

Identification of Proteins from Organisms with Unsequenced Genomes by Tandem Mass Spectrometry and Sequence-Similarity Database Searching Tools

Adam J. Liska and Andrej Shevchenko

I. INTRODUCTION

The analysis of proteomes by mass spectrometric methods that correlate peptide fragments from proteins with database entries *in silico* has been dependent on the sequencing of genomes. Mass spectrometry and database sequences have enabled the analysis of the human, mouse, and *Arabidopsis* proteomes, among others. Due to the high homology between living organisms at the molecular level, it is possible to use the available protein sequences accumulated in databases from a range of organisms as a reference for the identification of proteins from organisms with unsequenced genomes by sequence-similarity database searching. As research continues in organisms such as *Xenopus*, maize, cow, and others with limited database sequence resources, sequence-similarity searching is a powerful method for protein identification. This article focuses on MS BLAST (Shevchenko *et al.*, 2001) and MultiTag (Sunyaev *et al.*, 2003) as bioinformatic methods for the identification of proteins by the interpretation of tandem mass spectra of peptides and sequence-similarity searching.

II. MATERIALS AND INSTRUMENTATION

In analyses where MS BLAST is utilized for protein identification, tandem mass spectra of peptides can be acquired with any ionization source and mass spectrometer that enables *de novo* sequence prediction: nanoelectrospray, LC/MS/MS, o-MALDI (MALDI nanoelectrospray, LC/MS/MS, o-MALDI (MALDI quadrupole TOF), MALDI TOF-TOF, QqTOF, triple quad, PSD MALDI-TOF, and ion trap. Alternatively, in analyses where MultiTag is applied, tandem mass spectra of peptides can be acquired with any ionization source and tandem mass spectrometer that enables the creation of peptide sequence tags (Mann and Wilm, 1994): nanoelectrospray, LC/MS/MS, QqTOF, triple quad, or other novel system.

Software for *de novo* sequence prediction from tandem mass spectra is often included in the software packages associated with mass spectrometers: BioMultiview, BioAnalyst (both are from MDS Sciex, Canada), BioMassLynx (Micromass Ltd, UK), BioTools (Bruker Daltonics, Germany), and DeNovoX (ThermoFinnigan). The Lutefisk program (Johnson and Taylor, 2000) can be acquired from <http://www>.

hairyfatguy.com/Sherpa/. BioAnalyst with the ProBlast processing script can generate a complete MS BLAST query automatically from multiple-spectra files acquired by nanoelectrospray or LC/MS/MS (Nimkar and Loo, 2002). A web browser such as Internet Explorer or Netscape is also required to gain access to the MS BLAST web interface located at <http://dove.embl-heidelberg.de/Blast2/msblast.html>. An independent BLAST computer (Paracel BlastMachine system) may also be purchased and installed for rapid and private MS BLAST operation.

Sequence tags that contain a few confidently designated amino acid residues and mass values that lock the short sequence stretch into the length of the peptide can often be generated from tandem mass spectra with the software packages associated with mass spectrometers. Microsoft Excel or an alternative spreadsheet program is required for compiling of search results before submission to MultiTag. BioAnalyst software has an associated processing script that produces a list of database search results from a generated list of sequence tags to accelerate spectra processing with MultiTag.

III. PROCEDURES

A. Identification of Proteins by MS BLAST Database Searching

1. MS BLAST is a specialized BLAST-based tool for the identification of proteins by sequence-similarity searching that utilizes peptide sequences produced by the interpretation of tandem mass spectra (Shevchenko *et al.*, 2001). The algorithm and principles of BLAST sequence-similarity searching are reported in detail elsewhere (Altschul *et al.*, 1997). A useful list of BLAST servers accessible on the web is provided in Gaeta (2000).

2. Peptide sequences are generated from the interpretation of tandem mass spectra from the analysis of a single in-gel or in-solution digest of an unknown protein, edited and assembled into a query list for the MS BLAST search. If tandem mass spectra were interpreted by *de novo* sequencing software, disregard relative scores and use the entire list of candidate sequences (or some 50–100 top scoring sequence proposals per fragmented peptide precursor) (Fig. 1). Automated interpretation of tandem mass spectra often requires adjustment of parameters that affect the quality of predicted sequences. It is therefore advisable to test the settings in advance using digests of standard proteins and to adjust them if necessary. Note that the settings may depend on a charge state of the fragmented precursor ion. Use only the standard single-

letter symbols for amino acid residues. If the software introduces special symbols for modified amino acid residues, replace them with standard symbols.

3. When interpreting MS/MS spectra manually, try making the longest possible sequence stretches, although their accuracy may be compromised. For example, it is usually difficult to interpret unambiguously fragment ion series at the low *m/z* range because of abundant peaks of chemical noise and numerous fragment ions from other series. In this case, it is better to include many complete (albeit low confidence) sequence proposals into the query rather than using a single (although accurate) three or four amino acid sequence stretch deduced from a noise-free high *m/z* segment of the spectrum.

4. Gaps and ambiguities in peptide sequences can occur due to the fragmentary nature of tandem mass spectra of peptides. Some *de novo* sequencing programs may suggest a gap in the peptide sequence that can be filled with various isobaric combinations of amino acid residues. For example,

DTPS[...]HYNAR, [...] = [S, V] or [D, A]

If one or two combinations were suggested, include all variants into a searching string:

-DTPSSVHYNAR-DTPSVSHYNAR-
-DTPSDAHYNAR-DTPSADHYNAR-

If more combinations were possible, the symbol X can be used instead to fill the gap. Zero score is assigned to X symbol in PAM30MS scoring matrix and therefore it matches weakly any amino acid residue:

-DTPSXXHYNAR-

Note that MS BLAST is sensitive to the number of amino acid residues that are filling the gap. If the gap could be filled by a combination of two and three amino acid residues, consider both options in the query

-DTPSXXHYNAR-DTPSXXXHYNAR-

5. Isobaric amino acids need to be altered in the MS BLAST query. L stands for Leu (L) and Ile (I). Z stands for Gln (Q) and Lys (K), if undistinguishable in the spectrum. Use Q or K if the amino acid residue can be determined. The query string needs to be further altered for cleavage site specificity. If the proposed sequence is complete, a putative trypsin cleavage site symbol B is added prior to the peptide sequence:

...-BDTPSVDHYNAR-

It is often difficult to determine two amino acid residues located at the N terminus of the peptide. In this case, present them as

...-BXXPSVDHYNAR-...

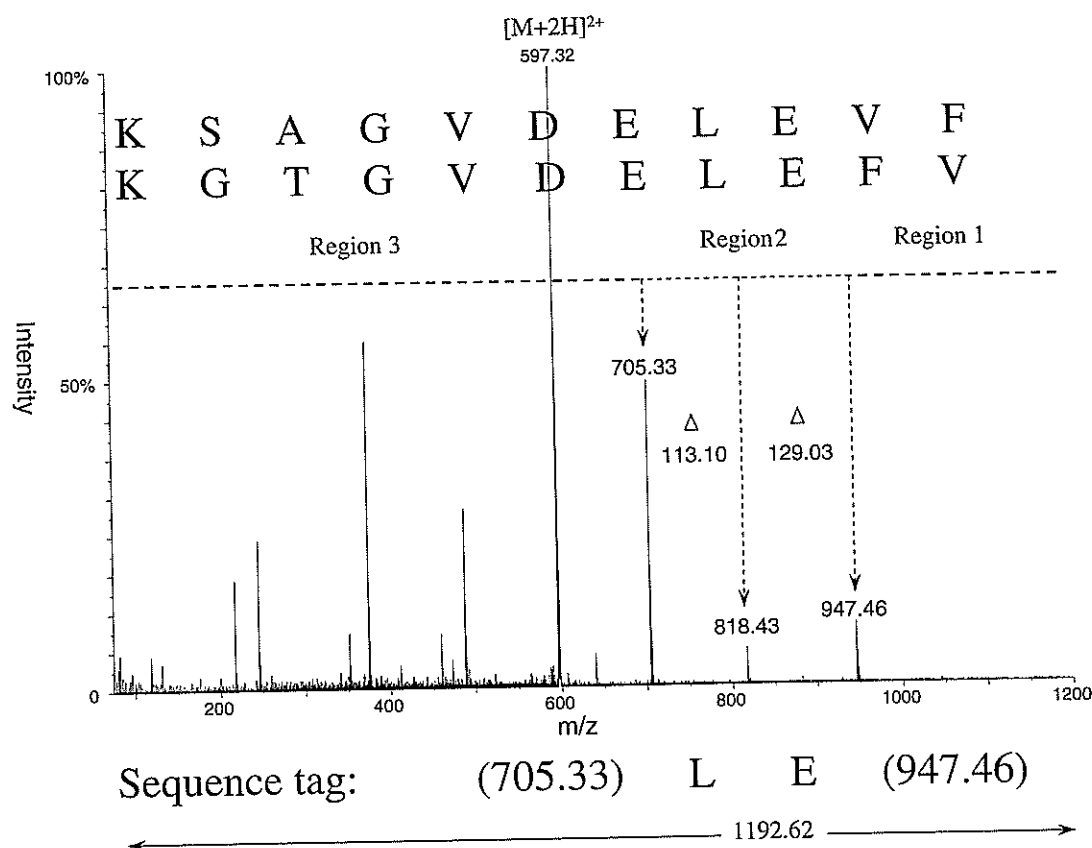


FIGURE 1 A spectrum with *de novo*-predicted amino acid sequences and a manually constructed sequence tag. Multiple candidate peptide sequences can be generated from a single spectrum for MS BLAST analysis, whereas MultiTag requires one sequence tag per spectrum.

MS BLAST will then consider BXX residues in possible sequence alignments.

6. The regular BLAST search must be altered in options and settings for an MS BLAST query:

NOGAP is absolutely essential, it turns off gapped alignment method so that only high-scoring pairs (HSPs) with no internal gaps are reported.

SPAN1 is absolutely essential, it identifies and fetches the best matching peptide sequence among similar peptide sequences in the query. Therefore the query may contain multiple partially redundant variants of the same peptide sequence without affecting the total score of the protein hit.

HSPMAX 100 limits the total number of reported HSPs to 100. Set it to a higher number (e.g., 200) if a large query is submitted and a complete list of protein hits (including low confidence hits) is required in the output.

SORT_BY_TOTALSCORE places the hits with multiple high scoring pairs to the top of the list.

Note that the total score is not displayed, but can be calculated, if necessary, by adding up scores of individual HSPs.

EXPECT: It is usually sufficient to set EXPECT at 100. Searching with higher EXPECT (as, 1000) will report many short low-scoring HSPs, thus increasing the sequence coverage by matching more fragmented peptides to the protein sequence. Note that low scoring HSPs do not increase statistical confidence of protein identification. The EXPECT setting also does not affect the scores of retrieved HSPs.

MATRIX: PAM30MS is a specifically modified scoring matrix. It is not used for conventional BLAST searching.

PROGRAM: blast2p.

DATABASE: nrdb95 are default settings of the MS BLAST interface.

FILTER: Filtering is set to "none" default. However, if the sequence query contains many repeating stretches (as ... EQEQEQ ...), filtering should be set to "default."

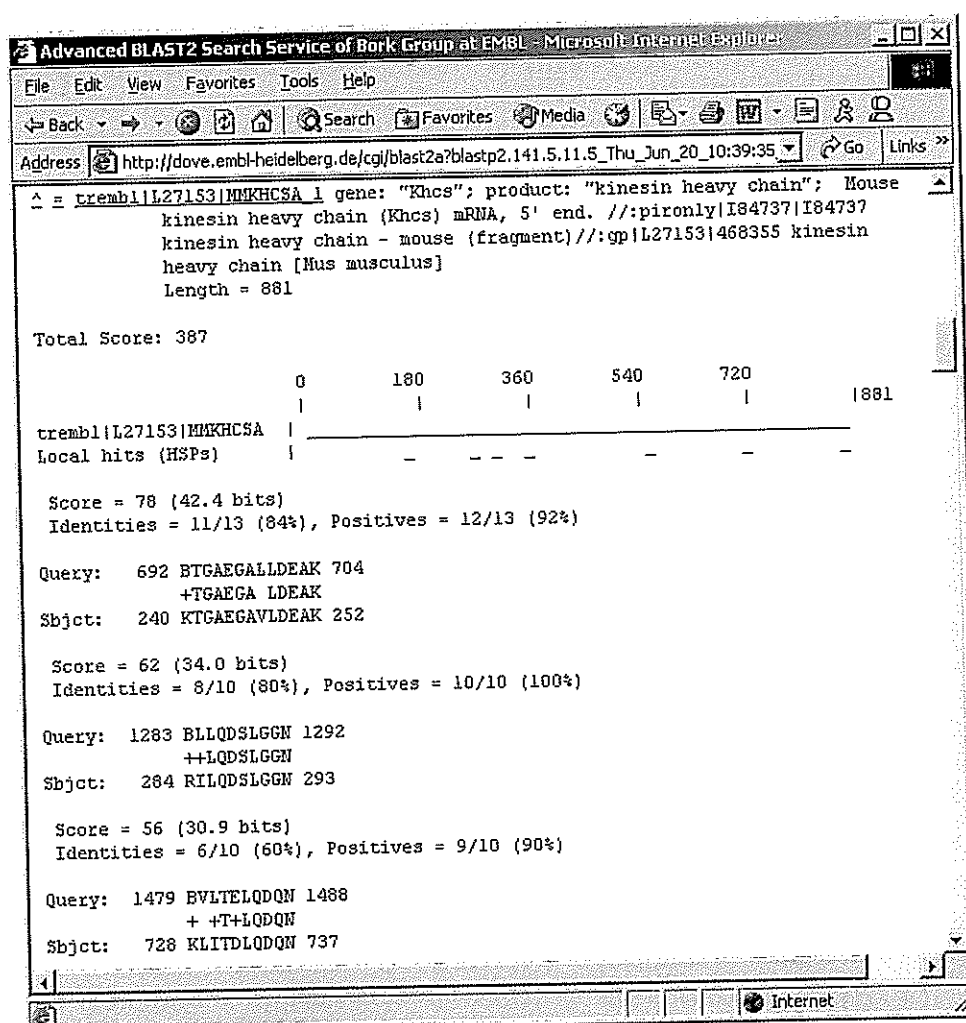


FIGURE 3 A section of the MS BLAST output where three of seven matched peptides to one database entry are shown. A list of matching database entries with the highest significance match at the top of the list is generated. Significant hits are color coded for easy data interpretation.

igned peptides from a *tblastn* search against specific genomic databases. This search enables the use of unannotated genomic sequences and makes it possible to identify novel genes in large genomes. However, in both EST and genomic searches, different scoring schemes would need to be developed and installed locally. To set up a local BLAST searching engine, VU-BLAST 2.0 can be acquired from <http://blast.vustl.edu/>.

3. Identification of Proteins by MultiTag Database Searching

1. MultiTag is a software program that sorts compiled results from database searches with partial and complete sequence tags and calculates the significance of matches that align multiple sequence tags (for a

complete description, see Sunyaev *et al.*, 2003). MultiTag is based on error-tolerant searching with multiple partial sequence tags. This technique enables the correlation of search results from multiple searches with sequence tags representative of numerous spectra and gauges the significance of those matches.

2. Sequence tags should be generated manually from tandem mass spectra acquired in the analysis of a single in-gel or in-solution digest of an unknown protein or proteins. Some mass spectrometer software enables the automatic prediction of sequence tags; however, for best results it is advisable to make sequence tags by manual interpretation or gauge the accuracy of the automatic prediction software with a standard protein prior to the analysis of unknown samples. One sequence tag per spectrum should be made using prominent Y-ions, usually larger

than the multiply charged precursor. Sequence tags made with two to four amino acids each from multiple MS/MS spectra should be compiled in a text file list that includes the tag followed by the parent mass.

(360.20)FLL(733.44)918.64

(561.27)LA(745.40)935.48

(866.41)DEA(1181.52)1422.59

3. Each sequence tag is used to search a protein database, and the results from four searches are compiled in a spreadsheet. For the most specific results, mass tolerances should be narrow, taking into consideration the best accuracy of the mass spectrometer employed. The database is first searched using the complete sequence tag:

(866.41)DEA(1181.52)

This search will only find proteins that contain peptides with exactly these amino acids, spaced with exact amino acids residues that give mass combinations to make up the gaps to the peptide's termini. The second search allows for one error within the amino acid representation itself:

(866.41)D?A(1181.52)

The third search allows for errors between the analyzed peptide and the database at the C terminus of the peptide (searching with regions 1 and 2 only):

DEA(1181.52)

The fourth search allows for errors between the analyzed peptide and the database at the N terminus of the peptide (searching with regions 2 and 3 only):

(866.41)DEA

In the first column of the results table, the parent mass should be followed by an "NC", "E", "N", or "C" for sequence tag search results from complete, one-error, regions 1 and 2 (N-terminal match), and regions 2 and 3 (C-terminal match), respectively. The second, third, fourth, fifth, and sixth columns should include the amino acid sequence of the peptide matched, the molecular weight of the protein the tag matched, the accession number of the entry, the name of the protein, and the species name, respectively (Fig. 4). Compile the search results in a spreadsheet from all searches with the sequence tags generated, and save results as a text file.

4. Submit the search results table to MultiTag. Designate the mass tolerance used for searching (in daltons), input the approximate number of entries in

the database searched, input the list of tags used to generate search results, and compute significance. Results will be sorted, with the database entry containing the most unlikely correlation event at the top of the list and probabilities are calculated (Fig. 5). Results can be evaluated based on the number of sequence tags matched, E-values and P-values [count (predicted)]. E-values lower than 1×10^{-3} and P-values lower than 1×10^{-4} can be considered significant matches. Final E-values are highly dependent on the number of tags submitted for database searching; more tags will tend to diminish the significance of the alignment of multiple tags to one database entry. Reported P-values are less affected by the number sequence tags in the query. P-values reflect an approximation of the probability that the tags that are aligned will match randomly to an entry in a database of a specific size and at a specific mass accuracy, while neglecting the query size (number of sequence tags). Low P-values (but with higher E-values) are good indicators of alignments having borderline significance that need further manual evaluation to conclude a confident identification. Three partial sequence tags are normally specific enough to identify one entry in a database of 1,000,000 entries, at a mass accuracy of 0.1 Da.

IV. COMMENTS

It is not known in advance if the sequence of the analyzed protein is already present in a database. Therefore, conventional database-searching routines based on stringent matching of peptide sequences should be applied first (Mann and Wilm, 1994; Perkin *et al.*, 1999). Only if the protein is unknown and no convincing cross-species matches can be obtained is it recommended to proceed with *de novo* interpretation of tandem mass spectra and sequence-similarity searching.

The success of MS BLAST and MultiTag identification depends on the size of a query and the corresponding database, the number of peptides aligned, the quality of peptide sequences or sequence tags, and the sequence similarity between the protein of interest and its homologues available in a database. On average, candidate sequences determined for five tryptic peptides should be submitted to MS BLAST or MultiTag searching to identify the protein by matching to a homologous sequence. With 10 sequences submitted and aligned, MS BLAST can identify 50% of homologues containing 50% sequence similarity, and MultiTag can identify 50% of homologues at 70% sequence similarity. However, because MultiTag can utilize less

MultiTag - [509.01.14input.txt]					
File View MultiTag Window					
Tag Mass	Sequence	Mass(kDa)	DB Accession	Protein name	Species
1173.628E	SAAKKVKNAEK	47.465796	gi 14210646	(AY033620) putative RNA-binding pro	Unknown
1173.628E	TGAELHLWLR	27.823135	gi 7019377	(NM_013393) cell division protein Ft	Human
1173.628E	TGAELHLWLR	27.487108	gi 13386002	(NM_026510) RIKEN cDNA Z310037B18 [Mouse
1173.628E	SAANKALNDKK	15.032396	gi 12744797	AF323725_1 (AF323725) PsaN precursor	Unknown
1173.628E	ASAEILSVDRV	47.31585	gi 15611313	(NC_000921) EXODEOXYRIBONUCLEASE LA	Helicobacter Pylori
1173.628E	ASAEILSVDRV	47.455069	gi 15644887	(NC_000915) exonuclease VII, large	Helicobacter Pylori
1173.628E	TGAETLVEEAK	12.263558	gi 401181	THGF_TOBAC FLOWER-SPECIFIC GAMMA-THIO	Unknown
1173.628E	SAAERKRQEK	79.528375	gi 18375979	(AL356173) conserved hypothetical p	Unknown
1173.628E	SAAERKRQEK	82.102377	gi 11359450	T49456 hypothetical protein B14D6.8	Unknown
1173.628E	SANEKKSINVK	143.453267	gi 17224297	AF218388_1 (AF218388) apoptotic pro	Rat
1173.628E	SANEKKSINVK	143.435205	gi 13027436	(NM_023979) apoptotic protease acti	Rat
1173.628E	SAAEAQATRGR	21.850893	gi 13162112	(AL512667) putative tetR-family tra	some Streptomyces
1173.628E	GTAEQRLFLVG	32.418744	gi 7294725	(AE003544) CG7547 gene product [Dros	Fruit Fly
1173.628E	SAAEQWQDL	74.867942	gi 17227422	(NC_003267) ORF_ID:all8048~unknown	Unknown
1173.628E	ASAEQRATQTI	36.747253	gi 17481280	(AB062896) vomeronasal receptor 1 A	Mouse
1173.628E	ASAEQRATQTI	35.05824	gi 3892596	(Y12724) pheromone receptor 2 [Mus m	Mouse
1173.628E	ASAEQRATQTI	35.84985	gi 17481276	(AB062895) vomeronasal receptor 1 A	Mouse
1173.628E	ASAEQRATQTI	36.402747	gi 18558569	(AY065464) vomeronasal receptor V1R	Mouse
1173.628E	ASAEQRATQTI	34.605305	gi 16716523	(NM_053218) vomeronasal 1 receptor,	Mouse
1173.628E	ASAEKGIASVRS	13.48812	gi 15802561	(NC_002655) orf, hypothetical prote	Escherichia Coli
1173.628E	SAAEQSGLDKNG	35.440047	gi 12620486	AF322012_67 (AF322013) ID142 [Brady	Unknown
1173.628E	ASAEKKRQATS	56.223343	gi 8570440	AC020622_1 (AC020622) Contains simil	Human
1173.628E	ASAEKKRQATS	64.846084	gi 15223502	(NM_100069) hypothetical protein [A	Mouse-Ear Cress
1173.628E	SAAEKLSEETL	272.279457	gi 4874311	AC006053_15 (AC006053) unknown prote	Mouse-Ear Cress
1173.628E	SAAEKLSEETL	60.547133	gi 15081785	(AY048285) At2g25730/F3N11.18 [Arab	Mouse-Ear Cress
1173.628E	SAAEKLSEETL	277.4555	gi 18400918	(NM_128132) unknown protein [Arabid	Mouse-Ear Cress
1173.628E	ASAEKKAESK	105.810124	gi 15900468	(NC_003028) translation initiation	Streptococcus Pyogenes
1173.628E	ASAEKYPHEF	31.186912	gi 17988633	(NC_003318) DIPEPTIDE TRANSPORT ATP	Unknown
1173.628E	GTAEKMPITSR	13.777599	gi 3204328	(AJ008500) gag protein [Human immuno	Human
1173.628E	GTAEKMPITSR	13.71449	gi 3204368	(AJ008521) gag protein [Human immuno	Human
1173.628E	TGAEKRSFVAD	80.784195	gi 7302767	(AE003803) CG4878 gene product [alt	Fruit Fly
1173.628E	SAAEKIVVSGG	50.711733	gi 17231331	(NC_003272) unknown protein [Nostoc	Unknown
1173.628E	SAAEKAVSAPPR	55.045677	gi 13471797	(NC_002678) ATP-binding protein of	Unknown
1173.628E	SAAEKFDVSMT	24.872427	gi 10956719	(NC_002490) conjugal transfer prote	Unknown
1173.628E	ASAEKEQIAQI	144.720391	gi 16555336	(AY056833) chitin synthase [Anophel	Unknown
1173.628E	GTAEQHGRNWK	46.230479	gi 16759094	(NC_003198) putative 15 element tra	Unknown
1173.628E	GTAEQHIKEGK	51.501952	gi 695769	(X84038) transposase [Xanthobacter au	Unknown
1173.628E	GTAEGGLAIGDT	86.792704	gi 8894820	(AL360055) putative ABC transport sy	some Streptomyces
1173.628E	SAAEKDGKQKE	10.64305	gi 16850306	(XM_103535) hypothetical protein XP	Human
1173.628E	TGAEKAPKSPSK	13.977534	gi 6009909	(AB018242) histone H2A-like protein	Unknown
1173.628E	ASAEQCGRQAGG	33.741525	gi 7798662	AF135145_1 (AF135145) class I chitin	Unknown
1173.628E	GTAEKMPITSR	13.580314	gi 3204322	(AJ008497) gag protein [Human immuno	Human
1173.628E	GTAEKMPITSR	13.697865	gi 3355417	(AJ011213) gag protein [Human immuno	Human
1173.628E	GTAEKMPITSR	13.807781	gi 3204271	(AJ008470) gag protein [Human immuno	Human
1173.628E	GTAEKMPITSR	13.637407	gi 3204303	(AJ008487) gag protein [Human immuno	Human
1173.628E	GTAEKMPITSR	13.552304	gi 3204338	(AJ008505) gag protein [Human immuno	Human

FIGURE 4 Compiled and formatted search results from sequence tag searching are input in the MultiTag software via opening a text-formatted results file.

tense and noisy spectra, it can outperform MS BLAST in many cases (for a more thorough discussion on analogue identification specificity, see Sunyaev *et al.*, 2003).

Both MS BLAST and MultiTag can identify proteins present in mixtures. Usually two or three components per sample can be identified easily. The sensitivity of both methods is determined primarily on the quality of the *de novo* sequences or sequence tags.

V. PITFALLS

1. Poor sample preparation can frequently deteriorate the quality of tandem mass spectra of peptides. The digestion of proteins with trypsin or other proteases should be carried out with chemicals of the

highest degree of purity available. Plasticware (pipette tips, gloves, dishes, etc.) may acquire a static charge and attract dust, thus leading to contamination of samples with human and sheep (wool) keratin during in-gel or liquid digestion. Any polymeric detergents (Tween, Triton) should not be used for cleaning the laboratory materials.

2. When generating *de novo* sequences or sequence tags, if the software automatically extrapolates the parent mass from the precursor isotope cluster in the MS/MS spectra or the proceeding survey scan in a LC/MS/MS run, it is advisable to manually calculate this value, as software may determine the parent mass incorrectly by designating an incorrect charge state or ^{12}C monoisotopic peak of the parent ion isotope cluster, thus disabling correct *de novo* sequence prediction and sequence tag prediction.

MultiTag - [509.01.14input.txt-output]

File View MultiTag Window

Nr	Tag Mass	Sequence	Mass (kDa)	DB Accession	Protein name	Species	Count (predicted)	E-value
0	1206.6986N;920.5164E;1091.6363N...	LYLVDLAGESEKV;STLMFGQ...	110.065...	gi 4758650	(NM_004522) kinesin family me...	Human	4.42795e-017	4.25215e-008
1	1206.6986N;920.5164N;1091.6363...	LYLVDLAGESEKV;STLMFGQ...	110.427...	gi 4758648	(NM_004521) kinesin family me...	Human	2.13484e-014	4.25215e-008
2	920.5164E;1091.6363N;1231.6488N...	STLMFGQR;LFVQDLTRV...	109.811...	gi 6680574	(NM_008449) kinesin family me...	Mouse	4.92744e-014	4.25215e-008
3	1206.6986N;920.5164E;1091.6363N...	LYLVDLAGESEKV;STLMFGQ...	117.923...	gi 481072	S37711 kinesin heavy chain - ...	Mouse	9.7062e-013	4.25215e-008
4	1206.6986N;920.5164E;1091.6363N...	LYLVDLAGESEKV;STLMFGQ...	117.889...	gi 6680570	(NM_008447) kinesin family me...	Mouse		
5	1206.6986N;920.5164N;1231.6488N...	LYLVDLAGESEKV;STLMFGQ...	43.572721	gi 14424665	AAH09353 (BC009353) Similar ...	Human	2.56829e-011	4.25215e-008
6	1206.6986N;920.5164E;1231.6488N...	LYLVDLAGESEKV;STLMFGQ...	36.815703	gi 3891936	Human Ubiquitous Kinesin Mot...	Human		
7	1206.6986N;920.5164E;1231.6488N...	LYLVDLAGESEKV;STLMFGQ...	18.178803	gi 2981494	(AF053473) kinesin heavy chai...	Mouse	1.16769e-009	7.21201e-007
8	920.5164N;1231.6488N	STLMFGQR;ILQDSLGGNCR	101.405...	gi 2119280	I84737 kinesin heavy chain - m...	Mouse	2.859e-008	9.96048e-006
9		STLMFGQR;ILQDSLGGNCR	80.580155	gi 13628366	(XM_005856) kinesin family me...	Human		
10		STLMFGQR;ILQDSLGGNCR	110.291...	gi 6680572	(NM_008448) kinesin family me...	Mouse		
11	920.5164E;1231.6488N	STLMFGQR;ILQDSLGGNCR	118.234...	gi 18579458	(XM_012156) kinesin family me...	Human	1.29941e-006	0.000244703
12		STLMFGQR;ILQDSLGGNCR	43.899682	gi 18579462	(XM_090306) hypothetical pro...	Human		
13		STLMFGQR;ILQDSLGGNCR	118.248...	gi 4826808	(NM_004984) kinesin family me...	Human		
14		STLMFGQR;ILQDSLGGNCR	35.831809	gi 9929983	(AB047624) hypothetical prote...	some ...		
15		STLMFGQR;ILQDSLGGNCR	13.397445	gi 3891777	B Chain B, Kinesin (Dimeric) Fr...	Rat		
16	920.5164N;1091.6363N	STLMFGQR;LFVQDLQNK	109.864...	gi 125415	KINH_1OLPE KINESIN HEAVY C...	Unknown	0.000122741	0.00747671
17	1173.628N;1401.7885E	GTAELKREVV;KLSVKNAA...	41.586238	gi 19114865	(NC_003424) hypothetical pro...	Fission...	0.00091398	0.03948
18	920.5164E;1231.6488N	STLMFGQR;ILQDSLGGNCR	11.324004	gi 3114354	B Chain B, Kinesin (Monomeric)...	Rat	0.00161464	0.0721143
19	1091.6363C;1630.936C	DFVQDVMK;TDDCEDFVQ...	25.268933	gi 15022431	(AB046578) orf [Treponema m...	Unknown	0.00185708	0.0844395
20		PSFVKGFLLR;TQNIAPSFV...	82.951548	gi 14043646	AAH07795 (BC007795) Similar ...	Baker...		
21		DFVKWSKGG;SYKSKDFVK...	28.849588	gi 5381159	(D49512) Chockroach lectin-lk...	Unknown		
22		DFVKWSKGG;SYKSKDFVK...	28.772455	gi 5381157	(D49511) Cockroach lectin-like ...	Unknown		
23		EGFVKMWVEK;TLQATEGFV...	49.410834	gi 15807360	(NC_001263) pyruvate dehydr...	Unknown		
24		PSFVKGFLLR;TQNIAPSFV...	86.740994	gi 18575674	(XM_034420) YME1 (S.cerevisi...	Baker...		
25		EFVQTLMLK;VFWSGEFVQ...	37.200563	gi 16330556	(NC_000911) unknown protein...	Unknown		
26		FFVKSRSKK;SSIKNFFVK...	121.038...	gi 12656113	AF229182_1 (AF229182) tran...	Unknown		
27		EFVKILPKL;AVFIPEFVKTL...	41.738234	gi 15894852	(NC_003030) NAD-dependent ...	some ...		
28		EFVKACVY;PKFWEFVK...	86.871085	gi 15231992	(NM_111730) hypothetical pro...	Mouse...		
29		DGVQSGKTGR;TYSTDG...	100.994...	gi 9294681	(AP001305) receptor-like prot...	Mouse...		
30		RFVKKAMKK;KSIARRFVK...	51.856202	gi 11499811	(NC_000917) covalytic acid a...	Mouse...		
31		PSFVKGFLLR;TQNIAPSFV...	86.789038	gi 14248493	AF151782_1 (AF151782) ATP...	Human		
32		PSFVKGFLLR;TQNIAPSFV...	80.061091	gi 7657689	(NM_014263) YME1 (S.cerevisi...	Baker...		
33		PSFVKGFLLR;AQNIAPSFV...	80.199582	gi 7305635	(NM_013771) YME1-like 1 (S.c...	Baker...		

Ready

FIGURE 5 A section of MultiTag output. Margins may be adjusted to see the full list of tags matched, as well as full peptide sequences aligned, names, and so on. Results may be saved as text files to be viewed in an appropriate spreadsheet application.

3. MultiTag is laborious. Without scripted sequence tag database searching and processing of search results, manual data processing can demand extended effort; however, in cases where conventional methods fail to identify analyzed proteins, positive identifications are of a high value to cell biological studies.

4. Poor queries tend to obscure protein identification by both MS BLAST and MultiTag. It is best to submit fewer higher quality sequences than numerous lower quality sequences to MS BLAST. MS BLAST is particularly susceptible to low-complexity glycine- and proline-rich sequences generated incorrectly by *de novo* software. These low-complexity sequences tend to mask correct alignments. MultiTag functions best with sequence tags containing multiple (2–4) amino acids that have a low prevalence, such as tryptophan

(W) or methionine (M), whereas common amino acids such as leucine (L) in the tag tend to be of less significance and are likely to produce more false positive. Sequence tags generated from larger peptides have more significance in a database search than those generated from smaller peptides.

References

- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, W., and Lipman, D. J. (1997). Gapped BLAST and FBLAST: A new generation of protein database search programs. *Nucleic Acids Res.* 25, 3389–3402.
- Gaeta, B. A. (2000). BLAST on the Web. *Biotechniques* 28, 436–440.
- Johnson, R. S., and Taylor, J. A. (2000). Searching sequence databases via *de novo* peptide sequencing by tandem mass spectrometry. *Methods Mol. Biol.* 146, 41–61.

- nn, M., and Wilm, M. (1994). Error-tolerant identification of peptides in sequence databases by peptide sequence tags. *Anal. Chem.* 66, 4390-4399.
- mkar, S., and Loo, J. A. (2002). Orlando FL. Application of a new algorithm for automated database searching of MS sequence data to identify proteins. Abstract 334.
- rkins, D. N., Pappin, D. J., Creasy, D. M., and Cottrell, J. S. (1999). Probability-based protein identification by searching sequence databases using mass spectrometry data. *Electrophoresis* 20, 3551-3567.
- Shevchenko, A., Sunyaev, S., Loboda, A., Bork, P., Ens, W., and Standing, K. G. (2001). Charting the proteomes of organisms with unsequenced genomes by MALDI-quadrupole time-of-flight mass spectrometry and BLAST homology searching. *Anal. Chem.* 73, 1917-1926.
- Sunyaev, S., Liska, A., Golod, A., Shevchenko, A., and Shevchenko, A. (2003). MultiTag: Multiple error-tolerant sequence tag search for the sequence-similarity identification of proteins by mass spectrometry. Submitted for publication.