

# Database Manager in Mascot 2.4

MASCOT

*{MATRIX}*  
*{SCIENCE}*

## Database Manager

### What does it do?

- Makes it easy to set up new fasta 'database' files for searching in Mascot
- Automate scheduled updates of fasta files

Database manager does something that sounds trivial. It makes it easy to setup a new database for searching with Mascot. It also makes it easy to automate regular updates of these databases.

## What's so difficult?

### More complex than it appears!

- Where to find the files - which may change
- Downloading very large files
- Unpacking large files
- No standard fasta format
- Formats keep changing
- Need additional taxonomy and unigene files
- Need to understand 'cron' or 'scheduled tasks'.

However, it's not quite as simple as that in practice for a number of reasons.

In most cases the databases are frequently updated and you want to get the latest version, but download locations and file names can change over time.

There are also issues with downloading very large files, so we need to check that the whole file has been downloaded. Running out of disk space or a lost connection are common problems. Also, it's quite possible for the file to be updated on the server while it is being downloaded which can result in a corrupt file.

Some tools, for example the built in unzip utility in Windows 7 will fail to unpack files larger than 4Gb

For each sequence, we need a unique identifier or accession but there is no standard fasta format.

Once we've got it right, it can change...

We also need a collection of taxonomy files so that we can filter searches based on taxonomy. Locations, names and formats of these files can change

It's also not trivial to automate regular updating of the fasta files. We suspect that many of our users just have very old versions of the databases.

## Predefined Database Definition

### Required information

- Location of files for download
- Rules to extract accession and description
- Rules to get reference information
- Getting the database version number
- AA or NA?
- Taxonomy and unigene files

### Predefined definitions stored on Matrix Science public server

**MASCOT** : Database Manager

© 2012 Matrix Science



We use the term Predefined Database Definition quite often in the software and I'd just like to explain what we mean by this.

Database manager needs this information to be able to configure a database:

- The location and the names of the files for download. If the files are compressed, it also needs to know how they are compressed.
- The fasta format isn't really a standard, and we need to be able to extract a unique identifier, or accession for each sequence. The predefined database definition includes rules to do this.
- It's also useful to be able to show reference information for the sequences. For SwissProt and TrEMBL, this can come from a local reference file or from an internet resource. Again, we need some rules for this.
- Most sequence database providers always give the same filename for the fasta file to make it easy to download with a fixed url. However, some provide a separate file with the version number in. Again, we need rules for this
- We need to know whether the fasta file is a protein database or a nucleic acid database. If you get that wrong, your searches clearly won't work.
- If the sequence database comes from multiple species, it's useful to be able to filter searches based on the taxonomy of your sample. For this to work, we need to know what taxonomy files to download.

- And finally, for the EST databases, it's really useful to be able to group the fragments using unigene indexes. Patrick will describe this in a later talk.

However, it's no good us just putting this information for each of the common databases on the CD we send out because the database providers often change some of these things and that prevents Mascot from working with the latest download. To cope with this, we now have a file on our public web server that contains all the rules. In future, when the NCBI or EBI change the format of a database, we can just change the rules and everyone's Mascot server will get updated automatically. Or, if someone produces a new database that is commonly used, we'll just add that to our predefined database definitions file.

## Database Manager - Create new

### 4 ways to make a database available for Mascot

- i. Using a predefined definition
- ii. Using a template, for example for a Uniprot species specific database
- iii. Creating a copy of an existing definition
- iv. Using a custom definition

When you start database manager for the first time it will import your existing configuration. After that you'll probably want to make new databases available and automate scheduled downloads for them.

There are four ways to make a new database available for Mascot.

I'll start by describing how to get NCBIInr online using a predefined definition.

Next, I'll describe how to use a template for a uniprot species specific database.

I'll just touch briefly on the last two options because they are fairly obvious.

**Database Manager**

- Databases (7)
- Parse rules (10)
- Tasks (0)
- Settings

**New database**

- Enable predefined definition
- Create new database
- Synchronize custom definitions

## Enable predefined database definition

Predefined database definitions are configuration entries for the most commonly used, publicly available databases. Configuration and FASTA files for predefined definitions will be automatically kept up to date as long as the Mascot Server machine is connected to the Internet.

Name						
contaminants	<input type="button" value="Enable"/>	IPI_human	<input type="button" value="Enable"/>	simple_AA_template	<input type="button" value="Enable"/>	
cRAP	<input type="button" value="Enable"/>	IPI_mouse	<input type="button" value="Enable"/>	simple_NA_template	<input type="button" value="Enable"/>	
Environmental_EST	<input type="button" value="Enable"/>	IPI_rat	<input type="button" value="Enable"/>	SwissProt_AC	<input type="button" value="Enable"/>	
EST_human	<input type="button" value="Enable"/>	IPI_zebrafish	<input type="button" value="Enable"/>	SwissProt_ID	<input type="button" value="Enable"/>	Already set up
EST_mouse	<input type="button" value="Enable"/>	Mammals_EST	<input type="button" value="Enable"/>	Trembl_AC	<input type="button" value="Enable"/>	
EST_others	<input type="button" value="Enable"/>	Mus_EST	<input type="button" value="Enable"/>	Trembl_ID	<input type="button" value="Enable"/>	
Fungi_EST	<input type="button" value="Enable"/>	NCBI_AA_template	<input type="button" value="Enable"/>	Unclassified_EST	<input type="button" value="Enable"/>	
Human_EST	<input type="button" value="Enable"/>	NCBI_NA_template	<input type="button" value="Enable"/>	UniProt_proteome_template	<input type="button" value="Enable"/>	
Invertebrates_EST	<input type="button" value="Enable"/>	<b>NCBIInr</b>	<input type="button" value="Enable"/>	UniRef100	<input type="button" value="Enable"/>	
IPI_arabidopsis	<input type="button" value="Enable"/>	Plants_EST	<input type="button" value="Enable"/>	Vertebrates_EST	<input type="button" value="Enable"/>	
IPI_bovine	<input type="button" value="Enable"/>	Prokaryotes_EST	<input type="button" value="Enable"/>			
IPI_chicken	<input type="button" value="Enable"/>	Rodents_EST	<input type="button" value="Enable"/>			

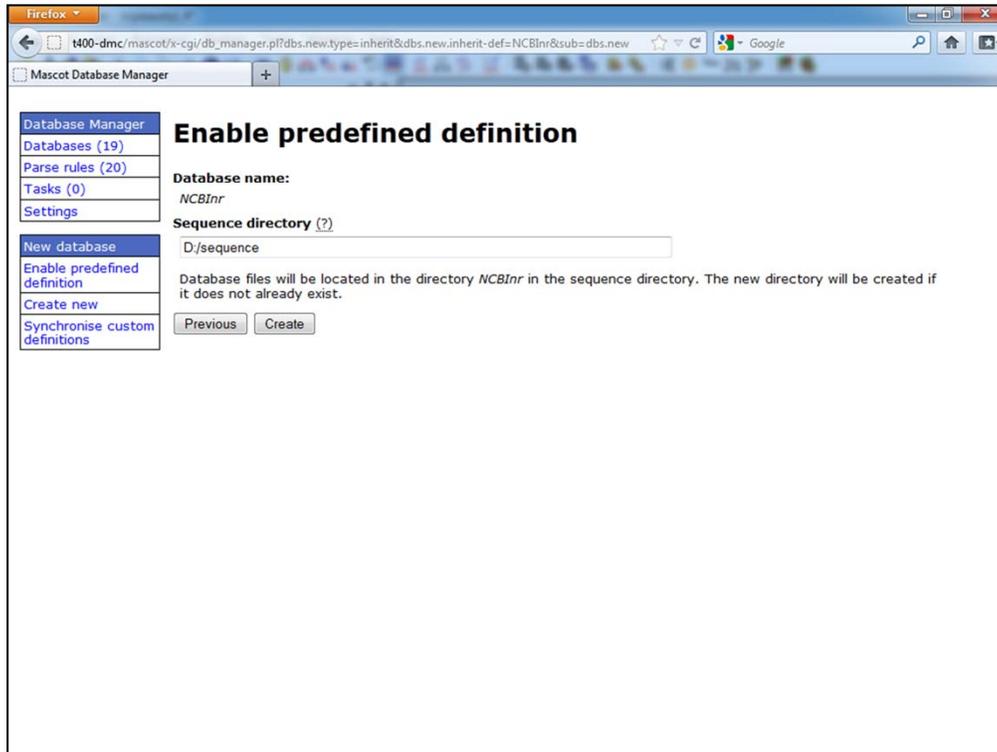
**MASCOT** : Database Manager

© 2012 Matrix Science

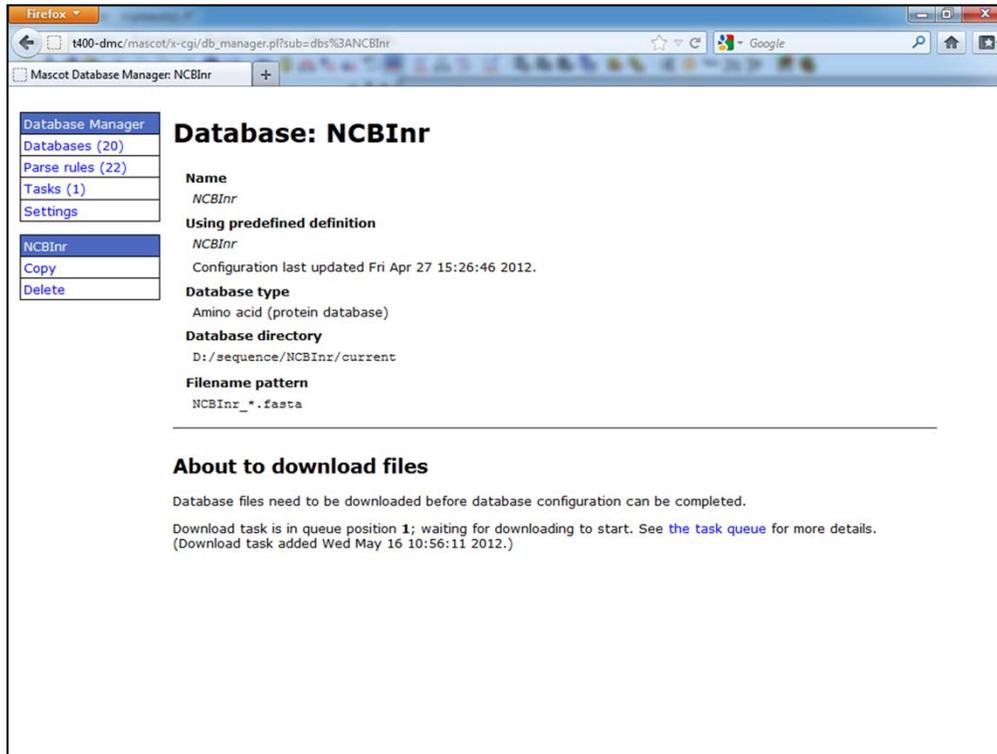
So, let's start with making NCBIInr available. We do this by enabling a predefined definition and this is ridiculously easy!

Simply click on the link on the left, and choose one of the available definitions. There are over 20 definitions available.

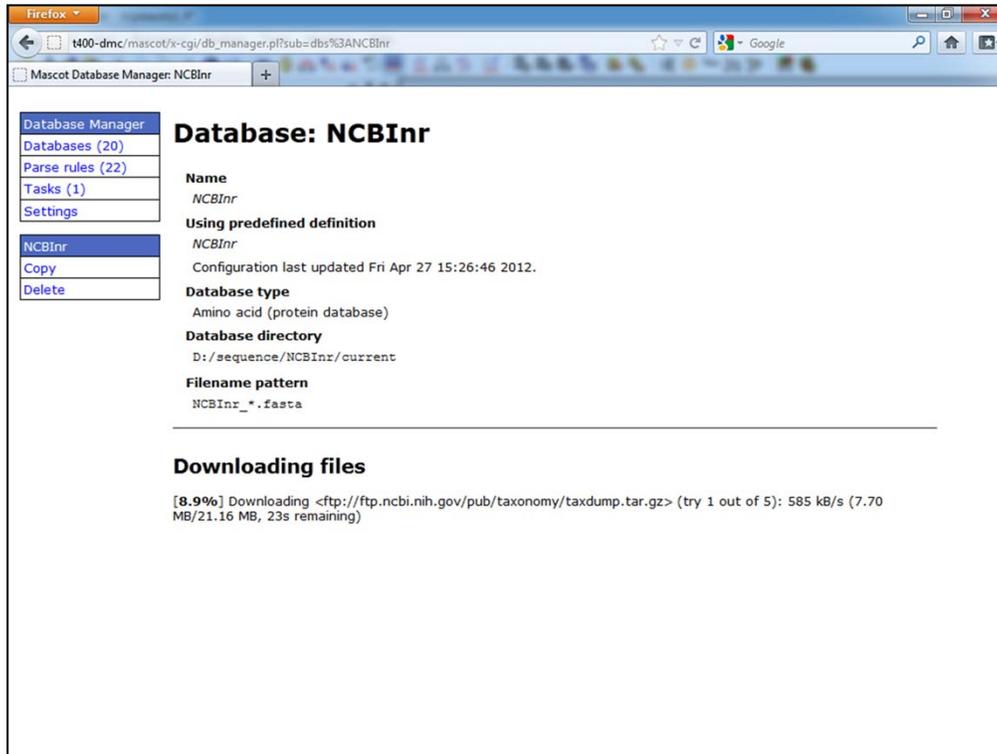
Just click on the enable button next to NCBIInr.



Database manager then asks you to confirm the sequence directory under which the NCBIInr directory will be created. You can change the default directory in the settings page, or just change it here. Mascot is installed on the C: drive on my computer, but I've decided to put all my databases on D: because there's plenty of spare space there. Click on create, and you could just go off and get a cup of coffee, or you could sit and watch the screen update. I don't recommend watching an update for EST\_others, which is a 10Gb download and over 20Gb uncompressed. I'll just show you what happens next.

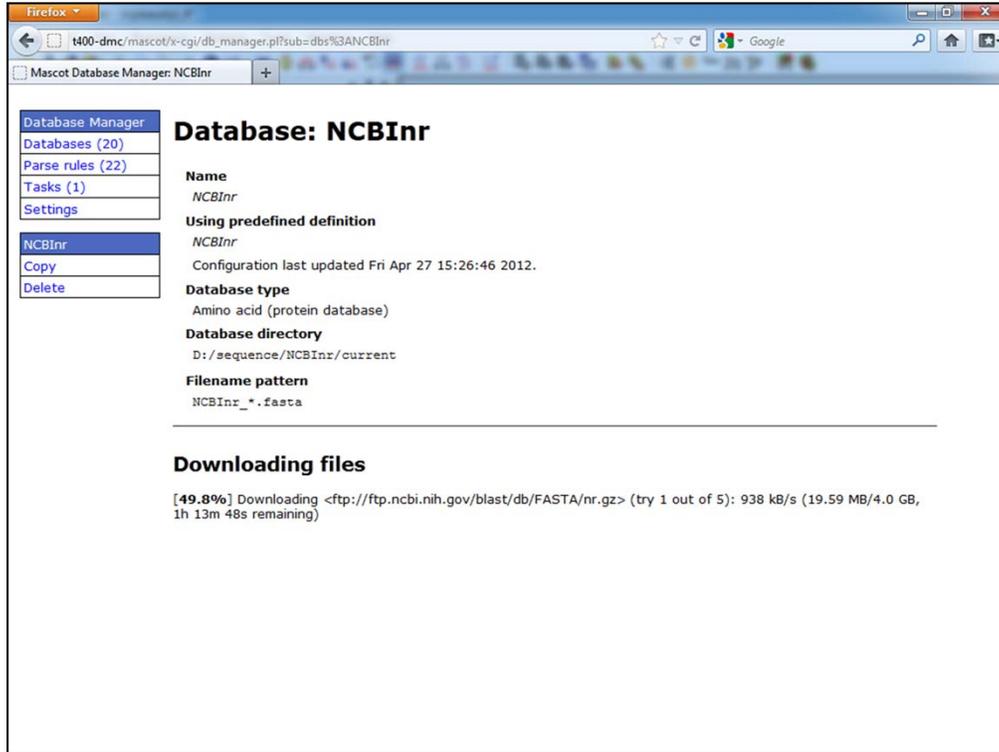


Database manager lets you know that it's going to start downloading files. You could click on the link to see the task queue if nothing happens in a minute or so,

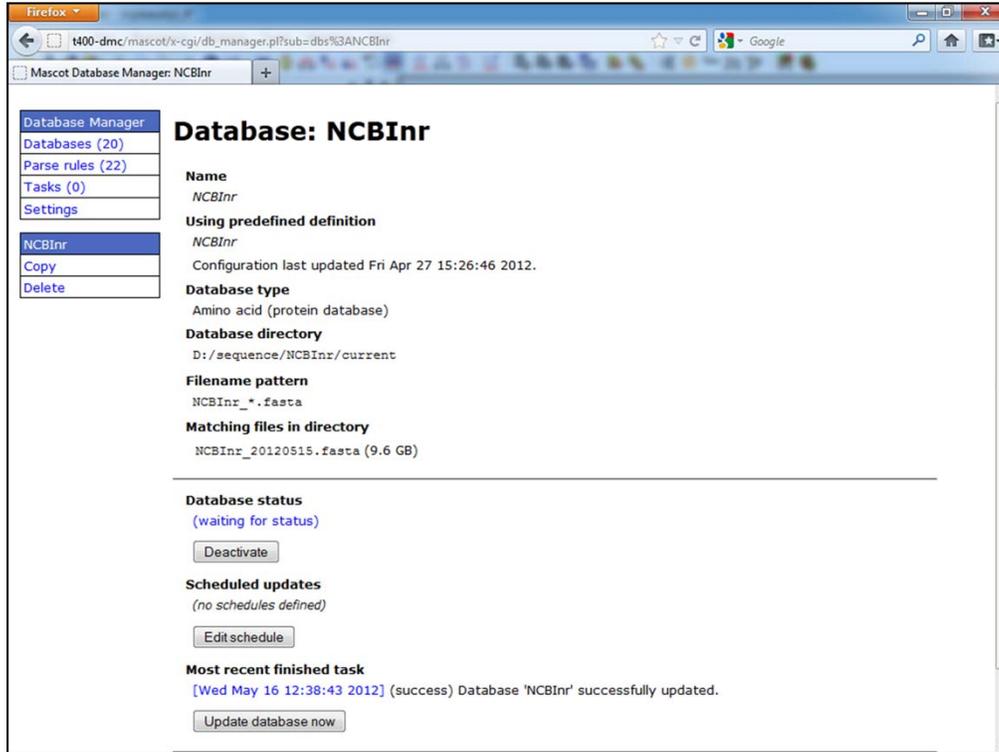


But most likely you'll see it starts to download pretty quickly.

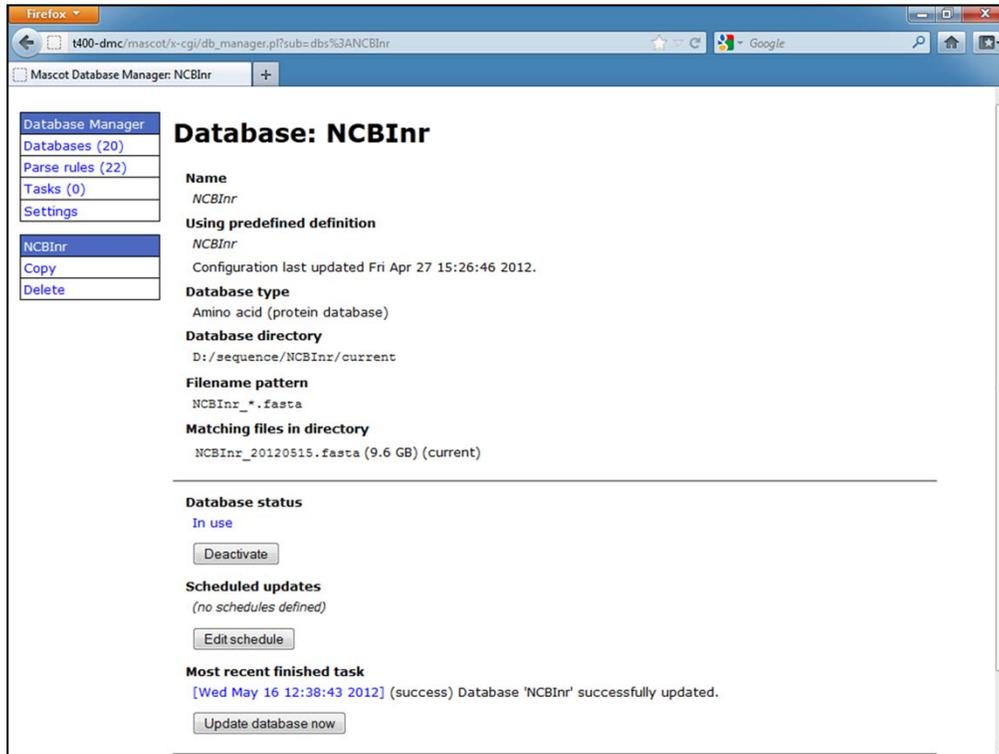
You can see here that it's started by downloading the taxonomy files that it needs



And it swiftly moves on to downloading the compressed fasta file which will probably take about an hour depending on the speed of your internet connection.

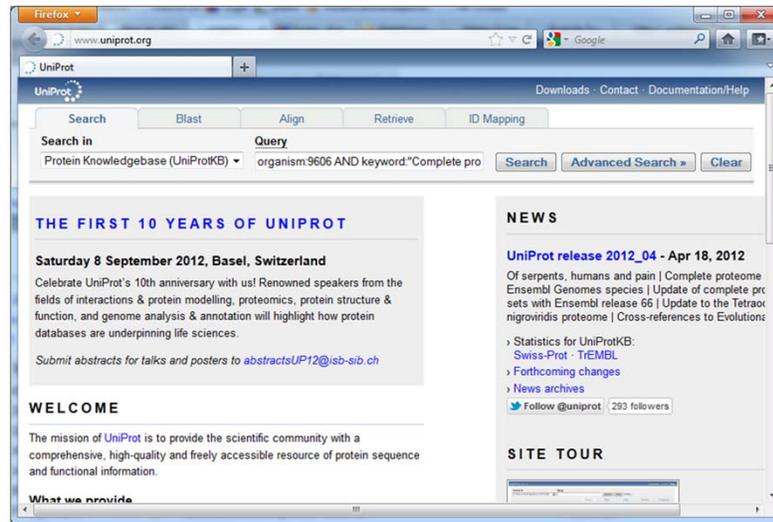


Once it's finished downloading and decompressing there will be a pause before it's available for searching, but you don't need to do anything.



And finally, it will come in use. If you want to automatically check for updates, say every week, you can do that here, but I'll come back to that later. There's also a manual "update database now" button if you want to just do this manually at any time.

## Species specific proteome from Uniprot



**MASCOT** : Database Manager

© 2012 Matrix Science

**MATRIX**  
**SCIENCE**

The next example is using a template, and in this case I'll generate a species specific proteome.

Last September saw the final update to the IPI databases which have become very popular in our community, and this is the recommended replacement. There are no downloadable fasta files, but rather a query interface and the option to produce fasta files for any species on the fly.

Firstly, I'll show you the uniprot 'interface' for producing the fasta files.

Go to [www.uniprot.org](http://www.uniprot.org), and in the search box enter the organism or taxonomy you are interested in and also the keyword "Complete proteome". I've typed organism:9606 and "Complete Proteome" to find all the human proteome proteins.

70,718 proteins

Download...

Or limit to SwissProt (20,243) AND reviewed:yes

Entry	Entry name	Status	Protein names	Gene names	Organism	Length	Keywords
Q2TB18	ASTE1_HUMAN	★	Protein asteroid homolog 1	ASTE1 HT001	Homo sapiens (Human)	679	Complete proteome; Reference proteome
Q8TF44	C2orf213_HUMAN	★	C2 calcium-dependent domain-containing protein C2orf213	C2CD4C FAM148C KIAA1957 NLF3	Homo sapiens (Human)	421	Complete proteome; Reference proteome
B7Z1M9	C2orf213_HUMAN	★	C2 calcium-dependent domain-containing protein C2orf213	C2CD4D FAM148D	Homo sapiens (Human)	353	Complete proteome; Reference proteome
Q5VUE5	C1orf53_HUMAN	★	Uncharacterized protein C1orf53	C1orf53	Homo sapiens (Human)	145	Complete proteome; Reference proteome
Q8NC38	C1orf213_HUMAN	★	Putative uncharacterized protein C1orf213	C1orf213	Homo sapiens (Human)	126	Complete proteome; Reference proteome
Q9UKZ1	C2orf29_HUMAN	★	UPF0760 protein C2orf29	C2orf29 C40	Homo sapiens (Human)	405	Complete proteome; Reference proteome

There are 70,718 proteins in SwissProt and TrEMBL, but I can limit it to just 20,243 proteins in SwissProt if I want, and this adds the text “reviewed:yes” to the URL

Clicking on the download link takes me to some choices:

For MS-MS, download with isoform sequences

```

>sp|P31946|1433B_HUMAN 14-3-3 pro...
MTMDKSELVQKAKLA . . .
>sp|P31946-2|1433B_HUMAN Isoform Short of 14-3-3 pro...
MDKSELVQKAK. . .

```

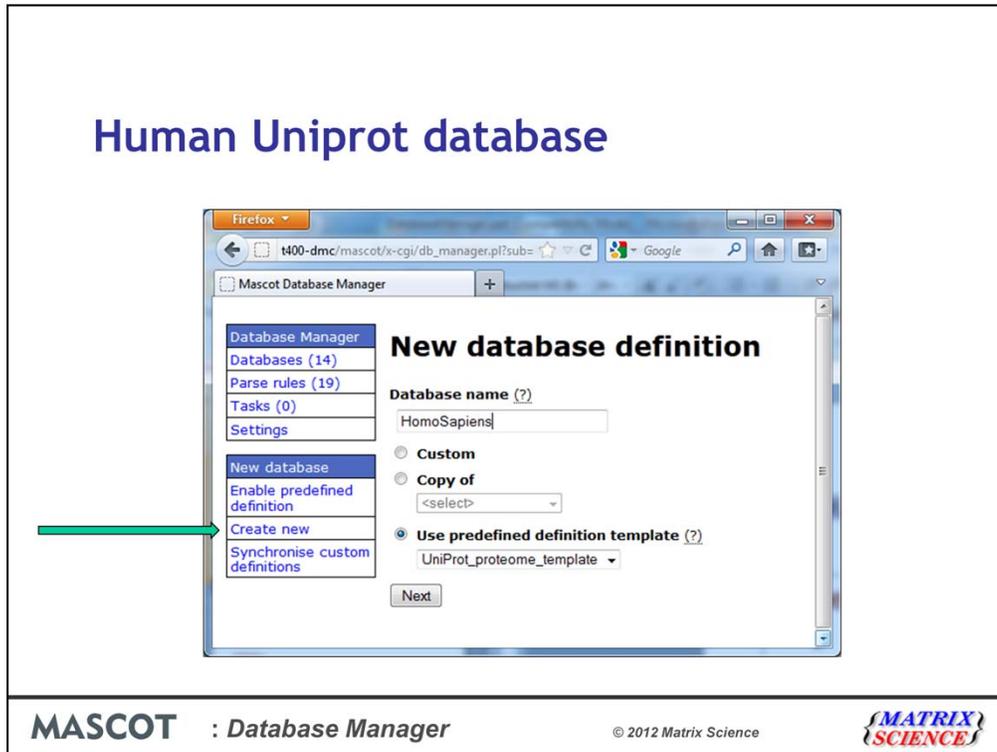
**MASCOT** : Database Manager © 2012 Matrix Science **MATRIX SCIENCE**

For searching ms-ms data with Mascot, we need to choose the fasta format and ideally should include with isoforms.

Isoforms means that multiple copies of some sequences are provided with the different isoforms. A -2, -3 -4 etc is appended to the accession. Of course, this means that the accession rather than the ID needs to be used for the Mascot accessions.

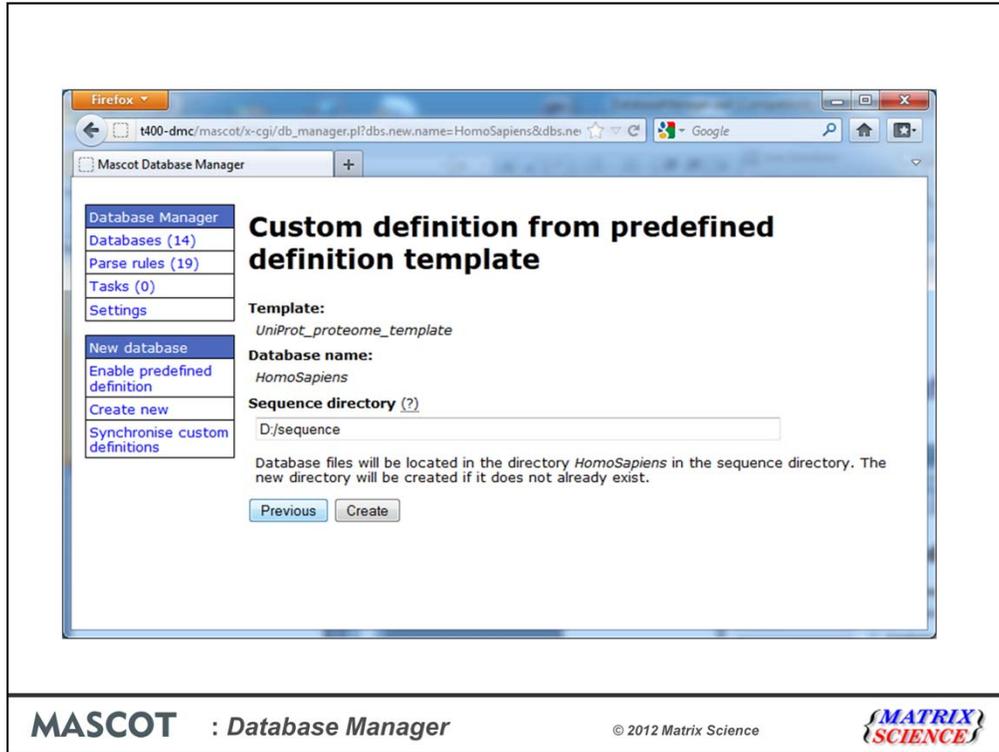
However, there's actually a better way of downloading if we want to keep the fasta file updated automatically.

## Human Uniprot database

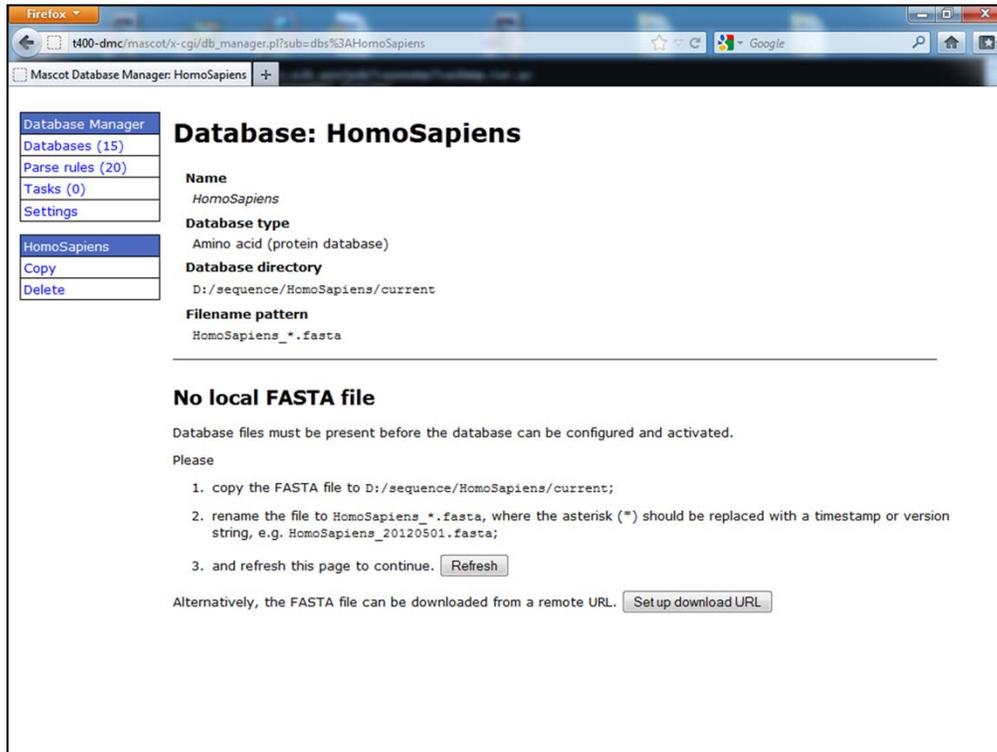


Back in Database Manager, click on Create new, and then enter the name and say that we want to use a template.

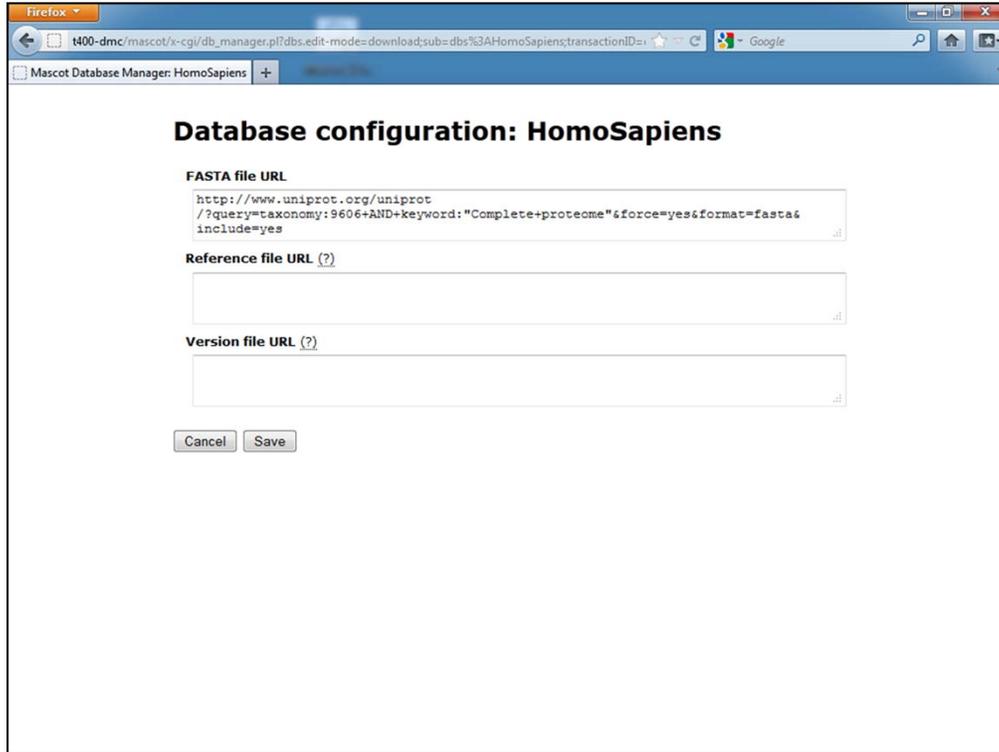
Next.



The next screen asks us to confirm the directory and we then click on Create



We now have a choice of downloading the file manually from uniprot.org as and copying it into the specified folder, or we can database manager to download it for us which is what I will do.



I don't expect you to remember this URL, but it's in the help for Database Manager on our public web site or on the uniprot web site.

The include=yes is for the isoforms.

If we use a URL here, we can come back to database manager later and just click a single button to force an update of the local database. Alternatively, we can set scheduled updates automatically. If we had chosen to download manually now, we would always need to download manually, so it's worth the extra effort here.

The screenshot shows a web browser window with the URL `t400-dmc/mascot/x-cgi/db_manager.pl?sub=db%3AHomoSapiens`. The page title is "Database: HomoSapiens". On the left, there is a navigation menu with options: "Database Manager", "Databases (15)", "Parse rules (20)", "Tasks (0)", "Settings", "HomoSapiens", "Copy", and "Delete". The main content area displays the following information:

- Name:** *HomoSapiens*
- Database type:** Amino acid (protein database)
- Database directory:** `D:/sequence/HomoSapiens/current`
- Filename pattern:** `HomoSapiens_*.fasta`

Below this, a section titled "No local FASTA file" states: "Database files must be present before the database can be configured and activated." It provides a "FASTA file URL": `http://www.uniprot.org/uniprot/?query=taxonomy:9606+AND+keyword:"Complete+proteome"&force=yes&format=fasta&include=yes`. There are two buttons: "Edit download URLs" and "Start downloading".

Alternatively, the page lists three steps:

1. copy the FASTA file to `D:/sequence/HomoSapiens/current`;
2. rename the file to `HomoSapiens_*.fasta`, where the asterisk (\*) should be replaced with a timestamp or version string, e.g. `HomoSapiens_20120501.fasta`;
3. and refresh this page to continue.

Click on "Start downloading".

The screenshot shows a Firefox browser window with the URL `t400-dmc/mascot/v-cgi/db_manager.pl?sub=db%3AHomoSapiens`. The page title is "Mascot Database Manager: HomoSapiens". On the left, there is a navigation menu with "Database Manager" selected, containing links for "Databases (15)", "Parse rules (20)", "Tasks (1)", and "Settings". Below this, the "HomoSapiens" section is active, with links for "Copy" and "Delete".

The main content area is titled "Database: HomoSapiens" and contains the following information:

- Name:** *HomoSapiens*
- Database type:** Amino acid (protein database)
- Database directory:** `D:/sequence/HomoSapiens/current`
- Filename pattern:** `HomoSapiens_*.fasta`

---

**Downloading files**

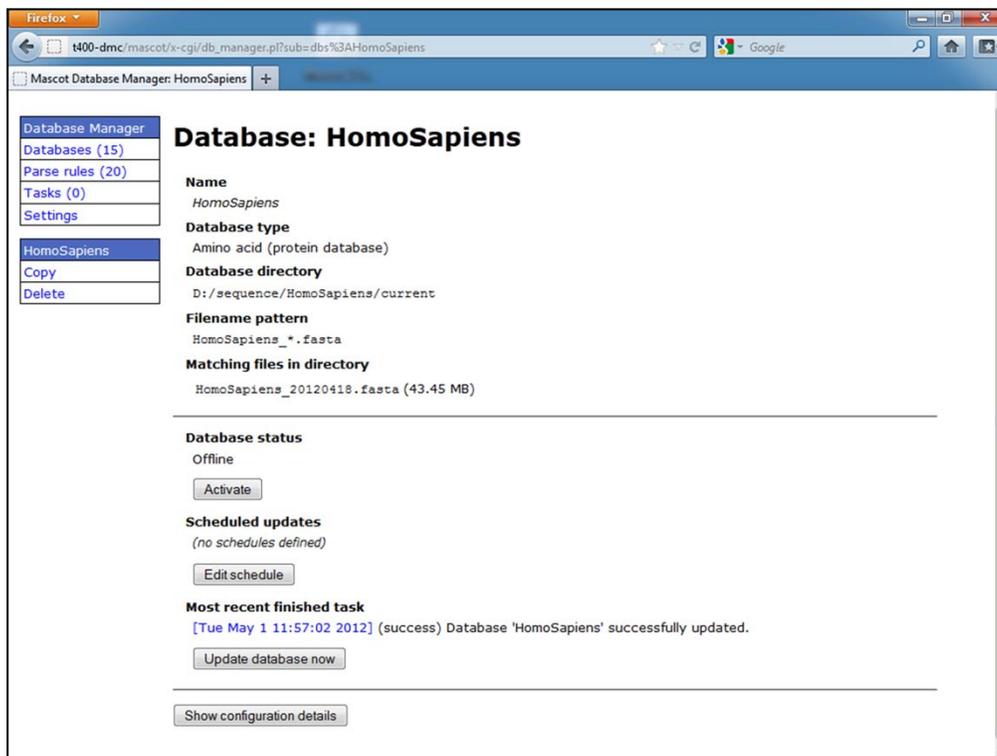
[10.1%] Downloading <<http://www.uniprot.org/uniprot/?query=taxonomy:9606+AND+keyword:%22Complete+proteome%22&force=yes&format=fasta&include=yes>> (try 1 out of 5): 585 kB/s (1.88 MB/unknown size)

**URLs to download**

**FASTA file URL**

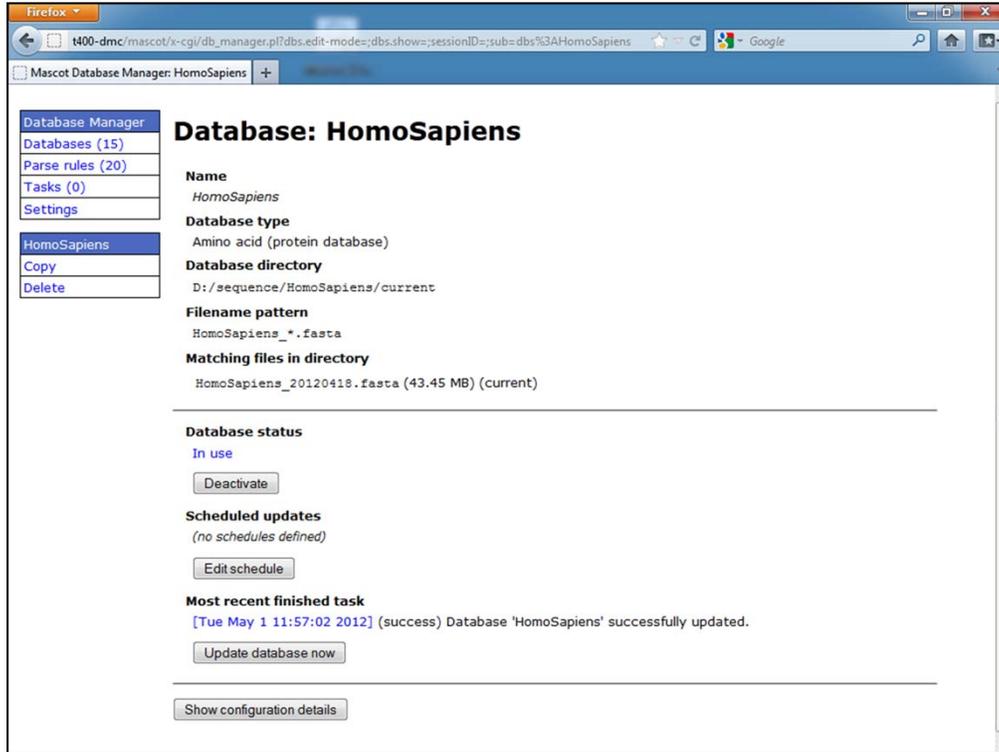
`http://www.uniprot.org/uniprot/?query=taxonomy:9606+AND+keyword:%22Complete+proteome%22&force=yes&format=fasta&include=yes`

And you'll see some progress



And finally you will need to click on Activate.

We don't need to worry about parse rules because it's all coped with in the predefined definition. If you remember, from a few screens ago, we needed to use the accession rather than the ID, and this is the default for this predefined definition.



And the database is now in use.

## Database Manager

### Create a copy of an existing database

- Copy of particular database release throughout a project, but also want the latest and greatest available separately
- A number of custom similar databases
- The configuration of the copy will not be kept up to date

### Custom definition

- Normally easier to use a template or copy.

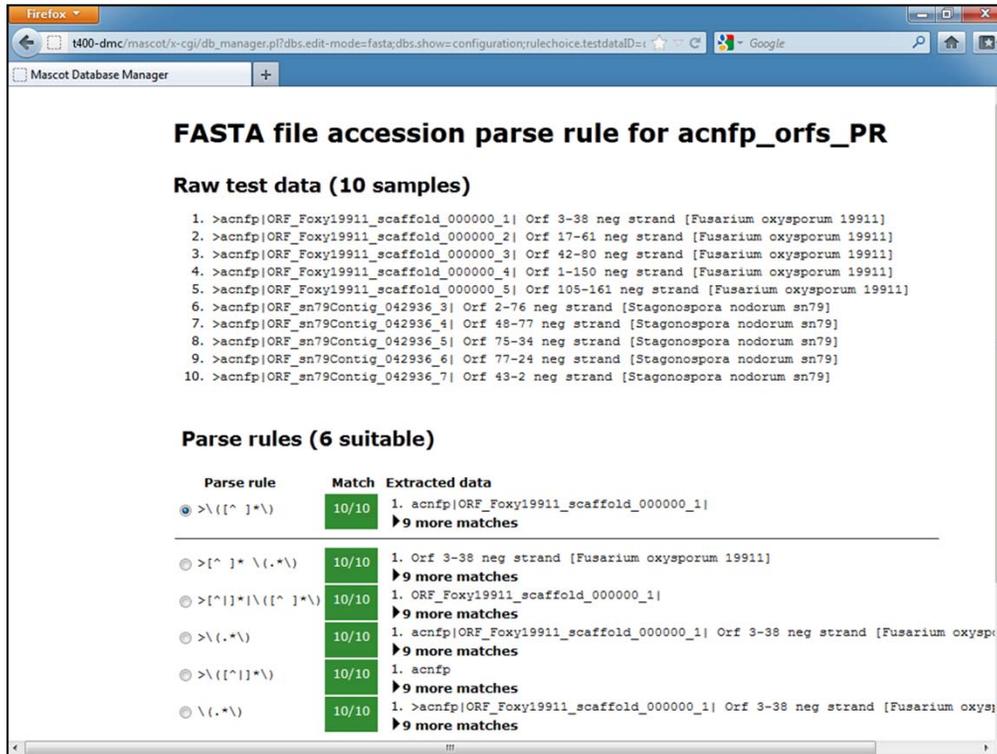
There are two other ways of installing a new database.

The first is to create a copy of an existing database. For example, you may wish to use a specific release of SwissProt throughout the duration of a lengthy project, but also want to have an up to date SwissProt available for searching. This is easy to do, and database manager will even make a copy of the files for you.

You may also have a number of databases with the same format. If you've set up one correctly, just create a copy of the configuration. With copies of any configuration, any changes to the parent database configuration won't trickle through to the copy.

Secondly, you can create a custom definition. To be honest, it's normally slightly easier to use a template such as the simple\_AA\_template or to make a copy of a similar working database.

In every case, you can still edit the rules (we call them parse rules) to extract the unique accession and also the sequence description. For most cases, this will follow the EBI or NCBI standards, but here's a different case:



Database manager searches all the rules that it knows about and finds which give any match at all and rejects the others. In this case, it found 6 possible rules. Database manager shows the first and last 5 sequences in the file. In this case, we can see that all sequences start with >acnfp, so it may make sense to exclude that.

The first suggestion includes everything up to the first space which is OK, but not ideal.

The second suggestion is the free form text. We'd be better off using that as a description.

The third suggestion excludes the acnfp and the pipe symbol, so is pretty much what we want

The fourth suggestion looks to be the whole description line which is too long

The fifth suggestion is not useful because all the accessions will be identical

The sixth suggestion is the whole line including the >

So, could choose the third suggestion. It's not quite perfect because it includes the trailing pipe symbol. Alternatively, we could define a new parse rule here which is much easier than in the older database maintenance utility.

## Database Manager - Scheduled Updates

Available for any database where the download URL has been specified

Can be daily, weekly, monthly

When using a predefined definition, a check is made on our public server for a change in the definition.

**MASCOT** : Database Manager

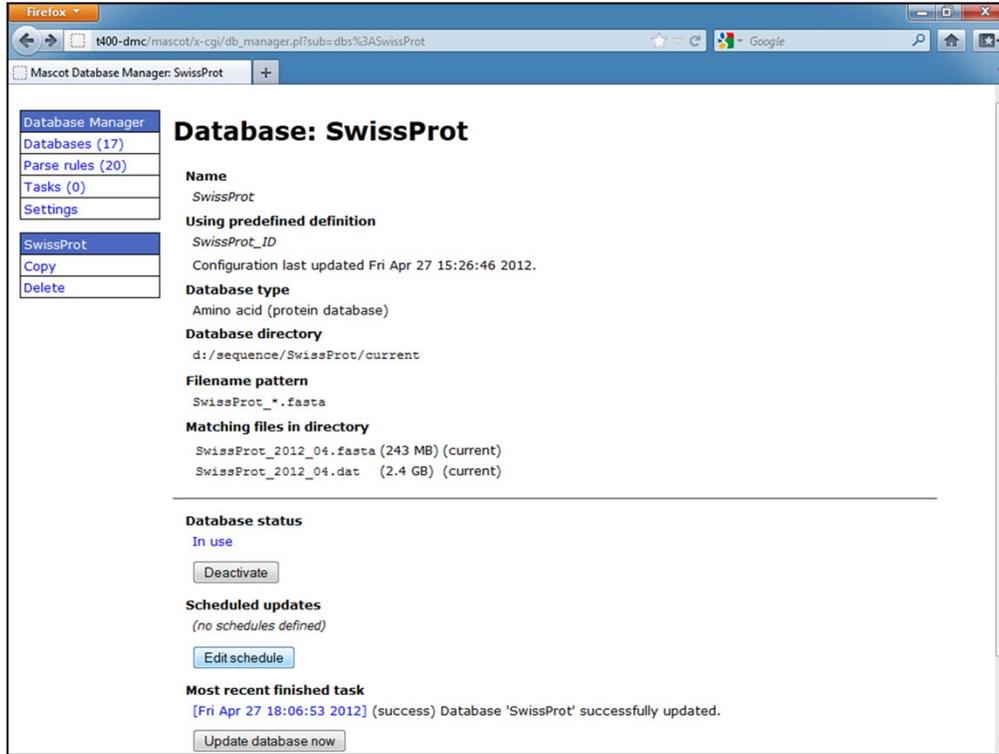
© 2012 Matrix Science

**MATRIX**  
**SCIENCE**

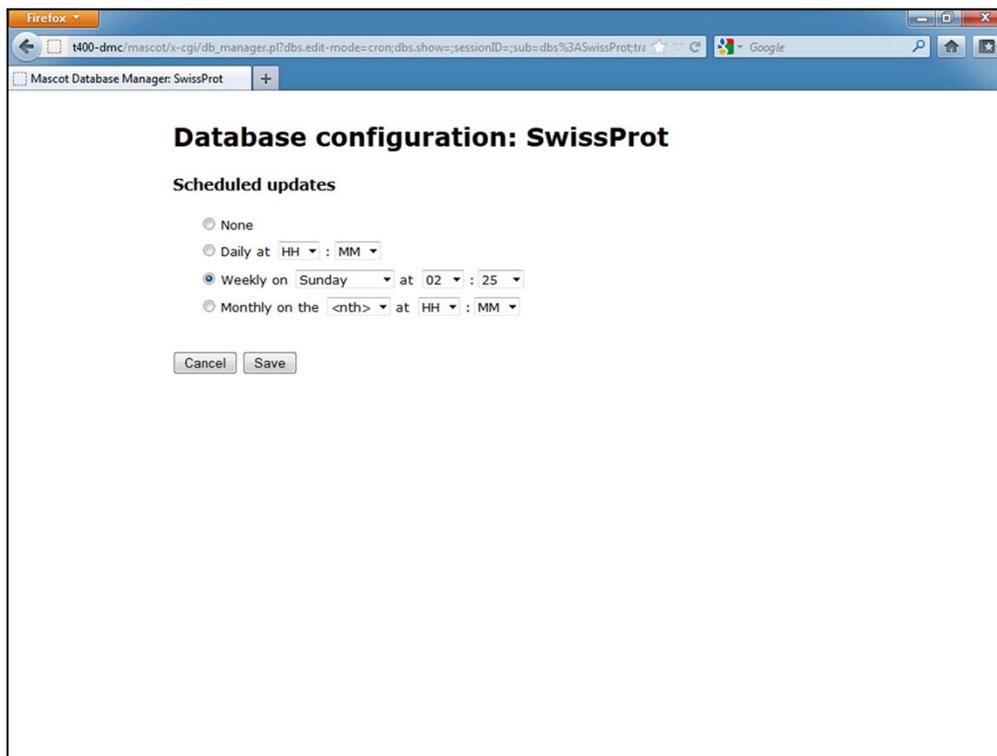
Most people, most of the time want to keep their sequence databases up to date. In previous versions of Mascot, we provided a utility called `db_update.pl` to do this. You needed to modify the script, test it at a command prompt and then run set up scheduled tasks or a cron job to automate it. Unfortunately, because this wasn't trivial, lots of people didn't bother. Also, for people who have set it up, they would sometimes find it stopped working because the database supplier moved the files, or changed the format.

It's much easier in Mascot 2.4:

- As long as you can specify a download URL, it can be automated.
- There's some simple choices for daily, weekly or monthly which should be sufficient for most cases.
- If the database supplier changes something, we will change our public definitions, and the updates should continue to work. I'll give an example of how to do it for SwissProt.



From the databases screen, click on SwissProt and then click on the Edit Schedule button



Here, I've decided to update it weekly on Sunday at 2:25 am.

In practice, SwissProt only gets updated monthly. However, database manager is pretty smart and will only download files if they have changed, so it doesn't matter that it's checking more often than it needs to.

## Database Manager

- Everything done in one place
- No need to edit perl scripts or running them at the command line
- No need to understand 'cron' or Windows 'scheduled tasks'
- All additional files such as taxonomy indexes and unigene indexes are kept up to date automatically
- Transparently deals with changes to database locations and formats.

**MASCOT** : Database Manager

© 2012 Matrix Science

**{MATRIX}**  
**{SCIENCE}**

I hope I've shown that it really is easy to use database manager and that this is a great improvement on the different utilities that we provided in previous versions. As you've seen, there's just one utility that does it all.

Also, there's no requirement to edit perl scripts, like the db\_update.pl script, and there's no need to run anything at a command line any more.

There's also no need to understand cron, or Windows scheduled tasks.

All additional files such as taxonomy indexes and unigene indexes are kept up to date automatically, so you don't need to worry about downloading them and making sure that they are unpacked and put into the correct directory.

And finally, when the NCBI or EBI change a database format or location, this should all be dealt with transparently.