# Turn-key Automation using Mascot Daemon
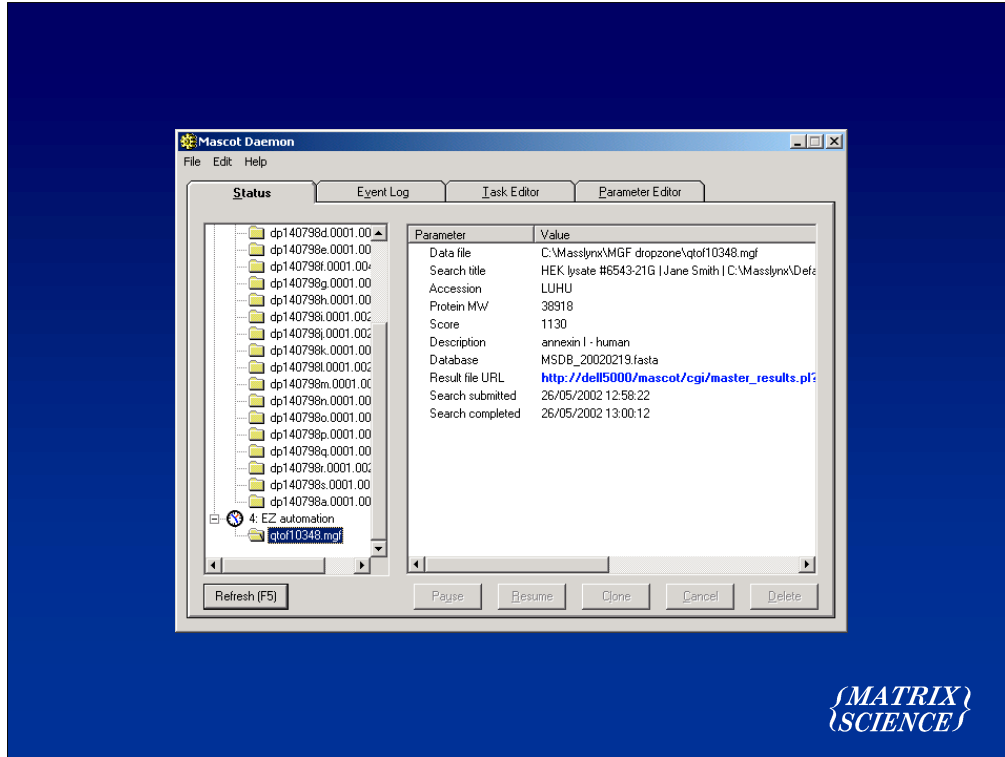
*MATRIX SCIENCE*

Mascot Daemon is our automation client. It's a part of the Mascot package that is only available when you have Mascot on an in-house server

## Mascot Daemon

**1. Batch task**
*A batch of data files to be searched immediately or at a defined time*

**2. Real-time monitor task**
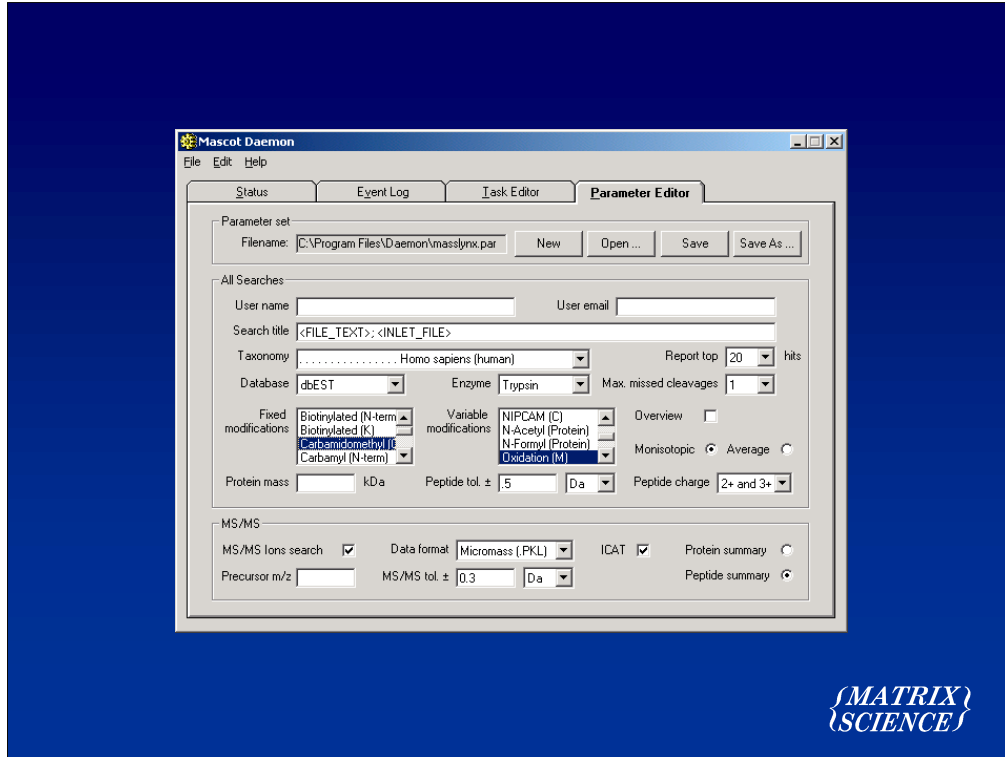*New files on a defined path are searched as they are created*

**3. Follow-up task**
*Accepts data files from another task. For example, to repeat a search against a different database*
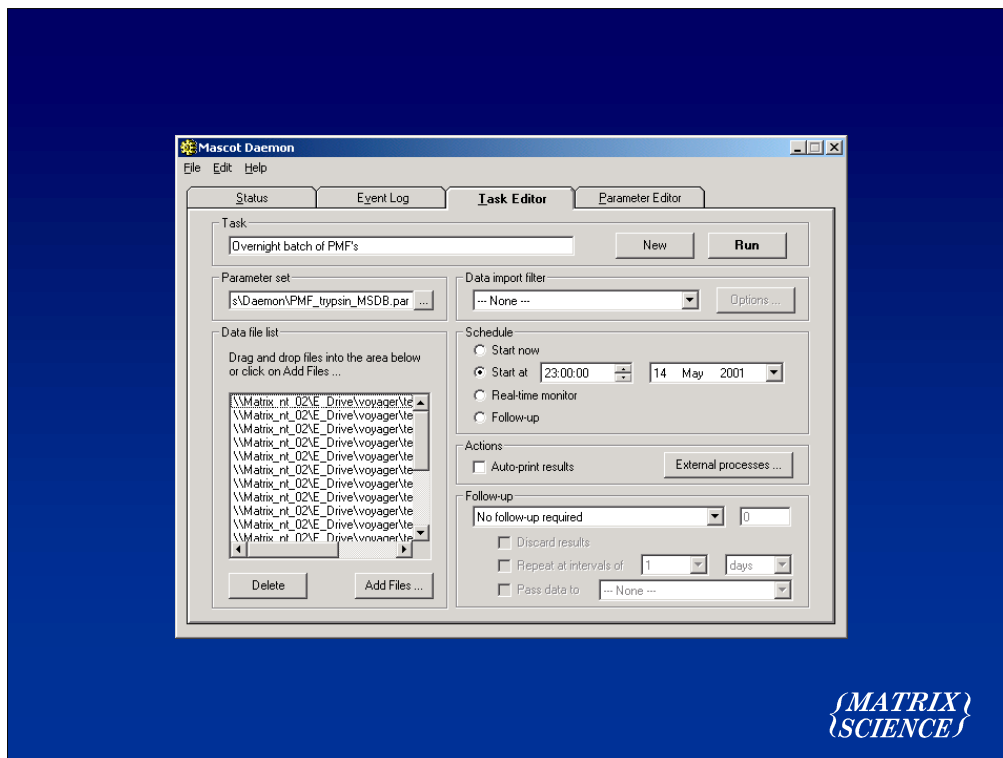
{MATRIX} {SCIENCE}

Mascot Daemon runs on any Win32 platform and supports three kinds of tasks.

The follow-up task is very powerful because it allows searches to be chained together to implement complex decision paths. For example, as batch of data files might be screened against a contaminants database containing entries for keratins, BSA, trypsin, etc. Those data files which fail to find a match can then be automatically searched against a non-redundant protein database. Spectra which are still unmatched can then be searched against a large EST database, etc., etc.
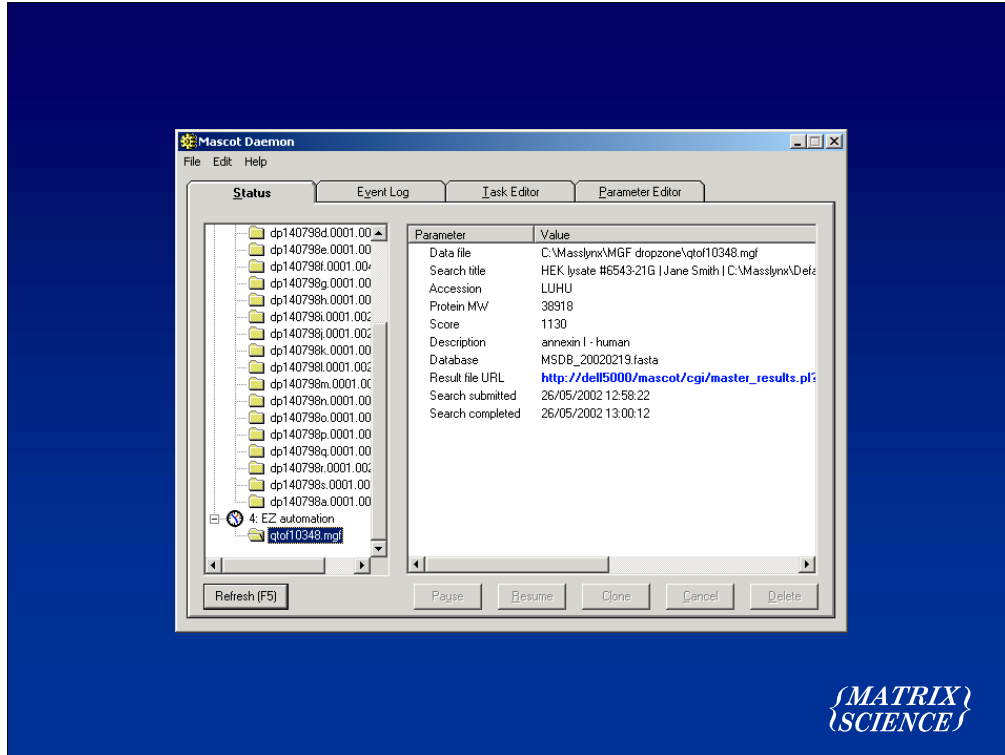
The parameter editor allows sets of search parameters to be defined and saved to disk, so that they can be used over and over again. The search parameters define *how* the data will be searched.

The Task Editor tab is used to define each task. A task defines *what* data will be searched and and *when* the search will take place.
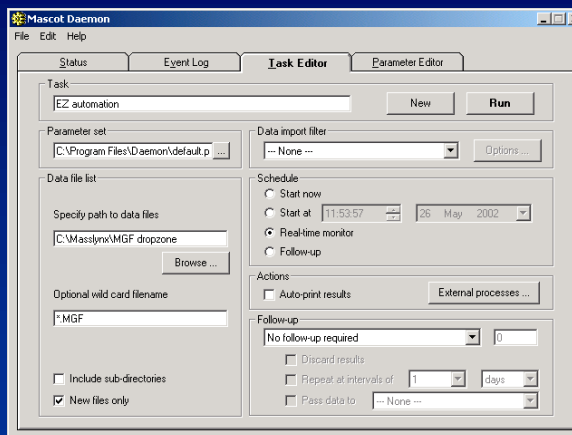
Here we have a simple batch task. A set of data files has been created, a parameter set has been chosen, and the task will run all the searches as a batch at a predefined time.

Each completed search is added to an Explorer-like tree. There is a minimum of summary information, and a hyperlink to the full result report

However, automation is more than just searching batches of files. This block diagram illustrates the workflow.

**Automation: zero coding**

**Use Mascot Daemon in real-time monitor mode to search files as they are created**

How far can you get in implementing this workflow with a minimum of coding? In fact, what can you do without any coding?

The starting point is to use Mascot Daemon in real-time monitor mode. This will pick up peak lists as they are created, and search them automatically.

For some file formats, we have data import filters. These allow Daemon to generate peak lists directly from the binary files created by the MS data system.

For example, the import filter for Sciex wiff files is fairly sophisticated.

We can also work directly from Xcalibur raw files and Applied Biosystems voyager data files

The MassLynx sample list import filter is a little different.

The aim is to extract filenames and sample information from the sample list. We still depend on Masslynx to generate PKL peak lists.

Tags in search title field can be substituted at run time by fields from the MassLynx sample list. For example, "HEK digest band 1", or any text from any of the columns, could appear at the top of the Mascot search report.

*Automation: more flexibility*

Use a script or macro to write peak lists in MGF format, including search parameters

```
1  # MGF format file generated by my VBA script
2  # 14:02:05 29 Nov 2001
3  #
4  SEARCH=MIS
5  USERNAME=Jane Smith
6  COM=HEK lysate #6543-21G | Jane Smith | C:\Masslynx\Default.pro\DATA\qtof10348.raw
7  DB=MSDB
8  CLE=Asp-N
9  TOL=0.1
10 TOLU=Da
11 MODS=Carbamidomethyl (C)
12 IT_MODS=Oxidation (M)
13 _sample_name=HEK lysate #6543-21G
14 _original_filename=C:\Masslynx\Default.pro\DATA\qtof10348.raw
15 _sample_barcode=4826659337524638329734239473216128138443
16
17 BEGIN IONS
18 TITLE=sum of scan(s) 521 - 528, RT 12:34
19 PEPMASS=1095.000000
20 CHARGE=2+
21 1093.439900   53107778.775450
22 1085.600000   58039365.132336
23 868.269559   11538842.201475
24 1321.475432   12294627.319613
25 1570.900000   13442700.087962
26 668.203086   7581443.399182
27 1521.100000   12143646.083781
```

*{MATRIX} {SCIENCE}*

That is about as far as we can get with zero coding. Some simple coding is required to achieve more flexibility. Maybe the MS data system supports a macro language, such as VBA. If so, you can write out peak lists in the Mascot Generic Format (MGF). This allows the search parameters to be embedded into the data file, avoiding the need to set up search parameters manually in Daemon. You can also add your own parameters, which will be passed through the search engine into the results file. Any parameter that starts with an underscore is a 'user' parameter.

It is also important to make full use of the search title field. This is displayed in Daemon

And also at the top of the results report. So, it is worth taking the trouble to include key sample tracking information in the title field

Especially when you need to find an old search from among the tens of thousands on the Mascot server
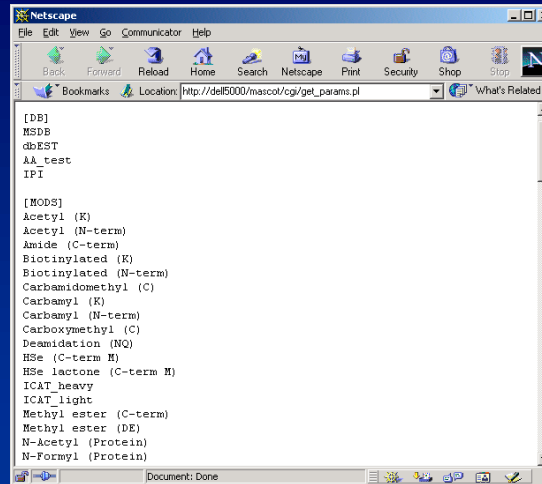
As mentioned earlier, any parameters with an underscore will be passed through to the result file. So you can include structured sample tracking information.

Another improvement to the workflow is to query the Mascot server to find out what databases are available,as well as the choice of enzymes, modifications, taxonomy, etc. This is achieved using a utility called get_params.pl

One thing we haven't addressed yet is how to tie up the result file to the input data.

This information can be found in the Mascot Daemon database. By default, this is an Access database.

However, the Mascot Daemon tables can live in any ODBC compliant database engine, such as Oracle or SQL Server. So, if you have an Oracle based LIMS, the Daemon tables can be inside the LIMS.

*Automation: parsing results*
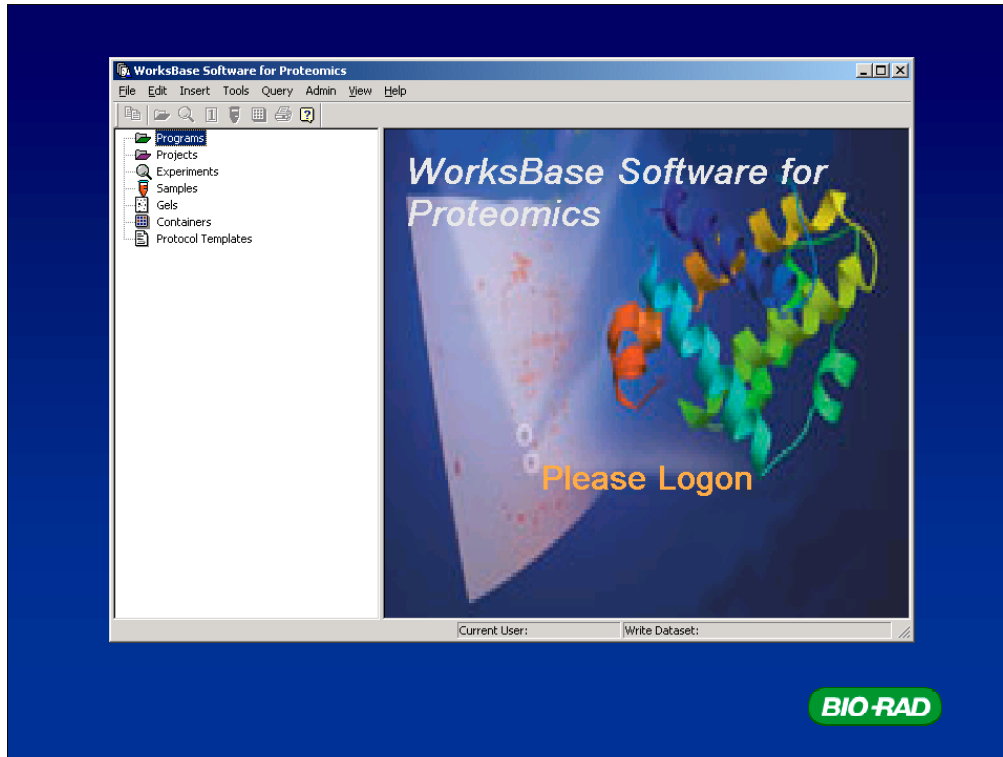
Mascot result files are simple text files. No need for screen scraping.

The next link in the chain is to export the results to a relational database, possibly the LIMS.
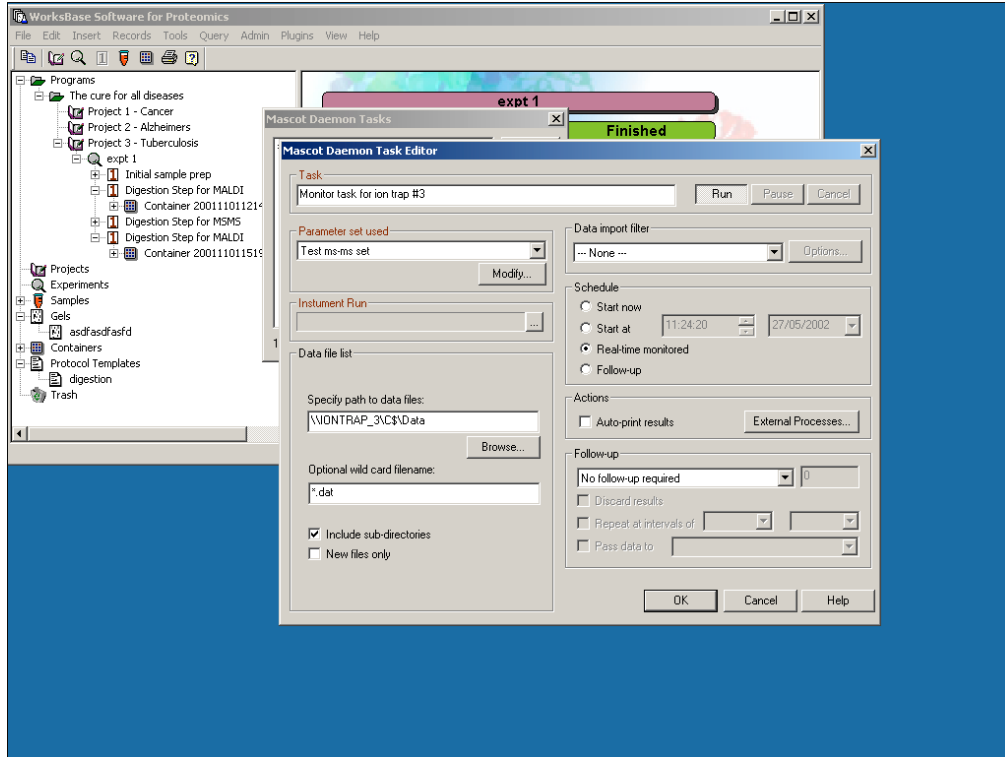
Mascot result files are structured text files, and it isn't difficult to write Perl scripts to fish out the information of interest
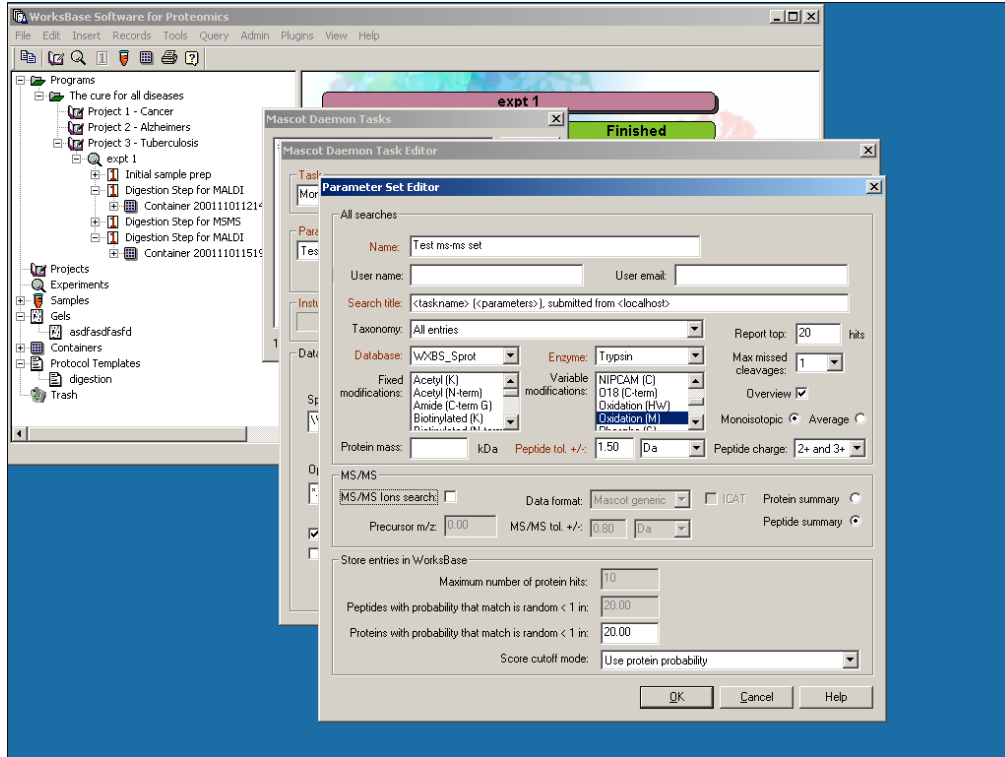
However, many programmers prefer to work in C++, so we are developing a class library that will support an object oriented approach to extracting the result information
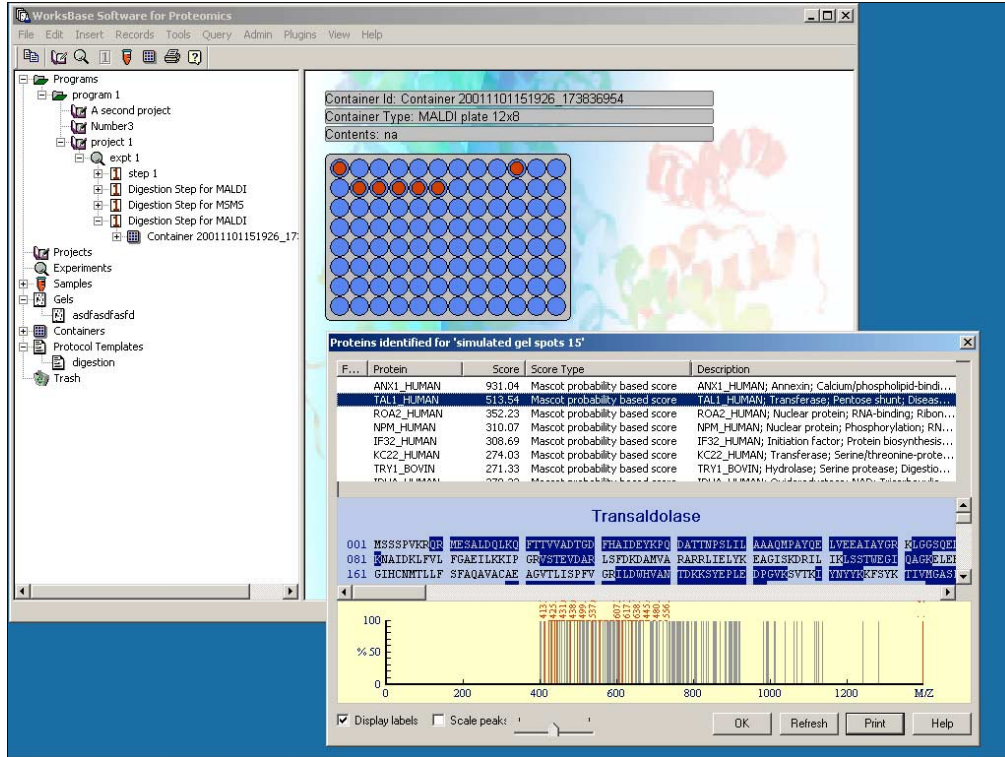
A good illustration of extensive integration between Mascot and a relational database is the WorksBase package from Bio-Rad.  WorksBase is a proteomics LIMS
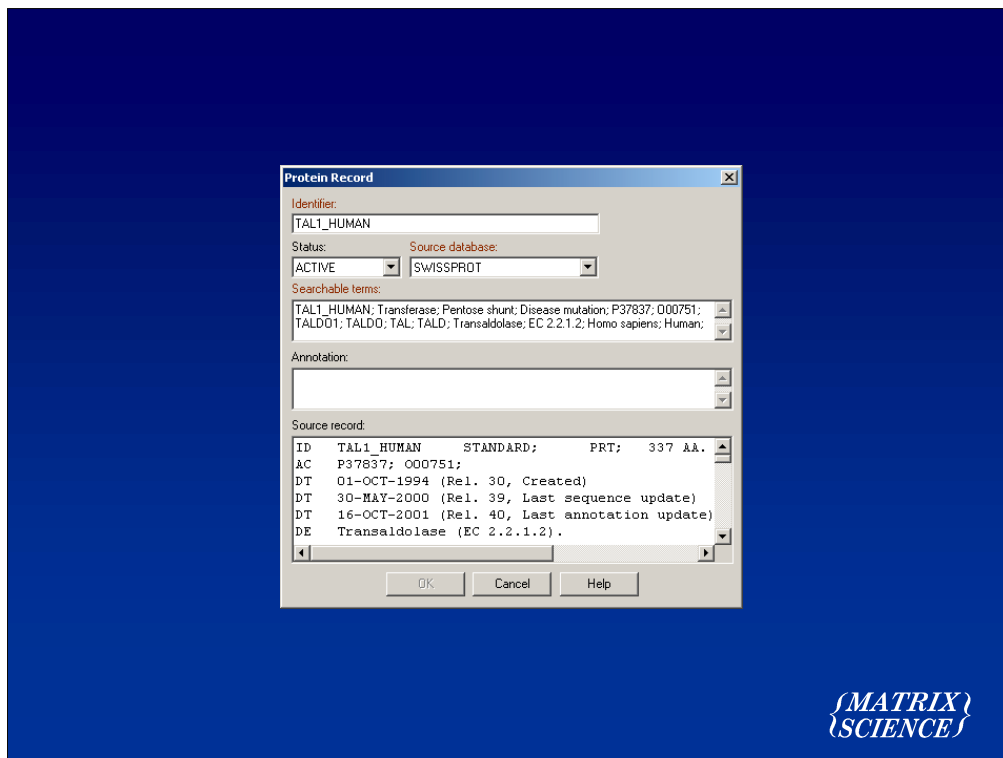
WorksBase provides the user interface for all aspects of Mascot searching. For example, tasks are defined in WorksBase dialogs and saved to WorksBase tables
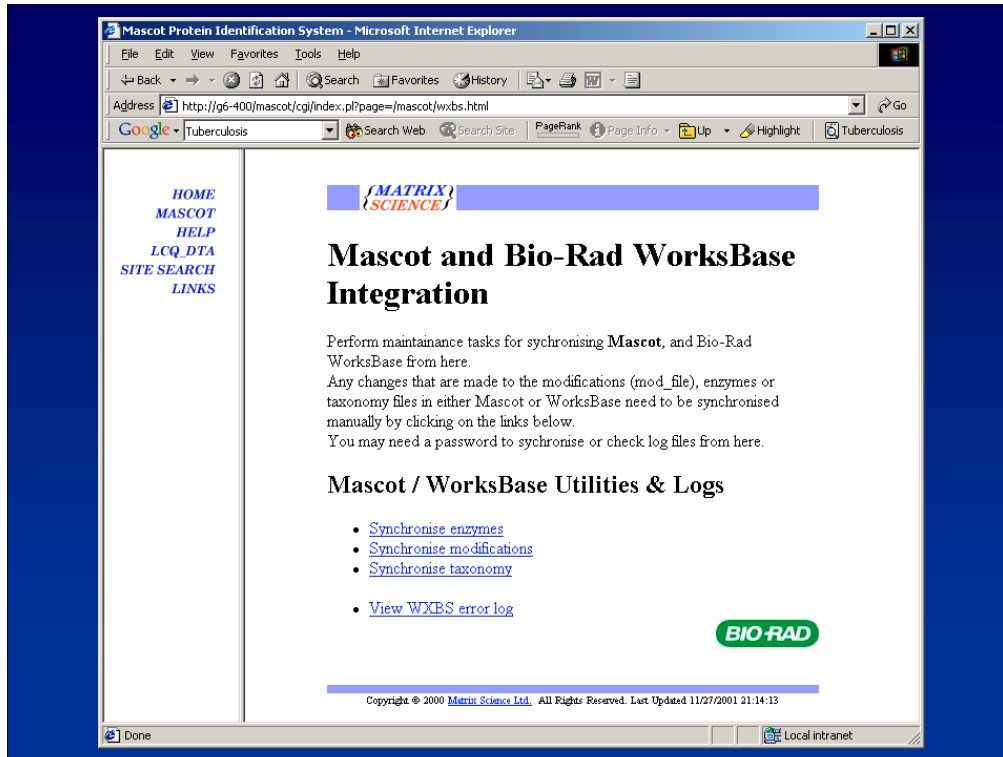
As are sets of search parameters

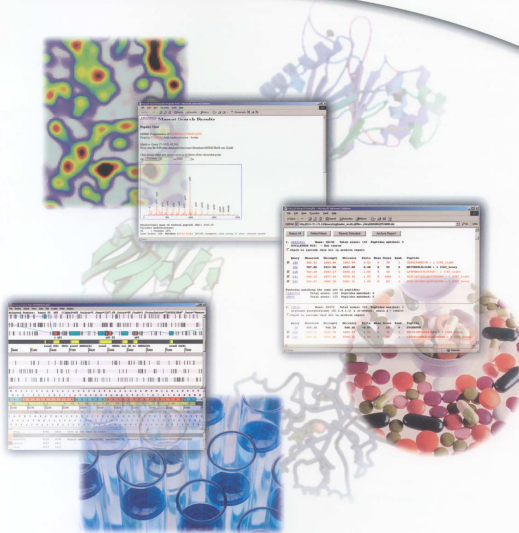The results are automatically exported to WorksBase, and reports can be generated by querying the database

Even the sequence database is contained within WorksBase tables

And there are web page links for administration tasks