Probability-Based Protein Identification for Post-Translational Modifications and Amino Acid Variants Using Peptide Mass Fingerprint Data

Tong WW, McComb ME, Perlman DH, Huang H, O'Connor PB, Costello CE, and Weng Z

Overview

- A web server that interprets peptide mass fingerprinting data by searching the spectrum against protein sequence databases using log likelihood ratio as scoring function.
- Log likelihood ratio is the best statistical model to distinguish correctly assigned peptides from incorrect assignments.
- Matching peaks with self-adjusting mass tolerance offers more flexibility and accuracy than the traditional mass window.
- Seamless pipeline for further processing of the result (internal calibration and PTM searches).
- BUPID outperforms conventional database search algorithms in sensitivity and specificity.

Introduction

A major goal of proteomics is to identify and characterize all proteins expressed in cells in various conditions. Mass spectrometry (MS) has become popular for identification of proteins in high-throughput proteomics research. Peptide mass fingerprint (PMF) with database search is a standard way to analyze high-throughput MS data of complex protein samples. This method compares the PMF spectrum against those expected for all possible proteins obtained from a sequence database. Each database protein is then assigned a probability score that reflects various aspects of the fit between spectrum and peptide. These scores help discriminate correct and incorrect protein assignments to the spectrum and therefore identify proteins in the sample.

BUPID

http://zlab.bu.edu/Amemee

We report a new method of database searching using MS data. The algorithm utilizes a log likelihood score to discriminate correctly assigned peaks from incorrectly assigned ones. The accumulative log likelihood score of all peaks that match with the protein represent the significance of the correlation between the protein and the spectrum. The raw score is normalized according to the size and mass of the protein in order to compare with other proteins against the same standard. The method was proven to have the best sensitivity and specificity in identifying proteins in a mixture. Currently BUPID only search for unmodified proteins. PTM search is under development and will be available to the public shortly.

The Principle of Protein Identification by Peptide Mass Fingerprinting



All database searching algorithms, although differ in scoring functions, share a similar approach. First, theoretical peptide mass values are created by applying simulative cleavage to sequences in the database. Then a set of peaks (matches) in the experimental spectrum are picked out and compared with calculated peptide mass values of the protein. Keep in mind that identifying matched peaks is as important as assessing their similarities.

Incentive to Use and Not to Use Peptide Mass Fingerprinting

- PMF is faster to acquire
- PMF is more sensitive
- PMF requires full (and correct) protein sequences in the databases

BUPID Uses Log Likelihood Ratio as the Target Function

$$score = \ln \frac{likelihood H_A}{likelihood H_0}$$
$$= \ln \frac{likelihood(spectrum data | predicted protein)}{likelihood(spectrum data | random background)}$$

Spectrum data:	All possible combination of peaks.
Predicted protein:	A protein considered to be in the sample.
Background:	- All proteins in the mixture.
	- All proteins in the database.

- All proteins in the universe.

Log likelihood ratio is considered the most powerful method to distinguish two hypotheses. In the database search problem, the null hypothesis is that a certain set of peaks in the spectrum is generated by random background noise. The alternative hypothesis is that the set of peaks is generated by the peptides corresponding to the predicted protein.

Calculation for the Likelihood of the Alternative Hypothesis



Automatically, the log likelihood of the alternative hypothesis is maximized if and only if all matched peaks and no miss-matched peaks are included in the calculation.

Calculation for Likelihood of Matches



BUPID assumes that the experimental error follows the Gaussian Distribution. Thus the likelihood of a match is the z score of the mass difference. However if an experimental peak is adjacent to several theoretical peaks (peptides from the same protein or different proteins), chances are it may not be the result of the closest one. Thus it is necessary to calculate the cumulative likelihood of the experimental peak given all possible peptides exist in the sample. The bottom line is that any peak can be the result of a peptide with any mass. An far-fetched match still stands a chance in the nature, although much more impossible than a hit right in the bull's eye.

Background probability



The background probability is generated by digesting all proteins the database. Assume that restriction sites appear randomly in the protein sequence with probability P. Length of peptides generated from the digestion would, therefore, be proportional to the exponential of (1-P). Peptide mass roughly follows the same distribution, as shown in the plot. Notice that miss cleavages generate more peptides and have larger peptide mass in average.

Lines in red/yellow, green and blue each correspond to no miss cleavage, one miss cleavage, and one or two miss cleavages.

Mass difference between peaks and their corresponding peptides



Peptide mass

Due to the non-uniform background, the log likelihood ratio of a match is affected by both the difference and the absolute mass of the peak-peptide pair. The threshold (the point where the likelihood that a peak is generated by the specific peptide is equal to the likelihood that the peak is generated by a random background noise) rises as peptide mass becomes larger.

Number of peaks in the spectrum

The more peaks in the spectrum, the easier for a protein to match and get a high score. Therefore the background probability is adjusted according to the number of peaks in the spectrum. This serves as a internal quality-control scheme.

Correction for larger and/or heavier proteins

The matching score of a protein is the sum of the log likelihood ratios of all matched peptides, adjusted for the size and mass of the protein.

Size of the protein

Proteins with more restriction sites in their sequences generate more peptides after the simulated digestion. They are in advantage against small proteins to match more experimental peaks in the spectrum. This tendency has to be adjusted before all scores could be compared on the same basis.

Mass of the protein

Once again, peaks with smaller mass are easier to match with peptides because there are so many of them. The fact that light-weighted peptides are more abundant shifts the peptide mass tolerance and favors large proteins in scoring. This is also normalized before the search results are ranked.

	M	ASCOT	SEARC	BUPID SEARCH RESULT							
Sample	HBA	HBB	HBG1	HBG2	HBA	HBB	HBG1	HBG2	HBD		
A12	2	1	6	4		2	1	7	3	4	
A02	2	1				2	1			3	
A17	2	1			57	2	1		15	3	
B1	2	1				1	2	23		4	
B4	2	1			3	3	1	19	2	4	
B6	3	1			117	2	1			3	
B8	1	2				1	2	59		7	
B18	3	2	4		1	2	3	8	1	4	
B21	2	1				2	1			4	
B22	1	2				1	2			5	
B23	2	1				1	2			4	
B24	2	1				3	1			2	
C1	3	1				2	1			3	
C2	2	1				2	1			3	
C3	2	1				2	1			3	
C4	2	1			79	1	2			3	
C5	2	4	3	1	5	2	3	5	1	4	
C6	2	1				2	1			3	
C7	2	1				2	1			4	
C8	2	1		58		1	2			3	
C12	3	2	4	1		2	3	5	1	4	
C13	3	2		22		2	1	41	4	3	
C14	3	2		4		3	1	28	2	4	
C15	3	2	88	89		2	1	9	5	3	
C18	2	1				2	1			3	
C19	3	1				3	1	40	76	2	
C21	3	1				3	1			2	
C21	2	1				2	1			3	
C22	1	2	92			1	2			3	
C23	2	1				2	1			3	

Result comparison: Human Blood Sample

- A random blood sample.
- In 6 cases, BUPID found all 5 chains within the top 20.
- Delta chain are found in all cases within top 10 predictions.
- If the researcher is willing to go down the list, he/she'll be able to find all five chains in BUPID's results.

Result Comparison: Artificial Spectra (5-Protein Mixture)



- A: Both MASCOT and BUPID return highly accurate protein identifications.
- B: MASCOT returns accurate result while BUPID gives near-random protein ID.
- C: BUPID returns accurate result while MASCOT gives near or worse than random protein ID.
- ALL SAMPLES are correctly interpreted by BUPID or MASCOT or both.

E2F1 Protein



Search Results [show all 7 hits]

	Click for	E-value	Raw	Matched			Sequence	Protein		Merge
	Detail		Score	Peptides	(total	effective)	Coverage	Sequence Name		Hits
1.	E2F1_HUMAN	0.341067	99.417922	17	165	94	41.7%	Transcription factor E2F1 (E2F-1) (Retinoblasto	Q01094	
2.	CDN7_HUMAN	0.403940	31.316942	5	70	41	25.3%	Cyclin-dependent kinase 4 inhibitor D (P19-INK4D).	P55273	
З.	TA6P_HUMAN	0.409788	6.313207	1	3	2	18.9%	TAP2-associated 6.5 kDa polypeptide.	Q9Y3F1	
4.	OBRG_HUMAN	0.429088	12.373958	2	30	13	17.6%	Leptin receptor gene-related protein (OB-R gene	015243	
5.	HV3S_HUMAN	0.430084	18.928037	3	45	34	35.7%	lg heavy chain ∀-III region JON.	P01780	
6.	UPA5_HUMAN	0.430411	5.144575	1	З	1	100.0%	Unknown protein from 2D-page of plasma (Spot 13	P30091	
7.	VIME_HUMAN	0.432666	102.438835	16	300	234	37.0%	Vimentin.	P08670	
									Trypsin	
										Plot

E2F1 Protein

						Result	as Summar	γ	ante a	
MOWS	E Ma	#/45 (9b) 9b asses Cov 7 atched	Mea TC Da	n Da To Da	ta MS-Digest I Index #	t Protein MM (Da)/pI	Accession	ⁿ Species Protein Name	Nutler of	5 -
1 2.685e	+05 10	5 (35) 35.03	5.6 -0.03	42 0.1	28	53652/5.1		HUMAN Vimentin		
2 1.298e	+05 1	5 (33) 40.0 3	3.3 -0.02	36 0.05	42	46920/4.8		Transcription factor E2F1 (E2F-1) (Retinoblastoma binding prot HUMAN 3) (RBBP-3) (RRB-binding protein E2F-1) (RBR3) (Retinoblastom secondated receives 1) (RBR-1)		o o o si sio sio sio sio Probalility Bured Nove Score
								Visual system homenhov 1 (Transmintion factor VSV1) (Patinal	Co	oncise Protein Summary Report
3 325	6	(13) 25.0 1	3.3 0.027	4 0.2	58	38432/9.0		HUMAN inner nuclear layer homeobox protein) (Homeodomain protein RINK)		Format As Concise Protein Summary 🛩 Help
4 269	8	(17) 19.01	7.8 -0.03	40 0.1	58	55950/9.4		HUMAN Protein KIAA0141		Significance threshold p< 0.06 Max. number of hits 20
5 187	6	(13) 9.0 1	3.3 -0.09	03 0.2	21	88931/8.7		HLMAN Cullin homolog 3 (CUL-3)	_	
6 174	8	(17) 17.01	7.8 -0.06	51 0.1	59	63795/9.4		HUMAN Zinc finger protein 382	Ľ	Re-Search All Search Unmatched
7 142	5	(11) 16.01	1.1 -0.07	99 0.1	50	45936/4.8		HUMAN Keratin, type I cuticular HA3-I (Hair keratin, type I HA3-I)	1.	Hixture 1 Total score: 217 Expect: 1e-17 Queries matched: 29
8 140	5	(11) 14.01	1.1 -0.02	94 0.2	79	57848/5.6		HUMAN Activin receptor type II precursor (ACTR-II) (ACTRIIA)		Components (only one family member shows for each component):
9 120	8	(17) 9.0 1	7.8 -0.007	17 0.2	24	97669/7.3		HUMAN Probable nucleolar complex protein 14		<u>001024-00-00</u> Mass: 46891 Score: 123 Expect: 2.6e-08 Queries matched: 15
10 113	7	(15) 22.01	5.6 -0.03	42 0.2	32	45887/8.7		HUMAN Hepatocyte nuclear factor 4-gamma (HNF-4-gamma)		(E27)_HURAN) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from Q01094 Transcription factor E271 (E27-1) (Retinoblastoma binding ; R08570 Base 51466 Source : 0.8 Funct: 0.8 ref. 6 marries mathematics : 0.1 (Section 2016) (Se
153222 . x	x	.x.x		. x x .	. x x x	. x x x .	x			(VIRE HURAN) Vimentin
131346	x	x x x		x	× · · · × · ·	x . x x x	. x x	. x		01004-00-00-00 Mars 4601 Server 19 Funct 2 de-08 America antidadi 15
142209	. x .				×	x . x .	× .			(E271_HUMAN) Splice isoform Displayed; Variant Displayed; Conflict Displayed; from Q01094 Transcription factor E2F1 (E2F-1) (Retinoblastoma binding ;
106972		.xxx		. x		×	×	х.		<u>001094-00-94-00</u> Mass: 46921 Score: 123 Expect : 2.6e-08 Queries matched: 15
149937		x	×>	.x.)		×				(E27] NUNAN) Splice isoform Displayed; Variant dbNP:3210176; Conflict Displayed; From Q01094 Transcription factor E2F1 (E2F-1) (Retinoblastoma bind: 021024-00-03-00 Mars: 46004 Search: doi:02.000400000000000000000000000000000000
47768 .x	x			. x	x . x . x	x	X	• •		(E2F1_NUMAN) Splice isoform Displayed; Variant dbSNP:3213174; Conflict Displayed; from Q01094 Transcription factor E2F1 (E2F-1) (Retinoblastoma bind;
40178 . x			x .		×			х.		<u>001094-00-92-00</u> Mass: 46923 Score: 110 Expect: 5.1e-07 Queries matched: 14
131176		x			x	.×××.				(E27)_HURAND Splice isoforms Displayed; Variant dBSWF3213173; Conflict Displayed; from Q01094 Transcription factor E291 (E27-1) (Retinoblastoma bind; 001094-00-011-00 Many 44872; Sonra: 90 Panetic 4, 2-0.6 Parvia matched; 13
01630		~ ~		I COM	~ ~		C ISSIE IN INT			(2) I UNDER Gales index links and links and links for a data for an analytic factor for the links and links an



MS-FIT

ExPASy Home page	Site Map S	earch ExPASy	Contact us	Proteomics tools	Swiss-Prot		ProF	ound	- Sear	ch Result Summary	The Rosler	uller U	sivenity Edition
	Search Switt-Prot/Tri	Mile 💌 for	60 0	lear			Protein	Cendidat	en for a	earch er/cell reported seasons (94820 sequer	ces searched	- 1	i lun la
	This page requir	es JavaScript enabled	to be fully function	al.			+1 1	.D++000	1.82	r gi12699110rd04P 005216.11 E2F transcription stimoblastoma-associated protein 1 [Flomo aspin	(furtor 1, rms)	E 4	18 4690 @
	Moreover you may n	ot see all the informat	ion available for thi	s page			+2 9	Da-009	0.65	gi4507895fredbiP_003371.1 vimentin [Homo av	epiens]	16 3	1 33.67 4
		(More information)				. 11	+3 3	De-018	- 1	#1136430-B-BAA11502 II KIAA0115 (Home	supiens]	11 9	3 210.03 4
Idente version 19/04/2005 Input	ummary Printable page H	elp Note: most headings a	re clickable, even if they d	on't appear as links.		Ald	+4 3	.6e-019	2.	r <u>g[Si21434;dbd[IAA331211]</u> actin bioding prot apiens]	eno HJ 001900 Alexan	d 3	3 620.29 6
Sample C Documents and Settings	mccomb AD My Documents	E2F1 test Mascot E2F1	ReCal 2 bit / Peaks	45 / Mass [1011.541 - 2402.26]	/ Intensity [1 - 1]	Sa	+5 3	.5e-019	- 3	r <u>gf7512967[piif[T02345</u> hypothetical protein KL fragment)	AA0334 - human	10 E	20 191 29 4
Date 29/05/2005 20:01:351 ITC			2.0000-0_0.01 / 00010	in / i and [control i closedo]	in more study [a a]	Da	+6 3	1=-019	- 1	gh4034829gh6AAD10838.1] kendein		1 3	4 376.32 4
Release Swiss-Prot Release 47.1 of 24 Proteins - Sequences 18500 / In rang	May-2005: 181821 entries 10982 / After digestion 417					Re Pre	*7 1	7=-019		r g21127713wf8F 006677.2 transcription slor renereption factor CA150, TATA box binding j isconsted factor, F8IA polymersse II, 5, 1504D, rotein-associated factor 25 (Homo regimes)	igation regulator 1; protein (TBP)- TATA box binding	4 8	123 88
 First Analysis on 417 seque Second Analysis on best 2 	nces : After Alignmen f first analysis : After Alignmen	t 5 / After pValue thresh t 2 / After pValue thresh	old 2 old 2 / Displayed 2				8 1	7e-019	-	r gj45033333re@HF 001371 [] dedicator of cytok of cyto-kinesis I [Homo supiens]	menir 1, dedicator	2 7	3 215.36 @
Peptides Generated 299836 / Matching	a peak 10133 / Average per pr	otein 27				Pe	*9 4	.7e-020	- 1	r <u>m627367partA45259</u> deemoyokin - Iranan (fra	(gnenta)	1 6	3 312.48 🖲
Random Generated 100000 / After Ali	nment 122 / Mean log(score) -	2.26 / Stdvar log(score)	0.63			Ra	10 4	4+-020		mi4506757frefNP_001026.1] symodime receptor	e 2 [Homo supiens]	2 3	7 364.41 8
New : Pr	tein with 50% sequence similar	ity are merged. (This se	tting is available in th	e Display tab of the form)		1	1. To man 2. Highly d Input So	niter protect	g unmator requests	ted manies, slick the symbol (1). er were given the same rate (10 user stick "4" to expan	disentariti	_	
Rank pValue Hits AC II	Descripti	on Mw.p	I Cov Taxon	TaxID Identified		Ra		Date & Ti	ID F2F	May 29 19 30 47 2003 UTC (Search Tune 1 17 se	(e)		
	40 first charac	ters kDa	% Simplified	Exact				Datah	ane MCD	Inr (2005/05/01)			
1 4.6e-9 7 001094 E2F1_F	JMAN Transcription factor E2F	1 (E2F-1) (Retin 47 4	.8 17 Homo sapier	is <u>9606</u> <u>Validate</u>			Taxes	any Caleg	ery Hon	o espirete (branat)			
2 1.1e-6 5 P08670 VIME_H	UMAN Vimentin.	54 5	.1 14 Homo sapler	ns <u>9606 Validate</u>		1	Protei	n Mass Ras	nge 0-3	000 kD s			
							Pre	tein yl Ru Secoch	age 0.0-	149 In sector other			
		Resubmit				- 11	Dig	eut Chemis	try Try	an brownin origi			
	(alternation					2	M	ax Missed	Cut 4				
araphical visualisation of the results :	owaraph					Gra	Modifications Note						
								Charge St	tate MH-	•			
							r	eptide Mas	545				

<u>Aldante</u>

ProFound

BUPID Web Interface



The standard interface of BUPID offers standard parameters.

Users also have access to more options one mouse-click away.

[Advanced Options]

Search Results [show 5 more hits] [show all 169 hits] Click for Matched Protein Raw Sequence Merae score Detail Sequence Name Hits Score Peptides (total effective) Coverage RS30 HUMAN 0.275882 16.704674 39 16 39.6% 40S ribosomal protein S30. Q05472 1. 3 K1CL HUMAN 2. 0.333126 75.270339 14 99 79 24.9% Keratin, type I cytoskeletal 9 (Cytokeratin 9) ... P35527 K1CJ HUMAN 0.343009 65.151715 22.3% P13645 З. 12 91 72 Keratin, type I cytoskeletal 10 (Cytokeratin 10 ... K2C1 HUMAN 0.344604 88.926094 33.4% Keratin, type II cytoskeletal 1 (Cytokeratin 1) P04264 4. 18 129 100 CAHA HUMAN 0.351723 38.177582 7 53 43 15.9% Carbonic anhydrase-related protein 10 (Carbonic ... Q9NS85 5. H13 HUMAN 6 57 18.6% P16402 6. 0.363524 27.715259 113 Histone H1.3 (Histone H1c). RL31 HUMAN 0.367282 21.808564 7. 4 61 36 32.8% 60S ribosomal protein L31. P12947 25.0% Zinc finger protein 125 (HZF-3) (Fragment). 8. Z125 HUMAN 0.367869 21.903299 4 33 24 P35274 RL44 HUMAN 0.370291 16.519316 31.3% 60S ribosomal protein L44 (L36a). P09896 9. 63 28 3 HV1C HUMAN 0.375541 17.414945 Ig heavy chain V-I region ND precursor (Fragmen ... 3 19 34.2% P01744 10. 21 P628 HUMAN 0.379520 11.876338 2 15 12 71.1% Protein PRO0628. Q9UI54 11. BAXC HUMAN 0.381664 10.699630 BAX protein, cytoplasmic isoform gamma. 2 11 10 57.5% Q07815 12. BUB3 HUMAN 0.383651 39.692537 Mitotic checkpoint protein BUB3. 7 21.6% 043684 13. 67 51 40S ribosomal protein S17. 14. RS17 HUMAN 0.384568 22.382246 5 53 31 22.4% P08708 S113 HUMAN 0.386341 17.076007 3 33 26.8% S100 calcium-binding protein A13. Q99584 15. 20 P09132 16. SR19 HUMAN 0.387692 22.615362 4 55 33 29.2% Signal recognition particle 19 kDa protein (SRP ... 17. ARGR HUMAN 0.389733 32.298601 6 107 55 26.1% Arginine-rich protein. P55145 18. K1CX HUMAN 0.390709 49.354253 11 93 75 16.0% Keratin, type I cytoskeletal 17 (Cytokeratin 17 P08779 IG2R HUMAN 0.391973 11.910041 Putative insulin-like growth factor II associat ... 19 2 15 12 51.6% P09565 TRIF HUMAN 0.393250 26.328247 5 79 55 19.9% Troponin I, fast skeletal muscle (Troponin I, f ... P48788 20. $\mathbf{\nabla}$ Trypsin Plot

Final results are ranked by the **score**. Click on the link in the first column to view each individual proteins. Check boxes in the last column and click **Plot** for more data-processing (see bellow).

[help]

Protein view provides

- A plot of spectrum and matched theoretical peaks
- Sequence and coverage
- Matched and miss matched peptides





Sequence

201	SPYYNTIDDL	KDQIVDLTVG	NNKTLLDIDN	TRMTLDDFRI	KFEMEQNLRQ	
151	NEKSTMQELN	SRLASYLDKV	QALEEANNDL	ENKIQDWYDK	KGPAAIQKNY	
101	GGGYSSSGGE	GGGFGGGSGG	GEGGGYGSGE	GGLGGFGGGA	GGGDGGILTA	
51	GYGGGSSRVC	GRGGGGSEGY	SYGGGSGGGF	SASSLGGGFG	GGSRGEGGAS	
1	MSCRQFSSSY	LISGGGGGGG	LGSGGSIRSS	YSRESSSGGR	GGGGRESSSS	

.... etc.

Matched Peptides

Start -	End	Length Peptide Mass	Peak Mass	Mass Error	Missed Cl.	eptide Sequence
1 -	4	4 495.193379			0 M S	
1 -	28	28 2622.191483			1 MS	SCRQFSSSYLTSGGGGGGGGGSGSIR
5 -	28	24 2145.008674			0 QF	FSSSYLTSGGGGGGGGGGSGSIR
5 -	33	29 2725,269201			1 QE	FSSSYLTSGGGGGGGGGGGSGSSSSSS
29-	33	5 598.271097			0 8 8	SYSR
29 -	40	12 1276.579632	1276.741975	- 0.162343	1 33	SYSRFSSSGGR
34-	40	7 696.319105			0 F S	
34-	45	12 1080.506071			1 F S	SSSGGRGGGGR
41-	45	5 402.197536				
41-	.58	18 1618.708416			1 GG	GGGRFSSSSGYGGGSSR
46-	58	13 1234.521450	1234.659975	- 0.138525	0 FS	SSSSGYGGGSSR
46-	62	17 1649.721628			1 F S	SSSSGYGGGSSRVCGR
- 59 -	62	4 433.210748			0 VC	CGR
59-	94	363119.354019			1 VC	CGRGGGGSFGYSYGGGSGGGFSASSLGGGFGGGSR
63 -	94	32 2704.153841	2704.437975	- 0.284134	0 GG	GGGSFGYSYGGGSGGGFSASSLGGGFGGGSR
63 -	153	91 7513,231551			1 GG	GGGSFGY SYGGGSGGGFSASSLGGGFGGGSRGFGGASGGG
95 -	153	59 4827.088280				FGGASGGGYSSSGGFGGGFGGGSGGGFGGGYGSGFGGLGG

Protein Mixture View

Spectrum coverage: 35.4%



Protein Mixture

		Raw	Matched			Sequence	Protein	
	E-value	Score	Peptides	(total	effective)	Coverage	Sequence Name	
RS30_HUMAN	0.275882	16.704674	3	39	16	39.6%	40S ribosomal protein S30.	Q05472
K2C1_HUMAN	0.344604	88.926094	18	129	100	33.4%	Keratin, type II cytoskeletal 1 (Cytokeratin 1)	P04264
CAHA_HUMAN	0.351723	38.177582	7	53	43	15.9%	Carbonic anhydrase-related protein 10 (Carbonic	Q9NS85
								Trypsin

- Protein mixture view plots all matched and miss matched theoretical peaks along with the spectrum. Users also have the choice to show peaks of self-digestion (peptides of the restriction enzyme).
- All matched peaks are also plotted back to back for users to identify shared peptides or peptides with similar masses. In popular search engines such as MASCOT, matched peaks are removed after each round. Thus peptides overlap with another peptide with a stronger signal may not be picked up in the database search. BUPID doesn't remove peaks and is free from such occasions.
- Mass errors of all peptides are plotted together. A trend line generated from least-square-fit helps to identify the trend of the error. If, for instance, that the researcher believes that the trend line resembles the system error of the machine, BUPID offers an internal calibration option to adjust for such error.

Search for Post-Translation Modifications



- BUPID can search for post-translation modifications within a set of user specified proteins.
- For each protein sequence, BUPID creates a database with post-translation modification and variations, against which the spectrum is searched again. Final results are ranked by the p-value of their log-likelihood score.



This project has been funded in whole or in part with Federal funds as part of the NHLBI Proteomics Initiative from the National Heart, Lung, and Blood Institute, National Institutes of Health, under Contract No. N01-HV-28178. We are also grateful for technological support from Bruker Daltonics.