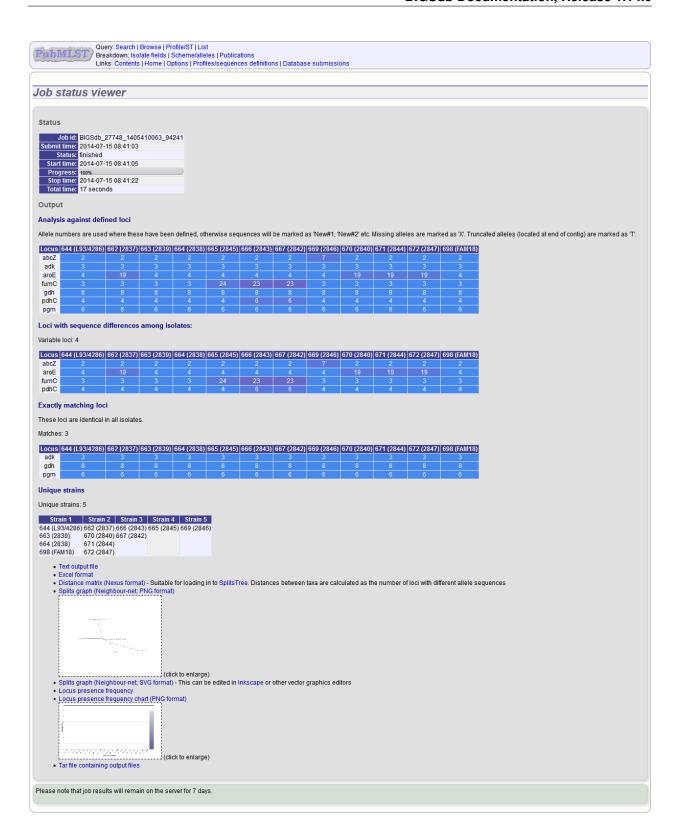
Select the isolate genomes that you wish to analyse and then either the loci from the list or a set of schemes. Press submit.



The job will be submitted to the job queue and will start running shortly. Click the link to follow the job progress and view the output.

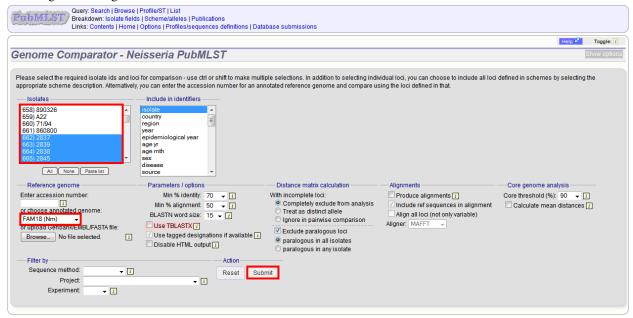


There will be a series of tables displaying variable loci, colour-coded to indicate allelic differences. Finally, there will be links to a distance matrix which can be loaded in to SplitsTree for further analysis and to a NeighborNet chart showing relatedness of isolates. Due to processing constraints on the web server, this NeighborNet is only calculated if 200 or fewer genomes are selected for analysis, but this can be generated in the stand-alone version of SplitsTree using the distance matrix if required.

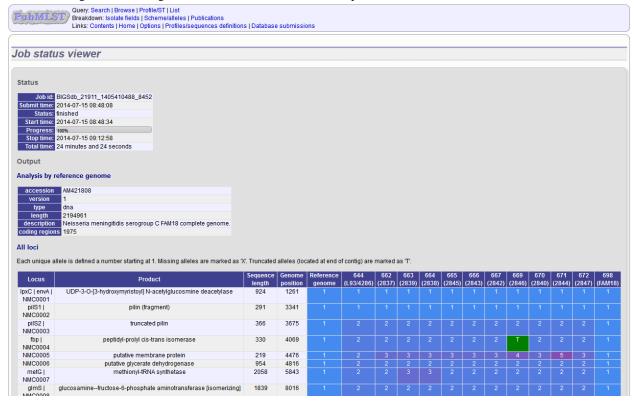


12.6.2 Analysis using annotated reference genome

Select the isolate genomes that you wish to analyse and then either enter a Genbank accession number for the reference genome, or select from the list of reference genomes (this list will only be present if the administrator has *set it up*). Selecting reference genomes will hide the locus and scheme selection forms.

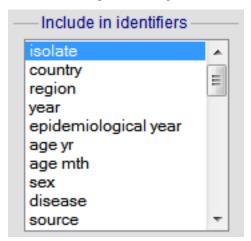


Output is similar to when comparing against defined loci, but this time every coding sequence in the annotated reference will be BLASTed against the selected genomes. Because allele designations are not defined, the allele found in the reference genome is designated allele 1, the next different sequence is allele 2 etc.



12.6.3 Include in identifiers fieldset

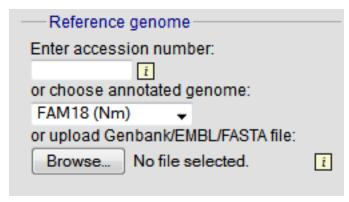
This selection box allows you to choose which isolate provenance fields will be included in the results table. This does not affect the output of the alignments as taxa names are limited in length by the alignment programs.



Multiple values can be selected by clicking while holding down Ctrl.

12.6.4 Reference genome fieldset

This section allows you to choose a reference genome to use as the source of comparator sequences.

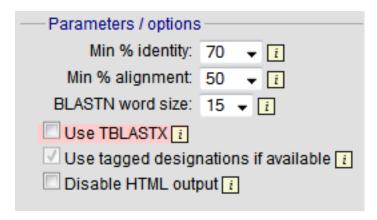


There are three possibilities here:

- 1. Enter accession number Enter a Genbank accession number of an annotated reference and Genome Comparator will automatically retrieve this from Genbank.
- 2. Select from list The administrator may have selected some genomes to offer for comparison. If these are present, simply select from the list.
- 3. Upload genome Click 'Browse' and upload your own reference. This can either be in Genbank, EMBL or FASTA format. Ensure that the filename ends in the appropriate file extension (.gb, .embl, .fas) so that it is recognized.

12.6.5 Parameters/options fieldset

This section allows you to modify BLAST parameters. This affects sensitivity and speed.



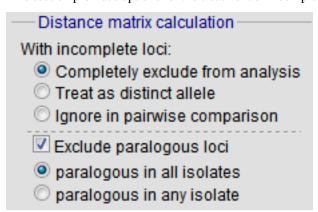
- Min % identity This sets the threshold identity that a matching sequence has to be in order to be considered (default: 70%). Only the best match is used.
- Min % alignment This sets the percentage of the length of reference allele sequence that the alignment has to cover in order to be considered (default: 50%).
- BLASTN word size This is the length of the initial identical match that BLAST requires before extending a match (default: 20). Increasing this value improves speed at the expense of sensitivity. The default value gives good results in most cases. The default setting used to be 15 but the new default of 20 is almost as good (there was 1 difference among 2000 loci in a test run) but the analysis runs twice as fast.
- Use TBLASTX This compares the six-frame translation of your nucleotide query sequence against the six-frame translation of the contig sequences. Sequences will be classed as identical if they result in the same translated sequence even if the nucleotide sequence is different. This is significantly slower than using BLASTN.

Additionally, two other options are available in this fieldset:

- Use tagged designations When analysing using defined loci, Genome Comparator can use the designations stored within the database (this is the default). This is much quicker since it doesn't need to run BLAST against these sequences. If a designation is missing, BLAST will be run for that locus anyway.
- Disable HTML output If running Genome Comparator against a large number of genomes, the resulting table may get so large that your web browser struggles to render it properly and may use up too much memory on your computer. Clicking this button prevents this output this output is not required for further analysis since everything present in it is also generated in Excel format at the end. HTML output is automatically disabled when more than 150 genomes are analysed.

12.6.6 Distance matrix calculation fieldset

This section provides options for the treatment of incomplete and paralogous loci when generating the distance matrix.



For incomplete loci, i.e. those that continue beyond the end of a contig so are incomplete you can:

- Completely exclude from analysis Any locus that is incomplete in at least one isolate will be removed from the analysis completely. Using this option means that if there is one bad genome with a lot of incomplete sequences in your analysis, a large proportion of the loci may not be used to calculate distances.
- Treat as a distinct allele This treats all incomplete sequences as a specific allele 'I'. This varies from any other allele, but all incomplete sequences will be treated as though they were identical.
- Ignore in pairwise comparison (default) This is probably the best option. In this case, incomplete alleles are only excluded from the analysis when comparing the particular isolate that has it. Other isolates with different alleles will be properly included. The effect of this option will be to shorten the distances of isolates with poorly sequenced genomes with the others.

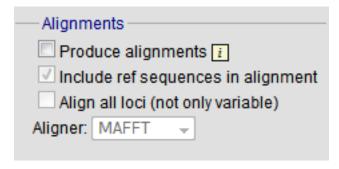
Paralogous loci, i.e. those with multiple good matches, can be excluded from the analysis (default). This is the safest option since there is no guarantee that differences seen between isolates at paralogous loci are real if the alternative matches are equally good. NB: Loci are also only classed as paralogous when the alternative matches identify different sequences, otherwise multiple contigs of the same sequence region would result in false positives.

When paralogous loci are excluded, there are two further options:

- Exclude when paralogous in all isolates (default). Loci are only classed as paralogous when there are multiple hits in every genome (except if a genome is missing the locus entirely, in which case that genome is ignored in the calculation). This is generally the option that you will want to use with the default BLAST parameters since you can often expect multiple hits even when loci are not paralogous if you have used relaxed thresholds.
- Exclude when paralogous in any isolate. Unless you use stringent BLAST thresholds, this is likely to overestimate the number of paralogous loci, but may be useful if you are specifically looking for them.

12.6.7 Alignments fieldset

This section enables you to choose to produce alignments of the sequences identified.



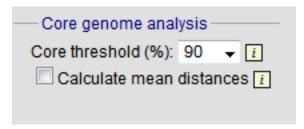
Available options are:

- Produce alignments Selecting this will produce the alignment files, as well as XMFA and FASTA outputs of aligned sequences. This will result in the analysis taking approximately twice as long to run.
- Include ref sequences in alignment When doing analysis using an annotated reference, selecting this will include the reference sequence in the alignment files.
- Align all loci By default, only loci that vary among the isolates are aligned. You may however wish to align all if you would like the resultant XMFA and FASTA files to include all coding sequences.
- Aligner There are currently two choices of alignment algorithm (provided they have both been installed)
 - MAFFT (default) This is the preferred option as it is significantly quicker than MUSCLE, uses less memory, and produces comparable results.

MUSCLE - This was originally the only choice. It is still included to enable previous analyses to be re-run
and compared but it is recommended that MAFFT issued otherwise.

12.6.8 Core genome analysis fieldset

This section enables you to modify the inclusion threshold used to calculate whether or not a locus is part of the core genome (of the dataset).

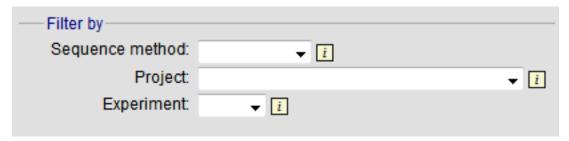


The default setting of 90% means that a locus is counted as core if it appears within 90% or more of the genomes in the dataset.

There is also an option to calculate the mean distance among sequences of the loci. Selecting this will also select the option to produce alignments.

12.6.9 Filter fieldset

This section allows you to further filter your collection of isolates and the contigs to include.



Available options are:

- Sequence method Choose to only analyse contigs that have been generated using a particular method. This depends on the method being set when the contigs were uploaded.
- Project Only include isolates belonging to the chosen project. This enables you to select all isolates and filter to a project.
- Experiment Contig files can belong to an experiment. How this is used can vary between databases, but this enables you to only include contigs from a particular experiment.

12.6.10 Understanding the output

Distance matrix

The distance matrix is simply a count of the number of loci that differ between each pair of isolates. It is generated in NEXUS format which can be used as the input file for SplitsTree. This can be used to generate NeighborNet, Split decomposition graphs and trees offline. If 200 isolates or fewer are included in the analysis, a Neighbor network is automatically generated from this distance matrix.

Unique strains

The table of unique strains is a list of isolates that are identical at every locus. Every isolate is likely to be classed as unique if a whole genome analysis is performed, but with a constrained set of loci, such as those for MLST, this will group isolates that are indistinguishable at that level of resolution.

12.7 BLAST

The BLAST plugin enables you to BLAST a sequence against any of the genomes in the database, displaying a table of matches and extracting matching sequences.

The function can be accessed by selecting the 'BLAST' link on the Analysis section of the main contents page.

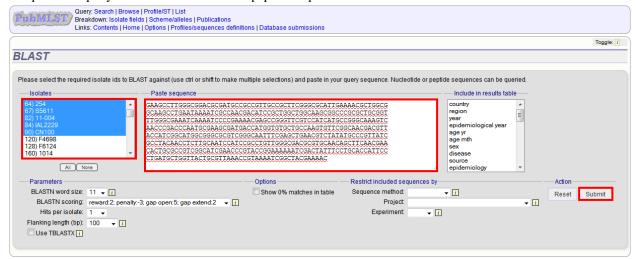


Alternatively, it can be accessed following a query by clicking the 'BLAST' button in the Analysis list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



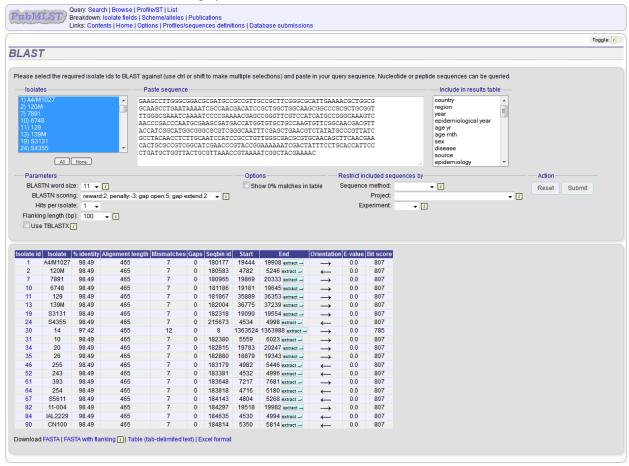
12.7. BLAST 285

Select the isolate records to analyse - these will be pre-selected if you accessed the plugin following a query. Paste in a sequence to query - this be either a DNA or peptide sequence.



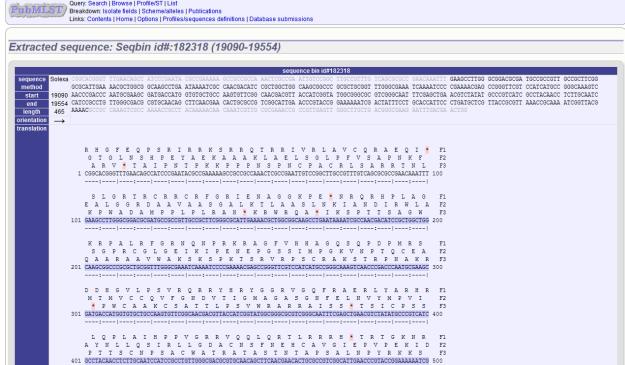
Click submit.

A table of BLAST results will be displayed.



Clicking any of the 'extract' buttons will display the matched sequence along with a translated sequence and flanking sequences.





At the bottom of the results table are links to export the matching sequences in FASTA format, (optionall) including flanking sequences. You can also export the table in tab-delimited text or Excel formats.



12.7.1 Include in results table fieldset

This selection box allows you to choose which isolate provenance fields will be included in the results table.

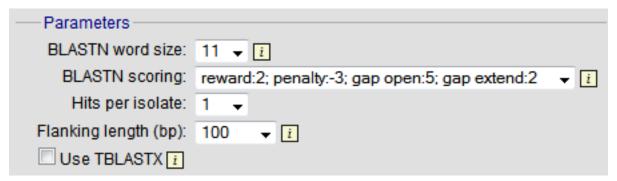
12.7. BLAST 287



Multiple values can be selected by clicking while holding down Ctrl.

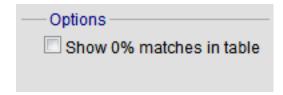
12.7.2 Parameters fieldset

This section allows you to modify BLAST parameters. This affects sensitivity and speed.



- BLASTN word size This is the length of the initial identical match that BLAST requires before extending a match (default: 11). Increasing this value improves speed at the expense of sensitivity.
- BLASTN scoring This is a dropdown box of combinations of identical base rewards; mismatch penalties; and
 gap open and extension penalties. BLASTN has a constrained list of allowed values which reflects the available
 options in the list.
- Hits per isolate By default, only the best match is shown. Increase this value to the number of hits you'd like to see per isolate.
- Flanking length Set the size of the upstream and downstream flanking sequences that you'd like to include.
- Use TBLASTX This compares the six-frame translation of your nucleotide query sequence against the six-frame translation of the contig sequences. This is significantly slower than using BLASTN.

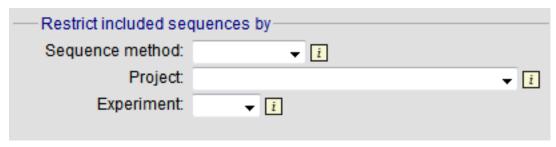
12.7.3 No matches



Click this option to create a row in the table indicating that a match was not found. This can be useful when screening a large number of isolates.

12.7.4 Filter fieldset

This section allows you to further filter your collection of isolates and the contig sequences to include.



Available options are:

- Sequence method Choose to only analyse contigs that have been generated using a particular method. This depends on the method being set when the contigs were uploaded.
- Project Only include isolates belonging to the chosen project. This enables you to select all isolates and filter to a project.
- Experiment Contig files can belong to an experiment. How this is used can vary between databases, but this enables you to only include contigs from a particular experiment.

12.8 BURST

BURST is an algorithm used to group MLST-type data based on a count of the number of profiles that match each other at specified numbers of loci. The analysis is available for both sequence definition database and isolate database schemes that have primary key fields set. The algorithm has to be *specifically enabled* by an administrator. Analysis is limited to 1000 or fewer records.

The plugin can be accessed following a query by clicking the 'BURST' button in the Analysis list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



If there multiple schemes that can be analysed, these can then be selected along with the group definition.

12.8. BURST 289



Modifying the group definition affects the size of groups and how they link together. By default, the definition is n-2 (where n is the number of loci), so for example on a 7 locus MLST scheme groups contain STs that match at 5 or more loci to any other member of the group.

Click Submit.

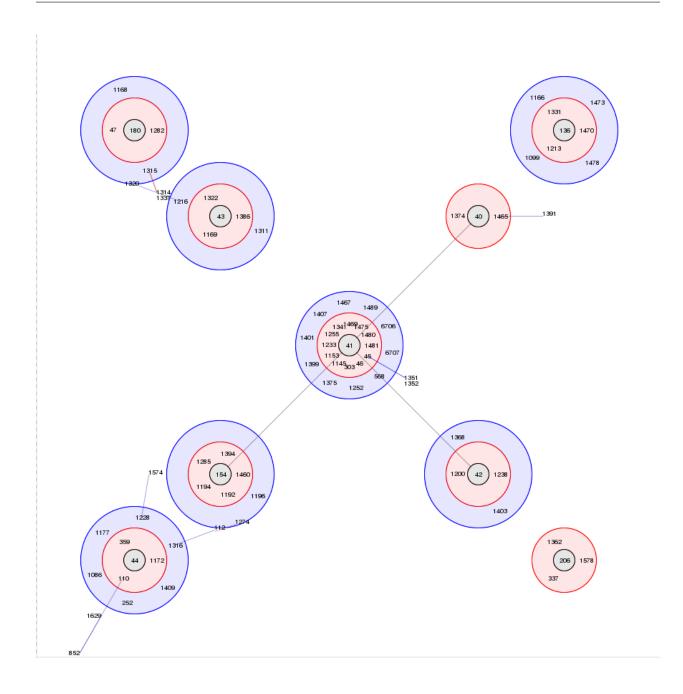
A series of tables will be displayed indicating the groups of profiles. Where one profile can be identified as a central genotype, i.e. the profile that has the greatest number of other profiles that are single locus variants (SLV), double locus variants (DLV) and so on, a graphical representation will be displayed. The central profile is indicated with an asterisk.

SLV profiles that match the central profile are shown within a red circle surrounding the central profile. Most distant profiles (triple locus variants) may be linked with a line. Larger groups may additionally have DLV profiles. These are shown in a blue circle.

12.8. BURST 291

	group:	6								
ST	Frequency	SLV	DLV	SAT						
32*	2	3	2							
230	1	1	3	1						
484	1	0	3	2						
1015	1	1	4							
1100	1	1	2	2						
1148	1	0	4	1						
1015 32 1100 1148										
	SVG file (right cli	ck to save	e)							

Groups can get very large, where linked profiles form sub-groups and an attempt is made to depict these.



12.9 Codon usage

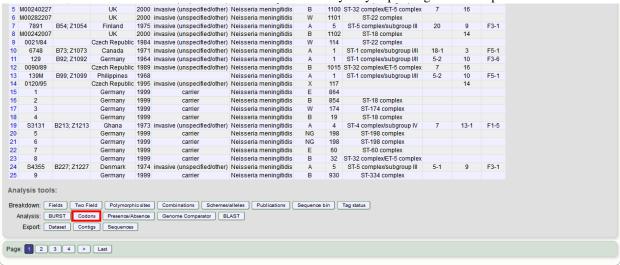
The codon usage plugin for isolate databases calculates the absolute and relative synonymous codon usage by isolate and by locus.

The function can be selected by clicking the 'Codon usage' link in the Analysis section of the main contents page.

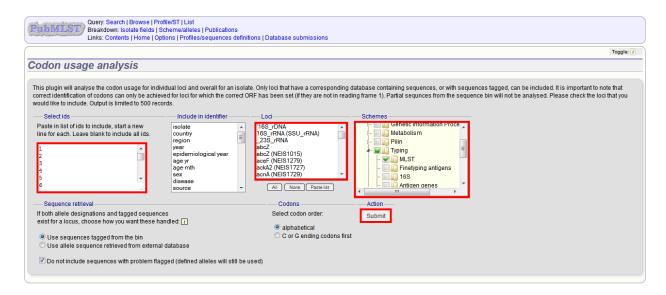
12.9. Codon usage 293



Alternatively, it can be accessed following a query by clicking the 'Codons' button in the Analysis list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



Enter the ids of the isolate records to analyse - these will be already entered if you accessed the plugin following a query. Select the loci you would like to analyse, either from the dropdown loci list, and/or by selecting one or more schemes.



Click submit. The job will be submitted to the queue and will start running shortly. Click the link to follow the job progress and view the output.



Four tab-delimited text files will be created.

- Absolute frequency of codon usage by isolate
- Absolute frequency of codon usage by locus
- · Relative synonymous codon usage by isolate
- Relative synonymous codon usage by locus

12.9. Codon usage 295



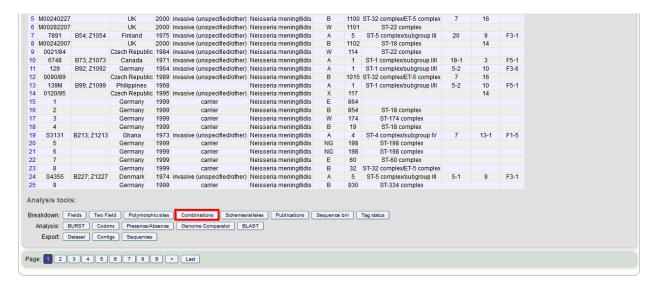
12.10 Unique combinations

The Unique Combinations plugin calculates the frequencies of unique file combinations within an isolate dataset. Provenance fields, composite fields, allele designations and scheme fields can be combined.

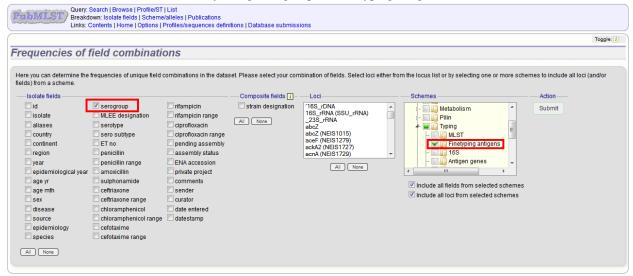
The function can be selected by clicking the 'Unique combinations' link in the Breakdown section of the main contents page. This will run the analysis on the entire database.



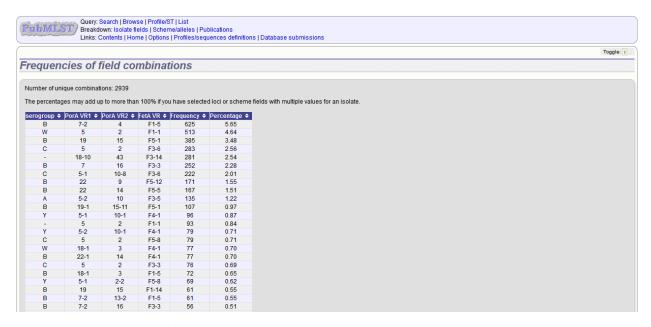
Alternatively, it can be accessed following a query by clicking the 'Combinations' button in the Breakdown list at the bottom of the results table. This will run the analysis on the dataset returned from the query. Please note that the list of functions here may vary depending on the setup of the database.



Select the combination of fields to analyse, e.g. serogroup and finetyping antigens.



Click submit. When the analysis has completed you will see a table showing the unique combinations of the selected fields along with the frequency and percentage of the combination.



The table can be downloaded in tab-delimited text or Excel formats by clicking the links at the bottom of the page.



12.11 Polymorphisms

The Polymorphisms plugin generates a *Locus Explorer* polymorphic site analysis on the alleles designated in an isolate dataset following a query.

The analysis is accessed by clicking the 'Polymorphic sites' button in the Breakdown list at the bottom of a results table following a query.

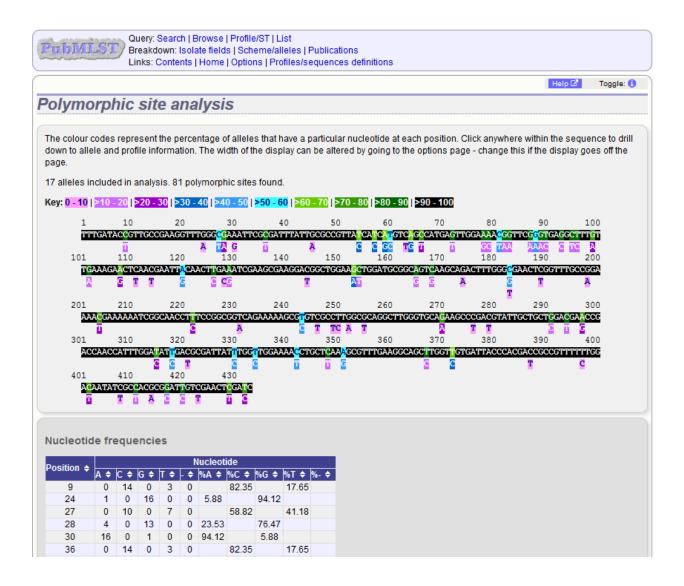


Select the locus that you would like to analyse from the list.



Click 'Analyse'.

A schematic of the locus is generated showing the polymorphic sites. A full description of this can be found in the *Locus Explorer polymorphic site analysis* section.



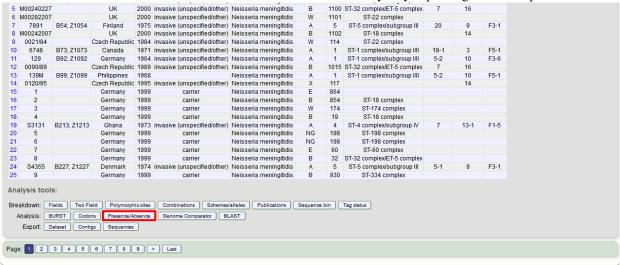
12.12 Presence/absence

This plugin displays the status of loci for isolate records. It will shown whether a locus has been designated with an allele name, has a sequence tag, or both.

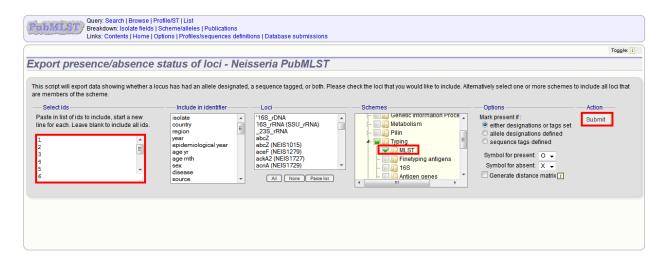
The function can be selected by clicking the 'Presence/absence status of loci' link in the 'Analysis' section of the main contents page.



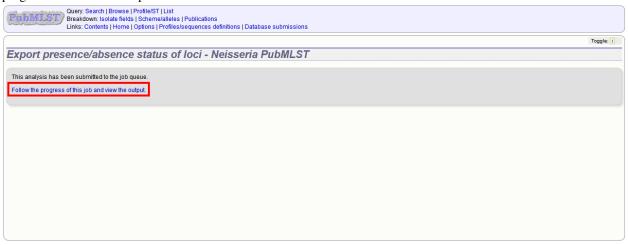
Alternatively, it can be accessed following a query by clicking the 'Presence/Absence' button in the Analysis list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



Enter the ids of the isolate records to analyse - these will be already entered if you accessed the plugin following a query. Select the loci you would like to analyse, either from the dropdown loci list, and/or by selecting one or more schemes.



Click submit. The job will be submitted to the queue and will start running shortly. Click the link to follow the job progress and view the output.



When complete, a single text file will have been generated.

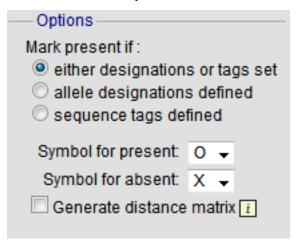


This is a tab-delimited text file that uses 'O' to represent presence and 'X' to represent a missing locus designation or tag.

id	pgm	adk	abcZ	pdhC	gdh	fumC	aroE
1	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0

12.12.1 Options

There are a number of options that can be selected to modify the output.



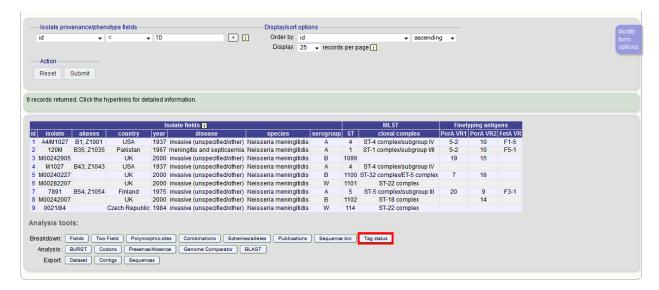
With these you can change the symbols used and whether designations, or tags, or both are counted.

You can also choose to generate a distance matrix based on presence/absence.

12.13 Tag status

The tag status plugin displays a graphical representation of the status of loci designations or tags for isolate data. It is accessed following a query by clicking the 'Tag status' button in the Breakdown section at the bottom of the results table.

12.13. Tag status 303



Select the loci you would like to analyse.



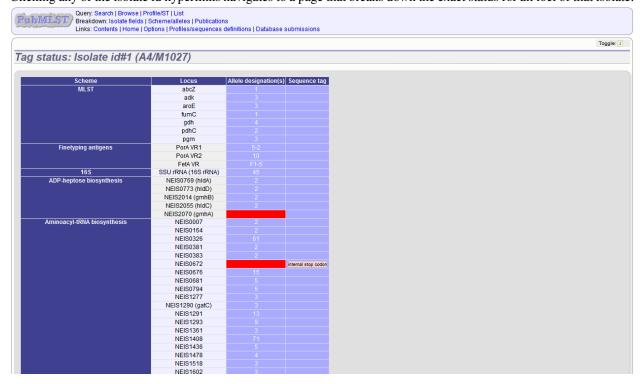
You should see a series of bars representing loci. The colour of these bars designates whether they have an allele designation only, a sequence tag only, both designations or tags, or whether they have flags set.



Hovering the mouse over the bars will indicate the scheme represented.

Note: Loci will be represented more than once if they are members of multiple selected schemes.

Clicking any of the isolate id hyperlinks navigates to a page that breaks down the exact status for all loci of that isolate.



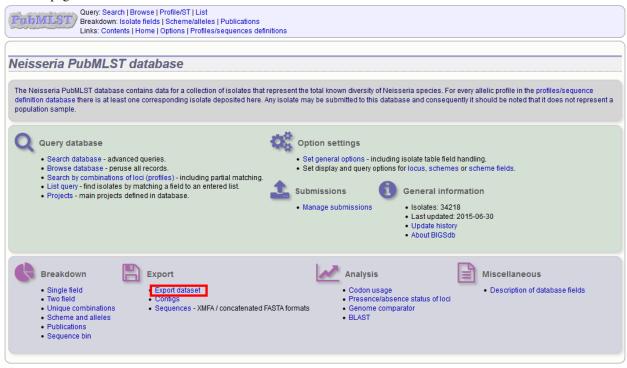
12.13. Tag status 305

There is a column each for allele designations and sequence tags. If an allele designation is defined, the allele identifier is displayed. Cells shaded in blue show that the designation or tag is present, whereas red indicates thet they are absent.

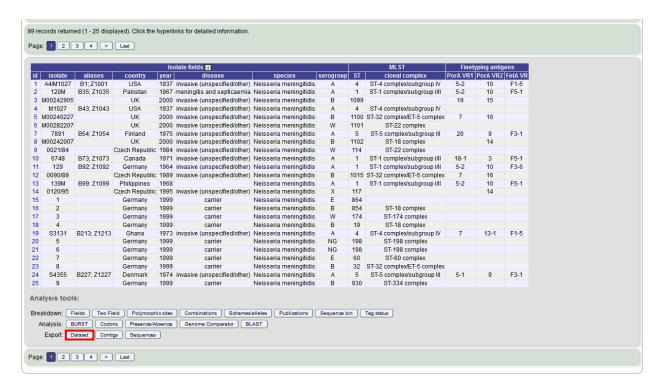
Data export plugins

13.1 Isolate record export

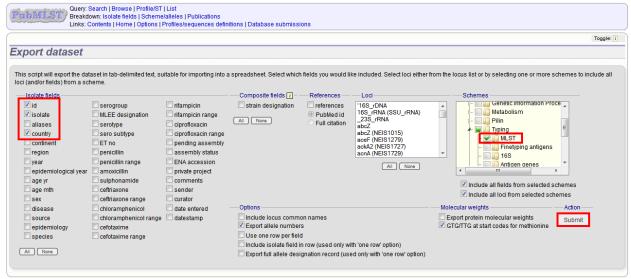
You can export the entire isolate recordset by clicking the 'Export dataset' link in the Export section of the main contents page.



Alternatively, you can export the recordsets of isolates returned from a database query by clicking the 'Dataset' button in the Export list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



Select the isolate fields and schemes to include.

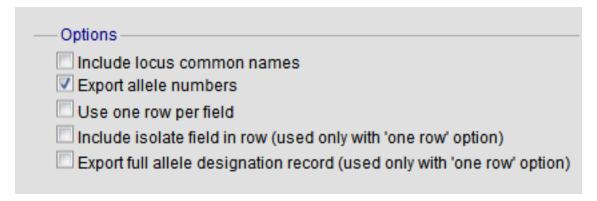


Click Submit.

You can then download the data in tab-delimited text or Excel formats.



13.1.1 Advanced options

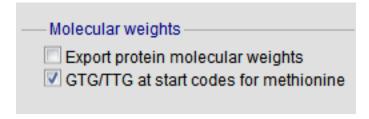


The options fieldset has the following options.

- Include locus common names any common name for the locus is displayed in parentheses following the primary name.
- Export allele numbers the allele designation is included for any locus included.
- Use one row per field this is an alternative output format where instead of each locus and field having a separate column, each field is export on a separate row.
- Include isolate field in row the name of the isolate is included as a separate column when exporting in 'one row per field' fomrmat.
- Export full allele designation record export sender, curator and datestamp information as separate rows when exporting allele designation data.

13.1.2 Molecular weight calculation

The plugin can also calculate the predicted molecular weight of the gene product of any allele designated in the dataset.

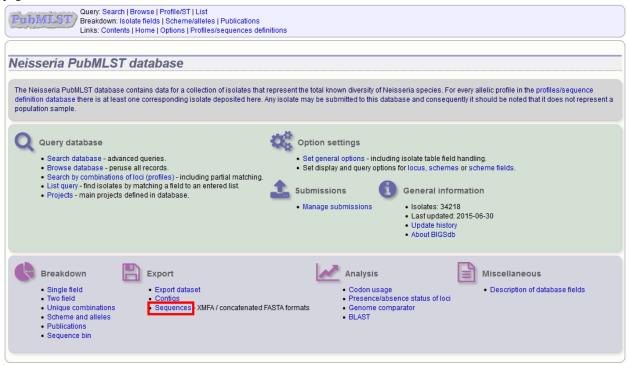


Click the 'Export protein molecular weight' checkbox. Additional columns (or rows depending on the output format) will be created to include the molecular weight data.

13.2 Sequence export

You can export the sequences for any set of loci designated in isolate records, or belonging to scheme profiles in the sequence definition database.

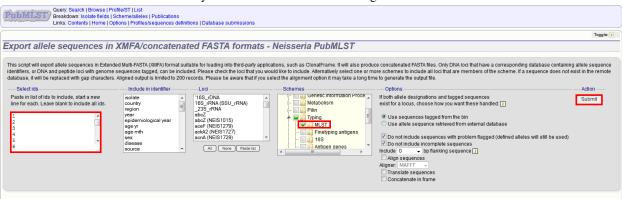
The sequence export function can be accessed by clicking the 'Sequences' link in the Export section of the contents page.



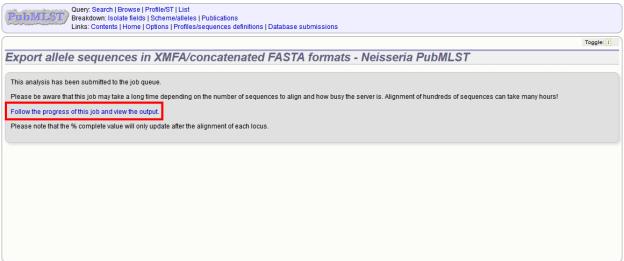
Alternatively, you can access this function by clicking the 'Sequences' button in the Export list at the bottom of the results table. Please note that the list of functions here may vary depending on the setup of the database.



Select the isolate or profile records to analyse - these will be pre-selected if you accessed the plugin following a query. Select the loci to include either directly within the loci list and/or using the schemes tree.



Click submit.

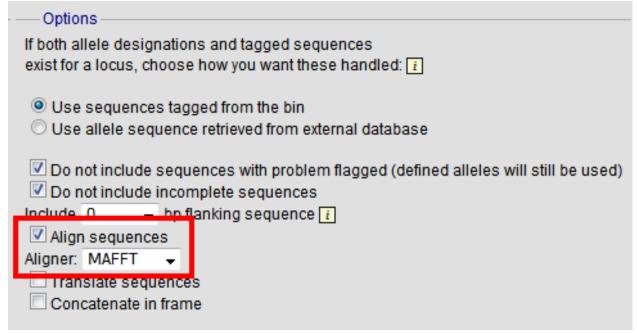


The job will be submitted to the job queue. Click the link to follow the progress and download the resulting files. Sequences will be export in XMFA and FASTA file formats.



13.2.1 Aligning sequences

By default, sequences will be exported unaligned - this is very quick since no processing is required. You can choose to align the sequences by checking the 'Align sequences' checkbox.



You can also choose to use MUSCLE or MAFFT as the aligner. MAFFT is the default choice and is usually much quicker than MUSCLE. Both produce comparable results.

13.3 Contig export

The contig export plugin can be accessed by clicking the 'Contigs' link in the Export section of the contents page of isolate databases.



Alternatively, it can be accessed following a query by clicking the 'Contigs' button in the Export section at the bottom of the results table.



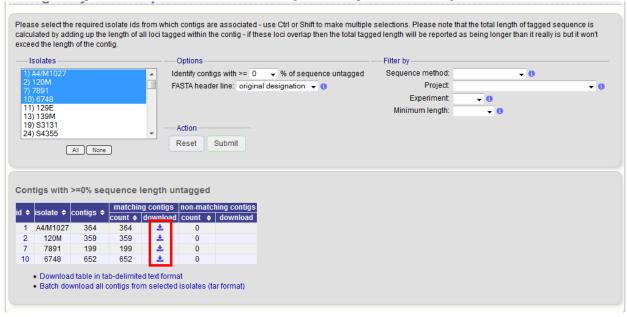
Select the isolates for which you wish to export contig data for. If the export function was accessed following a query, isolates returned in the query will be pre-selected.

13.3. Contig export 313

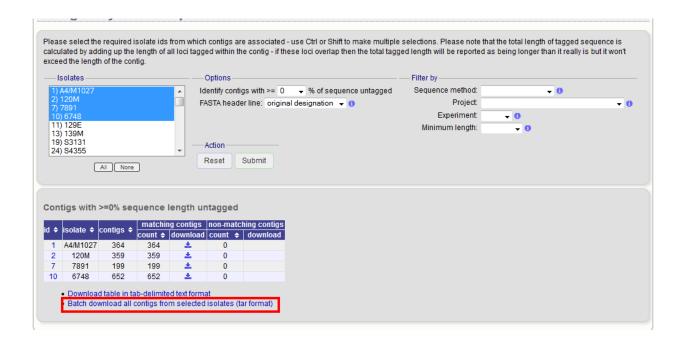


At its simplest, press submit.

A table will be produced with download links. Clicking these will produce the contigs in FASTA format.



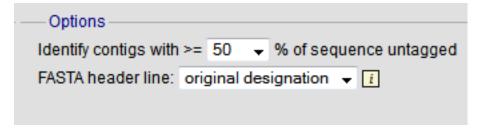
You can also download all the data in a tar file by clicking the 'Batch download' link.



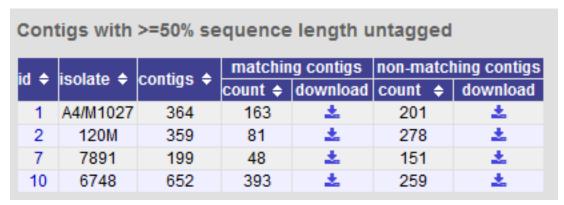
13.3.1 Filtering by tagged status of contigs

You can also export contigs based on the percentage of the sequence that has been tagged. This is useful to find sequences to target for gene discovery.

In order to export contigs where at least half the sequence has been tagged (and also the remaining contigs in a separate file), select '50' in the dropdown box for %untagged.



The resulting table has two download links for each isolate, one for contigs matching the condition, and one for contigs that don't match.

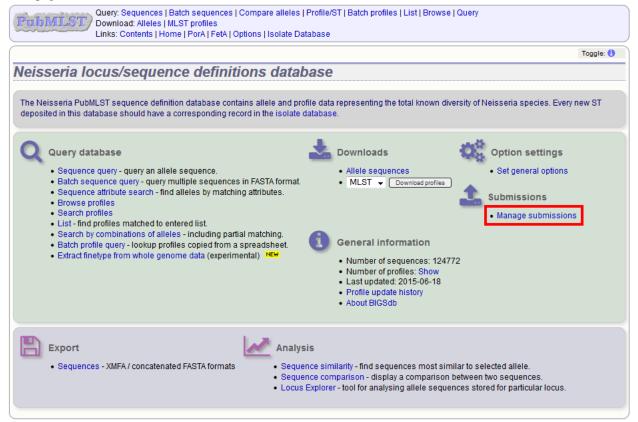


13.3. Contig export 315

Submitting data using the submission system

The automated submission system allows users to submit data (new alleles, profiles, or isolates) to the database curators for assignment and upload to the database. The submission system is enabled on a per-database basis so will not always be available.

If the system is enabled, new submissions can be made by clicking the 'Manage submissions' link on the database front page.



14.1 Registering a user account

You must have an account for the appropriate database in order to use the submission system. This will need to be set up by a curator, so contact them in the first instance.

14.2 Allele submission

New allele data can only be submitted from within the appropriate sequence definition database. Submissions consist of one or more new allele sequences for a single locus. You will need to create separate submissions for each locus - this is because different loci may be handled by different curators.

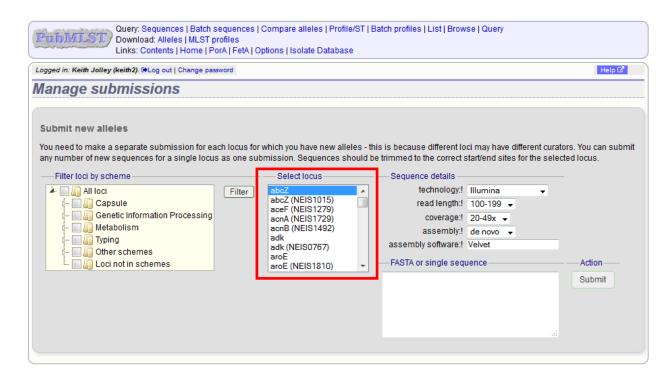
14.2.1 Start

Click the 'alleles' link under submission type on the submission management page.

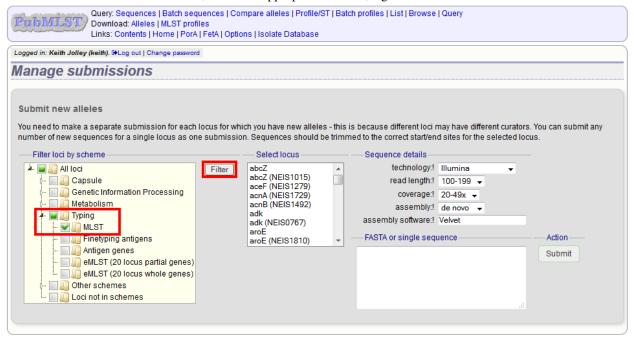


14.2.2 Select the submission locus

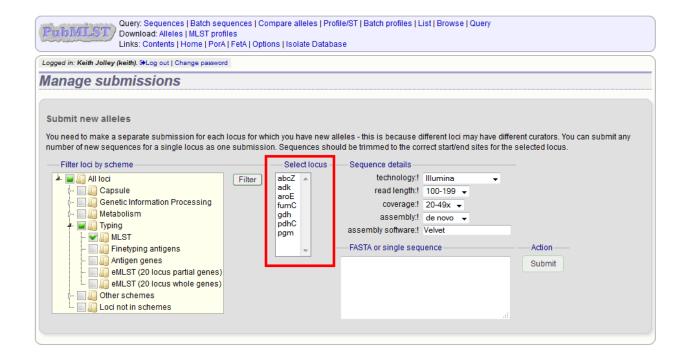
Select the locus from the locus list box:



The locus list may be very long in some databases. It may be possible to filter these to those belonging to specific schemes. If the scheme tree is shown, select the appropriate scheme, e.g. 'MLST' and click 'Filter'.



The locus list is now constrained making selection easier.



14.2.3 Enter details of sequencing method

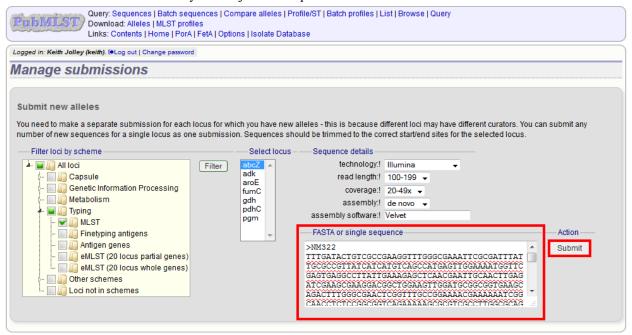
There are a number of fields that must be filled in so that the curator knows how the sequence was obtained:

- technology the sequncing platform used, allowed values are:
 - 454
 - Illumina
 - Ion Torrent
 - PacBio
 - Oxford Nanopore
 - Sanger
 - Solexa
 - SOLiD
 - other
 - unknown
- read length this is the length of sequencing reads. This is a required field for Illumina data, and not relevant to Sanger sequencing. Allowed values are:
 - <100
 - **-** 100-199
 - 200-299
 - 300-499
 - >500

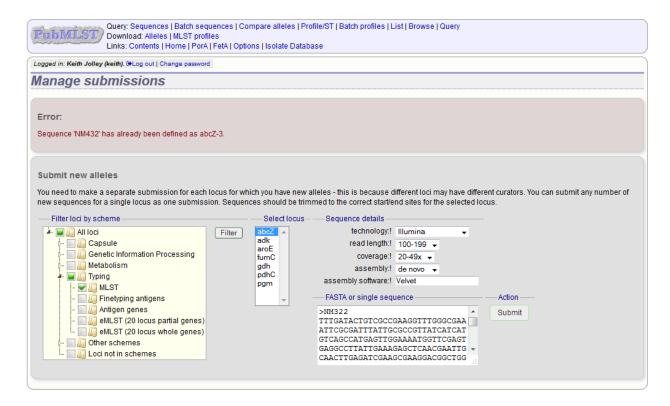
- coverage the mean number of reads covering each nucleotide position of the sequence. This is not relevant to Sanger sequencing, Allowed values are:
 - < 20x
 - -20-49x
 - 50-99x
 - > 100x
- assembly the means of generating the submitted sequence from the sequencing reads. Allowed values are:
 - de novo
 - mapped
- assembly software this is a free text field where you should enter the name of the software used to generate the submitted sequence.

14.2.4 Paste in sequence(s)

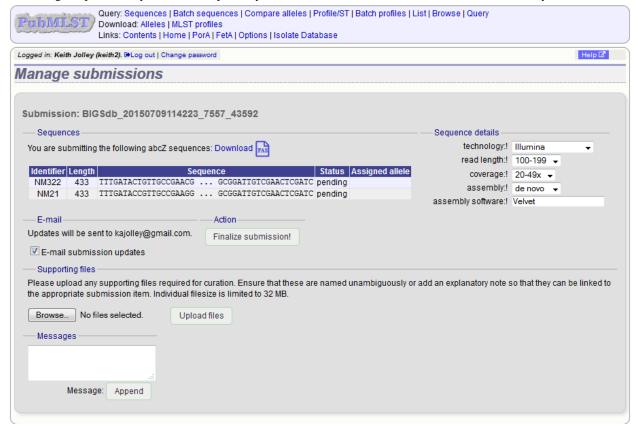
Paste in the new variant sequences to the box. This can either be a stand- alone sequence or multiple sequences in FASTA format. The sequences must be trimmed to the start and end points of the loci - check existing allele definitions if in doubt. The submission is likely to be rejected if sequences are not trimmed. Click submit.



The system will perform some basic checks on the submitted sequences. If any of the sequences have been defined previously they must be removed from the submission before you can proceed. Curators do not want to waste their time dealing with previously defined sequences.

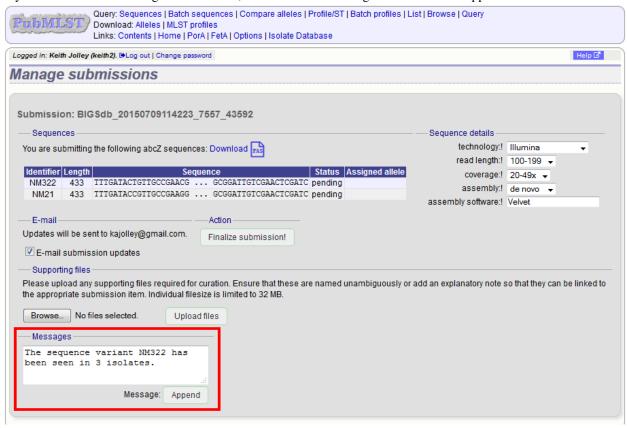


Assuming the preliminary checks have passed you will then be able to add additional information to your submission.

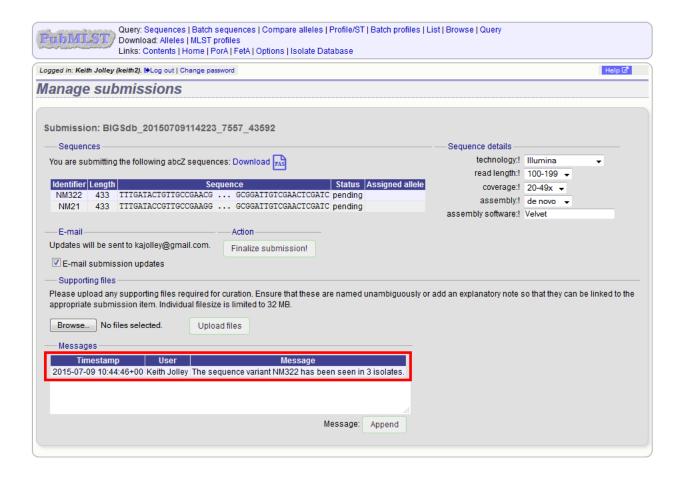


14.2.5 Add message to curator

If you wish to enter a message to the curator, enter this in the messages box and click 'Append'.



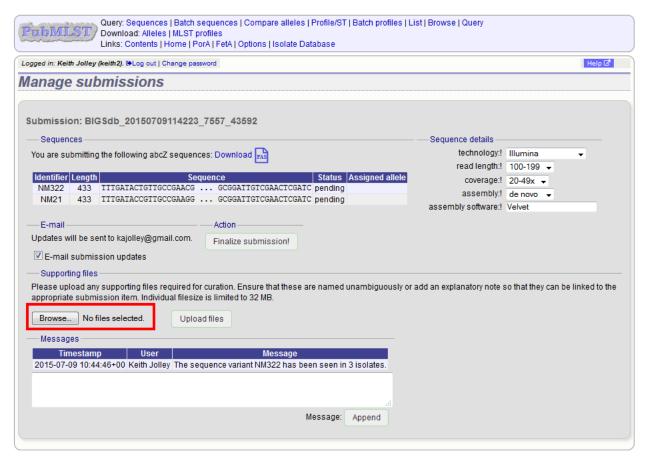
The message will be attached. A curator may respond to the message and attach their own, with the full conversation becoming part of the submission record.



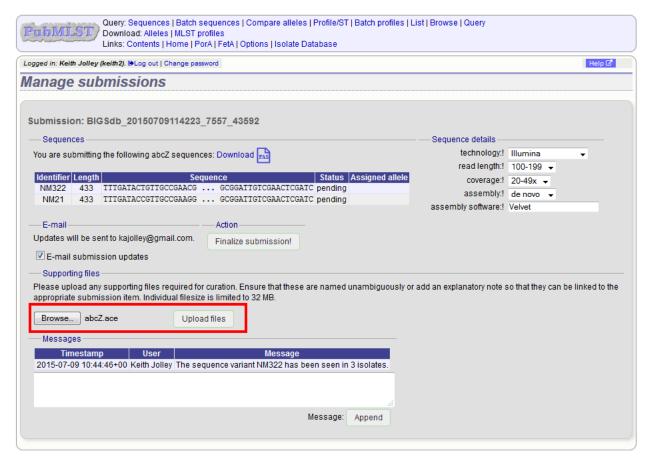
14.2.6 Add supporting files

Some submissions will require the attachment of supporting files. This will depend on the policies of the individual databases. Sequences determined by Sanger sequencing should normally have forward and reverse trace files attached.

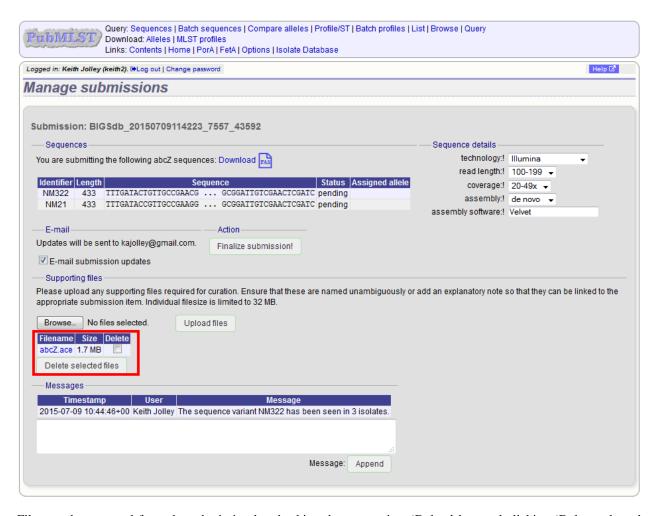
Files can be added to the submission by clicking the 'Browse' button in the 'Supporting files' section.



Select the file in the selection box, then click 'Upload files'.



The file will be uploaded and shown in a table.

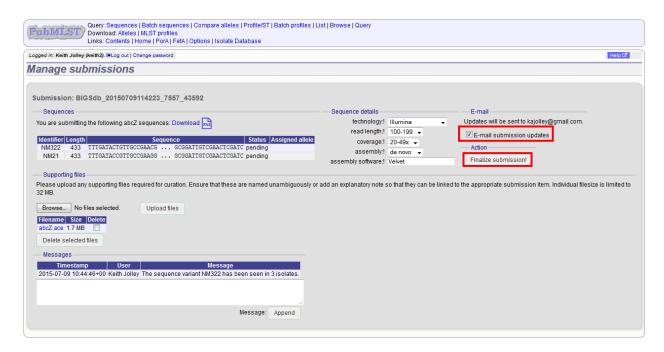


Files can be removed from the submission by checking the appropriate 'Delete' box and clicking 'Delete selected files'.

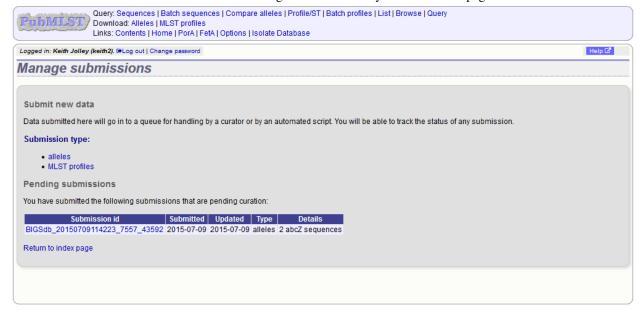
14.2.7 Finalize submission

Make sure the 'E-mail submission updates' box is checked if you wish to receive E-mail notification of the result of your submission. This setting is remembered between submissions.

Click 'Finalize submission!'.



Your submission will then be listed under 'Pending submissions' on your submission page.



14.3 Profile submission

14.3.1 Start

Note: Most MLST databases on PubMLST.org require you to submit an isolate record for each new ST that you wish to be defined. In these cases, you should add the isolate name to the id field of your profile submission and make a corresponding *isolate submission* containing the allelic profile.

Click the appropriate profiles link under submission type on the submission management page.



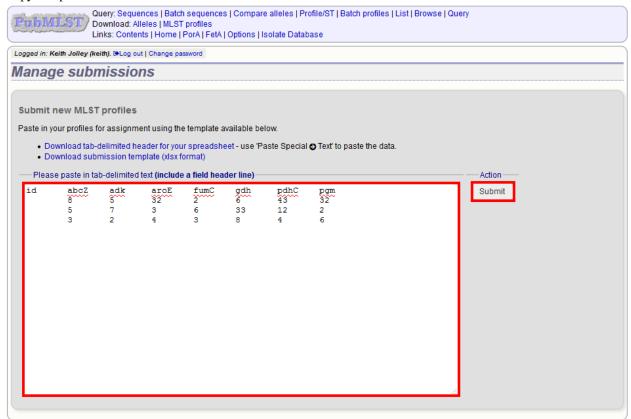
Download the Excel submission template.



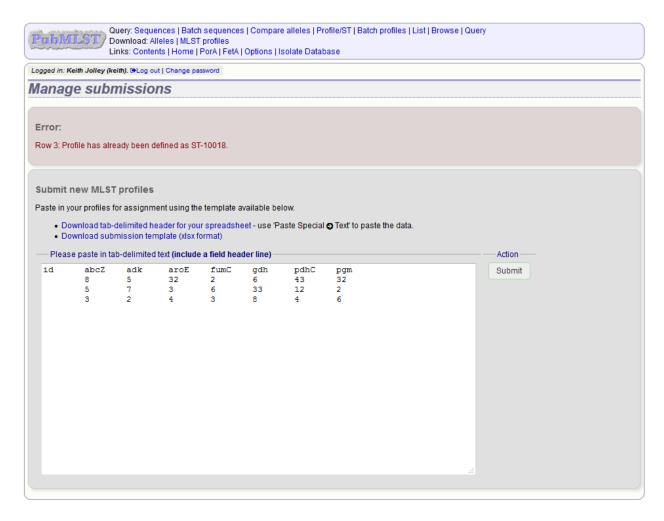
14.3.2 Paste in profile(s)

Fill in the template. The first column 'id' can be used to enter an identifier that is meaningful to you - it is used to report back the results but is not uploaded to the database. It can be left blank, or the entire column can be removed - in which case individual profiles will be identified by row number.

Copy and paste the entire contents of the submission worksheet. Click submit.



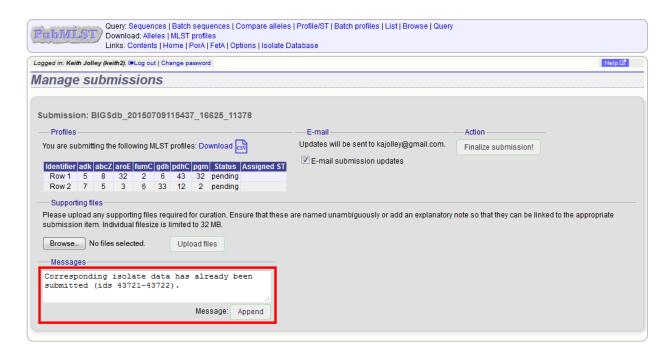
Some basic checks will be performed. These include whether the profile has already been assigned and whether each allele identifier exists. The submission cannot proceed if the checks fail.



Provided the checks pass, you will then be able to add additional information to your submission

14.3.3 Add message to curator

If you wish to enter a message to the curator, enter this in the messages box and click 'Append'.



The message will be attached. A curator may respond to the message and attach their own, with the full conversation becoming part of the submission record.



14.3.4 Add supporting files

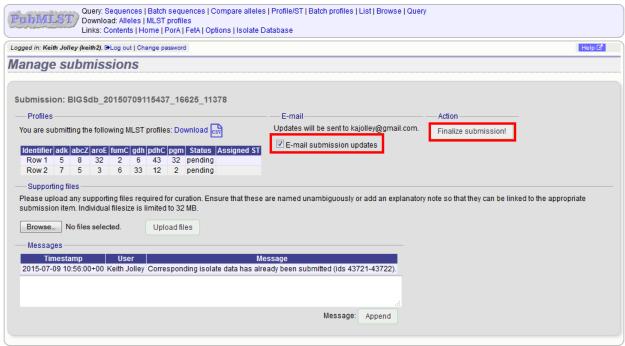
Some submissions may require the attachment of supporting files. These files can be added to the submission by clicking the 'Browse' button in the 'Supporting files' section.

Select the file in the selection box, then click 'Upload files'.

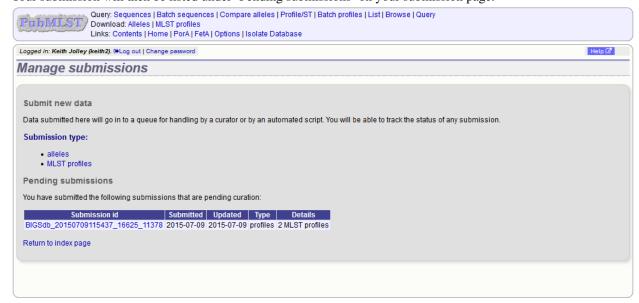
14.3.5 Finalize submission

Make sure the 'E-mail submission updates' box is checked if you wish to receive E-mail notification of the result of your submission. This setting is remembered between sessions.

Click 'Finalize submission!'.



Your submission will then be listed under 'Pending submissions' on your submission page.

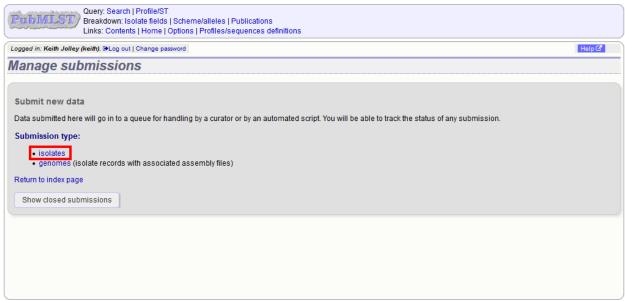


14.4 Isolate submission

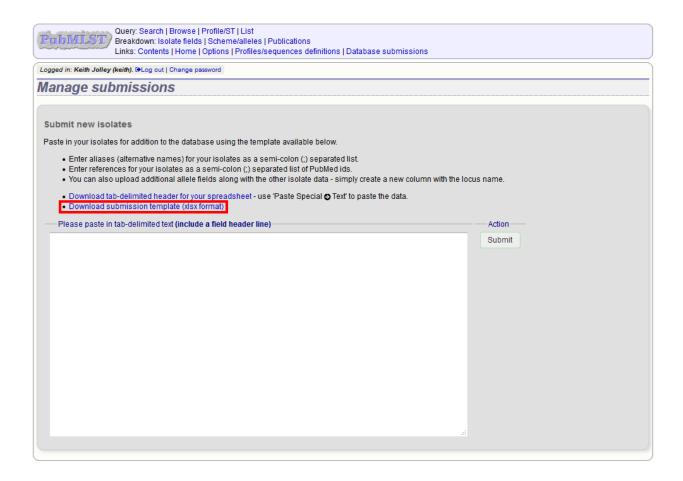
New isolate data can only be submitted from within the appropriate isolate database. You may be required to submit isolate data if you would like to get a new MLST sequence type defined, but this depends on individual database policy.

14.4.1 Start

Click the 'isolates' link under submission type on the submission management page.



Download the Excel submission template.

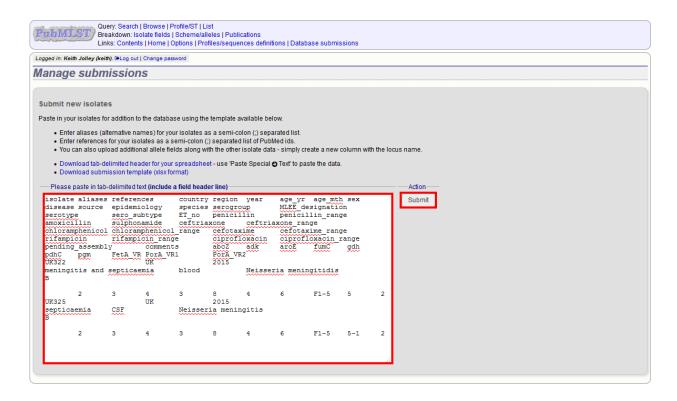


14.4.2 Paste in isolate data

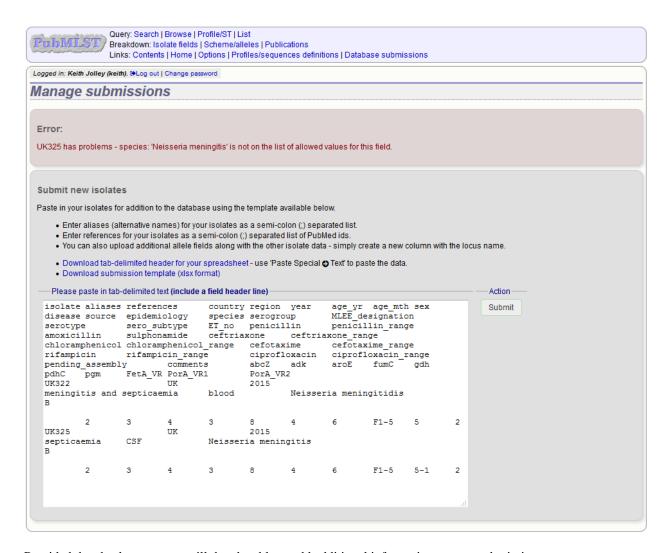
Fill in the template. Some fields are required and cannot be left blank. Check the 'Description of database fields' link on the database contents page to see a description of the fields and allowed values where these have been defined. Where allowed values have been set, the template will have dropdown boxes (although these require newer versions of Excel to work).

Some databases may have hundreds of loci defined, and most will not have a column in the template. You can add new columns for any loci that have been defined and for which you would like to include allelic information for. These locus names must be the primary locus identifier. A list of loci can be found in the 'allowed_loci' tab of the Excel submission template.

Copy and paste the entire contents of the submission worksheet. Click submit.



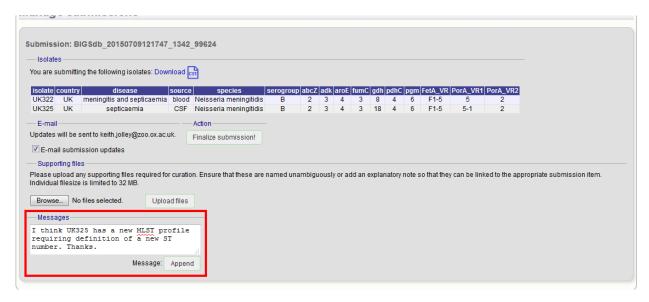
Some basic checks will be performed. These include checking all field values conform to allowed lists or data types. The submission cannot proceed if any checks fail.



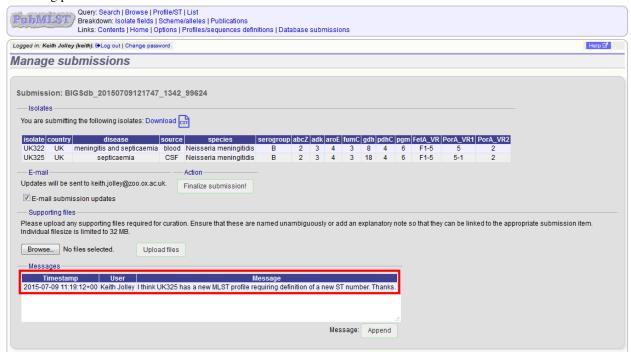
Provided the checks pass, you will then be able to add additional information to your submission.

14.4.3 Add message to curator

If you wish to enter a message to the curator, enter this in the messages box and click 'Append'.



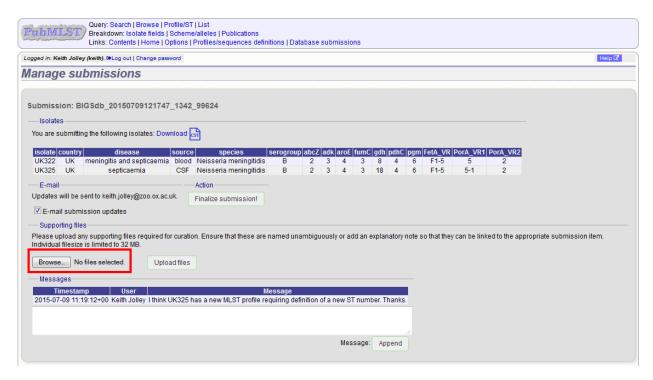
The message will be attached. A curator may respond to the message and attach their own, with the full conversation becoming part of the submission record.



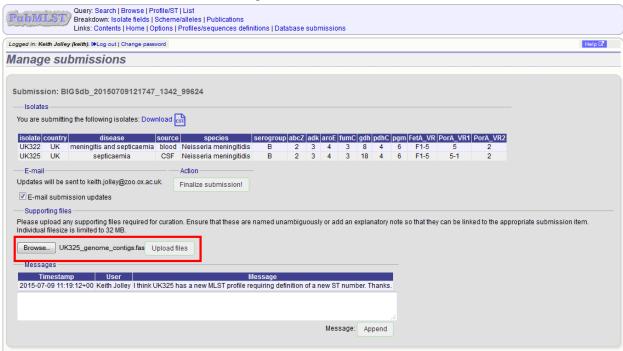
14.4.4 Add supporting files

You can add any files required to support the submission. You may, for example, wish to include a genome sequence for an isolate record (contigs in FASTA format). If you are doing this, make sure that the filename can be unambiguously linked to the appropriate isolate record and *add a message*.

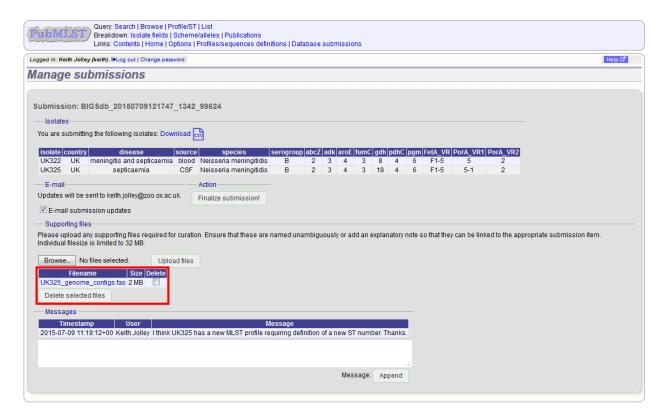
Files can be added to the submission by clicking the 'Browse' button in the 'Supporting files' section.



Select the file in the selection box, then click 'Upload files'.



The file will be uploaded and shown in a table.

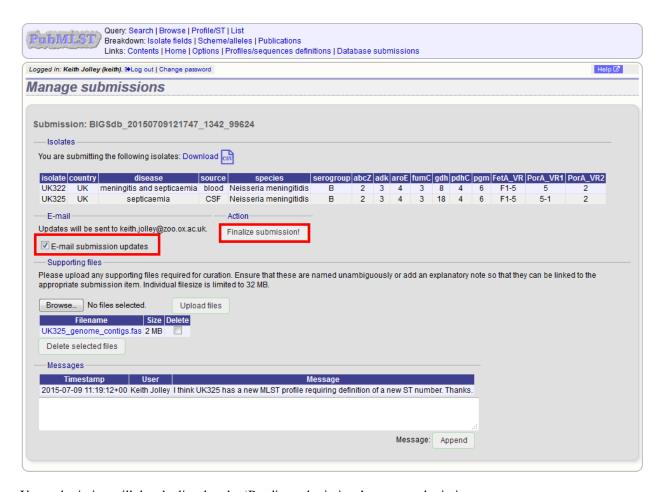


Files can be removed from the submission by checking the appropriate 'Delete' box and clicking 'Delete selected files'.

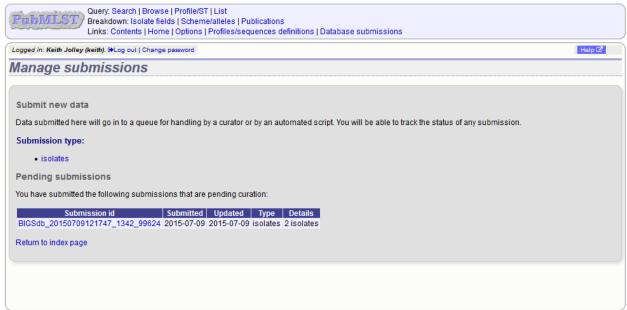
14.4.5 Finalize submission

Make sure the 'E-mail submission updates' box is checked if you wish to receive E-mail notification of the result of your submission. This setting is remembered between sessions.

Click 'Finalize submission!'.



Your submission will then be listed under 'Pending submissions' on your submission page.

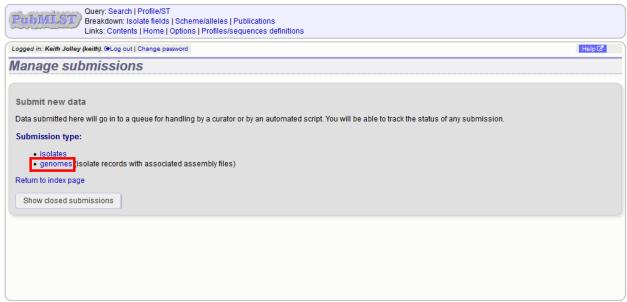


14.5 Genome submission

Submitting genomes uses the same process as standard *isolate submission*. The only difference is that there are a couple of extra required fields in the submission table:

- assembly_filename this is the name of the FASTA file containing the assembly contigs. This must be uploaded
 as a supporting file you will not be able to finalize the submission until every isolate record has a matching
 contig file.
- sequence_method the sequencing technology used to generate the sequences. The allowed values are listed on the submission page.

To start the submission, click the 'genomes' link under submission type on the submission management page.



Then follow the steps for *isolate submission*, uploading the contigs files as supporting files.

14.6 Removing submissions from your notification list

Once a submission has been closed by a curator, the results will be displayed in your 'Manage submissions' area. You can remove submissions once you have noted the result by clicking the 'Remove' link.



Alternatively, submissions will be removed automatically a specified period of time after closure. By default, this time is 90 days, but this can vary depending on the site configuration.

RESTful Application Programming Interface (API)

The REST API allows third-party applications to retrive data stored within BIGSdb databases or to send new submissions to database curators. To use the REST API, your application will make a HTTP request and parse the response. The response format is JSON (except for routes that request a FASTA or CSV file).

Access to protected resources, i.e. those requiring an account, can be accessed via the API using OAuth authentication.

15.1 Passing additional/optional parameters

If you are using a method called with GET, optional parameters can be passed as arguments to the query URL by adding a '?' followed by the first argument and its value (separated by a '='). Additional parameters are separated by a '&', e.g.

http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates?page=2&page_size=100

Methods called with POST require their arguments to be sent as JSON within the post body.

15.2 Resources

- GET / or /db List site resources
- GET /db/{database} List database resources
- GET /db/{database}/loci List loci
- GET /db/{database}/loci/{locus} Retrieve locus record
- GET /db/{database}/loci/{locus}/alleles Retrieve list of alleles defined for a locus
- GET /db/{database}/loci/{locus}/alleles_fasta Download alleles in FASTA format
- GET /db/{database}/loci/{locus}/alleles/{allele_id} Retrieve full allele information
- GET /db/{database}/schemes List schemes
- GET /db/{database}/schemes/{scheme id} Retrieve scheme information
- GET /db/{database}/schemes/{scheme_id}/fields/{field} Retrieve information about scheme field
- GET/db/{database}/schemes/{scheme_id}/profiles List allelic profiles defined for scheme
- GET /db/{database}/schemes/{scheme_id}/profiles_csv Download allelic profiles in CSV (tab-delimited) format

- GET /db/{database}/schemes/{scheme_id}/profiles/{profile_id} Retrieve allelic profile record
- GET /db/{database}/isolates Retrieve list of isolate records
- GET /db/{database}/isolates/{isolate_id} Retrieve isolate record
- GET/db/{database}/isolates/{isolate_id}/allele_designations Retrieve list of allele designations
- GET /db/{database}/isolates/{isolate_id}/allele_designations/{locus} Retrieve full allele designation record
- GET /db/{database}/isolates/{isolate id}/allele ids Retrieve allele identifiers
- GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_designations Retrieve scheme allele designation records
- GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_ids Retrieve list of scheme allele identifiers
- GET /db/{database}/isolates/{isolate_id}/contigs Retrieve list of contigs
- GET /db/{database}/isolates/{isolate_id}/contigs_fasta Download contigs in FASTA format
- GET /db/{database}/contigs/{contig_id} Retrieve contig record
- GET /db/{database}/fields Retrieve list of isolate provenance field descriptions
- GET /db/{database}/users/{user_id} Retrieve user information
- GET /db/{database}/projects Retrieve list of projects
- GET /db/{database}/projects/{project_id} Retrieve project information
- GET/db/{database}/projects/{project_id}/isolates Retrieve list of isolates belonging to a project
- GET /db/{database}/submissions Retrieve list of submissions
- POST/db/{database}/submissions Create new submission
- GET /db/{database}/submissions/{submission_id} Retrieve submission record
- DELETE /db/{database}/submissions/{submission_id} Delete submission record
- GET /db/{database}/submissions/{submission_id}/messages Retrieve submission correspondence
- POST /db/{database}/submissions/{submission_id}/messages Add submission correspondence
- GET/db/{database}/submissions/{submission_id}/files retrieve list of supporting files uploaded for submission
- POST/db/{database}/submissions/{submission_id}/files Upload submission supporting file
- $\bullet \ \textit{GET /db/{database}/submissions/{submission_id}/\textit{files/{filename}} Download \ submission \ supporting \ file \ \textit{files/fil$
- DELETE /db/{database}/submissions/{submission id}/files/{filename} Delete submission supporting file

15.2.1 GET / or /db - List site resources

Required route parameters: None **Optional query parameters:** None

Example request URI: http://rest.pubmlst.org/

Response: List of resource groupings (ordered by name). Groups may consist of paired databases for sequence definitions and isolate data, or any set of related resources. Each group contains:

- name [string] short name (usually a single word)
- · description [string] fuller description

- databases [array] list of database objects, each consists of three key/value pairs:
 - name [string] name of database config
 - description [string] short description of resource
 - href [string] URI to access resource

15.2.2 GET /db/{database} - List database resources

These will vary depending on whether the resource is an isolate or a sequence definition database.

Required route parameter: database [string] - Database configuration name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates

Response: Object containing a subset of the following key/value pairs:

- fields [string] URI to isolate provenance field information
- isolates [string] URI to isolate records
- schemes [string] URI to list of schemes
- loci [string] URI to list of loci
- projects [string] URI to list of projects

15.2.3 GET /db/{database}/loci - List loci

Required route parameter: database [string] - Database configuration name

Optional parameters:

- page [integer]
- page size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/loci

Response: Object containing:

- records [int] Number of loci.
- loci [array] List of *URIs to defined locus records*. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2. Resources 347

15.2.4 GET /db/{database}/loci/{locus} - Retrieve locus record

Provides information about a locus, including links to allele sequences (in seqdef databases).

Required route parameters:

- database [string] Database configuration name
- locus [string] Locus name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/loci/abcZ

Response: Object containing a subset of the following key/value pairs:

- id [string] locus name
- data_type [string] 'DNA' or 'peptide'
- allele_id_format [string] 'integer' or 'text'
- allele_id_regex [string] regular expression constraining allele ids
- common_name [string]
- aliases [array] list of alternative names of the locus
- length_varies [boolean]
- length [integer] length if alleles are of a fixed length
- coding_sequence [boolean]
- orf [integer] 1-6
- schemes [array] list of scheme objects, each consisting of:
 - scheme [string] URI to scheme information
 - description [string]
- min_length [integer] (seqdef databases) minimum length for variable length loci
- max length [integer] (seqdef databases) maximum length for variable length loci
- alleles [string] (seqdef databases) URI to list of allele records
- alleles_fasta [string] (seqdef databases) URI to FASTA file of all alleles of locus
- curators [array] (seqdef databases) list of URIs to user records of curators of the locus
- publications [array] (segdef databases) list of PubMed id numbers of papers describing the locus
- full_name [string] (seqdef databases)
- product [string] (seqdef databases)
- description [string] (seqdef databases)
- extended_attributes [array] (seqdef databases) list of extended attribute objects. Each consists of a subset of the following fields:
 - field [string] field name
 - value_format [string] 'integer', 'text', or 'boolean'
 - value_regex [string] regular expression constraining value
 - description [string] description of field

- length [integer] maximum length of field
- required [boolean]
- allowed_values [array] list of allowed values
- genome_position [integer] (isolate databases)

15.2.5 GET /db/{database}/loci/{locus}/alleles - Retrieve list of alleles defined for a locus

Required route parameters:

- database [string] Database configuration name
- locus [string] Locus name

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.
- added_after [date] Include only alleles added after specified date (ISO 8601 format).
- updated_after [date] Include only alleles last modified after specified date (ISO 8601 format).

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/loci/abcZ/alleles

Response: Object containing:

- records [int] Number of alleles
- alleles [array] List of *URIs to defined allele records*. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.6 GET /db/{database}/loci/{locus}/alleles_fasta - Download alleles in FASTA format

Required route parameters:

- database [string] Database configuration name
- locus [string] Locus name

Optional parameters:

- added_after [date] Include only alleles added after specified date (ISO 8601 format).
- updated_after [date] Include only alleles last modified after specified date (ISO 8601 format).

15.2. Resources 349

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/loci/abcZ/alleles_fasta

Response: FASTA format file of allele sequences

15.2.7 GET /db/{database}/loci/{locus}/alleles/{allele_id} - Retrieve full allele information

Required route parameters:

• database [string] - Database configuration name

- locus [string] Locus name
- allele_id [string] Allele identifier

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/loci/abcZ/alleles/5

Response: Object containing the following key/value pairs:

- locus [string] URI to locus description
- allele_id [string] allele identifier
- sequence [string] sequence
- status [string] either 'Sanger trace checked', 'WGS: manual extract', 'WGS: automated extract', or 'unchecked'
- sender [string] URI to user details of sender
- curator [string] URI to user details of curator
- date_entered [string] record creation date (ISO 8601 format)
- datestamp [string] last updated date (ISO 8601 format)

15.2.8 GET /db/{database}/schemes - List schemes

Required route parameter: database [string] - Database configuration name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes

Response:

- records [integer] Number of schemes
- schemes [array] list of scheme objects, each containing:
 - scheme [string] URI to scheme information
 - description [string]

15.2.9 GET /db/{database}/schemes/{scheme_id} - Retrieve scheme information

Includes links to allelic profiles (in seqdef databases, if appropriate). Required route parameters:

- database [string] Database configuration name
- scheme_id [integer] Scheme id number

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes/1

Response: Object containing a subset of the following key/value pairs:

- id [integer]
- description [string]
- locus count [integer] number of loci belonging to scheme
- loci [array] list of URIs to locus descriptions
- has_primary_key_field [boolean]
- fields [array] list of URIs to scheme field descriptions
- primary_key_field [string] URI to primary key field description
- profiles [string] URI to list of profile definitions (only seqdef databases)
- profiles_csv [string] URI to tab-delimited file of all scheme profiles
- curators [array] (seqdef databases) list of URIs to user records of curators of the scheme

15.2.10 GET /db/{database}/schemes/{scheme_id}/fields/{field} - Retrieve information about scheme field

Required route parameters:

- database [string] Database configuration name
- scheme_id [integer] Scheme id number
- field [string] Field name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes/1/fields/ST

Response: Object containing the following key/value pairs:

- field [string] field name
- type [string] data type of field (integer or text)
- primary_key [boolean] true if field is the scheme primary key

15.2.11 GET /db/{database}/schemes/{scheme_id}/profiles - List allelic profiles defined for scheme

Required route parameters:

- database [string] Database configuration name
- scheme_id [integer] Scheme id

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

- added_after [date] Include only profiles added after specified date (ISO 8601 format).
- updated_after [date] Include only profiles last modified after specified date (ISO 8601 format).

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes/1/profiles

Response: Object containing:

- records [int] Number of profiles
- profiles [array] List of URIs to defined profile records. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.12 GET /db/{database}/schemes/{scheme_id}/profiles_csv - Download allelic profiles in CSV (tab-delimited) format

Required route parameters:

- database [string] Database configuration name
- scheme_id [integer] Scheme id

Optional parameters:

- added_after [date] Include only profiles added after specified date (ISO 8601 format).
- updated_after [date] Include only profiles last modified after specified date (ISO 8601 format).

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes/1/profiles_csv

Response: Tab-delimited text file of allelic profiles

15.2.13 GET /db/{database}/schemes/{scheme_id}/profiles/{profile_id} - Retrieve allelic profile record

Required route parameters:

- database [string] Database configuration name
- scheme_id [integer] Scheme id
- profile_id [string/integer] Profile id

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/schemes/1/profiles/11

Response: Object containing the following key/value pairs:

- primary_key_term [string/integer] The field name is the primary key, e.g. ST. The value is the primary key value (primary_id used as an argument).
- alleles [object] list of URIs to allele descriptions

- other_scheme_fields [string/integer] Each scheme field will have its own value if defined. The field name is the name of the field.
- sender [string] URI to user details of sender
- curator [string] URI to user details of curator
- date_entered [string] record creation date (ISO 8601 format)
- datestamp [string] last updated date (ISO 8601 format)

15.2.14 GET /db/{database}/isolates - Retrieve list of isolate records

Required route parameter: database [string] - Database configuration name

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.
- added_after [date] Include only isolates added after specified date (ISO 8601 format).
- updated_after [date] Include only isolates last modified after specified date (ISO 8601 format).

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates

Response: Object containing:

- records [int] Number of isolates
- isolates [array] List of URIs to isolate records. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.15 GET /db/{database}/isolates/{isolate_id} - Retrieve isolate record

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1

Response: Object containing some or all of the following key/value pairs:

• provenance [object] - set of key/value pairs. Keys are defined by calling the /fields route. The fields will vary by database but will always contain the following:

- id [integer]

- sender [string] URI to user details of sender
- curator [string] URI to user details of curator
- date_entered [string] record creation date (ISO 8601 format)
- datestamp [string] last updated date (ISO 8601 format)
- publications [array] (seqdef databases) list of PubMed id numbers of papers that refer to the isolate
- sequence_bin [object] consists of the following key/value pairs:
 - contigs_fasta [string] URI to FASTA file containing all the contigs belonging to this isolate
 - contigs [string] URI to list of contig records
 - contig_count [integer] number of contigs
 - total_length [integer] total length of contigs
- allele_designations [object] consists of the following key/value pairs:
 - allele_ids URI to list of all allele_id values defined for the isolate
 - designation_count number of allele designations defined for the isolate
 - full_designations URI to list of full allele designation records
- schemes [array] list of scheme objects, each containing the following:
 - description [string] description of scheme
 - loci_designated_count [integer] number of loci within scheme that have an allele designated for this isolate.
 - allele_ids [string] URI to list of all allele_id values defined for this scheme for this isolate
 - full_designations [string] URI to list of full allele designation records for this isolate
 - fields [object] consisting of key/value pairs where the key is the name of each scheme field
- projects [array] list of project objects, each containing the following:
 - id [string] URI to project information
 - description [string] description of project
- new_version [string] URI to newer version of record
- old_version [string] URI to older version of record

15.2.16 GET /db/{database}/isolates/{isolate_id}/allele_designations - Retrieve list of allele designation records

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1/allele_designations

Response: Object containing:

- records [int] Number of allele designations
- allele_designations [array] List of *URIs to allele designation records*. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.17 GET /db/{database}/isolates/{isolate_id}/allele_designations/{locus} - Retrieve full allele designation record

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier
- locus [string] Locus name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1/allele_designations/BACT000065

Response: List of allele_designation objects (there may be multiple designations for the same locus), each containing:

- locus [string] URI to locus description
- allele_id [string]
- method [string] either 'manual' or 'automatic'
- status [string] either 'confirmed' or 'provisional'
- comments [string]
- sender [string] URI to user details of sender
- curator [string] URI to user details of curator
- datestamp [string] last updated date (ISO 8601 format)

15.2.18 GET /db/{database}/isolates/{isolate_id}/allele_ids - Retrieve allele identifiers

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1/allele_ids

Response: Object containing:

- records [int] Number of allele id objects
- allele_ids [array] List of allele id objects, each consisting of a key/value pair where the key is the locus name. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.19 GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_designations - Retrieve scheme allele designation records

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier
- scheme_id [integer] Scheme identifier

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1/schemes/1/allele_designations

Response:

- records [int] Number of allele designation objects
- allele_designations [array] List of *allele designation objects* for each locus in the specified scheme that has been designated.

15.2.20 GET /db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele_ids - Retrieve list of scheme allele identifiers

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier
- scheme_id [integer] Scheme identifier

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/isolates/1/schemes/1/allele_ids

Response:

- records [int] Number of allele id objects
- allele_ids [array] List containing allele id objects for each locus in the specified scheme that has been designated. Each allele_id object contains a key which is the name of the locus with a value that may be either a string, integer or array of strings or integers (required where there are multiple designations for a locus). The data type depends on the allele_id_format set for the specific locus.

15.2.21 GET /db/{database}/isolates/{isolate id}/contigs - Retrieve list of contigs

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst neisseria isolates/isolates/1/contigs

Response: Object containing:

- records [int] Number of contigs
- contigs [array] List of *URIs to contig records* Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.22 GET /db/{database}/isolates/{isolate_id}/contigs_fasta - Download contigs in FASTA format

Required route parameters:

- database [string] Database configuration name
- isolate_id [integer] Isolate identifier

Optional parameter:

• header [string] - either 'original_designation' or 'id' (default is 'id'). This selects whether the FASTA header lines contain the originally uploaded FASTA headers or the sequence bin id numbers.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/1/contigs_fasta?header=original_designation

Response: FASTA format file of isolate contig sequences

15.2.23 GET /db/{database}/contigs/{contig id} - Retrieve contig record

Required route parameters:

- database [string] Database configuration name
- contig_id [integer] Contig identifier

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/contigs/180062

Response: Contig object consisting of the following key/value pairs:

- id [integer] contig identifier
- isolate_id [integer] isolate identifier
- sequence [string] contig sequence
- · length [integer] length of contig sequence
- method [string] sequencing method
- sender [string] URI to user details of sender
- curator [string] URI to user details of curator
- date_entered [string] record creation date (ISO 8601 format)
- datestamp [string] last updated date (ISO 8601 format)

15.2.24 GET /db/{database}/fields - Retrieve list of isolate provenance field descriptions

Required route parameters:

• database [string] - Database configuration name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/fields

Response: Array of field objects, each consisting of some or all of the following key/value pairs:

- name [string] name of field
- type [string] data type (int, text, date, float)
- length [integer] maximum length of field
- required [boolean] true if field value is required
- min [integer] minimum value for integer values
- max [integer] maximum value for integer values
- regex [string] regular expression that constrains the allowed value of the field
- comments [string]
- allowed values [array] list of allowed values for the field [string]

15.2.25 GET /db/{database}/users/{user id} - Retrieve user information

Users may be data submitters or curators.

Required route parameters:

- database [string] Database configuration name
- user_id [integer] User id number

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/users/2

Response: Object containing the following key/value pairs:

- id [integer] user id number
- first_name [string]
- surname [string]
- affiliation [string] institutional affiliation
- email [string] E-mail address

15.2.26 GET /db/{database}/projects - Retrieve list of projects

Required route parameter: database [string] - Database configuration name

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/projects

Response:

- projects [array] List of project objects, each containing:
 - project [string] URI to project information
 - description [string]
 - isolate_count [integer] number of isolates in project

15.2.27 GET /db/{database}/projects/{project_id} - Retrieve project information

Required route parameters:

- database [string] Database configuration name
- project_id [integer] Project id number

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/projects/3

Response: Object containing a subset of the following key/value pairs:

- id [integer]
- description [string]
- isolates [string] URI to list of URIs of member isolate records.

15.2.28 GET /db/{database}/projects/{project_id}/isolates - Retrieve list of isolates belonging to a project

Required route parameter:

- database [string] Database configuration name
- project_id [integer] Project id number

Optional parameters:

- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/projects/3/isolates

Response: Object containing:

- records [int] Number of isolates in the project
- isolates [array] List of URIs to isolate records. Pages are 100 records by default. Page size can be modified using the page_size parameter.
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results
 - first URI to first page of results
 - last URI to last page of results
 - return_all URI to page containing all results (paging disabled)

15.2.29 GET /db/{database}/submissions - retrieve list of submissions

Required route parameter: database [string] - Database configuration name

Optional parameters:

- type [string] either 'alleles', 'profiles' or 'isolates'
- status [string] either 'closed' or 'pending'
- page [integer]
- page_size [integer]
- return_all [integer] Set to non-zero value to disable paging.

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_isolates/submissions

Response: Object containing:

- records [int] Number of submissions
- submissions [array] List of URIs to submission records
- paging [object] Some or all of the following:
 - previous URI to previous page of results
 - next URI to next page of results

- first URI to first page of results
- last URI to last page of results
- return_all URI to page containing all results (paging disabled)

15.2.30 POST /db/{database}/submissions - create new submission

Required route parameter: database [string] - Database configuration name

Required additional parameters:

- type [string] either:
 - alleles (sequence definition databases only)
 - profiles (sequence definition databases only)
 - isolates (isolate databases only)
 - genomes (isolate databases only)

The following are required with the specified database type:

Allele submissions

- locus [string] name of locus
- technology [string] name of sequencing technology: either '454', 'Illumina', 'Ion Torrent', 'PacBio', 'Oxford Nanopore', 'Sanger', 'Solexa', 'SOLiD', or 'other'
- read_length [string] read length of sequencing: either '<100', '100-199', '200-299', '300-499', '>500', or any positive integer (only required for Illumina)
- coverage [string] mean coverage of sequencing: either '<20x', '20-49x', '50-99x', '>100x', or any positive integer (only required for Illumina)
- · assembly [string] assembly method: either 'de novo' or 'mapped'
- software [string] name of assembly software
- sequences [string] either single raw sequence or multiple sequences in FASTA format

Profile submissions

- scheme id [int] scheme id number
- profiles [string] tab-delimited profile data this should include a header line containing the name of each locus

Isolate submissions

 isolates [string] - tab-delimited isolate data - this should include a header line containing each field or locus included

Genome submissions

• isolates [string] - tab-delimited isolate data - this should include a header line containing each field or locus included as well as for 'assembly_filename' and 'sequence_method'. The 'sequence_method' should be either '454', 'Illumina', 'Ion Torrent', 'PacBio', 'Oxford Nanopore', 'Sanger', 'Solexa', 'SOLiD', or 'other'. Following submission, contig files should be uploaded with the same names as set for 'assembly_filename'. This can be done using the *file upload route*.

Optional parameters:

- message [string] correspondence to the curator
- email [int] set to 1 to enable E-mail updates (E-mails will be sent to the registered user account address).

Response: Object containing:

• submission - URI to submission record

For genome submissions, the response object will also contain:

- missing_files [array] List of filenames that need to be uploaded to complete the submission. These filenames
 are defined in the 'assembly_filename' field of the isolate record upload. The files should contain the contig
 assemblies.
- message [string] 'Please upload missing contig files to complete submission.'

15.2.31 GET /db/{database}/submissions/{submission_id} - Retrieve submission record

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/submissions/BIGSdb_20151013081836_14559_14740

Response: Object containing some of the following:

- id [string] Submission id
- type [string] Either 'alleles', 'profiles', 'isolates'
- date_submitted [string] Submission date (ISO 8601 format)
- datestamp [string] Last updated date (ISO 8601 format)
- submitter [string] URI to user details of submitter
- curator [string] URI to user details of curator
- status [string] either 'started', 'pending', or 'closed'
- outcome [string] either 'good' (data uploaded), 'bad' (data rejected), or 'mixed' (parts of submission accepted)
- correspondence [array] List of correspondence objects in time order. Each contains:
 - user [string] URI to user details of user
 - timestamp [string]
 - message [string]

Allele submissions

- locus [string] name of locus
- technology [string] name of sequencing technology: either '454', 'Illumina', 'Ion Torrent', 'PacBio', 'Oxford Nanopore', 'Sanger', 'Solexa', 'SOLiD', or 'other'
- read_length [string] read length of sequencing: either '<100', '100-199', '200-299', '300-499', '>500', or any positive integer (only required for Illumina)
- coverage [string] mean coverage of sequencing: either '<20x', '20-49x', '50-99x', '>100x', or any positive integer (only required for Illumina)
- assembly [string] assembly method: either 'de novo' or 'mapped'
- software [string] name of assembly software

- seqs [array] List of sequence objects each containing:
 - seq_id [string] Sequence identifier
 - assigned_id [string] Allele identifier if uploaded to the database (otherwise undefined)
 - status [string] Either 'pending', 'assigned', or 'rejected'
 - sequence [string]

Profile submissions

- scheme [string] URI to scheme information
- profiles [array] List of profile record objects. Each contains:
 - profile_id [string] Record identifier
 - assigned_id [string] Profile identifier if uploaded to the database (otherwise undefined)
 - status [string] Either 'pending', 'assigned', or 'rejected'
 - designations [object] containing key/value pairs for each locus containing the allele identifier

Isolate submissions

isolates [array] - List of isolate record objects. Each contains key/value pairs for included fields.

15.2.32 DELETE /db/{database}/submissions/{submission_id} - Delete submission record

You must be the owner and the record must be closed.

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/submissions/BIGSdb_20151013081836_14559_14740

Response: message [string] - 'Submission deleted.'

15.2.33 GET /db/{database}/submissions/{submission_id}/messages - Retrieve submission correspondence

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/submissions/BIGSdb_20151013081836_14559_14740/me

Response: Array of correspondence objects in time order. Each contains:

- user [string] URI to user details of user
- timestamp [string]
- message [string]

15.2.34 POST /db/{database}/submissions/{submission_id}/messages - Add submission correspondence

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Required additional parameter:

• message [string] - Message text

Optional parameters: None

Response: message [string] - 'Message added.'

15.2.35 GET /db/{database}/submissions/{submission_id}/files - Retrieve list of supporting files uploaded for submission

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Optional parameters: None

Example request URI: http://rest.pubmlst.org/db/pubmlst_neisseria_seqdef/submissions/BIGSdb_20151013081836_14559_14740/file

Response: Array of URIs to files

15.2.36 POST /db/{database}/submissions/{submission_id}/files - Upload submission supporting file

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id

Required additional parameters:

- filename [string] Name of file to store within submission
- upload [base64 encoded data] Raw file data

Optional parameters: None

Response: message [string] - 'File uploaded.'

15.2.37 GET /db/{database}/submissions/{submission_id}/files/{filename} - Download submission supporting file

Required route parameters:

- database [string] Database configuration name
- submission_id [string] Submission id
- filename [string] Name of file

Optional parameters: None Response: File download

15.2.38 DELETE /db/{database}/submissions/{submission_id}/files/{filename} - Delete submission supporting file

Required route parameters:

- database [string] Database configuration name
- submission id [string] Submission id
- filename [string] Name of file

Optional parameters: None

Response: message [string] - 'File deleted.'

15.3 Authentication

Protected resources, i.e. those requiring a user to log in, can be accessed via the API using OAuth (1.0A) authentication (see IETF RFC5849 for details). Third-party client software has to be registered with the BIGSdb site before they can access authenticated resources. The overall three-legged flow works as follows:

- 1. Developer signs up and gets a consumer key and consumer secret specific to their application.
- 2. Application gets a request token and directs user to authorization page on BIGSdb.
- 3. BIGSdb *asks user for authorization* for application to access specific resource using their credentials. A verifier code is provided.
- 4. The application exchanges the request token and OAuth verifier code for an *access token and secret* (these do not expire but may be revoked by the user or site admin).
- 5. Application uses access token/secret to request session token (this is valid for 12 hours).
- 6. All calls to access protected resources are signed using the session token/secret and consumer key/secret.

It is recommended that application developers use an OAuth library to generate and sign requests.

15.3.1 Developer sign up to get a consumer key

Application developers should apply to the site administrator of the site running BIGSdb. The administrator can *generate a key and secret* using a script - both of these will need to be used by the application to sign requests.

The client id is usually a 24 character alphanumeric string. The secret is usually a 42 character alphanumeric (including punctuation) string, e.g.

- client_id: efKXmqp2D0EBlMBkZaGC2lPf
- client_secret: F\$M)_+fQ2AFFB2YBDfF9fpHF^qSWJdmmN%L4Fxf5Gur3

15.3. Authentication 365

15.3.2 Getting a request token

- **Relative URL:** /db/{database}/oauth/get_request_token
- Supported method: GET

The application uses the consumer key to obtain a request token. The request token is a temporary token used to initiate user authorization for the application and will expire in 60 minutes. The request needs to contain the following parameters and to be signed using the consumer secret:

- · oauth_consumer_key
- oauth_request_method ('GET')
- oauth request url (request URL)
- oauth signature method ('HMAC-SHA1')
- · oauth_signature
- oauth_timestamp (UNIX timestamp seconds since Jan 1 1970) this must be within 600 seconds of the current time.
- oauth callback ('oob' for desktop applications)
- oauth_nonce (random string)
- oauth_version ('1.0')

If the application has been registered and has been granted permission to access the specific resource, a JSON response will be returned containing the following parameters:

· oauth_token

- This is the request token. It is usually a 32 character alphanumeric string.
- e.g. fKFm0WNhCfbEX8zQm6qhDA8K23FOWDGE

oauth_token_secret

- This is the secret associated with the request token. It is usually a 32 character alphanumeric string.
- e.g. aZ0fncP7i5w5jlebdK5zyQ4vrRRVcdnv

· oauth_callback_confirmed

- This parameter is always set to true.

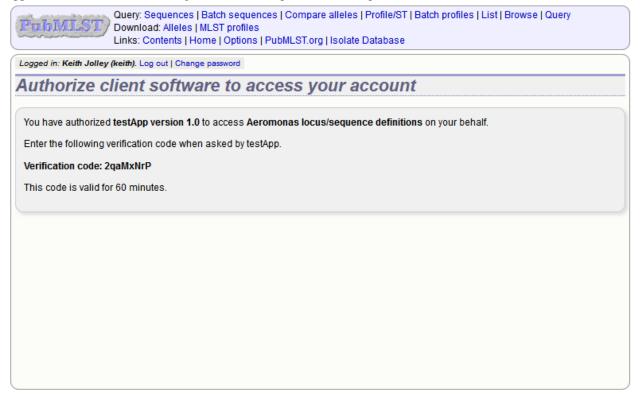
15.3.3 Getting user authorization

Once a request token has been obtained, this can be used by the end user to grant permission to access a specific resource to the application. The application should direct the user to the client authorization page (authorizeClient) specific to a database within BIGSdb, e.g. http://pubmlst.org/cgi-bin/bigsdb/bigsdb.pl?db=pubmlst_neisseria_seqdef&page=authorizeClient&oauth_token=fKFm0WNhCfbEX8zQm6qhDA8K23FOWI

The user will be asked if they wish to grant access to the application on their behalf:



If they authorize the access, they will be presented with a verifier code. This should be entered in to the client application which will use this together with the request token to request an access token.



The verifier code is valid for 60 minutes.

15.3. Authentication 367

15.3.4 Getting an access token

- **Relative URL:** /db/{database}/oauth/get_access_token
- Supported method: GET

The application uses the request token, verifier code and its consumer key to obtain an access token. The access token does not expire but can be revoked by both the end user or the site admininstrator. The request needs to contain the following parameters and to be signed using the consumer secret and request token secret:

- · oauth_consumer_key
- oauth_request_method ('GET')
- oauth request url (request URL)
- oauth signature method ('HMAC-SHA1')
- · oauth_signature
- oauth_token (request token)
- oauth_timestamp (UNIX timestamp seconds since Jan 1 1970) this must be within 600 seconds of the current time.
- oauth_nonce (random string)
- oauth_version ('1.0')

If the application has been registered and has been granted permission to access the specific resource, a JSON response will be returned containing the following parameters:

- · oauth_token
 - This is the access token. It is usually a 32 character alphanumeric string.
 - e.g. SDrC74ZV15SYSqY8lWZqrRxnyDnNGVFO
- oauth_token_secret
 - This is the secret associated with the access token. It is usually a 32 character alphanumeric string.
 - e.g. tYI2SPzgiO02IRVzW4JR1ez6Vvm4gVyv

15.3.5 Getting a session token

- **Relative URL:** /db/{database}/oauth/get_session_token
- Supported method: GET

The application uses the access token and its consumer key to obtain a session token. The session token is valid for 12 hours before it expires. The request needs to contain the following parameters and to be signed using the consumer secret and access token secret:

- · oauth consumer key
- oauth_request_method ('GET')
- oauth_request_url (request URL)
- oauth_signature_method ('HMAC-SHA1')
- · oauth_signature
- oauth_token (access token)

- oauth_timestamp (UNIX timestamp seconds since Jan 1 1970) this must be within 600 seconds of the current time.
- oauth_nonce (random string)
- oauth_version ('1.0')

If the application has been registered and has been granted permission to access the specific resource, a JSON response will be returned containing the following parameters:

· oauth token

- This is the session token. It is usually a 32 character alphanumeric string.
- e.g. H8CjIS8Ikq6hwCUqUfF114pTaCY18Ljw

· oauth_token_secret

- This is the secret associated with the session token. It is usually a 32 character alphanumeric string.
- e.g. RfponbaNPO7tkZ2miHFISk0pMndePNfJ

15.3.6 Accessing protected resources

• Supported method: GET

The application uses the session token and its consumer key to access a protected resource. The request needs to contain the following parameters and to be signed using the consumer secret and session token secret:

- oauth_consumer_key
- oauth_request_method ('GET')
- oauth_request_url (request URL)
- oauth_signature_method ('HMAC-SHA1')
- · oauth_signature
- oauth_token (session token)
- oauth_timestamp (UNIX timestamp seconds since Jan 1 1970) this must be within 600 seconds of the current time.
- oauth_nonce (random string)
- oauth_version ('1.0')

15.3. Authentication 369

Frequently asked questions (FAQs)

16.1 General

1. What is the minimum specification of hardware required to run BIGSdb?

The software will run on very modest hardware - a number of PubMLST mirrors have been set up on virtual machines with 1 processor core and 4 GB RAM. This should be considered an absolute minimum specification though. For an installation with only local users, the following minimum is recommended:

- 4 processor cores
- 16 GB RAM
- 50 GB partition for temporary files
- 100 GB partition for databases

As usual, the more RAM that is available the better. Ideally you would want enough RAM that the whole database(s) can reside in memory (an approximation is roughly twice the total size of your contigs), although this is not absolutely required.

Offline jobs, such as *Genome Comparator* will use a processor core each, so if you want to run multiple jobs in parallel then you may want more cores (and memory). Tagging of new genomes using the offline *autotagger* can be run in multi-threaded mode so the more cores available the faster this will be.

As a comparison, the PubMLST site is run on two machines - separate web and database servers. All offline jobs and tagging of genomes is performed on the database server. These have the following specification:

- web server: 16 cores, 64GB RAM
- database server: 64 cores, 1TB RAM, 3TB RAID 10 local storage
- 2. Why might icons be missing when using Internet Explorer?

This can occur if you have Compatibility Mode enabled. BIGSdb generates valid HTML5 and Compatibility Mode should not be used. Please ensure this is not enabled in the Internet Explorer tools section.

16.2 Installation

- 1. BIGSdb is accumulating files in various temp directories is this normal and how do I clean them out? See: *Periodically delete temporary files*.
- 2. BIGSdb is complaining of an invalid script path what does this mean?

In your database config.xml file system tag are two attributes - script_path_includes and curate_path_includes. These contain regexes that the web url to your script (bigsdb.pl and bigscurate.pl respectively) must match. This prevents somebody from accessing a private database using an instance of bigsdb.pl that is not in a protected directory if you're using apache authentication.

So, if you access the script from http://localhost/cgi-bin/bigsdb/bigsdb.pl then you can set script_path_includes to something like "/bigsdb/" (which is the default), or "/cgi-bin/" or just "/" if you don't care about this check.

16.3 Administration

1. How can I make some isolates public but not others?

The easiest way to do this is to set up two or more separate configuration directories that refer to the database. The URLs to access these will differ by the value of the 'db' attribute, which refers to the name of the configuration directory (in /etc/bigsdb/dbases/). The database view accessed by each of these configurations can be different as can the access restrictions.

Example:

We have a database 'bigsdb_test' that contains data, only some of which we wish to make publicly available. The isolates to make public are all members of a project. First we can make a view of the isolates table that contains only isolates within this project.

For isolates in project id 3, create a database view by logging in to psql as the postgresql user. We will name this view 'public'.:

```
sudo su postgres
psql bigsdb_test

CREATE VIEW public AS SELECT * FROM isolates WHERE id IN (SELECT isolate_id
   FROM project_members WHERE project_id=3);
GRANT SELECT ON public TO apache;
```

Create a private configuration that can access everything in the database in /etc/bigsdb/dbases/test_private. This will be accessible from http://IP_ADDRESS/cgi-bin/bigsdb/bigsdb.pl?db=test_private.

The important attributes to set in the system tag of the config.xml file in this directory are::

```
view="isolates"
read_access="authenticated_users"
```

This means that anyone with an account can log in and view all the isolates (because the view is set to the isolates table).

Now create a public configuration in /etc/bigsdb/dbases/test_public. This will be accessible from http://IP_ADDRESS/cgi-bin/bigsdb/bigsdb.pl?db=test_public. It is better to create a symlink to the private config.xml and then override the attributes that are different. So create a symlink to the private config file:

```
cd /etc/bigsdb/dbases/test_public
sudo ln -s ../test_private/config.xml .
```

You can now override the view and access settings. Within /etc/bigsdb/dbases/test_public, create a file called system.overrides and add the following:

```
view="public"
read_access="public"
```

See also Restricting particular configurations to specific user accounts.

16.3. Administration 373

Appendix

17.1 Query operators

Various query forms have operators for use with field values. Available operators are:

- =
- Exact match (case-insensitive).
- · contains
 - Match to a partial string (case-insensitive), e.g. searching for clonal complex 'contains' st-11 would return all STs belonging to the ST-11 complex.
- · starts with
 - Match to values that start with the search term (case-insensitive).
- · ends with
 - Match to values that end with the search term (case-sensitive).
- >
- Greater than the search term.
- <
- Less than the search term.
- NOT
 - Match to values that do not equal the search term (case-insensitive).
- · NOT contain
 - Match to values that do not contain the search term (case-insensitive).

17.2 Sequence tag flags

Sequences tagged in the sequence bin can have features indicated by specific flags. The presence of these flags can be queried. These are a superset of *flags available for allele sequences*. Available flags are:

- · ambiguous read
 - Genome sequence contains ambiguous nucleotides in coding sequence.

- · apparent misassembly
 - Sequence has a region of very high identity to existing allele in one region but looks completely different in another.
- · atypical
 - Catch-all term for a sequence that is unusual compared to other alleles of locus.
- · contains IS element
 - Coding sequence is interrupted by insertion sequence.
- · downstream fusion
 - No stop codon present resulting in translation continuing.
- · frameshift
 - Frameshift in sequence relative to other alleles, not resulting in internal stop codon.
- internal stop codon
 - Frameshift in sequence relative to other alleles, resulting in internal stop codon.
- · no start codon
 - No apparent start codon in immediate vicinity of usual start.
- no stop codon
 - No stop codon in immediate vicinity of usual stop.
- phase variable: off
 - Coding sequence has a homopolymeric run with a frameshift resulting in a stop codon preventing complete translation.
- truncated
 - Coding sequence is unusually short resulting in a truncated protein (not the same as running off the end of a contig).
- · upstream fusion
 - No apparent start codon in immediate vicinity of usual start, likely due to a gene fusion (sequence is transcribed together with upstream coding sequence).

17.3 Allele sequence flags

Sequences can be flagged with specific attributes - these are searchable when doing a sequence attribute query. These are a subset of *flags available for tagged sequences*. These are mainly for use with whole genome MLST type data. Multiple flags can be selected by Ctrl-clicking the list. Available flags are:

- · atypical
 - Catch-all term for a sequence that is unusual compared to other alleles of locus.
- contains IS element
 - Coding sequence is interrupted by insertion sequence.
- · downstream fusion
 - No stop codon present resulting in translation continuing.

- frameshift
 - Frameshift in sequence relative to other alleles, not resulting in internal stop codon.
- · internal stop codon
 - Frameshift in sequence relative to other alleles, resulting in internal stop codon.
- · no start codon
 - No apparent start codon in immediate vicinity of usual start.
- no stop codon
 - No stop codon in immediate vicinity of usual stop.
- · phase variable: off
 - Coding sequence has a homopolymeric run with a frameshift resulting in a stop codon preventing complete translation.
- truncated
 - Coding sequence is unusually short resulting in a truncated protein (not the same as running off the end of a contig).
- · upstream fusion
 - No apparent start codon in immediate vicinity of usual start, likely due to a gene fusion (sequence is transcribed together with upstream coding sequence).

CHAPTER 18

Database schema

- Sequence definition database
- Isolate database

A	download alleles in FASTA format, 349
access	download allelic profiles in CSV (tab-delimited) for-
control lists, 36	mat, 352
restricting, 36	download contigs in FASTA format, 357
adding	download submission supporting file, 364
classification groups, 80	GET /, 346
isolates, 132	GET /db, 346
locus, 45, 48, 49, 55, 84	GET /db/{database}, 347
MLST scheme, 78	GET /db/{database}/contigs/{contig_id}, 357
schemes, 60	GET /db/{database}/fields, 358
allele	GET /db/{database}/isolates, 353
breakdown, 268	GET /db/{database}/isolates/{isolate_id}, 353
allele definition	GET/db/{database}/isolates/{isolate_id}/allele_designations,
records, 201	354
allele designations	<pre>GET/db/{database}/isolates/{isolate_id}/allele_designations/{locus},</pre>
count, 230	355
query, 229	GET /db/{database}/isolates/{isolate_id}/allele_ids,
status, 232	355
allele sequence	GET /db/{database}/isolates/{isolate_id}/contigs,
identify, 205	357
allele sequences	GET/db/{database}/isolates/{isolate_id}/contigs_fasta,
adding, 110	357
alleles, 3	GET/db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele
list query, 213	356
API authentication	GET/db/{database}/isolates/{isolate_id}/schemes/{scheme_id}/allele
access token, 367	356
accessing protected resources, 369	GET /db/{database}/loci, 347
consumer key, 365	GET /db/{database}/loci/{locus}, 347
request token, 365	GET /db/{database}/loci/{locus}/alleles, 349
session token, 368	GET/db/{database}/loci/{locus}/alleles/{allele_id},
user authorization, 366	350
API resources	GET /db/{database}/loci/{locus}/alleles_fasta, 349
add submission correspondence, 363	GET /db/{database}/projects, 359
create new submission, 361	GET /db/{database}/projects/{project_id}, 359
DELETE/db/{database}/submissions/{submission_id},	GET /db/{database}/projects/{project_id}/isolates,
363	359
DELETE/db/{database}/submissions/{submission_id}/file	es GFT (db/{database}/schemes, 350
365	GE1/do/{database}/schemes/{scheme_id}, 330
delete submission record, 363	GET/db/{database}/schemes/{scheme_id}/fields/{field},
delete submission supporting file, 365	351

auto

GET /db/{database}/schemes/{scheme_id}/profiles,	stop, 188
351	automated assignment
GET/db/{database}/schemes/{scheme_id}/profiles/{p	=
352	autotagger, 183
GET/db/{database}/schemes/{scheme_id}/profiles_cs	sv, stop, 188
352	В
GET /db/{database}/submissions, 360	_
GET /db/{database}/submissions/{submission_id},	BLAST, 285
362	breakdown
GET/db/{database}/submissions/{submission_id}/file	es, allele, 268
364	provenance field, 261
GET/db/{database}/submissions/{submission_id}/file	
364	sequence bin, 274
GET/db/{database}/submissions/{submission_id}/me	essage _{wo-field} , 265
363	browse
GET /db/{database}/users/{user_id}, 358	scheme profiles, 215
list allelic profiles defined for scheme, 351	BURST, 289
list database resources, 347	
list loci, 347	C
list schemes, 350	caching
list site resources, 346	schemes, 39
POST /db/{database}/submissions, 361	classification groups, 4
POST/db/{database}/submissions/{submission_id}/fi	les, adding, 80
364	-1:4411:4:
POST/db/{database}/submissions/{submission_id}/m	nessage ESTful interface, 107
363	client databases, 70
retrieve allele identifiers, 355	clustering
retrieve contig record, 357	core genome, 82
retrieve full allele designation record, 355	codon usage, 293
retrieve full allele information, 350	Locus Explorer, 259
retrieve information about scheme field, 351	composite fields, 96
retrieve isolate record, 353	configuration
retrieve list of allele designations, 354	export, 105
retrieve list of alleles defined for a locus, 349	configuration settings
retrieve list of contigs, 357	validation, 104
retrieve list of isolate provenance field descriptions,	core genome
358	clustering, 82
retrieve list of isolate records, 353	=
retrieve list of isolates belonging to a project, 359	count
retrieve list of projects, 359	allele designations, 230 sequence tags, 232
retrieve list of scheme allele identifiers, 356	sequence tags, 232
retrieve list of submissions, 360	D
retrieve list of supporting files uploaded for submis-	
sion, 364	defining
retrieve locus record, 347	exemplar alleles, 185
retrieve project information, 359	E
retrieve scheme allele designation records, 356	E
retrieve scheme information, 350	exemplar alleles
retrieve specific allelic profile record, 352	defining, 185
retrieve submission correspondence, 363	export
retrieve submission record, 362	configuration, 105
retrieve user information, 358	extended attributes
upload submission supporting file, 364	locus, 57
allele definer, 186	provenance fields, 99

382 Index

F	modifying display
filters, 236	loci, 252 schemes, 252
G	
Genome Comparator, 276 genome filtering, 87 in silico hybridization, 90 in silico PCR, 87 groups scheme, 66 H hosts	offline curation auto allele definer, 186 autotagger, 183 options, 244 isolate record, 247 main results table, 246 provenance fields, 249 query, 250
mapping, 39	Р
 identify	partitioning sets, 41
allele sequence, 205 sequence type, 207 in silico hybridization, 90	passwords setting, 36 first user, 37
in silico PCR, 87 isolate records, 197	performance caching schemes, 39 mod_perl, 39
isolate record options, 247 isolates	permissions, 31 locus curation, 34 scheme curation, 34
adding, 132	plugins enabling, 37
L	polymorphic sites Locus Explorer, 255
list query alleles, 213	polymorphisms, 298 profile
loci, 3 modifying display, 252	records, 203 profiles, 4
adding, 45, 48, 49, 55, 84 copying existing record, 53 extended attributes, 57	projects, 154 provenance field breakdown, 261 provenance fields
Locus Explorer, 255 codon usage, 259 polymorphic sites, 255 translated sequences, 260	extended attributes, 99 options, 249 publications, 144
locus positions	Q
setting, 91	query allele designations, 229
main results table options, 246 mapping	options, 250 scheme profiles, 216 sequence tags, 234
hosts, 39 MLST, 3	R
MLST scheme adding, 78 mod_perl, 39	records allele definition, 201 isolate, 197
1110u DC11, 37	

Index 383

```
profile, 203
    sequence bin, 203
     sequence tag, 202
RESTful interface
     client authorization, 107
rule-based queries, 76
S
scheme
    breakdown, 268
    groups, 66
scheme profiles
     automated assignment, 79
     browse, 215
     query, 216
schemes, 3
    adding, 60
     caching, 39
     modifying display, 252
sequence bin
    breakdown, 274
     records, 203
sequence similarity
    determining, 220
sequence tag
    records, 202
    status, 303
sequence tags, 4
    count, 232
     query, 234
sequence type
    identify, 207
sets, 4
     partitioning, 41
stop
     auto allele definer, 188
     autotagger, 188
Τ
translated sequences
    Locus Explorer, 260
two-field
    breakdown, 265
unique combinations, 296
updates
     disabling, 38
user groups, 31
user types, 31
users
     adding, 109
```

384 Index