



# Faster, More Sensitive Peptide ID by Sequence DB Compression

**Nathan Edwards**

**Center for Bioinformatics and Computational Biology**



# MS/MS Search Engines

- Fail when peptides are missing from sequence database
- Protein sequence databases serve many masters
  - Full length protein sequences not needed for MS/MS
  - Explicit variant enumeration is needed for MS/MS
- Much peptide sequence information is lost, inaccessible, or not integrated
  - Protein isoforms, Sequence variants, SNPs, alternate splice forms, ESTs
- Some peptides are more interesting than others
  - Protein identification is only part of the story

# Swiss-Prot

**NiceProt View of Swiss-Prot:**  
**P13746**

[Printer-friendly view](#) [Submit update](#) [Quick BlastP search](#)

[\[Entry info\]](#) [\[Name and origin\]](#) [\[References\]](#) [\[Comments\]](#) [\[Cross-references\]](#) [\[Keywords\]](#) [\[Features\]](#) [\[Sequence\]](#) [\[Tools\]](#)

*Note: most headings are clickable, even if they don't appear as links. They link to the [user manual](#) or [other documents](#).*

Entry information	
Entry name	1A11_HUMAN
Primary accession number	P13746
Secondary accession numbers	O19605 O19606 Q29747 Q29835 Q9BCN0 Q9MYI5 Q9TQE9 Q9TQP6 Q9TQP7
Entered in Swiss-Prot in	Release 13, January 1990
Sequence was last modified in	Release 13, January 1990
Annotations were last modified in	Release 42, October 2003
Name and origin of the protein	
Protein name	HLA class I histocompatibility antigen, A-11 alpha chain [Precursor]
Synonym	MHC class I antigen A*11
Gene name	HLA-A or HLAA
From	<a href="#">Homo sapiens (Human)</a> [TaxID: 9606]
Taxonomy	<a href="#">Eukaryota</a> ; <a href="#">Metazoa</a> ; <a href="#">Chordata</a> ; <a href="#">Craniata</a> ; <a href="#">Vertebrata</a> ; <a href="#">Euteleostomi</a> ; <a href="#">Mammalia</a> ; <a href="#">Eutheria</a> ; <a href="#">Primates</a> ; <a href="#">Catarrhini</a> ; <a href="#">Hominidae</a> ; <a href="#">Homo</a> .

# Swiss-Prot Variant Annotations

**Comments**

- **FUNCTION:** Involved in the presentation of foreign antigens to the immune system.
- **SUBUNIT:** Heterodimer of an alpha chain and a beta chain (beta-2-microglobulin).
- **SUBCELLULAR LOCATION:** Type I membrane protein.
- **ALTERNATIVE PRODUCTS:**
  - Alternative splicing [2 named forms] [Display all isoform sequences in Fasta format](#)

Name	1
Isoform ID	P13746-1
This is the isoform sequence <a href="#">displayed in this entry</a> .	

Name	2
Synonyms	Long
Isoform ID	<a href="#">P13746-2</a>
<i>Note:</i> Only produced by allele A*1103.	
Features which should be applied to build the isoform sequence: <a href="#">VSP_008099</a> .	

- **POLYMORPHISM:** The following alleles of A-11 are known: A\*1101 (A-11E), A\*1102 (A-11K), A\*1103, A\*1104, A\*1105 and A\*1107. The sequence shown is that of A\*1101.

**Copyright**

This Swiss-Prot entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL outstation - the European Bioinformatics Institute. There are no restrictions on its use by non-profit institutions as long as its content is in no way modified and this statement is not removed. Usage by and for commercial entities requires a license agreement (See <http://www.isb-sib.ch/announce/> or send an email to [license@isb-sib.ch](mailto:license@isb-sib.ch))

**Cross-references**

X13111; CAA31503.1; -	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>
X13112; CAA31504.1; -	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>
D16841; BAA04117.1; -	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>
D16842; BAA04118.1; -	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>
M16010; AAA65449.1; -	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>
M16007; AAA65449.1; JOINED	<a href="#">[EMBL / GenBank / DDBJ]</a> <a href="#">[CoDingSequence]</a>

# Swiss-Prot Variant Annotations

NiceProt View of Swiss-Prot: P13746 - Microsoft Internet Explorer

Address: [http://us.expasy.org/cgi-bin/niceprot.pl?1A11\\_HUMAN](http://us.expasy.org/cgi-bin/niceprot.pl?1A11_HUMAN)

Features

[Feature table viewer](#) [Feature aligner](#)

Key	From	To	Length	Description	FTId
SIGNAL	<a href="#">1</a>	<a href="#">24</a>	24		
CHAIN	<a href="#">25</a>	<a href="#">365</a>	341	HLA class I histocompatibility antigen, A-11 alpha chain.	
DOMAIN	<a href="#">25</a>	<a href="#">114</a>	90	Extracellular alpha-1.	
DOMAIN	<a href="#">115</a>	<a href="#">206</a>	92	Extracellular alpha-2.	
DOMAIN	<a href="#">207</a>	<a href="#">298</a>	92	Extracellular alpha-3.	
DOMAIN	<a href="#">299</a>	<a href="#">308</a>	10	Connecting peptide.	
TRANSMEM	<a href="#">309</a>	<a href="#">332</a>	24		
DOMAIN	<a href="#">333</a>	<a href="#">365</a>	33	Cytoplasmic tail.	
CARBOHYD	<a href="#">110</a>	<a href="#">110</a>		N-linked (GlcNAc...) (By similarity).	
DISULFID	<a href="#">125</a>	<a href="#">188</a>		By similarity.	
DISULFID	<a href="#">227</a>	<a href="#">283</a>		By similarity.	
VARSPPLIC	<a href="#">337</a>	<a href="#">337</a>		S -> SGGEGVK (in <a href="#">isoform 2</a> ).	VSP_008099
VARIANT	<a href="#">43</a>	<a href="#">43</a>	*	E -> K (in allele A*1102).	<a href="#">VAR_004353</a>
VARIANT	<a href="#">133</a>	<a href="#">133</a>	*	F -> L (in allele A*1107).	<a href="#">VAR_016731</a>
VARIANT	<a href="#">168</a>	<a href="#">168</a>	*	K -> E (in allele A*1105).	<a href="#">VAR_016732</a>
VARIANT	<a href="#">175</a>	<a href="#">175</a>	*	H -> R (in allele A*1103).	<a href="#">VAR_016733</a>
VARIANT	<a href="#">176</a>	<a href="#">176</a>	*	A -> E (in allele A*1103).	<a href="#">VAR_016734</a>
VARIANT	<a href="#">187</a>	<a href="#">187</a>	*	R -> T (in allele A*1104).	<a href="#">VAR_016735</a>
VARIANT	<a href="#">345</a>	<a href="#">345</a>	*	T -> S (in allele A*1105).	<a href="#">VAR_016736</a>

Sequence information

Length: 365 AA [This is the length of the unprocessed precursor]      Molecular weight: 40937 Da [This is the MW of the unprocessed precursor]      CRC64: FE449CE2D4BF6CC5 [This is a checksum on the sequence]

10      20      30      40      50      60

MANHARDTLLIIGCALALTECHLACGSRVYVYCHSRDRCORREELANGVIRDTQSFRE

# Swiss-Prot Sequence

**Sequence information**

Length: 365 AA [This is the length of the unprocessed precursor]      Molecular weight: 40937 Da [This is the MW of the unprocessed precursor]      CRC64: FE449CE2D4BF6CC5 [This is a checksum on the sequence]

```

      10      20      30      40      50      60
      |      |      |      |      |      |
MAVMAPRTLL LLLSGALALT QTWAGSHSMR YFYTSVSRPG RGEPRFIAVG YVDDTQFVRF

      70      80      90     100     110     120
      |      |      |      |      |      |
DSDAASQRME PRAPWIEQEG PEYWDQETRN VKAQSQTRV DLGTLRGYYN QSEDEGSHTIQ

      130     140     150     160     170     180
      |      |      |      |      |      |
IMYGCDVGPD GRFLRGYRQD AYDGKDYIAL NEDLRSWTA DMAAQITKRK WEAHAHAEQQ

      190     200     210     220     230     240
      |      |      |      |      |      |
RAYLEGRCVE WLRRYLENGK ETLQRTDPPK THMTHHPISD HEATLRCWAL GFYPAEITLT

      250     260     270     280     290     300
      |      |      |      |      |      |
WQRDGEDQTQ DTELVETRPA GDGTFQKWA VVVPSEGEQR YTCHVQHEGL PKPLTLRWEL

      310     320     330     340     350     360
      |      |      |      |      |      |
SSQPTPIVIG ILAGLVLLGA VITGAVVA AV MWRKSSDRK GGSYTQAASS DSAQGSQVSL

TACKV

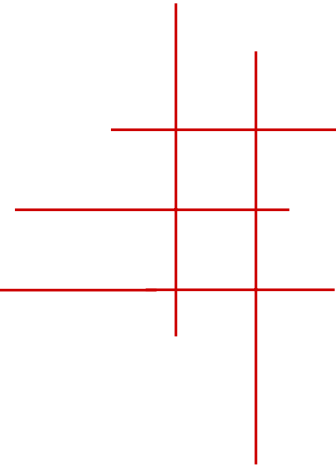
```

P13746 in [FASTA format](#)

[View entry in original Swiss-Prot format](#)  
[View entry in raw text format \(no links\)](#)

# Swiss-Prot

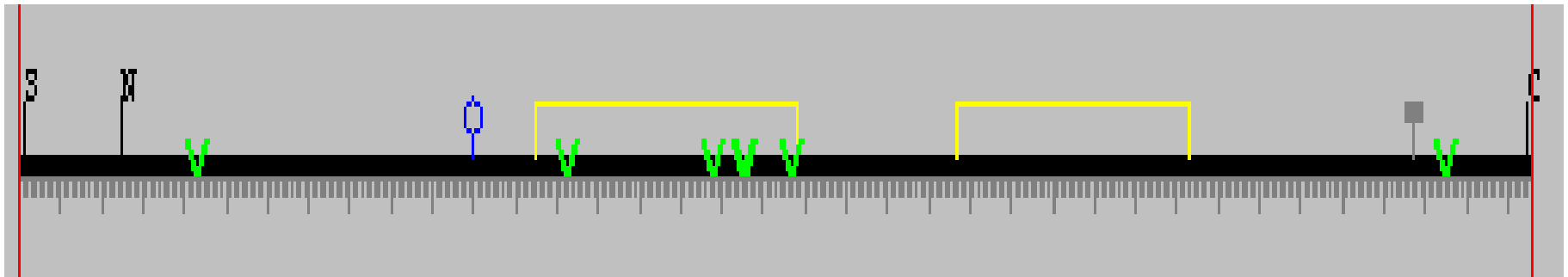
---



- VarSplic enumerates all variants, conflicts, isoforms
- Swiss-Prot sequence size:
  - 60 Mb
- VarSplic sequence size:
  - 95 Mb
- How many more peptide candidates?

# Swiss-Prot Variant Annotations

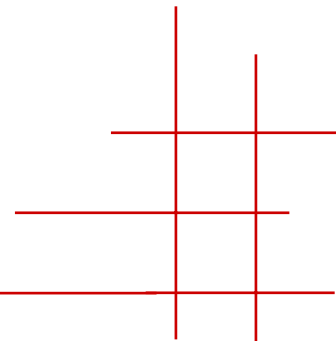
Feature viewer



Variants

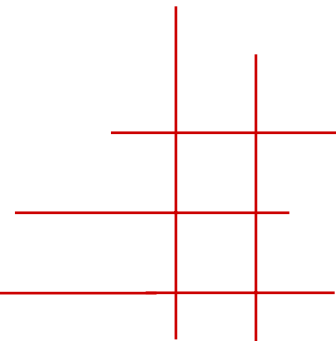


# Swiss-Prot VarSplic Output



```
P13746-00-01-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-01-01-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-00-00-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-00-03-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-01-03-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-00-04-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G K P R F I A V G Y V D D T Q F V R F
P13746-01-04-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G K P R F I A V G Y V D D T Q F V R F
P13746-00-05-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-01-05-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-01-00-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-00-02-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
P13746-01-02-00      MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F Y T S V S R P G R G E P R F I A V G Y V D D T Q F V R F
***** : *****
```

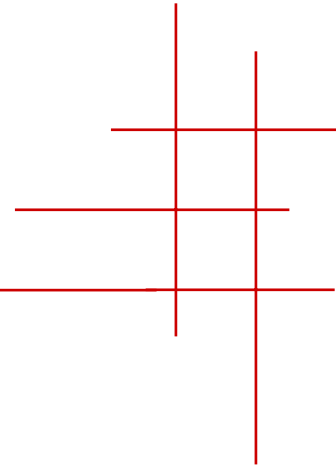
# Swiss-Prot VarSplic Output



```
P13746-00-01-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYTQAASSDSAQ
P13746-01-01-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYTQAASSDSAQ
P13746-00-00-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYTQAASSDSAQ
P13746-00-03-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYTQAASSDSAQ
P13746-01-03-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYTQAASSDSAQ
P13746-00-04-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYTQAASSDSAQ
P13746-01-04-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYTQAASSDSAQ
P13746-00-05-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYTQAASSDSAQ
P13746-01-05-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYTQAASSDSAQ
P13746-01-00-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYTQAASSDSAQ
P13746-00-02-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSS-----DRKGGSYSQAASSDSAQ
P13746-01-02-00      SSQPTIPIVGIIAGLVLLGAVITGAVVAAVMWRKSSGGEGVKDRKGGSYSQAASSDSAQ
*****                *****:*****
```

# Peptide Candidates

---



- Parent ion
  - Typically < 3000 Da
- Tryptic Peptides
  - Cut at K or R
- Search engines
  - Don't handle > 4+ well
  - Long peptides don't fragment well
- # of distinct 30-mers upper bounds total peptide content

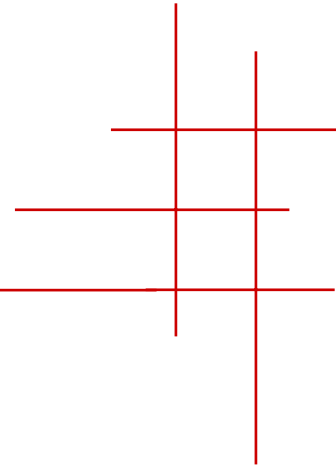
# Peptide Candidates

- At most 1% additional peptides in ~ 1.6 times as much sequence

Sequence Database	Swiss-Prot	VarSplic
Size	60 Mb	95 Mb
30-mers ( $N_{30}$ )	46 Mb	47 Mb
Overhead	29%	101%

# Sequence Database Compression

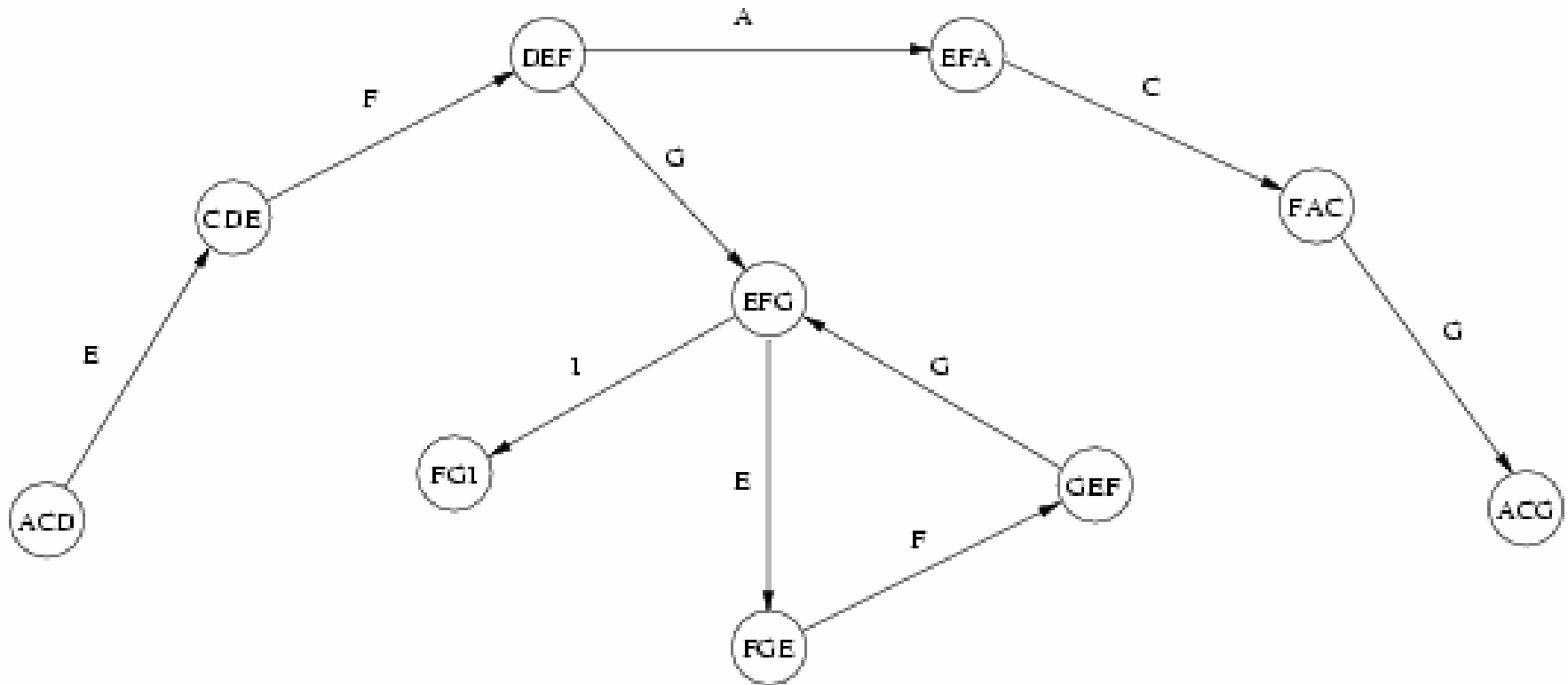
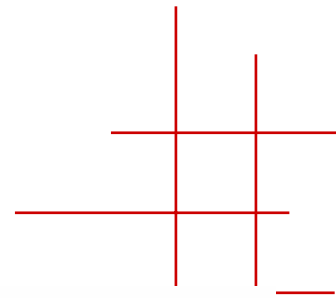
---



Construct sequence database that is

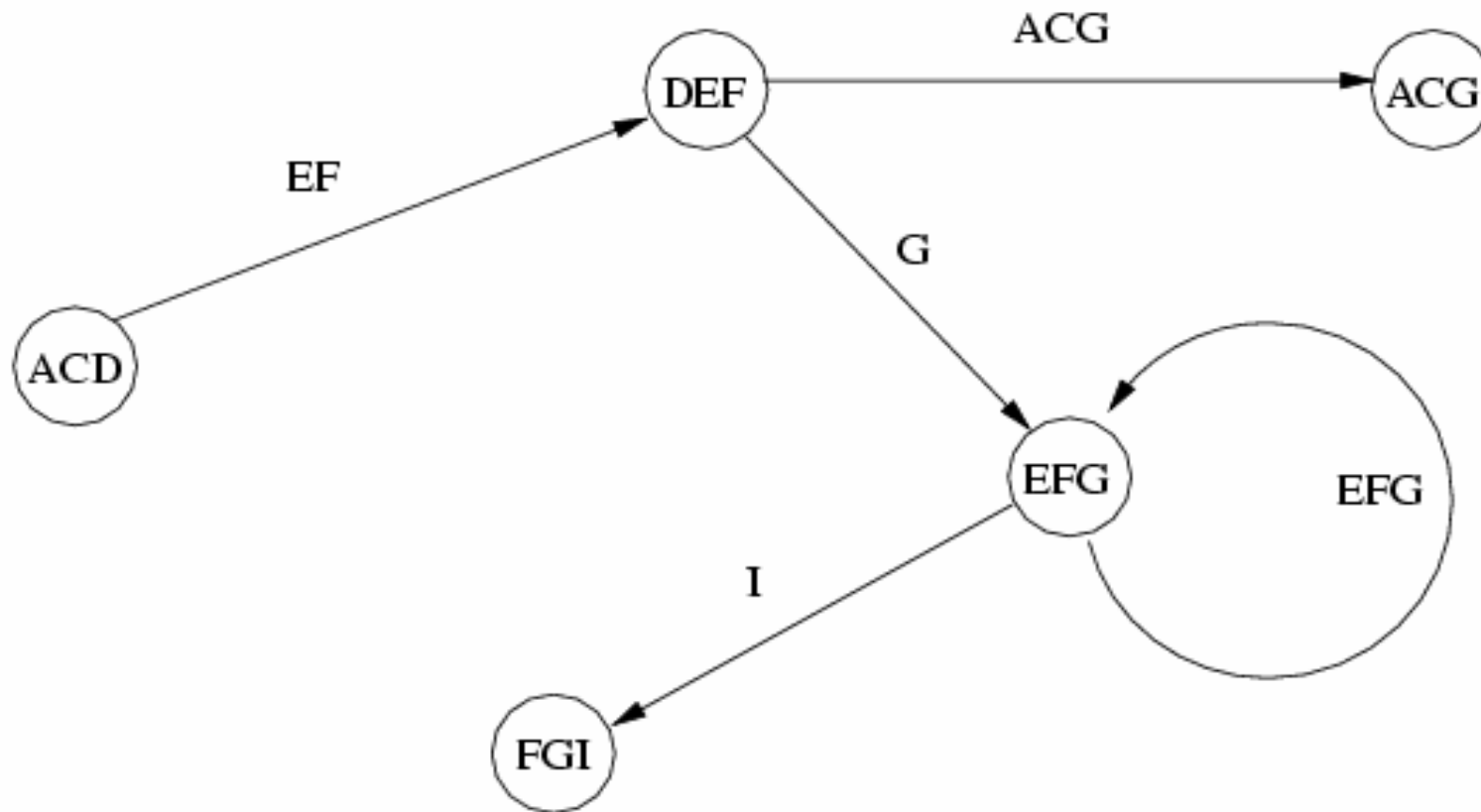
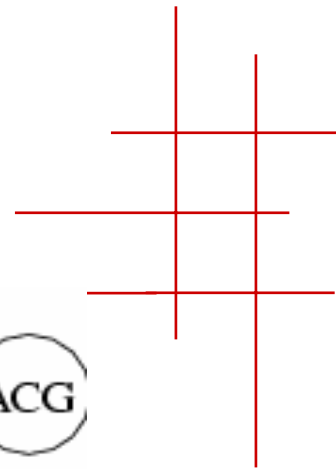
- **Complete**
  - All 30-mers are present
- **Correct**
  - No other 30-mers are present
- **Compact**
  - No 30-mer is present more than once

# SBH-graph



ACDEFGI, ACDEFACG, DEFGEFGI

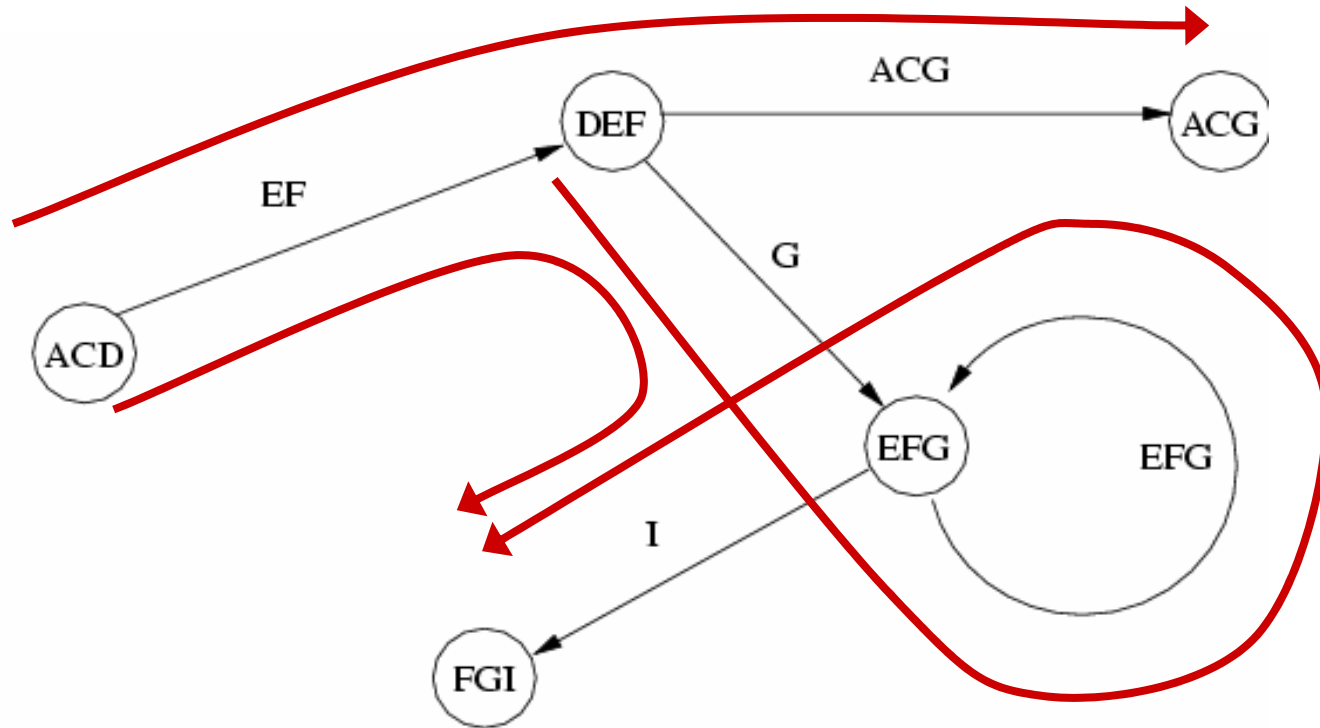
# Compressed SBH-graph



ACDEFGI, ACDEFACG, DEFGEFGI

# Sequence Databases & CSBH-graphs

- Sequences correspond to paths

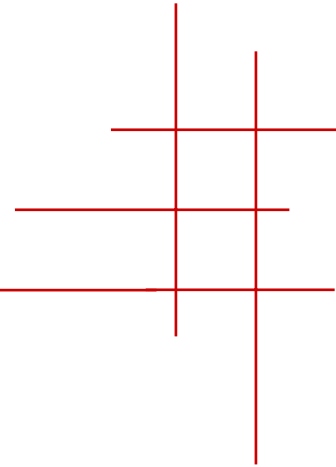


ACDEFGI, ACDEFACG, DEFGEFGI



# Sequence Databases & CSBH-graphs

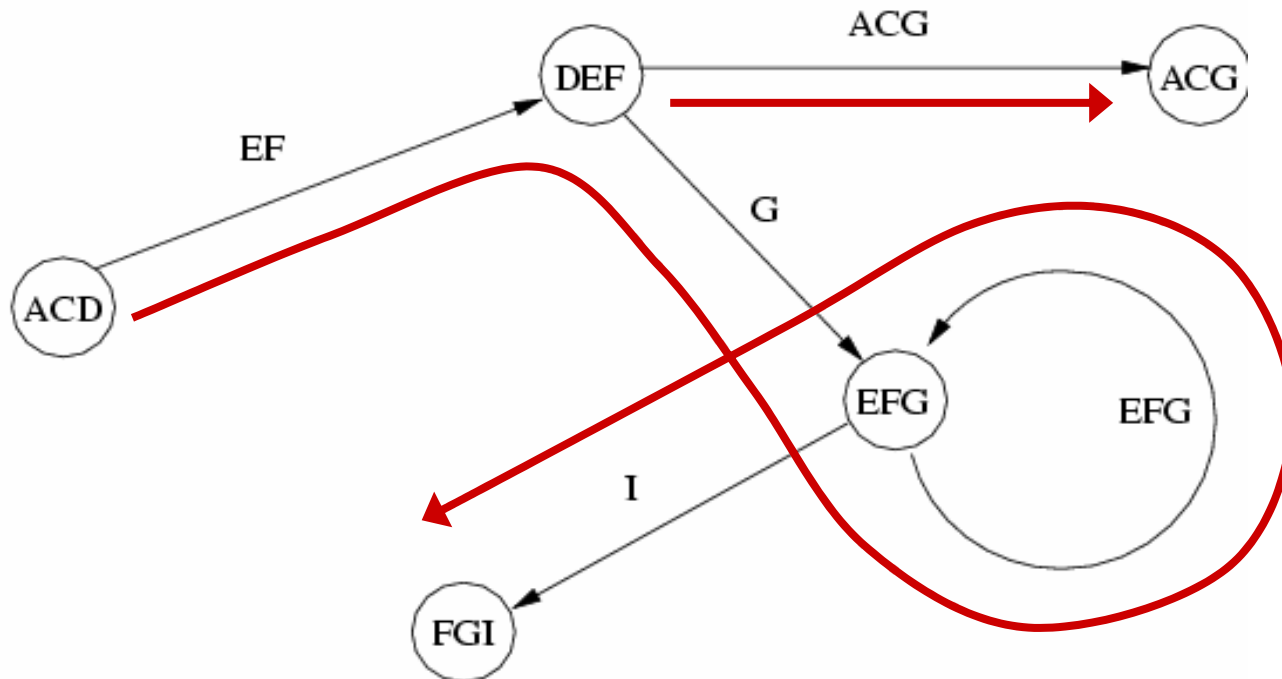
---



- **Complete**
  - All edges are on some path
- **Correct**
  - Output path sequence only
- **Compact**
  - No edge is used more than once
- **C<sup>3</sup> Path Set** uses all edges exactly once.

# Sequence Databases & CSBH-graphs

- Use each edge exactly once



ACDEFGGEFGI, DEFACG

# AA Sequence Databases

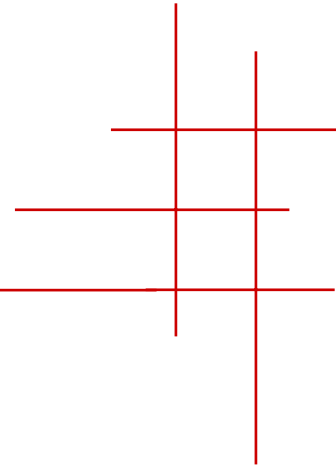
Sequence Database	Sequence Length	Distinct 30-mers	Overhead
IPI-HUMAN	20358846	12115520	68%
IPI	54145883	29769766	81%
Swiss-Prot	56454588	44374286	27%
Swiss-Prot-VS	89541275	45307827	97%
UniProt	472581860	274510105	72%
UniProt-VS	506796094	275391669	84%
MSDB	481919777	276523755	74%
NRP	495502241	283160529	75%
NCBI-nr	619132252	378721915	63%
UnionNR	674700840	385369671	75%
Union	2157353500	385369671	460%

# Minimum Size C<sup>3</sup> Sequence Database

Sequence Database	C <sup>3</sup> 30-mer Enumeration	Overhead	Compression	Compression Bound
IPI-HUMAN	13854679	14.35%	68.05%	59.51%
IPI	37961385	27.52%	70.11%	54.98%
Swiss-Prot	52662145	18.68%	93.28%	78.60%
Swiss-Prot-VS	54534356	20.36%	60.90%	50.60%
UniProt	337119564	22.81%	71.34%	58.09%
UniProt-VS	338890778	23.06%	66.87%	54.34%
MSDB	342924164	24.01%	71.16%	57.38%
NRP	351600578	24.17%	70.96%	57.15%
NCBI-nr	463517034	22.39%	74.87%	61.17%
UnionNR	473665310	22.91%	70.20%	57.12%
Union	473665310	22.91%	21.96%	17.86%

# Implementation

---



- Suitable for use by Mascot, SEQUEST, X-Tandem, etc...
  - FASTA format
  - Use special “restart symbol”!
- All connection to protein context is lost
  - Must do exact string search to find peptides in original database

# Implementation

>1  
ACDEFGI

>2  
ACDEFACG

>3  
DEFGEFGI

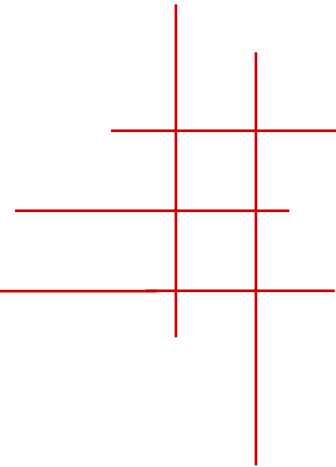


>A  
ACDEFGEFGI

>B  
DEFACG

- Can't guarantee that "restarts" won't create "false" tryptic peptides

# Implementation



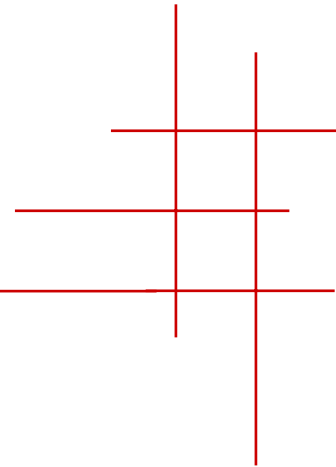
- Mark original entries' start and end
- “Restart” with unusable symbol
  - **Weight(J) = 10kDa**
- Forbid tryptic digestion at J
  - **Mascot cofiguration:**

**Title:Trypsin**  
**Cleavage:KR**  
**Restrict:PJ**  
**Cterm**

- What about “short” proteins?
- **C<sup>3</sup> guarantee for tryptic peptides of length  $\leq (k-2)$**

# Faster Peptide ID

---



## Peptide identification workflow

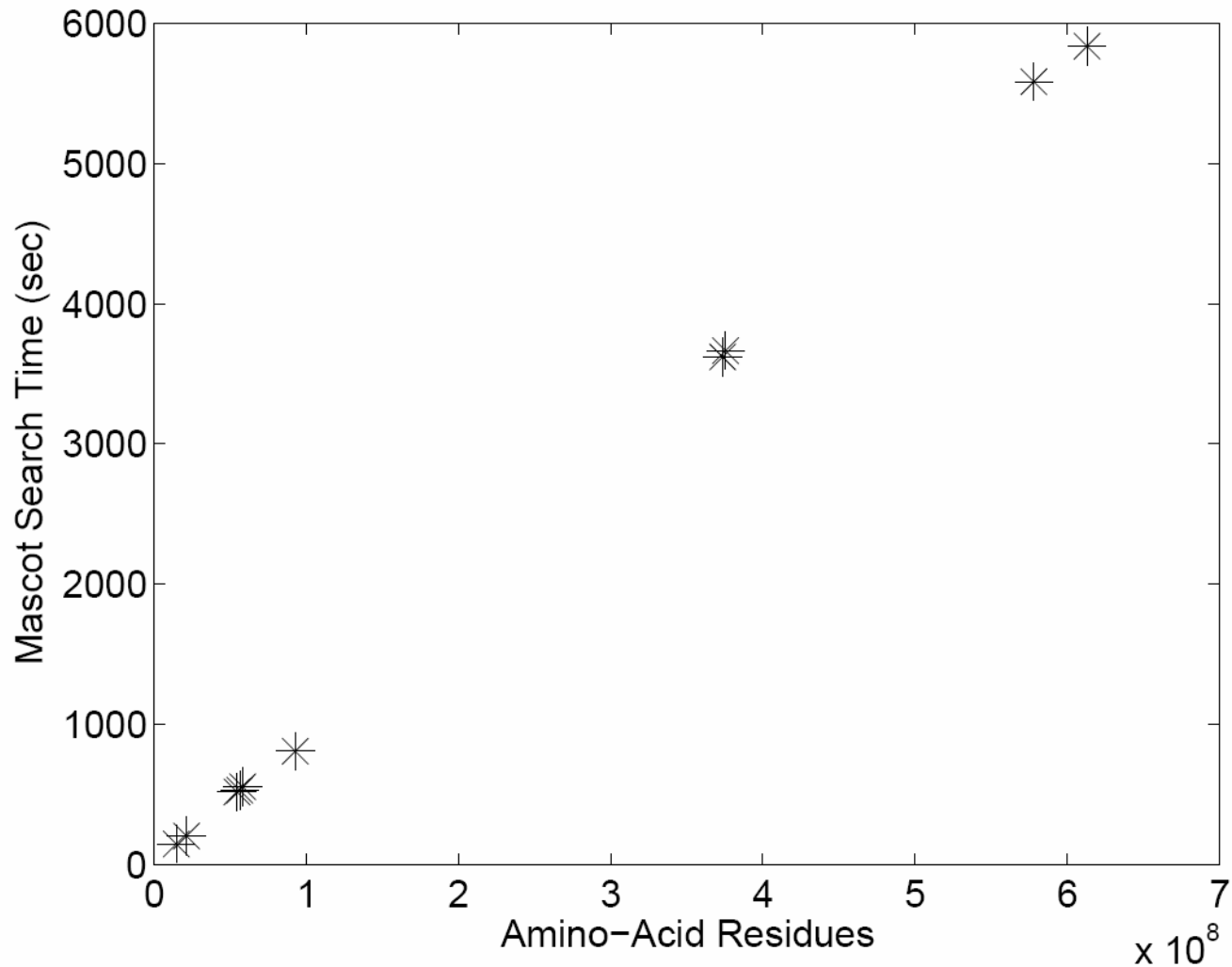
- Search spectra against C<sup>3</sup> sequence
- Insert protein context using exact string search
  - Parse search output
  - Search original sequence using exact peptide sequence
  - Re-construct search output



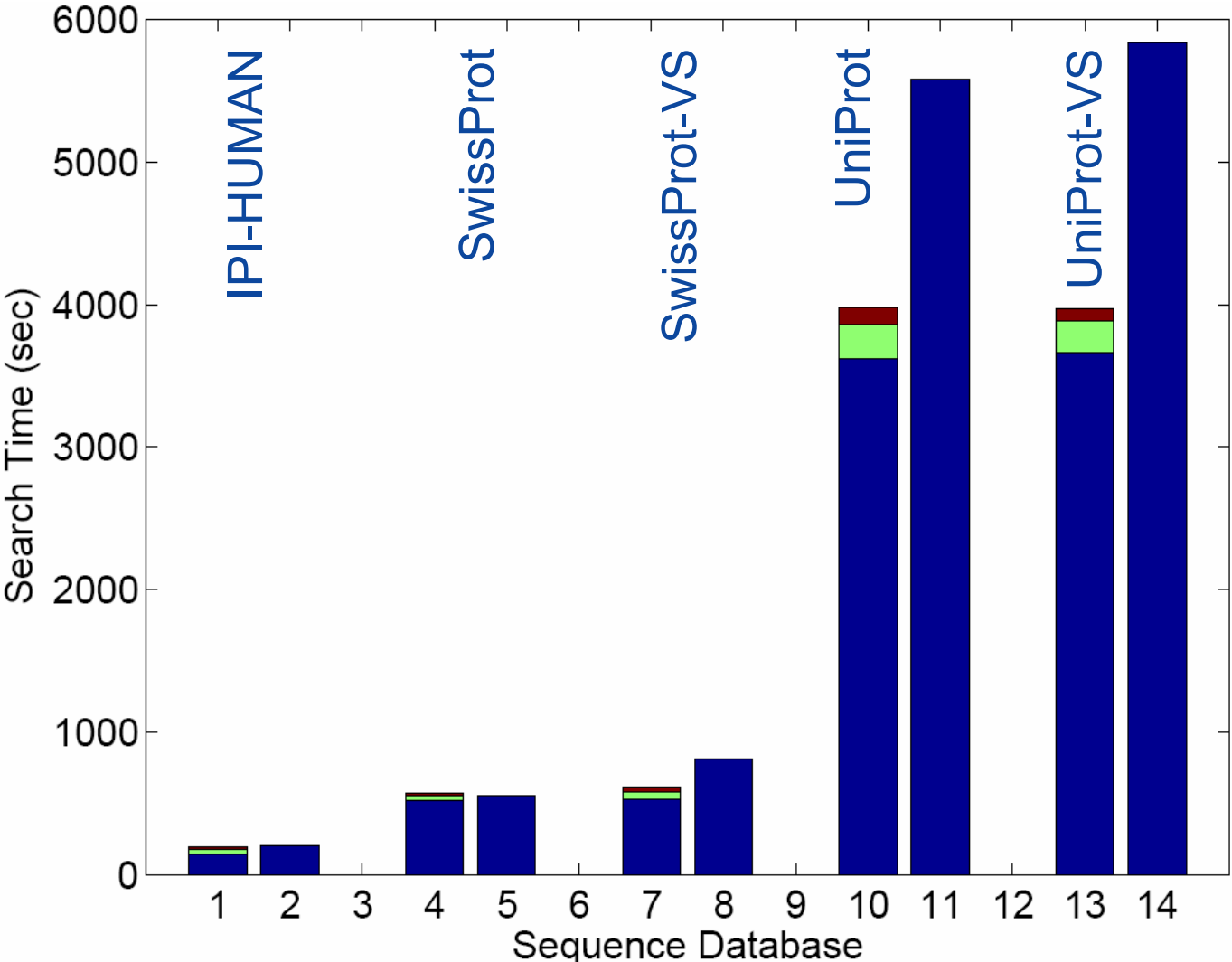
# Faster Peptide ID

- ISB 17 Protein Mix
  - LC/MS/MS, 2023 MS/MS spectra
- Mascot v2.0
  - Dell PC w/ 512 Mb RAM
  - Parent Tolerance 2Da
  - Fragment Tolerance 0.15Da
  - Up to 2 Missed Cleavages
- IPI-HUMAN, SwissProt, UniProt
  - plus SwissProt-VS UniProt-VS versions

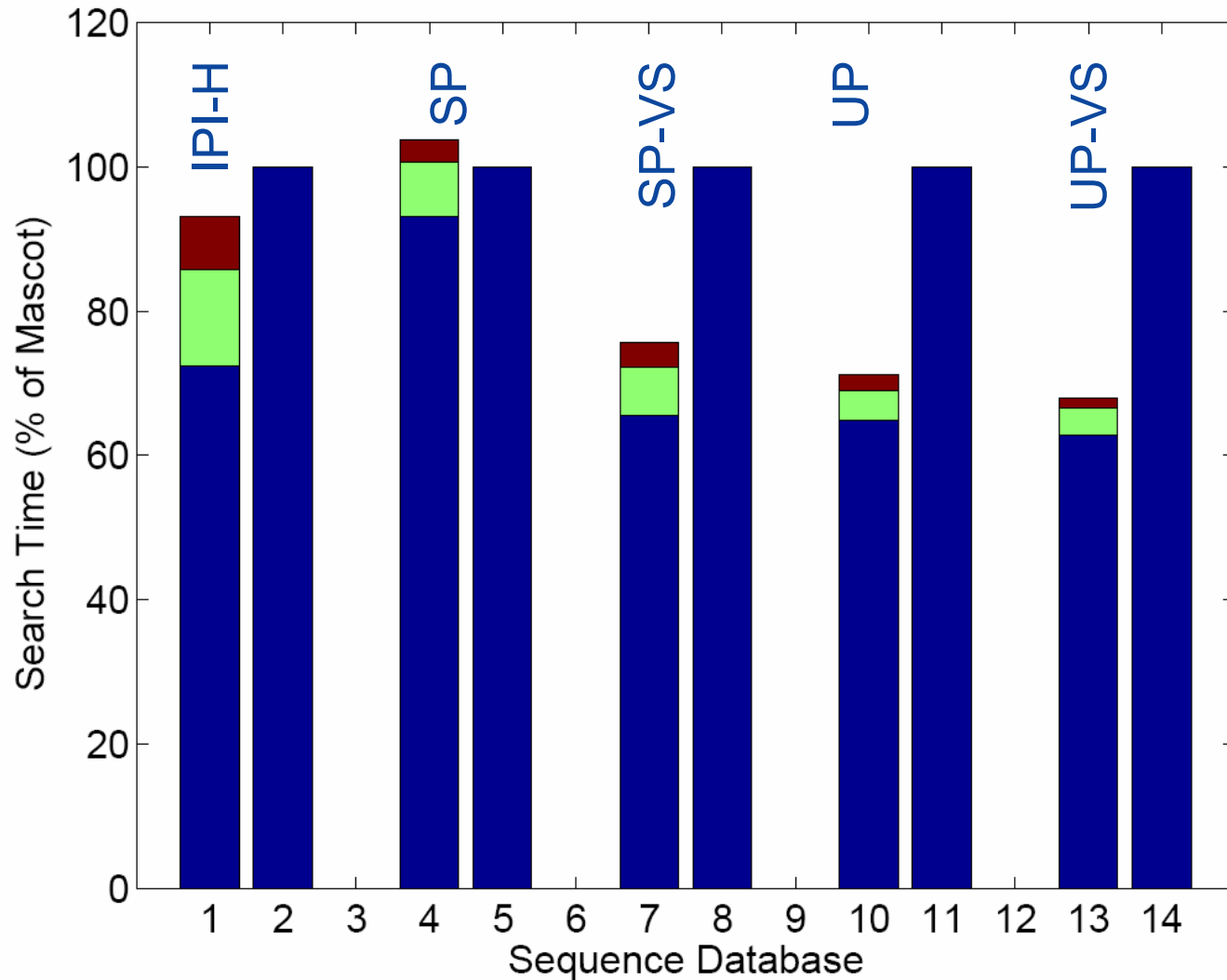
# Mascot Running Time



# Total Search Time



# Relative Search Time





# Mascot Results: C<sup>3</sup> DB

## Peptide Summary Report

[Switch to Protein Summary Report](#)

To create a bookmark for this report, right click this link: [Peptide Summary Report \( Conversion of 17mix\\_test2.xml to mascot\\_generic by msxml2other\)](#)

Select All

Select None

Search Selected

Error tolerant

Archive Report

1. [AAAAAAHJAX](#) Mass: 116278 Score: 762 Peptides matched: 38

AAAAAAHJAX

Check to include this hit in error tolerant search or archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Expect	Rank	Peptide
<input checked="" type="checkbox"/> <a href="#">587</a>	435.70	869.39	869.42	-0.03	0	(24)	4.3	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">588</a>	435.70	869.39	869.42	-0.03	0	25	3.6	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">589</a>	435.70	869.39	869.42	-0.03	0	(22)	6.2	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">641</a>	450.67	899.32	899.38	-0.06	0	39	0.19	1	FNDDFS
<input checked="" type="checkbox"/> <a href="#">764</a>	495.22	988.42	988.48	-0.06	0	39	0.2	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">765</a>	495.22	988.42	988.48	-0.06	0	(36)	0.36	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">767</a>	495.22	988.42	988.48	-0.06	0	(32)	0.94	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">769</a>	495.22	988.42	988.48	-0.06	0	(38)	0.25	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">845</a>	534.23	1066.44	1066.48	-0.05	0	31	0.72	1	WVGYGQDSR
<input checked="" type="checkbox"/> <a href="#">873</a>	550.23	1098.45	1098.55	-0.10	0	(72)	6.5e-005	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">874</a>	550.23	1098.45	1098.55	-0.10	0	73	5.9e-005	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">1020</a>	626.80	1251.58	1251.65	-0.07	0	36	0.24	1	LAHPFFASWR
<input checked="" type="checkbox"/> <a href="#">1021</a>	626.80	1251.58	1251.65	-0.07	0	(33)	0.45	1	LAHPFFASWR
<input checked="" type="checkbox"/> <a href="#">1037</a>	633.26	1264.51	1264.61	-0.10	0	(24)	4.5	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1038</a>	633.26	1264.51	1264.61	-0.10	0	37	0.23	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1072</a>	650.28	1298.55	1298.62	-0.07	0	47	0.023	1	ELNYGPHQWR
<input checked="" type="checkbox"/> <a href="#">1132</a>	671.31	1340.60	1340.66	-0.06	0	(63)	0.00049	1	VDEDQFFPAVPK
<input checked="" type="checkbox"/> <a href="#">1133</a>	671.31	1340.60	1340.66	-0.06	0	(57)	0.0023	1	VDEDQFFPAVPK
<input checked="" type="checkbox"/> <a href="#">1134</a>	671.31	1340.60	1340.66	-0.06	0	70	0.00011	1	VDEDQFFPAVPK
<input checked="" type="checkbox"/> <a href="#">1135</a>	671.31	1340.60	1340.66	-0.06	0	(38)	0.15	1	VDEDQFFPAVPK
<input checked="" type="checkbox"/> <a href="#">1172</a>	681.33	1360.65	1360.71	-0.06	0	38	0.16	1	LWSAEIPNLYR
<input checked="" type="checkbox"/> <a href="#">1173</a>	681.33	1360.65	1360.71	-0.06	0	(29)	1.5	1	LWSAEIPNLYR

# Mascot Results: C<sup>3</sup> DB ++

## *{MATRIX}* Mascot Search Results *{SCIENCE}*

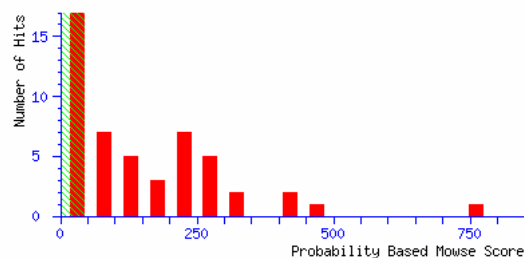
User : Search for Compression Paper  
Email :  
Search title : Conversion of 17mix\_test2.xml to mascot generic by msxml2other  
MS data file : C:\Documents and Settings\edwardnj\Desktop\17prot\_mix.mgf  
Database : Varsplic (203425 sequences; 92944773 residues)  
Timestamp : 5 Oct 2004 at 14:09:37 GMT  
Significant hits: [P00722](#) (BGAL\_ECOLI) Beta-galactosidase (EC 3.2.1.23) (Lactase)  
[P00634-00-00-00](#) (PPB\_ECOLI) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from P00634 Alkaline phosphatase precursor (  
[P00489](#) (PHS2\_RABIT) Glycogen phosphorylase, muscle form (EC 2.4.1.1) (Myophosphorylase)  
[P00921-00-01-00](#) (CAH2\_BOVIN) Splice isoform Displayed; Variant one of the major forms; Conflict Displayed; from P00921 Carbonic anhydrase  
[Q29443](#) (TRFE\_BOVIN) Serotransferrin precursor (Transferrin) (Siderophilin) (Beta-1-metal binding globulin)  
[P02769-00-00-00](#) (ALBU\_BOVIN) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from P02769 Serum albumin precursor (Allerg  
[P21948](#) (PPB\_ESCFE) Alkaline phosphatase precursor (EC 3.1.3.1) (APase)  
[Q28641](#) (MYH4\_RABIT) Myosin heavy chain, skeletal muscle, juvenile  
[P02608](#) (MLRS\_RABIT) Myosin regulatory light chain 2, skeletal muscle isoform type 2 (G2) (DTNB) (MLC-2)  
[P11217-00-05-00](#) (PHS2\_HUMAN) Splice isoform Displayed; Variant GSD-V-VAR\_014005; Conflict Displayed; from P11217 Glycogen phosphorylase, :  
[P11217-00-12-00](#) (PHS2\_HUMAN) Splice isoform Displayed; Variant GSD-V-VAR\_014009; Conflict Displayed; from P11217 Glycogen phosphorylase, :  
[P49064](#) (ALBU\_FELCA) Serum albumin precursor (Allergen Fel d 2)  
[P12882](#) (MYH1\_HUMAN) Myosin heavy chain, skeletal muscle, adult 1 (Myosin heavy chain IIx/d) (MyHC-IIx/d)  
[P07724](#) (ALBU\_MOUSE) Serum albumin precursor  
[Q9Y623](#) (MYH4\_HUMAN) Myosin heavy chain, skeletal muscle, fetal (Myosin heavy chain IIb) (MyHC-IIb)  
[P13538](#) (MYSS\_CHICK) Myosin heavy chain, skeletal muscle, adult  
[P11055](#) (MYH3\_HUMAN) Myosin heavy chain, fast skeletal muscle, embryonic (Muscle embryonic myosin heavy chain) (SMHCE)  
[P00432](#) (CATA\_BOVIN) Catalase (EC 1.11.1.6)  
[P06278](#) (AMY\_BACLI) Alpha-amylase precursor (EC 3.2.1.1) (1,4-alpha-D-glucan glucanohydrolase) (BLA)  
[P09812](#) (PHS2\_RAT) Glycogen phosphorylase, muscle form (EC 2.4.1.1) (Myophosphorylase)

### Probability Based Mowse Score

Ions score is  $-10 \cdot \log(P)$ , where P is the probability that the observed match is a random event.

Individual ions scores  $> 43$  indicate identity or extensive homology ( $p < 0.05$ ).

Protein scores are derived from ions scores as a non-probabilistic basis for ranking protein hits.



# Mascot Results: Original DB

## *{MATRIX}* Mascot Search Results *{SCIENCE}*

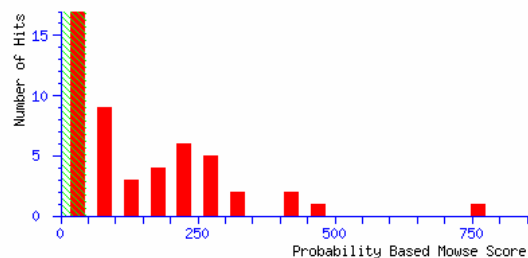
User : Search for Compression Paper  
Email :  
Search title : Conversion of 17mix\_test2.xml to mascot generic by msxml2other  
MS data file : C:\Documents and Settings\edwardnj\Desktop\17prot\_mix.mgf  
Database : Varsplic (203425 sequences; 92944773 residues)  
Timestamp : 5 Oct 2004 at 13:55:52 GMT  
Significant hits: [P00722](#) (BGAL\_ECOLI) Beta-galactosidase (EC 3.2.1.23) (Lactase)  
[P00634-00-00-00](#) (PPB\_ECOLI) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from P00634 Alkaline p  
[P00921-00-01-00](#) (CAH2\_BOVIN) Splice isoform Displayed; Variant one of the major forms; Conflict Displayed; from P00  
[P00489](#) (PHS2\_RABIT) Glycogen phosphorylase, muscle form (EC 2.4.1.1) (Myophosphorylase)  
[Q29443](#) (TRFE\_BOVIN) Serotransferrin precursor (Transferrin) (Siderophilin) (Beta-1-metal binding globulin)  
[P02769-00-00-00](#) (ALBU\_BOVIN) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from P02769 Serum alb  
[P21948](#) (PPB\_ESCFE) Alkaline phosphatase precursor (EC 3.1.3.1) (APase)  
[P02608](#) (MLRS\_RABIT) Myosin regulatory light chain 2, skeletal muscle isoform type 2 (G2) (DTNB) (MLC-2)  
[Q28641](#) (MYH4\_RABIT) Myosin heavy chain, skeletal muscle, juvenile  
[P11217-00-05-00](#) (PHS2\_HUMAN) Splice isoform Displayed; Variant GSD-V-VAR\_014005; Conflict Displayed; from P11217 G1  
[P11217-00-12-00](#) (PHS2\_HUMAN) Splice isoform Displayed; Variant GSD-V-VAR\_014009; Conflict Displayed; from P11217 G1  
[P49064](#) (ALBU\_FELCA) Serum albumin precursor (Allergen Fel d 2)  
[P07724](#) (MYH4\_HUMAN) Myosin heavy chain, skeletal muscle, fetal (Myosin heavy chain I1b) (MyHC-I1b)  
[Q9Y623](#) (MYH1\_HUMAN) Myosin heavy chain, skeletal muscle, adult 1 (Myosin heavy chain I1x/d) (MyHC-I1x/d)  
[P12882](#) (MYSS\_CHICK) Myosin heavy chain, skeletal muscle, adult  
[P13538](#) (CATA\_BOVIN) Catalase (EC 1.11.1.6)  
[P00432](#) (MYH3\_HUMAN) Myosin heavy chain, fast skeletal muscle, embryonic (Muscle embryonic myosin heavy cha  
[P11055](#) (AMY\_BACLI) Alpha-amylase precursor (EC 3.2.1.1) (1,4-alpha-D-glucan glucanohydrolase) (BLA)  
[P06278](#) (PHS2\_RAT) Glycogen phosphorylase, muscle form (EC 2.4.1.1) (Myophosphorylase)  
[P09812](#)

### Probability Based Mowse Score

Ions score is  $-10 \cdot \log(P)$ , where P is the probability that the observed match is a random event.

Individual ions scores  $> 45$  indicate identity or extensive homology ( $p < 0.05$ ).

Protein scores are derived from ions scores as a non-probabilistic basis for ranking protein hits.





# Mascot Results: C<sup>3</sup> DB ++

## Peptide Summary Report

[Switch to Protein Summary Report](#)

To create a bookmark for this report, right click this link: [Peptide Summary Report \( Conversion of 17mix\\_test2.xml to mascot generic by msxml2other\)](#)

Error tolerant

1. [P00722](#) **Mass:** 116278 **Score:** 762 **Peptides matched:** 38

(BGAL\_ECOLI) Beta-galactosidase (EC 3.2.1.23) (Lactase)

Check to include this hit in error tolerant search or archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Expect	Rank	Peptide
<input checked="" type="checkbox"/> <a href="#">587</a>	435.70	869.39	869.42	-0.03	0	(24)	4.3	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">588</a>	435.70	869.39	869.42	-0.03	0	25	3.6	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">589</a>	435.70	869.39	869.42	-0.03	0	(22)	6.2	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">641</a>	450.67	899.32	899.38	-0.06	0	39	0.19	1	FNDDFS
<input checked="" type="checkbox"/> <a href="#">764</a>	495.22	988.42	988.48	-0.06	0	39	0.2	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">765</a>	495.22	988.42	988.48	-0.06	0	(36)	0.36	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">767</a>	495.22	988.42	988.48	-0.06	0	(32)	0.94	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">769</a>	495.22	988.42	988.48	-0.06	0	(38)	0.25	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">845</a>	534.23	1066.44	1066.48	-0.05	0	31	0.72	1	WVGYGQDSR
<input checked="" type="checkbox"/> <a href="#">873</a>	550.23	1098.45	1098.55	-0.10	0	(72)	6.5e-005	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">874</a>	550.23	1098.45	1098.55	-0.10	0	73	5.9e-005	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">1020</a>	626.80	1251.58	1251.65	-0.07	0	36	0.24	1	LAHPPFASWR
<input checked="" type="checkbox"/> <a href="#">1021</a>	626.80	1251.58	1251.65	-0.07	0	(33)	0.45	1	LAHPPFASWR
<input checked="" type="checkbox"/> <a href="#">1037</a>	633.26	1264.51	1264.61	-0.10	0	(24)	4.5	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1038</a>	633.26	1264.51	1264.61	-0.10	0	37	0.23	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1072</a>	650.28	1298.55	1298.62	-0.07	0	47	0.023	1	ELNYGPHQWR
<input checked="" type="checkbox"/> <a href="#">1132</a>	671.31	1340.60	1340.66	-0.06	0	(63)	0.00049	1	VDEDQPPFAVPK
<input checked="" type="checkbox"/> <a href="#">1133</a>	671.31	1340.60	1340.66	-0.06	0	(57)	0.0023	1	VDEDQPPFAVPK
<input checked="" type="checkbox"/> <a href="#">1134</a>	671.31	1340.60	1340.66	-0.06	0	70	0.00011	1	VDEDQPPFAVPK
<input checked="" type="checkbox"/> <a href="#">1135</a>	671.31	1340.60	1340.66	-0.06	0	(38)	0.15	1	VDEDQPPFAVPK
<input checked="" type="checkbox"/> <a href="#">1172</a>	681.33	1360.65	1360.71	-0.06	0	38	0.16	1	LWSAETPNLYR

# Mascot Results: Original DB

## Peptide Summary Report

[Switch to Protein Summary Report](#)

To create a bookmark for this report, right click this link: [Peptide Summary Report \( Conversion of 17mix\\_test2.xml to mascot generic by msxml2other\)](#)

Select All

Select None

Search Selected

Error tolerant

Archive Report

1. [P00722](#) Mass: 116278 Score: **760** Peptides matched: 38

(BGAL\_ECOLI) Beta-galactosidase (EC 3.2.1.23) (Lactase)

Check to include this hit in error tolerant search or archive report

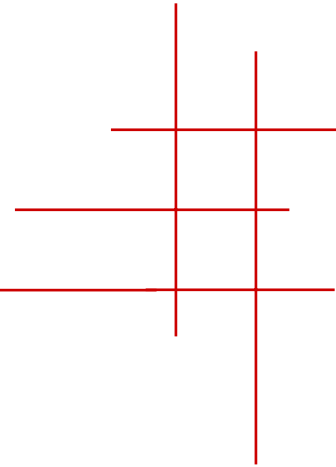
Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Expect	Rank	Peptide
<input checked="" type="checkbox"/> <a href="#">587</a>	435.70	869.39	869.42	-0.03	0	(24)	6.7	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">588</a>	435.70	869.39	869.42	-0.03	0	25	5.7	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">589</a>	435.70	869.39	869.42	-0.03	0	(22)	9.6	1	YWQAFR
<input checked="" type="checkbox"/> <a href="#">641</a>	450.67	899.32	899.38	-0.06	0	39	0.29	1	FNDDFS
<input checked="" type="checkbox"/> <a href="#">764</a>	495.22	988.42	988.48	-0.06	0	39	0.35	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">765</a>	495.22	988.42	988.48	-0.06	0	(36)	0.61	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">767</a>	495.22	988.42	988.48	-0.06	0	(32)	1.6	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">769</a>	495.22	988.42	988.48	-0.06	0	(38)	0.43	1	WLPAMSER
<input checked="" type="checkbox"/> <a href="#">845</a>	534.23	1066.44	1066.48	-0.05	0	31	1.2	1	WVGYGQDSR
<input checked="" type="checkbox"/> <a href="#">873</a>	550.23	1098.45	1098.55	-0.10	0	(72)	0.00011	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">874</a>	550.23	1098.45	1098.55	-0.10	0	73	9.8e-005	1	IDPNAWVER
<input checked="" type="checkbox"/> <a href="#">1020</a>	626.80	1251.58	1251.65	-0.07	0	36	0.42	1	LAHPPFASWR
<input checked="" type="checkbox"/> <a href="#">1021</a>	626.80	1251.58	1251.65	-0.07	0	(33)	0.78	1	LAHPPFASWR
<input checked="" type="checkbox"/> <a href="#">1037</a>	633.26	1264.51	1264.61	-0.10	0	(24)	8	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1038</a>	633.26	1264.51	1264.61	-0.10	0	37	0.4	1	HQQQFFQFR
<input checked="" type="checkbox"/> <a href="#">1072</a>	650.28	1298.55	1298.62	-0.07	0	47	0.039	1	ELNYGPHQWR
<input checked="" type="checkbox"/> <a href="#">1132</a>	671.31	1340.60	1340.66	-0.06	0	(63)	0.00082	1	VDEDPFFAVPK
<input checked="" type="checkbox"/> <a href="#">1133</a>	671.31	1340.60	1340.66	-0.06	0	(57)	0.0038	1	VDEDPFFAVPK
<input checked="" type="checkbox"/> <a href="#">1134</a>	671.31	1340.60	1340.66	-0.06	0	70	0.00018	1	VDEDPFFAVPK
<input checked="" type="checkbox"/> <a href="#">1135</a>	671.31	1340.60	1340.66	-0.06	0	(38)	0.24	1	VDEDPFFAVPK
<input checked="" type="checkbox"/> <a href="#">1172</a>	681.33	1360.65	1360.71	-0.06	0	38	0.26	1	LWSAEIPNLYR

# More Sensitive Peptide ID

- Significances,  $p$ -values, Expect values
  - Normalize for number of trials
  - Blast:
    - Size of sequence database
  - Mascot etc.:
    - Peptides with appropriate parent mass
  - Redundant peptide sequences increase the number of trials, artificially.
    - Trials are not independent!
  - Less redundancy results in a better significance estimate

# Available for Download

---



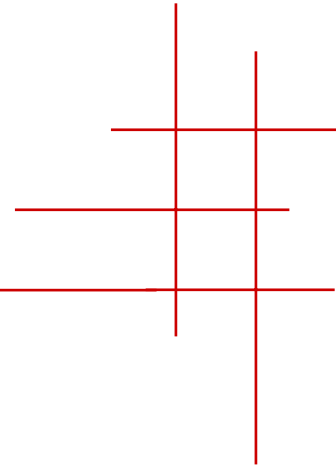
## Available for download:

- C<sup>3</sup> Sequence Databases
- Exact peptide search code
- Mascot parse + rewrite scripts

Other search engines and sequence databases on request

# Genomic Peptide Sequences

---



- Many putative peptide sequences never become “protein” sequences
  - Genomic DNA,
  - Refseq mRNA, ESTs
  - SNP/Polymorphism databases
  - Variant records in SwissProt
- Genomic annotation seeks “full length” genes and proteins

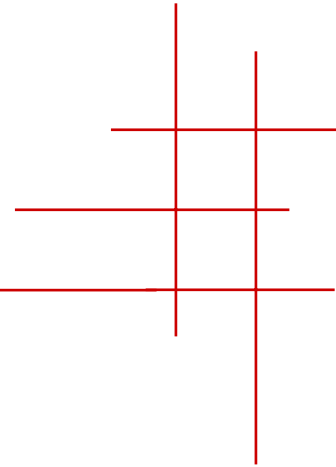
# Genomic Peptide Sequences

- Genomic DNA
  - Exons & introns, 6 frames, large (6 \* 3Gb)
- Refseq mRNA
  - No introns, 3 frames, small (36Mb)
  - Most protein sequences already represented in sequence databases
- ESTs
  - No introns, 3 frames, large (3 \* 3Gb)
  - Used by gene & alternative splicing pipelines
  - Highly redundant, nucleotide error rate ~ 1%

# Compressed EST Peptides

- 3 frame translation
- Break at non-amino-acids
  - stop codons + X
- Discard AA sequence < 50 AA long
- Result: 1.1Gb
- C3 Compress
  - Discard 30-mers that occur only once
- Result: 83Mb ( < 2 \* SwissProt )

# Compressed EST Peptides



- Genomic Peptides
  - Refseq mRNA & Proteins, ESTs
  - Size: 96Mb
- Genomic + IPI-HUMAN Peptides
  - Size: 99Mb
  - 30-mers:
    - 11.2Mb 30-mers in common
    - 1.6Mb in IPI-HUMAN alone
    - 37.2Mb in putative genomic peptides alone



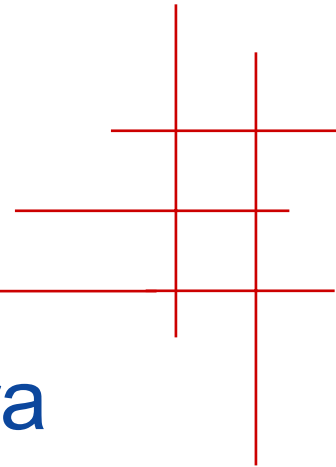
# Extensions and Futher Work



- Better compression
  - Relax compactness constraint
  - Enumerate tryptic peptides only
  - Relax correctness constraint
- Decouple peptide ID from protein ID
- Genomic peptides from GenScan exons

# Thanks

---



- Informatics Research @ ABI & Celera
  - Ross Lippert, Clark Mobarry, Bjarni Halldorsson
- UMIACS @ University of Maryland, CP
  - V.S. Subrahmanian, Fritz McCall, Doan Pham
- Fenselau Lab @ UM, CP