

Search engine that uses mass spectrometry data to identify proteins from primary sequence databases

MOWSE II, involve amino acid sequence and composition qualifiers
Uses indexed molecular weight databases

Drawback

Database had to be built for each new enzyme and for each set of amino acid residue masses.

Difficult for searching proteins in which residues had been chemically or post-translationally modified.

MOWSE III New algorithms
Facility to specify selected MS/MS fragment ion masses as an "ions" qualifier to a peptide mass value
Now available as mascot from www.matrixscience.com

SEQUEST

Developed by Jimmy Eng/John Yates

Correlates uninterpreted tandem mass spectra of peptides with amino acid sequences from protein and nucleotide databases.

SEQUEST will determine the amino acid sequence and thus the protein(s) and organism(s) that correspond to the mass spectrum being analyzed.

ProteinProspector
<http://prospector.ucsf.edu/>

PeptideSearch
<http://www.narrador.embl-heidelberg.de/GroupPages/Homepage.html>

X! Tandem

The Global Proteome Machine, Advanced Search Page (GPM)
X! Tandem open source is software that can match tandem mass spectra with peptide sequences
http://gpm.igm.jhmi.edu/tandem/therpm_tandem_a.html

OMSSA

The Open Mass Spectrometry Search Algorithm [OMSSA] is a faster search engine for identifying MS/MS peptide spectra by searching libraries of known protein sequences.

Classic probability score using an explicit model for matching experimental spectra to sequences

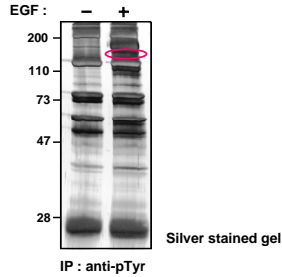
**The Myth of Kozak's Consensus Sequence:
Translation Initiation Codon**

- **CCACCATGG**
- **Most upstream ATG used for translation initiation**
- **Biologists look for this sequence and annotate any ATG near the 5' end of the clone as the initiator methionine**

N-terminal Acetylation

- **Perhaps the most common co-translational modification (60-85% of proteins in yeast)**
- **Usually, aminopeptidases cleave one or two N-terminal amino acids followed by acetylation of the 'mature' protein**
- **So, if you find an N-acetylated peptide, the initiation methionine can be established.**

**MS-Based Identification of a 130 kDa Protein
in the EGF Receptor Signaling Pathway**



**Assignment of the initiator methionine in a
cDNA 'fragment' based on an N-terminal peptide**

>KIAA0229 (1180 residues) FRAGMENT

SWGKREGVVSPAGLGGALPGDGKFGSPSRLGCSLGEGVQRVAALGMGKEQ
ELLRAARTGHLPAVEKLLSGKRLSSGFGGGGGGGGGGGGGGGGGGGGLGS
SSHPLSSLLSMWRGPNVNCVDSTGYTPLHHAALNGHRRSSSSRSQDSAEQD
DQVPEQFSGLLHGSSPVCVGDQPFQLLCTAGQSHPDGSPQQGACHKASM
QLEETGVHAPGASQPSALDQSKRVGYLTGLPTTNSRSHPETLHTASPHPGGA
EEGDRSGAR

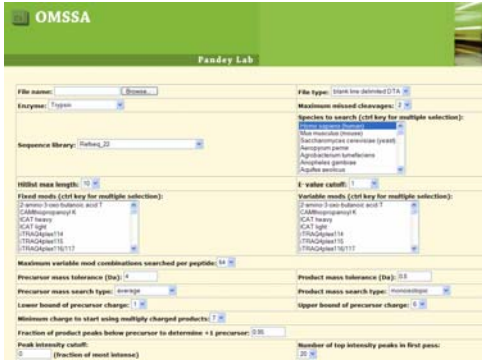
**Assignment of the initiator methionine in a
cDNA 'fragment' based on an N-terminal peptide**



>KIAA0229 (1180 residues) FRAGMENT

SWGKREGVVSPAGLGGALPGDGKFGSPSRLGCSLGEGVQRVAALGMGKEQ
LLRAARTGHLPAVEKLLSGKRLSSGFGGGGGGGGGGGGGGGGGGGGLGS
SHPLSSLLSMWRGPNVNCVDSTGYTPLHHAALNGHRRSSSSRSQDSAEQD
DQVPEQFSGLLHGSSPVCVGDQPFQLLCTAGQSHPDGSPQQGACHKASML
EETGVHAPGASQPSALDQSKRVGYLTGLPTTNSRSHPETLHTASPHPGGAEE
GDRSGAR

Open Mass Spectrometry Search Algorithm [OMSSA]



Protein Databases

- Swiss-Prot
- nr (non-redundant protein database)
- RefSeq
- IPI (International Protein Index)

Swiss-Prot

<http://us.expasy.org/sprot/>

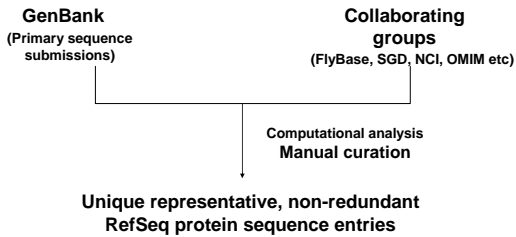
- Swiss-prot is part of the ExPASy (Expert Protein Analysis System) proteomics server of the Swiss Institute of Bioinformatics.
- A highly curated protein sequence database with minimal redundancy
- Swiss-Prot currently contains 172,000 protein sequences representing 8,859 species
- 13,500 Human protein sequences

RefSeq (Reference Sequence) database

<http://www.ncbi.nlm.nih.gov/RefSeq/>

- RefSeq database is a result of collaborative effort of NCBI and other groups and databases like TIGR, FlyBase, WormBase etc.
- A comprehensive, integrated and highly non-redundant curated protein sequence database
- 29,500 Human protein sequences
- Contains protein sequences from all major research organisms
- Alternate splice forms listed individually
- Also contains predicted proteins translated from predicted transcripts (designated as XP_ entries)

RefSeq (Reference Sequence) database

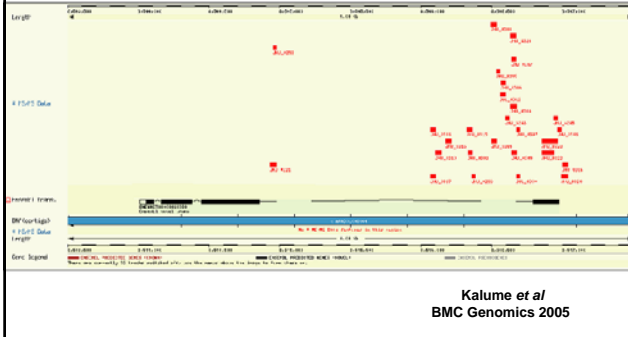


Ensembl database

<http://www.ebi.ac.uk/ensembl/>

- Ensembl is a joint project between the EMBL-EBI and the Wellcome Trust Sanger Institute that aims at developing a system that maintains automatic annotation of large eukaryotic genomes. database is a result of collaborative effort of NCBI and other groups and databases like TIGR, FlyBase, WormBase etc.
- It is a comprehensive source of stable annotation with confirmed gene predictions that have been integrated from external data sources.

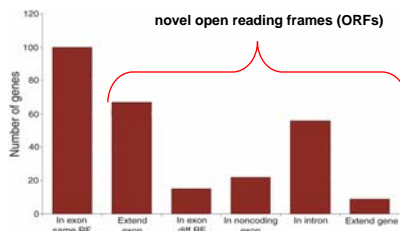
Correction of a Predicted Transcript



Mapping of Open Reading Frames (ORFs)

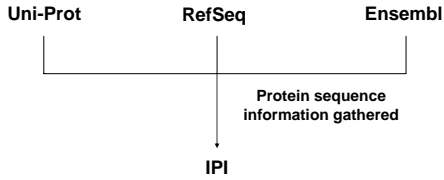


Statistics of ORFs Mapped to Genome



IPI (International Protein Index) database

<http://www.ebi.ac.uk/IPI/>



IPI (International Protein Index) database

- IPI is a protein database from the European Bioinformatics Institute
- Has protein sequence information from Human, Mouse, Rat, Zebra fish and Arabidopsis species only
- 58,000 Human protein sequences
- A redundant database
- Has information on protein isoforms
- The sequence identifiers and sequence entries are not generally stable

EMBL-EBI
European Bioinformatics Institute

International Protein Index

IPI provides a log-level guide to the cross references that describe the annotation of higher eukaryotic organisms. IPI

1. efficiently maintains a database of cross references between the primary data sources
2. provides nonredundant yet maximally complete sets of proteins for biological systems (one sequence per transcript)
3. maintains stable identifiers with chronological numbering to allow the tracking of sequences as IPI between IPI releases

IPI is updated monthly in accordance with the latest data released by the primary data sources.

UniProt
the universal protein resource

Ensembl
Protein and genome annotation database for vertebrate genomes

RefSeq
The Reference Sequence (RefSeq) collection provides a curated, non-redundant set of sequences.

Human Genome Annotation

Genome Annotation by Mass Spectrometry: What Can We Gain?

- Assigning start codons
- Proteins isoforms (alternative splicing, novel exons)
- Novel genes (proteins less than 100 amino acids not predicted by programs)
- cSNPs
- Correction of incorrect gene predictions (50% of the genes in human are predicted)
- Validation of gene predictions

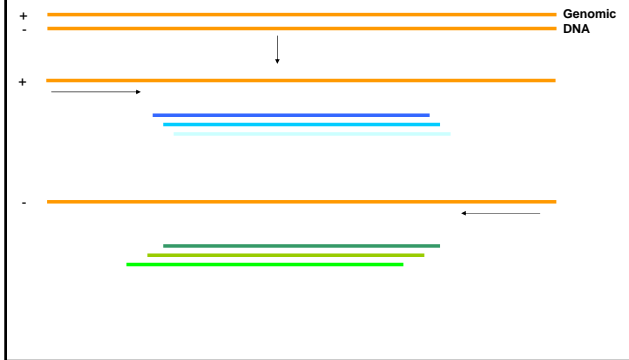
When is a peptide not identified from a database search?

- Protein not described (i.e. novel protein)
- Polymorphisms
- Alternative splice forms
- Novel exon
- Wrong annotation

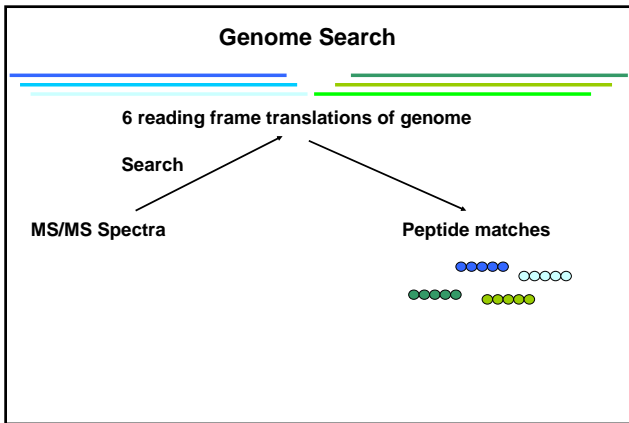
How do you identify such events?

- For novel genes and novel exons use the human genome sequence
- For polymorphisms and alternate splice forms, use a computational strategy

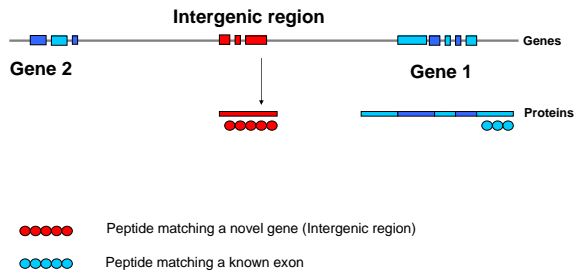
Genome Search



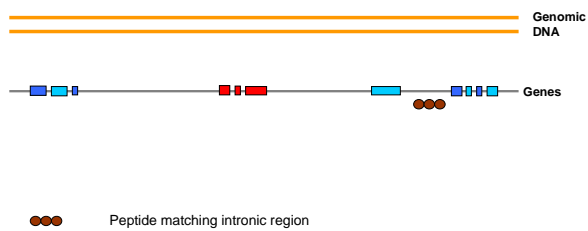
Genome Search



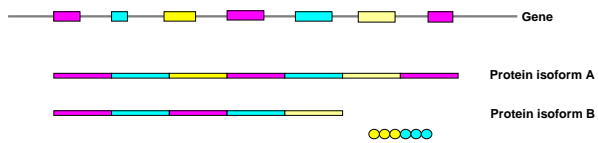
Peptide mapping onto the genome – Identifying a novel gene



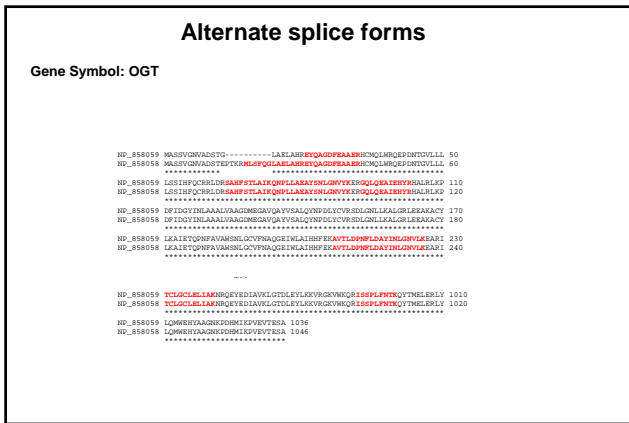
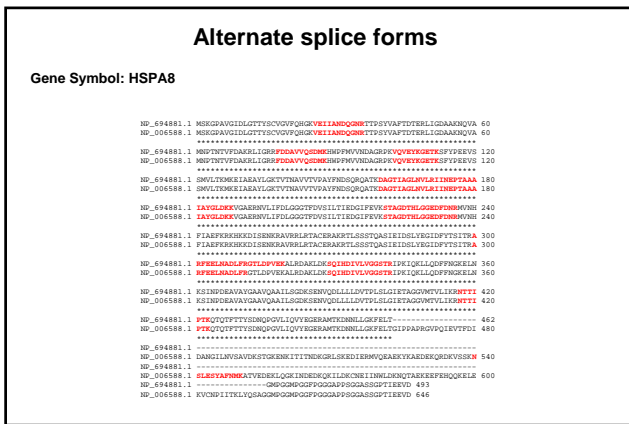
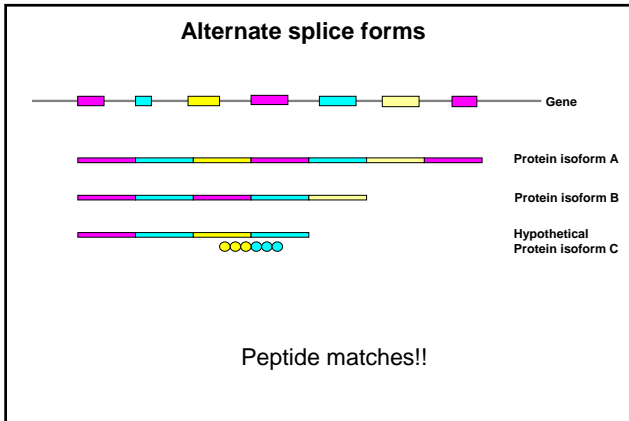
Peptide mapping onto the genome – Identifying a novel exon



Alternate splice forms



No Match!!



Genome Annotation

Genome Annotation by Mass Spectrometry: What Can We Gain?

- Assigning start codons
- Proteins isoforms (alternative splicing, novel exons)
- Novel genes (proteins less than 100 amino acids not predicted by programs)
- cSNPs
- Correction of incorrect gene predictions (~50% of the genes in human are predicted)
- Validation of gene predictions

Use of Ensembl Distributed Annotation System to Validate a Known Transcript

