

# Automated integration of Mass-Spectrometry based Proteomics evidence for improvement of gene annotations

Ritesh Krishna  
Institute of Integrative Biology  
University of Liverpool

# Objective

- Use of mass-spectrometry data to confirm that predicted gene models are translated into proteins
- Use the dataset to examine whether there are supporting evidences for alternative gene predictions at particular loci
- To check whether the same data can be used to identify novel genes

## Data Summary

## News

- 18 November 2010 AmoebaDB 1.3 Released
- 18 November 2010 PlasmoDB 7.1 Released
- 18 November 2010 TriTrypDB 2.5 Released



We are pleased to announce our 2011 EuPathDB Workshop, June 12-15, 2011 in Athens, GA, USA. For more information and to apply click here. **Application deadline is February 1, 2011.**

EuPathDB Bioinformatics Resource Center for Biodefense and Emerging/Re-emerging Infectious Diseases is a portal for accessing genomic-scale datasets associated with the eukaryotic pathogens (*Cryptosporidium*, *Encephalitozoon*, *Entamoeba*, *Enterocytozoon*, *Giardia*, *Leishmania*, *Neospora*, *Plasmodium*, *Toxoplasma*, *Trichomonas* and *Trypanosoma*).



AmoebaDB



CryptoDB



GiardiaDB



MicrosporidiaDB



PlasmoDB



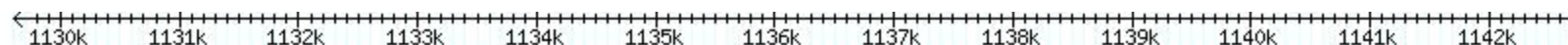
ToxoDB



TrichDB



TriTrypDB



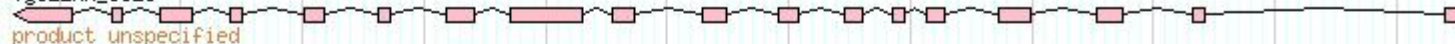
### Annotated Genes (with UTRs in gray when available)

TGME49\_026960



### GLEAM Gene Models

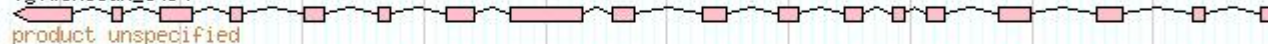
TgGLEAM\_6626



product unspecified

### TwinScan Gene Models

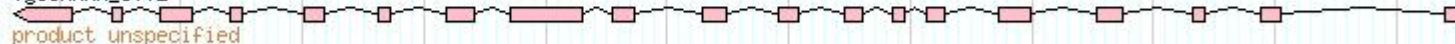
TgTwinScan\_5494



product unspecified

### GlimmerHMM Gene Models

TgGlimHMM\_3001



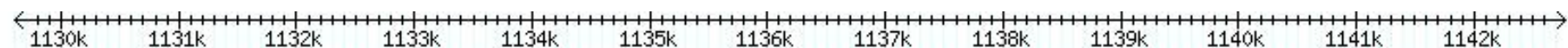
product unspecified

### TigrScan Gene Models

TgTigrScan\_6172

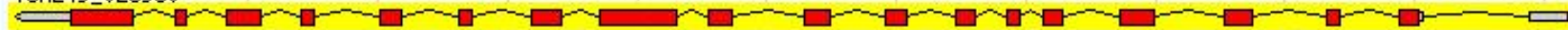


product unspecified



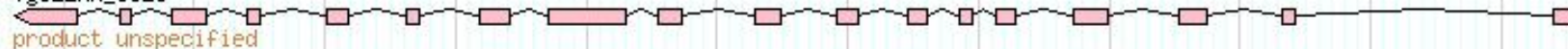
**Annotated Genes (with UTRs in gray when available)**

TGME49\_026960



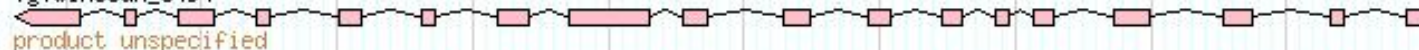
**GLEAN Gene Models**

TgGLEAN\_6626



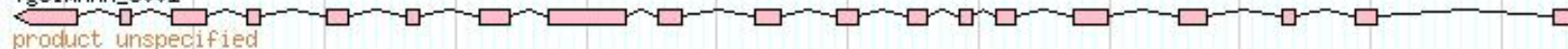
**TwinScan Gene Models**

TgTwinScan\_5494



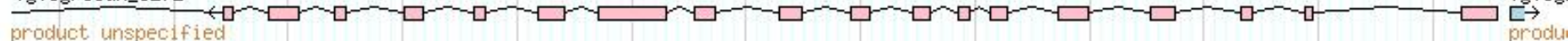
**GlimmerHMM Gene Models**

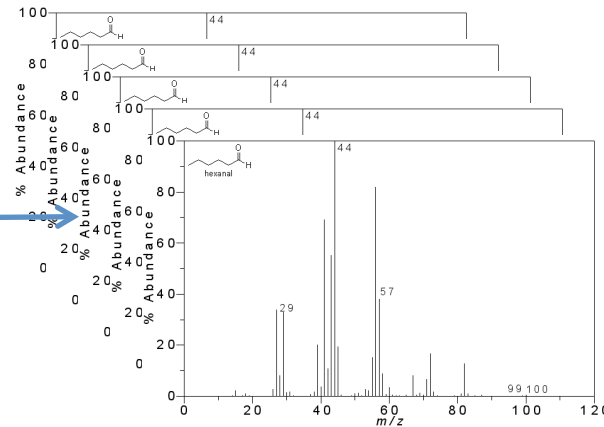
TgGlmHMM\_3001



**TigrScan Gene Models**

TgTigrScan\_6172

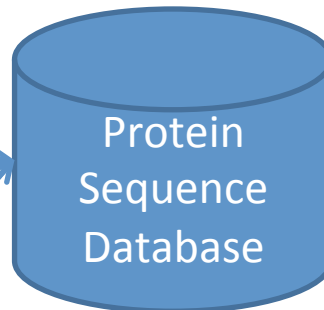




## Gene Models



... Any protein  
sequence database of  
your choice



## Search Engines

Mascot



Ommsa

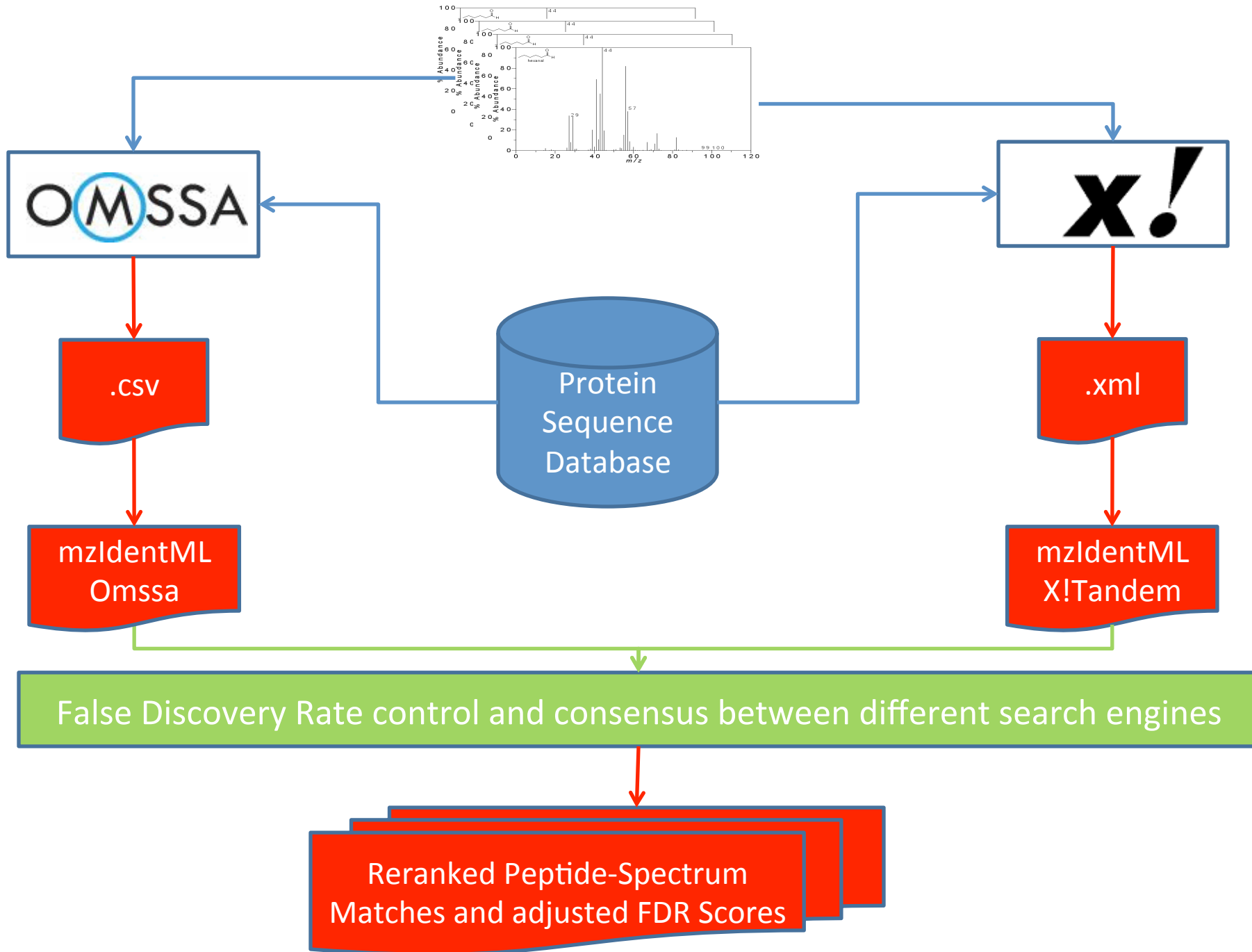


X!Tandem



...etc.

Peptide-Spectrum matches and Protein Identification



# Modes of operation



Single input file



Whole directory of input files



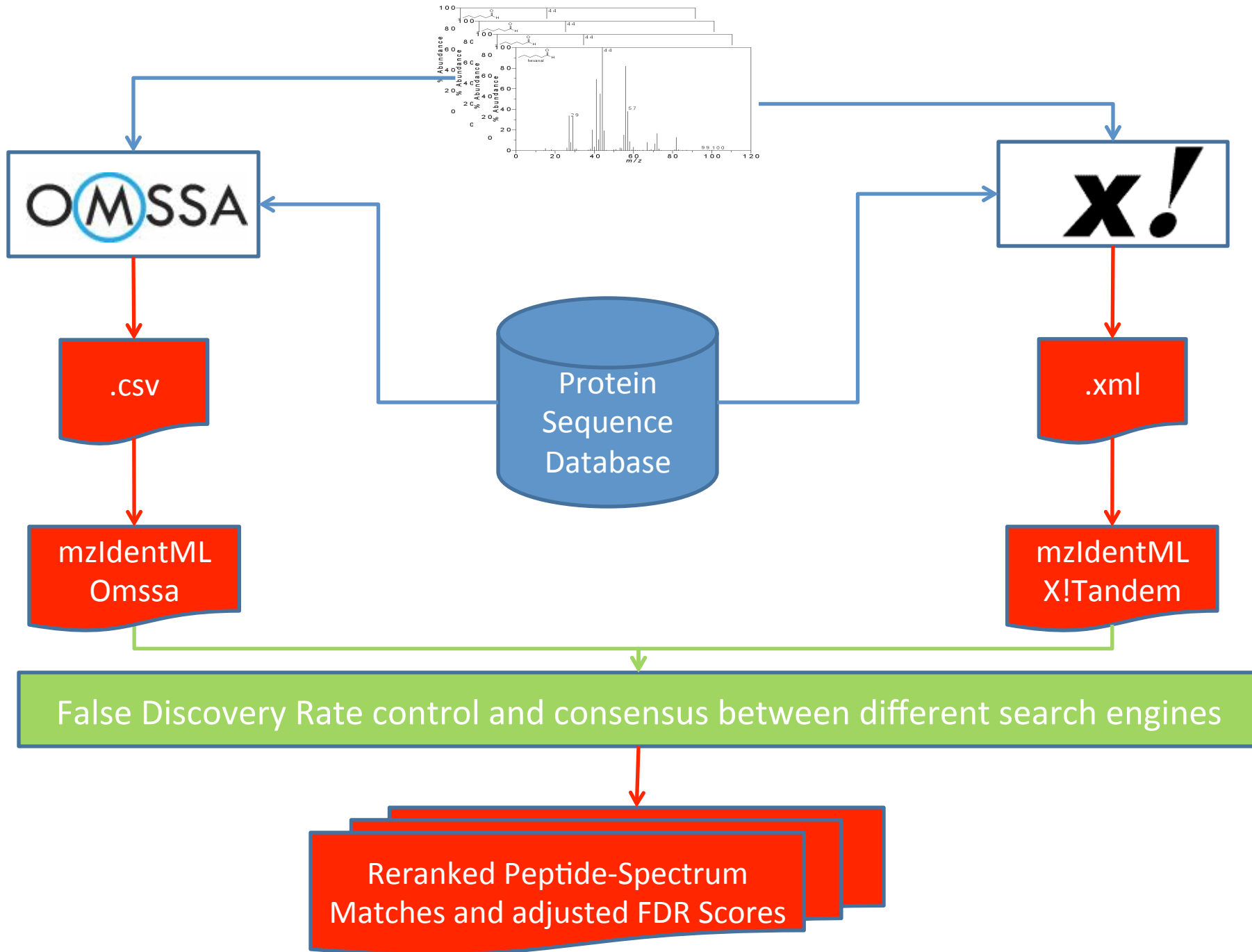
Single machine



A cluster of machines

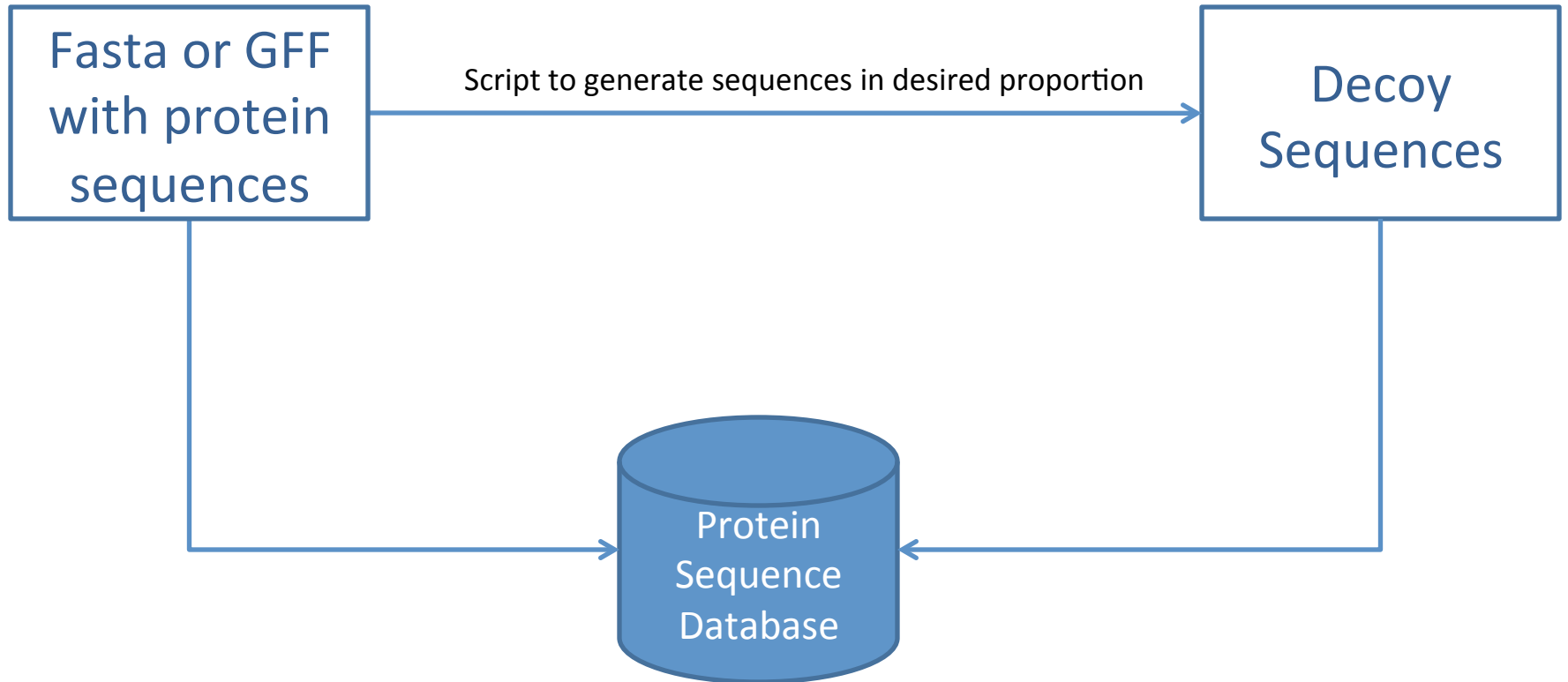
FASTA

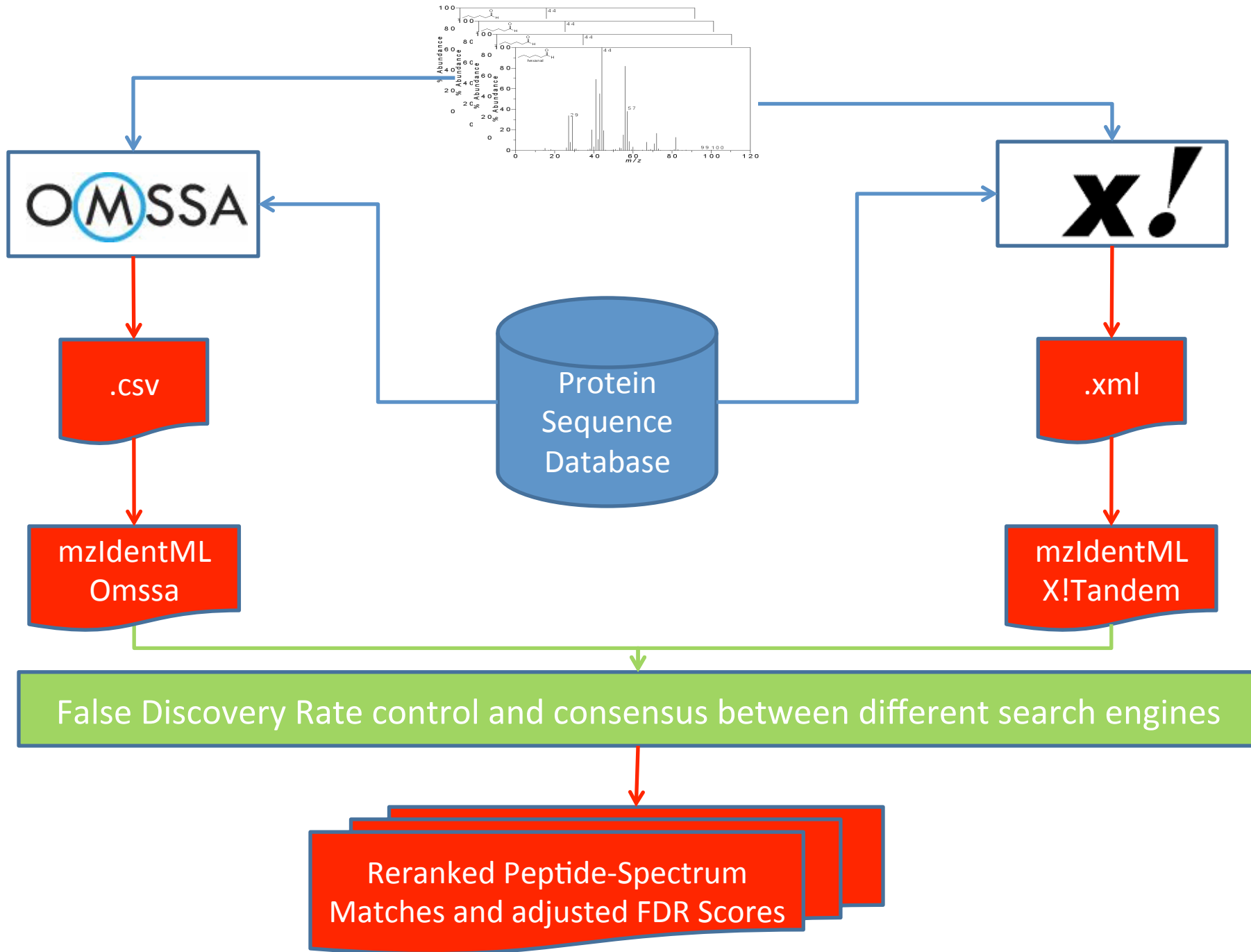
GFF





# Protein Sequence Database

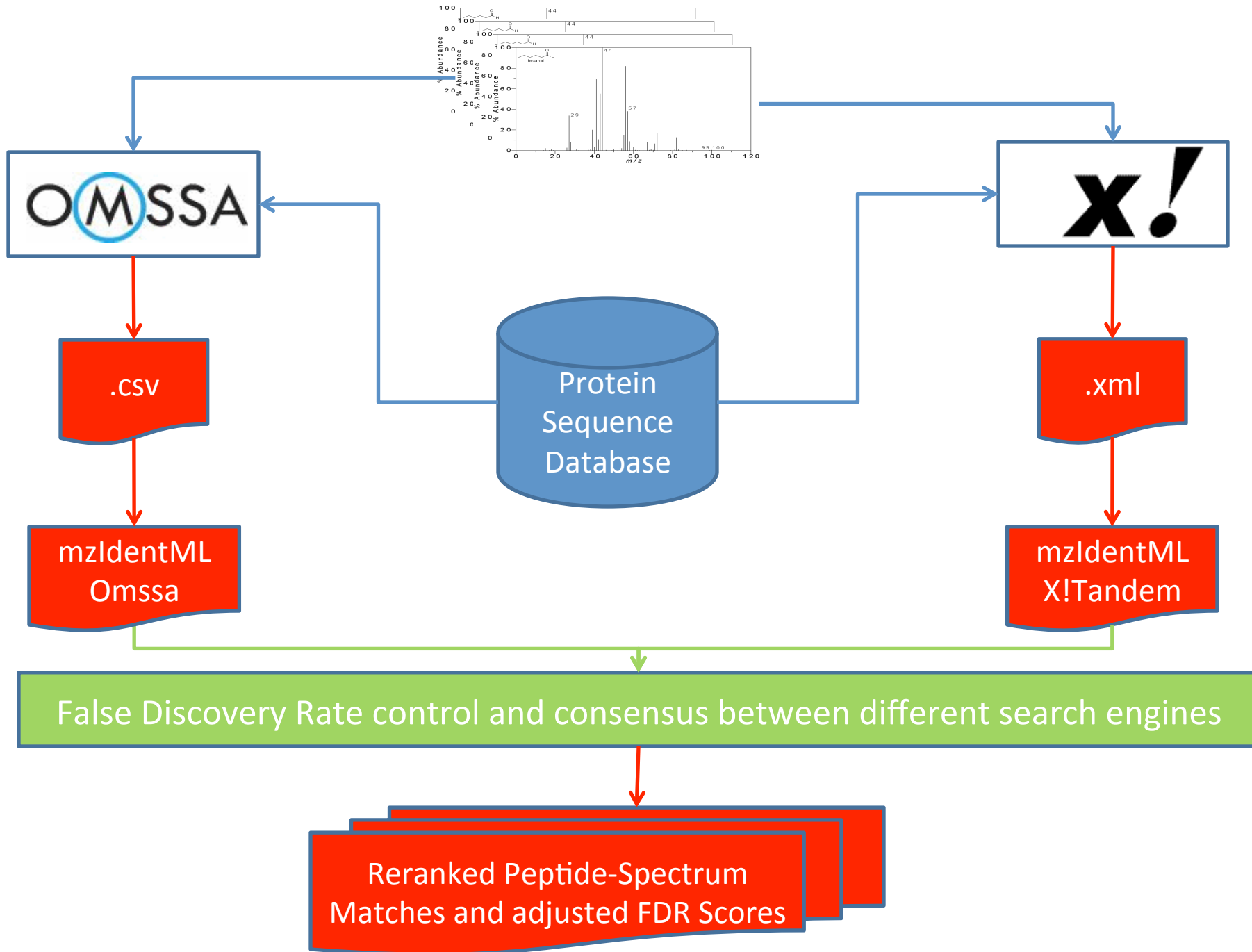






Controlling the search engines –

- Specification of digestive enzyme used
- Specification of parent and fragment tolerances
- Specification of maximum missed cleavage allowed
- Specification of Post translational modifications



- Increase confidence by creating a consensus between Omssa and X!Tandem results
- Use of Decoy database to compute False Discovery Rate
- Details of the algorithm available at the following publications

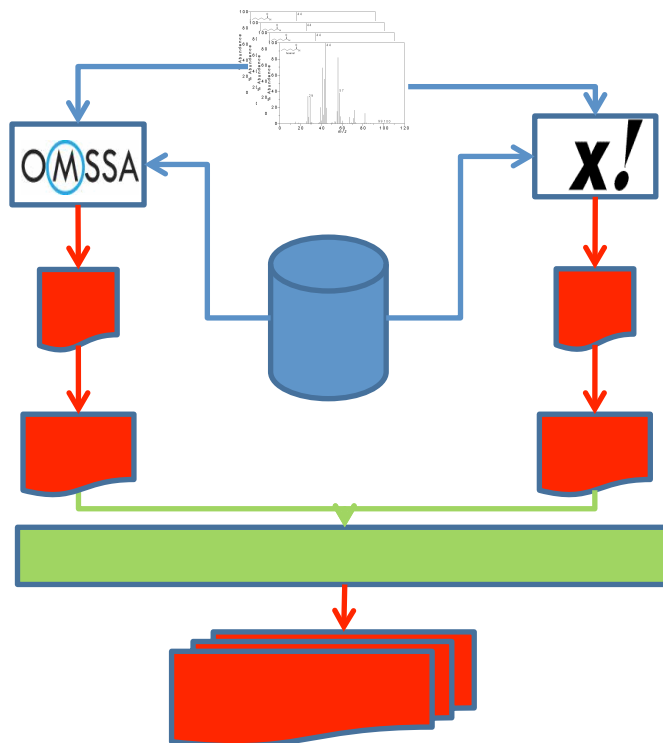
Jones, A. R., Siepen, J. A., Hubbard, S. J. and Paton, N. W. (2009), Improving sensitivity in proteome studies by analysis of false discovery rates for multiple search engines. *PROTEOMICS*, 9: 1220–1229.

Wedge, D. C., Krishna, R., Blackhurst, P., Siepen, J.A., Jones, A.R., Hubbard, S.J. (2011), **FDRAnalysis: A Tool for the Integrated Analysis of Tandem Mass Spectrometry Identification Results from Multiple Search Engines.** *Journal of Proteome Research*, 10 (4), 2088-2094

# FASTA



# Text Files



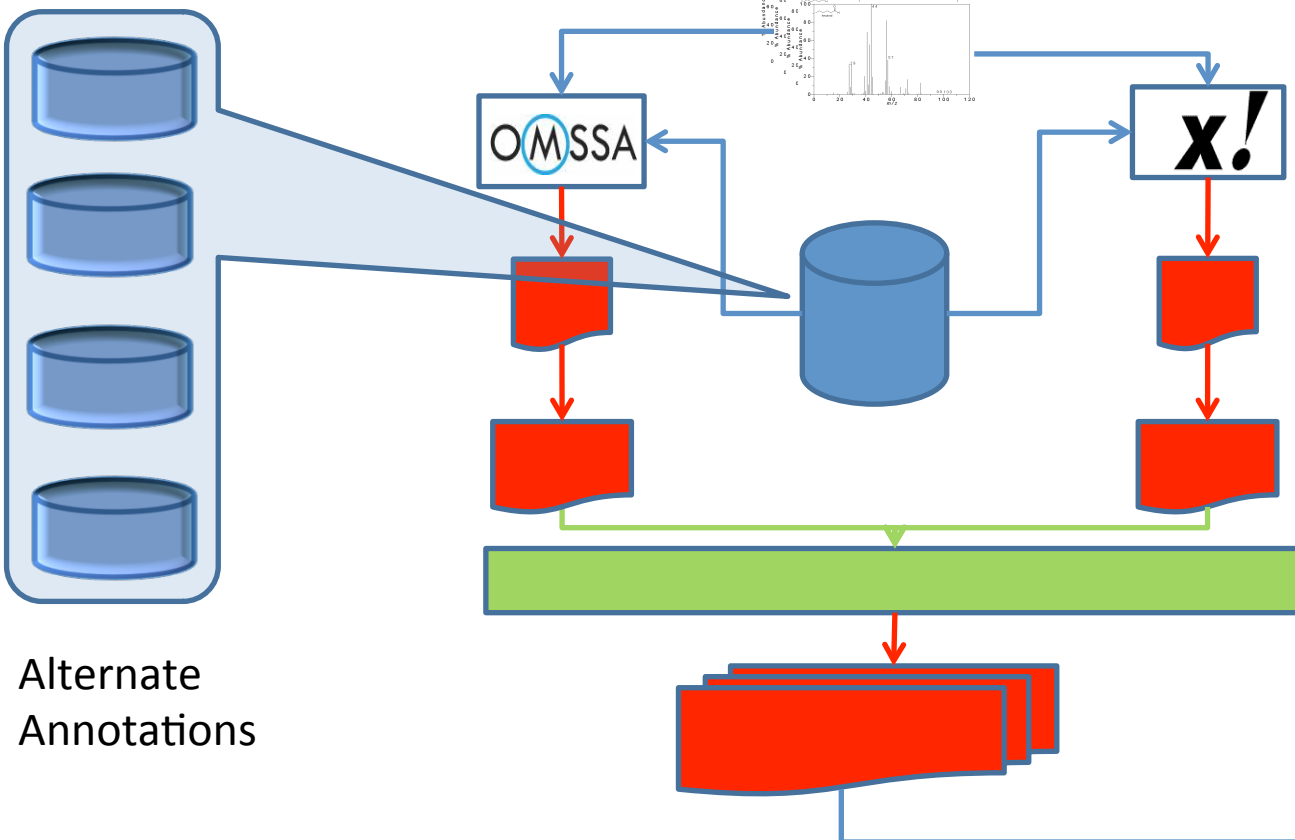
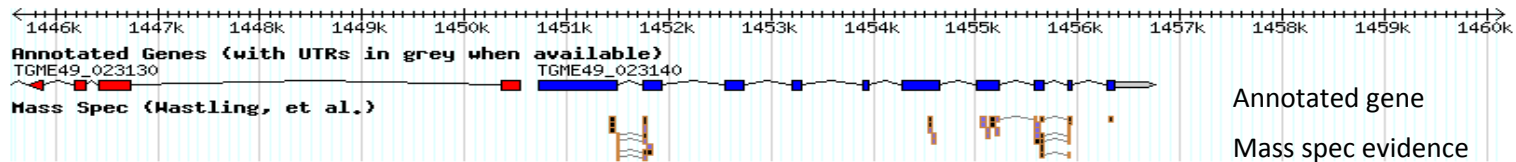
# GFF



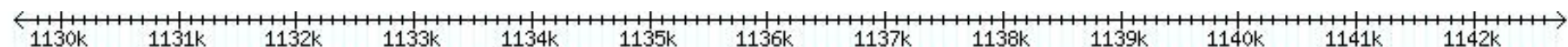
# GFF

Output for <u>each</u> input peak file
Omssa output in native CSV format
Omssa output converted to mzIdentML format
X!Tandem output in native XML format
X!Tandem output converted to mzIdentML format
Summary file reporting the consensus between Omssa and X! Tandem results – Re-ranked peptide-spectrum matches and adjusted FDR scores in a tab delimited text file

Output for <u>complete</u> dataset
Each input peak file has its own output directory with above listed files
A text file summarizing all the identified proteins, respective peptides and relevant statistics
A GFF3 file if operating in the GFF mode







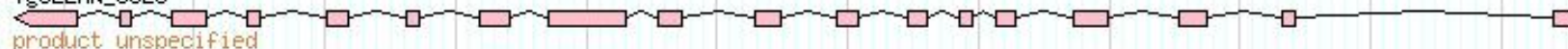
**Annotated Genes (with UTRs in gray when available)**

TGME49\_026960



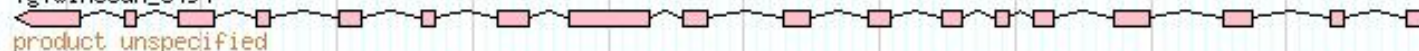
**GLEAN Gene Models**

TgGLEAN\_6626



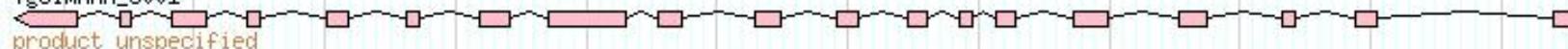
**TwinScan Gene Models**

TgTwinScan\_5494



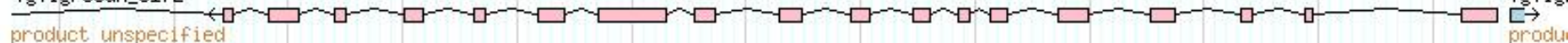
**GlimmerHMM Gene Models**

TgGlmHMM\_3001



**TigrScan Gene Models**

TgTigrScan\_6172



**Mass Spec (Wastling, et al.)**



**Mass Spec (Hu, et al.)**



**Mass Spec (Einstein)**



**Mass Spec (Carruthers, et al.)**



**Mass Spec (Moreno, et al.)**



**Mass Spec (Alternative Models)**



# Current status

## Test organisms



## Developers resources



# Acknowledgements



UNIVERSITY OF  
**LIVERPOOL**

Dr. Andy Jones

Prof. Jonathan Wastling

Dr. Dong Xia



**Penn**  
UNIVERSITY of PENNSYLVANIA

Prof. David Roos

Dr. Brian Brunk

Dr. Omar Harb

Contact e-mail – [Ritesh.Krishna@liverpool.ac.uk](mailto:Ritesh.Krishna@liverpool.ac.uk)