# Sequence Database Administration
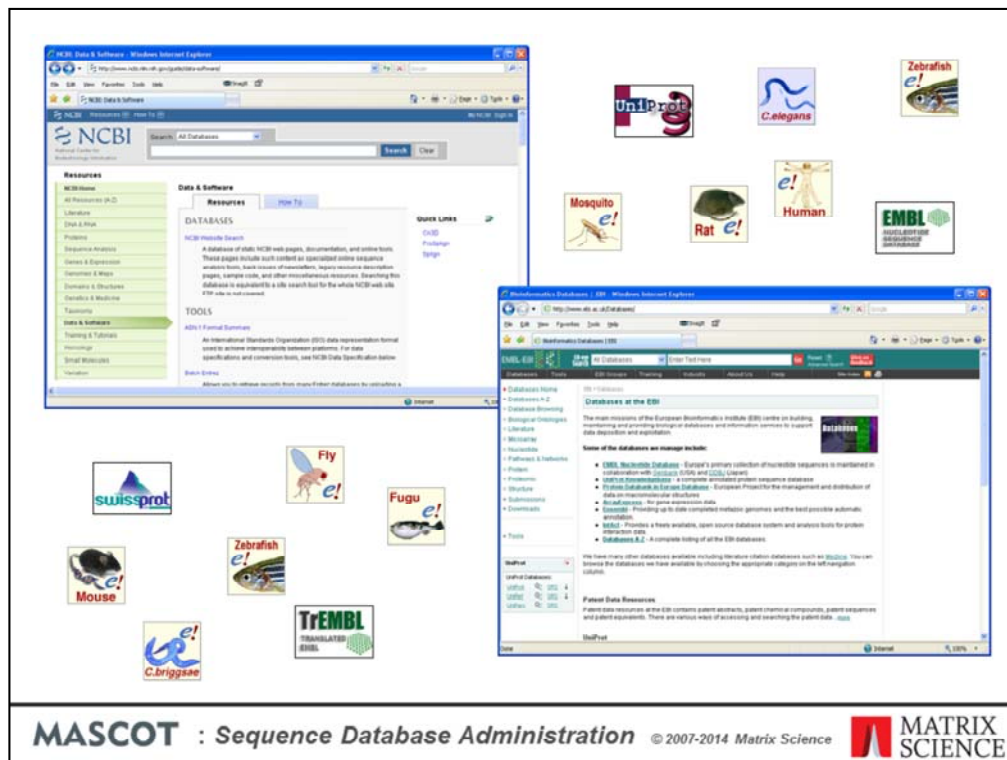
MASCOT

MATRIX SCIENCE

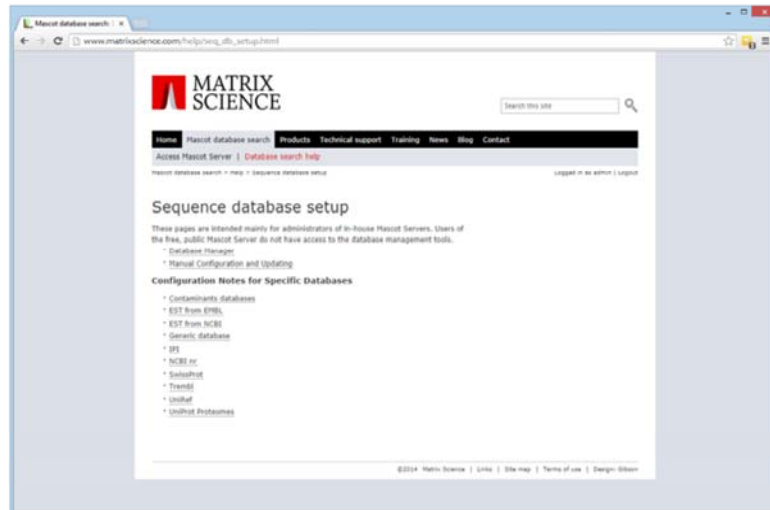MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

When you install Mascot, it includes a copy of the Swiss-Prot protein database. However, it is almost certain that you and your colleagues will want to search other databases as well. There are very many to choose from, and Mascot allows you to have as many databases as you wish on-line for searching at any one time.

Matrix Science doesn't supply sequence databases. Most databases are public domain, and there are a few sites that provide comprehensive database repositories. Two of the best known are NCBI and EBI. Here, you can download nr, Swiss-Prot, EMBL, Trembl, UniRef100, etc.

For specialised databases, such as individual genomes, you may have to track down the FTP site of the group that is doing the sequencing.

This topic described the general procedure for adding a new database to Mascot and keeping it up-to-date

For the latest information about the major public databases, refer to the help pages on the Matrix Science public web site. The help pages in your in-house copy of Mascot are similar, but become progressively out-of-date.

## Sequence Database Requirements

**Mascot can search any database available in Fasta format**
- Amino acid
- Nucleic acid
   Genomic DNA, EST's, ORF's, mRNA, etc

**Must have local Fasta file**
- (Mascot streams through the database during each search)

**Other files are optional**
- Taxonomy indexes
- Full text annotations.

MASCOT : *Sequence Database Administration* © 2007-2014 Matrix Science  MATRIX SCIENCE

To perform Mascot searches against a database, at a minimum, we need a FASTA file.

If the database contains nucleic acid sequences, there is no need to pre-translate the sequences. Mascot performs a 6 frame translation during each search. Nucleic acid databases come in several flavours. They may be described as genomic DNA, Expressed Sequence Tags, Open Reading Frames, messenger RNA, etc. As far as Mascot is concerned, the main differences are the quality and length of the individual entries. The relative merits of searching protein, EST and DNA sequences are discussed in Choudhary *et. al. Matching peptide mass spectra to EST and genomic DNA databases*. Trends in Biotechnology, 19, S17-S22 (2001)

If the database contains entries from multiple organisms, and you want to be able to filter searches by taxonomy, this will require some additional files, which vary from database to database

Some databases, such as SwissProt, also come with 'full text' files, containing comprehensive annotations.

**FASTA Format**

```
>Title text
SEQUENCESEQUENCESEQUENCESEQUENCESEQUENCESEQUENCE
SEQUENCESEQUENCESEQUENCESEQUENCESEQUENCESEQUENCE
SEQUENCESEQUENCESEQUENCE
>Next title
NEXTSEQUENCE …

>gi|6|bgi|Contig1.seq_7|2412 3299 [+3 L= 888] [Delayed
>20021010.2.1   1112073F09.y1 1112091F10.y1 1112073F0
>IPI:IPI00140097.1|REFSEQ_XP:XP_168061 Tax_Id=9606
>CCRB cytochrome c [validated] – rabbit
>gi|129249|sp|P02820|OSTC_BOVIN Osteocalcin precursor
>"ORF5 | start 2178-1309 | frame -1 | length=870 |
```

MASCOT : *Sequence Database Administration* © 2007-2014 Matrix Science    MATRIX SCIENCE

Perhaps this is a good moment to clarify exactly what we mean by a FASTA file.

FASTA is a very popular standard because it is so simple. On the down-side, it isn't much of a standard ... almost anything goes.

FASTA specifies that there will be a title line, starting with a 'greater than' character, followed by one or more lines containing the sequence in 1 letter code.

The problem is the lack of a well defined syntax for the title line. Here are a handful of examples of FASTA title lines. As you can see, there isn't much similarity. For a Mascot search, we need to find a short, unique identifier or accession string for each sequence. As you can see from these examples, the position of the identifier and the delimiters (e.g. spaces, pipe symbols, commas) varies considerably

## Parse Rules

**Parse rules are Basic Regular Expressions**

```
>IPI:IPI00043251.2|REFSEQ_XP:XP_064505
 Tax_Id=9606 similar to keratin 18,
 cytoskeletal - human (fragment)
```

Accession from Fasta title:          `">IPI:\([^| .]*\)"`

Description from Fasta title:          `">[^ ]* \(.*\)"`

MASCOT : *Sequence Database Administration* © *2007-2014 Matrix Science* MATRIX SCIENCE

The way Mascot handles this is to use regular expressions to describe how to parse information from the title lines in any particular database. Regular expressions will be familiar to anyone with a Unix background, but there may be a bit of a learning curve for someone with more of a Windows or Mac background.

Here, for example, we have a title line from the IPI human database. Let's say that we want to use IPI00043251 as the unique accession string and everything after the first space should be treated as the description.

The regular expressions, or parse rules, used to extract this information look like this.

The string we want to extract is always within back-slashed parentheses. For the accession, we show the first few characters as literal text. We then say that we want to take all the following characters, stopping when we hit either a pipe symbol, a space, or a period. In fact, it is the period which applies in this example. The contents of the square brackets are known as a character class, and the circumflex at the beginning means 'not'. The asterisk means 'as many as available'.

For the description, we discard everything up to and including the first space. This is done using a character class of 'not a space' followed by one literal space. Then, we use back-slashed parentheses, take everything to the end of the title. The period matches to any character, so .* matches to all the remaining text.

In Mascot 2.4, we introduced Database Manager, which handles both database configuration and the downloading of files from external servers. This replaced two utilities in Mascot 2.3 and earlier: the browser-based Database Maintenance, used for configuration, and the command line Database Update, used for downloading. If you are still using Mascot 2.3 or earlier, you will probably find the archived, 2.3 version of this presentation provides more practical information.

The file formats and download locations of sequence databases change from time to time. One of the smart features of Database Manager is that database configurations for many public databases are updated automatically, by downloading configuration data from the Matrix Science web site.

# Key Concepts

**Predefined Database Definition**
- Configuration information for the most popular public databases is kept up-to-date on the Matrix Science web site, and downloaded as required by Database Manager

**Custom Database Definition**
- If you want to search a database that is not included in the list of Predefined Database Definitions, or if you want to configure one of these databases in some non-standard way, you create a Custom Database Definition

**Synchronisation**
- If a custom definition is very similar to a predefined definition, it can be converted into a predefined definition by being synchronised

**Update Schedule**
- An schedule can be created to update all the files associated with a database automatically.

MASCOT : *Sequence Database Administration* © 2007-2014 Matrix Science — MATRIX SCIENCE

Let's review a few of the important terms used in Database Manager:

A Predefined Database Definition is one in which the configuration information is kept up-to-date on the Matrix Science web site, and downloaded as required by Database Manager. You don't need to know file URLs or worry about parse rules, etc. for a Predefined Database.

If you want to search a database that is not included in the list of Predefined Database Definitions, or if you want to configure one of these databases in some non-standard way, you create a Custom Database Definition.

If a Custom Database Definition is very similar to a Predefined Database Definition, it can be converted into a predefined definition by being synchronised. The advantage of doing this is that the configuration will then be kept up-to-date automatically.

An Update Schedule can be created to update all the files associated with a database automatically. Maybe once each week or each month. Files will only be downloaded if a new version is available.

Database Manager must be allowed underline exclusive control of database configuration. Editing mascot.dat outside of Database Manager will just cause confusion because Database Manager re-writes mascot.dat whenever a configuration changes. If you prefer to configure sequence databases manually, by editing mascot.dat, never run Database Manager.

The first time Database Manager is run, it tries to match existing database definitions against predefined definitions and reports the quality of the match as none, poor, good, or perfect. For poor or good matches, the differences can be inspected. Usually, these arise because the existing definition is out-of-date in some respect. You can choose whether to synchronise an existing definition, making it predefined, or keep it as a custom definition.

If the Mascot Server is not allowed to access the Internet, choose *Keep as Custom* because synchronisation of any definition requires the database files to be updated.

Having made your selections, choose *Import* to proceed. The list of Databases will be displayed, with status information for those that have been synchronised and are being updated.

## Adding a New Database

**Enable predefined definition**
- Apart from confirming a location for the downloaded files, everything will be handled automatically.

**Create New; Custom**
- Create a new custom database definition from scratch.

**Create New; Copy Of**
- Create a new custom database definition by copying an existing definition.

**Create New; Use predefined definition template**
- Create a new custom database definition by starting from a predefined definition.

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science    MATRIX SCIENCE

You can add new databases in four different ways:

1. Enable predefined definition: Apart from confirming a location for the downloaded files, everything will be handled automatically. Only one instance of each predefined definition can be enabled at any one time, as database names must be unique. If you want something similar to a predefined database, but with configuration changes, choose the final option: Use predefined definition template.

2. Create New; Custom: Create a new custom database definition from scratch.

3. Create New; Copy Of: Create a new custom database definition by copying an existing definition. You will be required to enter a new database name and given the choice of copying the existing database files.

4. Create New; Use predefined definition template: Create a new custom database definition by starting from a predefined definition. The differences between this and enabling a predefined definition are (i) you can make changes to the configuration, (ii) the definition will not be kept up-to-date automatically.

Let's look at the first of these options: enabling a predefined definition, using NCBInr as the example.

Scroll down to NCBInr and choose Enable

Enable Predefined Definition

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

The default location for the local copies of the sequence database files is specified in Database Manager settings. You can also change the location for the new database here. New directories will be created automatically, unless they already exist. For each database, there is a directory with the same name as the database. Under this directory are three sub-directories. The incoming directory provides a workspace for downloading and processing a new database file. The current directory contains the active database, and this is where Mascot Monitor creates the compressed files that will be memory mapped. The old directory is where the immediate past database files are archived … just in case.

Choose Create

All the files will now be downloaded automatically. For NCBInr, this is the Fasta file plus the files that are needed to create a taxonomy index. The lower part of the page is updated with status information. You don't have to leave this page open; you can close the browser and return later.

Some database files are very large, and downloads can fail for all sorts of reasons. Database Manager tries each download 5 times before giving up. If you have persistent problems, check the support page on our web site to see if there are any known issues.

Assuming the download is successful, this page will be displayed as the new files are compressed and the database is brought online. As soon as the new database shows as 'In Use', it is ready for searching.

To setup automatic updates of the database files, choose Edit Schedule from here or via the Scheduled Updates side menu

# Scheduled Updates

**MASCOT** : *Sequence Database Administration* © 2007-2014 Matrix Science

It is usually best to download the files at a quiet time, like the middle of the night or at the weekend.

Note that keeping the definition up-to-date and keeping the database files up-to-date are two different things. A predefined database definition is kept up to date automatically while a custom database definition is not. The only requirement for keeping the files up to date is that the definition includes URLs for downloading the required files. Files are not updated by default; you have to save a schedule for the database specifying how often to look for new files. If no new Fasta file is available at the scheduled time, nothing will be downloaded.

Creating a new database by starting from a copy of an existing database is usually more convenient than starting from scratch. It is also a good way to preserve a copy of a particular version of a database. Imagine you have SwissProt configured to be automatically updated every month or so. If you want to keep a copy of the current version, so that it can be used for all searches during a year long project, choose SwissProt from the drop down list and give it a suitable name so that it won't be confused with the 'live' version

You are given the option to copy the existing files.

Create New – Copy of

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

Unless you wish to make some change to the configuration, all you need to do now is to choose Activate. When submitting searches for the year-long project, you choose SP-frozen rather than SwissProt

To illustrate how a predefined definition can be used as a template, we'll set up a database for the Uniprot proteome of rice. In a browser, go to the Uniprot web site, www.uniprot.org, and follow the 'Complete Proteome' links. You could search on rice or, if you remember the latin name, oryza

There is a choice of three (one hit is a virus that infects rice). We'll choose Oryza sativa subsp. Japonica.

You could download the file manually. If so, click on the proteome ID link, then a download button, and select all proteins in Fasta format.

A better option is to set up a download URL, because this will allow us to configure automatic updating of the database files. First, make a note of the taxonomy ID for this strain - 39947

In Database Manager, choose Create new. Enter a suitable name and select the uniprot proteome predefined definition as a template.

Choose Create

If we had chosen to download the Fasta manually, we could follow the instructions to upload or copy the Fasta file to the target directory and rename it to match the Filename pattern. Since we want to schedule automatic file updates, we choose instead to download from a remote URL and, in the next page, choose Setup download URL

The format for the URL can be found on the Mascot help page for UniProt proteomes under Sequence database setup. You just need to change the taxonomy ID to the one for rice that you noted earlier. There is no reference file to download … we'll link out to Uniprot to get annotation text for the Mascot Protein View report. There is no version file either, so each update will be identified using an ISO datestamp

Save, and we're ready to start downloading

## Create New – Predefined as Template

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

Once the download is complete, the page is updated. Assuming we don't want to inspect or modify the configuration, just two things left to do. Make the database active in Mascot and create an update schedule.

## Create New – Custom

### Custom is rarely required
#### Fasta from Uniprot
UniProt_proteome_template
#### Fasta from Genbank
NCBI_AA_template
NCBI_NA_template
#### Most other cases
simple_AA_template
simple_NA_template

MASCOT : *Sequence Database Administration* © 2007-2014 Matrix Science    MATRIX SCIENCE

In most cases, it is faster to start from one of the predefined definitions, and modify it, than choose custom

But, to illustrate, lets configure a protein Fasta for the recently completed tomato genome. Choose a file and download it

The first thing to do with any unknown Fasta file is open it in a text editor that can handle large files and take a look at the syntax of the title lines. If you don't have a suitable text editor, you can use more at a command prompt. Make sure you look at different places in the file because it may be a merge of entries from different sources

**Create New – Custom**

**Look at through the the Fasta file to get a sense of what aspects of the titles are constant**

>Solyc00g005000.2.1 Aspartic proteinase nepenthesin I  IPR001461 Peptidase A1

>Solyc04g047700.1.1 Unknown Protein

>Solyc09g008500.1.1 Non-specific lipid-transfer protein  IPR013770 Plant lipid transfer protein and hydrophobic protein, helical

>Solyc12g100360.1.1 Calpain-like protein  IPR001300  Peptidase C2, calpain

**Best to use the very simple rules (simple_AA_template)**

">\([^ ]*\)"
">[^ ]* \(.*\)"

MASCOT : *Sequence Database Administration* © 2007-2014 Matrix Science | MATRIX SCIENCE

As is often the case, a simple rule that takes everything between the ">" symbol and the first space as the accession is the safest choice. Everything after the first space can be treated as the description. These rules are pre-defined in Database Manager and you could set up this database by using simple_AA_template as the template. But, to illustrate the flexibility of Database Manager, we'll follow the custom definition route, and even create a new parse rule.

So, we Create new, enter a suitable name, and choose Custom

This is an Amino Acid database and we already have the Fasta file. Set these options and choose Create

We've chosen to rename the Fasta file and upload it via the web browser

Once the file is uploaded and the name matches the wild card pattern, it will be recognized, and we can proceed to edit the configuration. If Database Manager doesn't recognize the presence of the Fasta, maybe there is a typo in the database name or maybe you have the file permissions / security settings set so that a CGI process cannot read the file.

First thing we have to do is choose parse rules for accession and description.

All the existing parse rules are tested against ten of the title lines, five from the start of the file and five from the end. Four parse rules give matches to all ten. We need to study the matches and choose the one that pulls out a suitable accession string. I think its fairly obvious that the one we've selected is the most suitable

When it comes to the description parse rule, the first one listed would be fine but, as a challenge, we'll devise a new rule to extract only the Interpro reference and the text that follows. We click on 'Create new parse rule'

There is a brief reminder of how basic regular expressions work, and you can use trial and error; pressing test until your rule succeeds. This screen shot shows a rule that appears to work. It locks on to the text IPR followed by 6 digits, then takes everything from that point to the end of the line. This is actually not a good rule, because not all titles include an Interpro reference, but let's not worry about that right now. We confirm our choices and move on to the remainder of the configuration

42

Create New – Custom

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

The next step deals with taxonomy and annotations.

If this was a comprehensive database, containing entries from many different organisms, we might want to select a rule for determining the taxonomy of each entry. There is a drop down list of rules for the major public databases. Most of these require additional files, which can be downloaded automatically by Database Manager whenever the Fasta file is updated. It is extremely unlikely that you will need to create a new taxonomy definition. But, for completeness, the syntax is described in Chapter 9 of the Installation & Setup manual.

Annotation text usually comes from some type of web service. Mascot submits a request using an accession string and the text is returned and embedded in the Protein View report. This form allows you to select from existing URLs or create a new one. In rare cases, the full text for the entire database is downloaded as a local file. The only common examples of this are the SwissProt and Trembl DAT files

For the tomato genome, as is often the case with simple, single organism databases, we don't need to worry about taxonomy and there isn't a suitable source for full-text annotation reports.

At the bottom of this page, we can choose 'Save and Finish' then, on the next page, 'Activate'.

Create New – Custom

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science

All being well, a short time later, our new database will show as 'In Use'. You'll notice that there is no option to update 'ketchup' because we copied the file manually.

If there are problems, and the database fails to reach 'In use', you'll need to follow the Status link to Database Status

Database status provides an overview of all the active databases. It also provides links to other pages of useful information.

Initially, there will be a single information block for each database on this page. When a database is updated, a second information block is added. One is for the new or incoming database, the other is for the old or outgoing. If all is well, one of the pair will have the status of "In use", and the other "Not in use".

If there is a problem, the status will be an error message and a 'compression warning' link added to the relevant error messages.

The database statistics are very useful for diagnosing problems and checking up on the health of a database

For example, does the number of entries look about right? Sometimes, a download may be truncated and the problem go undetected

Are there any invalid characters in the sequences? If there are, this should definitely be investigated

Mascot has a parameter, MaxSequenceLen, to set the length of the longest sequence. The default is 50,000. The higher this value, the more memory Mascot uses, so it should not be set to a ridiculously high value. If any sequences are "too long", then you need to increase MaxSequenceLen to something a little greater than the length of the longest sequence. If you are trying to search an assembled genome, you might want to consider searching shorter sequences instead, such as a database of contigs.

If your protein database seems to be composed entirely of A, C, G, and T, then it may be worth double checking that you downloaded the correct file..

Although it is rarely possible to achieve 100% accuracy for taxonomy, you certainly want the accuracy to be better than 99%. Otherwise, the results could be misleading. Near the bottom of the stats file is a list of the number of entries with 0, 1, 2, etc., taxonomy identifiers. From time to time, check that the number of entries with no taxonomy identifiers is well below 1% of the database. Here, it is a very healthy 0.03%

If you look at the configuration details for any database, there is a section for Performance Settings

# Configuration - Performance

**Database configuration: NCBInr**

**Performance settings**

Number of threads [?]
-1

☑ Use memory mapping [?]
☐ Lock to memory [?]

[Cancel] [Save]

MASCOT : *Sequence Database Administration* © 2007-2014 *Matrix Science*    MATRIX SCIENCE

A Mascot search can use multiple threads, so as to make use of all the logical processors covered by the licence. Usually, it is best to leave threads set to -1, which means automatic. If you want to restrict the number of threads on a non-cluster (SMP) system, you can do so by setting a value of 1 or more. Each CPU in the Mascot licence allows use of up to 4 cores, which requires 8 threads for a hyperthreaded processor or 4 otherwise. On a cluster system, the number of threads is set for each search node in a separate configuration file, nodelist.txt.

Database files should always be memory mapped because this gives the fastest access times. Memory mapped files can be locked in memory, but only if the computer has sufficient RAM. Having a database locked in memory means that it can never be swapped out to disk, ensuring there is never a lag as the database is read from disk. If you try to lock databases into RAM when there isn't room, this will not be a major problem. The locking will fail, generate an error message, and Mascot will carry on regardless. A more serious problem is when there is just sufficient RAM to lock the databases, but none left over for searches or other applications. In this case, the whole system will slow down and the hard disk will be observed to be "thrashing". Eventually, the system is likely to hang or crash. In general, it is better to let the operating system manage which files are held in memory and not lock any databases into memory.

### Database Tips

**Check the statistics file from time to time**
**Always memory map databases**
**Be selective when locking databases into memory**
- Only the small databases, which are searched frequently, should be locked in memory

**Can place sequence databases on any local drive**
**Don't download files onto a Windows desktop**
- They will get very restricted security settings

**Don't create a sequence database with inconsistent title syntax**
- Must be able to extract a unique identifier (accession) from all entries with a single parse rule

**Use predefined databases where available**
- Configuration kept up-to-date automatically.

MASCOT : Sequence Database Administration © 2007-2014 Matrix Science — MATRIX SCIENCE

This slide recaps some important tips.

Check the statistics file from time to time, particularly after configuring a new database

Always memory map database files to make access as fast as possible, but be selective about locking databases into memory. Only the smaller databases, which are searched regularly, should be locked in memory

You can place sequence database files on any local drive. Under Unix, you can use NFS mounted drives as long as the connection is fast and stable

If you download files manually, don't download to your Windows desktop. Chances are that Database Manager won't be able to see the file because it becomes private to your Windows login. If this happens, add the local Users group to the security settings for each file and give the Users group full control

If you create your own Fasta file, use a consistent title syntax. It must be possible to extract a unique identifier (accession) from all entries with a single parse rule

Use predefined databases where available because this means the configuration is kept up-to-date automatically.