

ABRF-PRG05: De Novo Peptide Sequence Determination

C.W. Turck¹, A.M. Falick², J.A. Kowalak³, W.S. Lane⁴, T.A. Neubert⁵, B.S. Phinney⁶, S.T. Weintraub⁷, and K.A. West⁸

¹Max Planck Institute of Psychiatry, Munich, Germany; ²HHMI, University of California, Berkeley CA; ³National Institutes of Health (NIMH), Bethesda MD; ⁴Harvard University, Cambridge MA; ⁵New York University School of Medicine, New York, NY; ⁶Michigan State University, East Lansing MI; ⁷The University of Texas Health Science Center at San Antonio, TX; ⁸Galson Laboratories, East Syracuse NY

GOAL

The 2005 Proteomics Research Group (PRG) study was designed to assist ABRF members to evaluate their proficiency in sequencing unknown peptides that are not included in any published database. From a comparison of the peptides obtained by different strategies, participating laboratories will be able to gauge their own capabilities and establish realistic expectations for the approaches that were used.

ABSTRACT

A common request of proteomics core facilities is protein identification. However, in some instances primary sequence information for the protein in question is not present in public databases. In other cases, the amino acid sequence of a protein may differ in some way from the predicted sequence in the database as a result of gene mutation, gene splicing, and/or multiple posttranslational modifications. Thus, it may be necessary to determine the sequence of one or more peptides *de novo* to identify and/or adequately characterize the protein of interest. The primary goal of this study was to give participating laboratories an opportunity to evaluate their peptide sequencing capabilities. Samples containing 3-6 pmol each of five synthetic peptides with amino acid sequences not present in public databases were sent to 106 participating laboratories. At least one non-standard amino acid was present. A detailed analysis of the sequencing results is presented below.

INTRODUCTION

Proteomics core laboratories are often presented with unknown proteins to be identified. Sometimes, these are not identifiable by commonly-used strategies that involve tryptic digestion and database searching. There are several reasons why a peptide might not be identified by means of a standard database search using mass spectrometric data. It might be modified in some way that is not being considered by the database search program being used; it might not have a required sequence characteristic (e.g., no C-terminal K or R from a tryptic digest), or it might come from an organism for which the primary sequence is not known. Sometimes a homologous protein can be identified, but this requires that the sequences have a high degree of identity. For example, if an unknown protein is 95% identical to a known one, there is a greater than 60% probability that a 20-residue peptide from the unknown protein will have at least one substitution compared to the corresponding known peptide. A direct sequencing approach may then be the only viable approach.

The primary goal of the 2005 PRG study was to give each participating laboratory a chance to evaluate its capabilities and practices in the following areas:

Peptide sequencing

Methods for the identification of *unknown amino acids*

Use of software to assist in the interpretation of the *de novo* sequence data

The sequences of the peptides synthesized for this study are shown in Table 1. No specific approaches were recommended, although it was anticipated that tandem mass spectrometry and possibly Edman sequencing would be employed. Each of the laboratories that requested a sample was provided with a mixture containing 3-6 pmol each of five synthetic peptides which had solutions that were not present in any public database. The sample was supplied as a dried pellet that could be dissolved in most common aqueous solutions; one peptide (A1) proved somewhat difficult to dissolve. As with any "real-life" sample, there were minor contaminants present. There was either a K or an R at the C-terminus of each peptide; analogous to tryptic peptides; one peptide had a double "missed cleavage" and another contained hydroxyproline. Participants were asked to return experimental evidence for each sequence they determined, along with a completed web-based questionnaire.

METHODS

Synthesis The peptides were synthesized and purified at the following locations: A1, A2 and A3 at the HHMI Mass Spectrometry Laboratory at UC Berkeley; T50 at the NYU Protein Chemistry Laboratory; and J1 at the Macromolecular Structure Facility, Michigan State University. The synthetic peptides were analyzed by reversed-phase HPLC and MALDI-TOF mass spectrometry to verify purity.

Composition analysis - Amino acid analysis was conducted on small portions of A2, A3 and T50, individually dissolved in the appropriate volume of water to yield 1 mg/mL stock solutions. For each of these three peptides, 3 µL of the stock solution was added to an amino acid analysis tube. The blank contained 3 µL of 1% acetic acid. The samples were dried in a vacuum centrifuge, sealed and analyzed in duplicate for amino acid content using a Waters AccQ-Quant precolumn AAA in conjunction with a Waters 2690 HPLC equipped with a Waters 2475 fluorometer.

For distribution to requesting laboratories, 3-6 pmol of each peptide was added to a 0.5-mL, polystyrene tube and the peptide mixtures were dried in a vacuum centrifuge. Dried samples were sent to 76 laboratories in North America, 20 in Europe and 10 in other countries.

RESULTS

Sequence data were submitted by 48 laboratories. A summary of the study results, organized according to instrument configuration and ionization method, is shown in Table 2. A compilation of all results received is shown in Table 3. The following approaches were used: MS alone (42); Edman degradation (1); Edman degradation plus MS (5).

The great majority of laboratories reported the correct nominal peptide masses; peptide A2 was often found to contain an oxidized Met. Differences in sample preparation and use of derivatization prior to analysis did not influence the success rate for sequencing, although one group used a variety of derivatization strategies and obtained the correct sequence for 4 of the 5 peptides. Static neopentyl worked as well as on-line prefractionation by capillary HPLC.

Laboratories using a TIOF generally had a slightly higher success rate in obtaining the correct sequences for these peptides. These instruments typically use MALDI ionization; for this study it was not possible to assess the relative importance of ionization mode versus instrument type as related to the TIOF results. In addition, the scores for laboratories reporting use of both an ion trap and another type of instrument were notably higher than those using a trap alone. But, it is important to remember that there is a wide range of capabilities and levels of expertise among the participating laboratories, and the total number of responses is not very large.

Manual verification was used by all laboratories; software alone did not appear to be sufficient to provide complete sequences.

The success rates for sequencing the individual peptides varied (Fig. 1). This is most likely due to differences in the sequences. The internal Lys residues combined with the multiple Leu and Ile (scored as 0.5 if not distinguished) undoubtedly contributed to the low scores for peptide T50. Peptide A1 was the longest and, therefore, expected to be the most difficult.

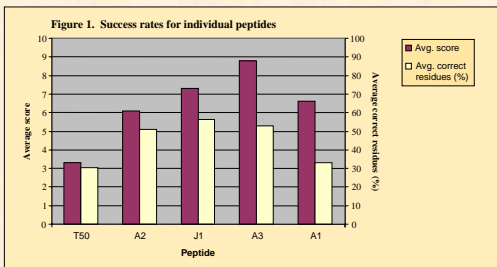


Figure 1. Success rates for individual peptides

DISCUSSION

The purpose of this study was to evaluate the capabilities of core laboratories to determine the sequences of peptides not found in any published database. Overall, the results show that this is an area that needs improvement in many core laboratories. A sufficient amount of each of the peptides was supplied such that sample quantity should not have been a limitation (although solubility issues might have caused problems for sequencing of peptide A1). Peptides T50 and A1 were the most difficult, probably due to specific sequence features of those peptides. In general, laboratories that reported using more than one type of instrument did slightly better than those that only used a single instrument. It is possible that facilities with multiple instruments might have a larger staff with more overall expertise. Too few cases in which Edman sequencing was used were reported to draw any conclusions. However, quantity limitations and time constraints make it generally not possible to separate the peptides in a digest sufficient for Edman analysis.

Although software is available that is designed to perform *de novo* sequencing, these programs do not as yet appear to be totally reliable. The peptides used in this study were by design not naturally occurring sequences. In many "real" cases, a partial sequence obtained by mass spectrometry, even with errors, can be linked to a protein by a BLAST search. But that would require that a protein of sufficient homology be present in a published database. While that strategy would not be successful for the synthetic peptides provided in this study, it should be routinely considered. It is clear that manual interpretation was necessary for the peptides analyzed in this study. Commercially available instruments can usually provide sufficient information to determine the sequences of most unknown peptides. However, it is critically important not only to acquire the spectra with the requisite mass accuracy and resolution, but also to be skilled in data interpretation. For example, there are two Hyp residues in peptide J1. The residue mass of Hyp (113.04768) is 36.4-mu lower than Ht. (113.08406). Using some commercial instruments, it is possible to measure CID fragment masses with sufficient accuracy to distinguish between these residues.

Finally, expertise in *de novo* sequencing is clearly essential, regardless of whether the data are acquired by mass spectrometry or Edman analysis or both. Whereas proteins that are present in a published database can be identified on a routine basis by scientists who are not experts in interpretation of mass spectra, the same cannot be said for proteins for which sequences are not included in any database. The results of this study provide excellent justification for core laboratories to have not only state-of-the-art instrumentation but also personnel with expertise in instrument operation and data analysis.

CONCLUSIONS

- Most core laboratories need to improve their capabilities for *de novo* sequencing. The average success rate in this study was relatively low. Note that this study addresses issues that are very different from identifying a protein that is in a database.
- MALDI ionization and TIOF instrumentation appeared to be more successful than the alternatives, but the numbers of participating laboratories are too small to make this finding statistically reliable.
- No sample preparation or derivatization strategy was notably more successful.
- Laboratories that used more than one type of instrument were slightly more successful than those that only used a single type of instrument.
- Software for *de novo* sequencing alone was not sufficient for successful sequence analysis of the test peptides.
- Expertise in MS data acquisition and manual interpretation is essential for success.

ACKNOWLEDGMENTS

We thank David S. King of the HHMI Mass Spectrometry Laboratory at the University of California, Berkeley, for synthesis and purification of peptides A1, A2, and A3, Joe Leykam at the Macromolecular Structure Facility at Michigan State University for synthesizing peptide J1 and for the amino acid analyses, Ron Beavis and Janet Brostrovian at the NYU Protein Chemistry Laboratory for the synthesis of peptide T50, Vieve Shetty, Chongfeng Xu and Yun Lu of the NYU Protein Analysis Facility for mass spectrometry analysis of the samples, Dawn Maynard of the NIMH at the National Institutes of Health for mailing and receiving correspondence that the participants remained anonymous, and Debra Diana of the NYU Skirball Institute of Biomedical Medicine for receiving confirmatory data.

Table 1. Amino acid sequences of the five peptides present in the unknown sample

Peptide	No.	M _r (Da)	MH ⁺	Sequence
T50	1	1192.8276	1193.8349	LGAIILKLIKLPK
A2	2	1395.6610	1396.6683	AYTFNMGQHSLSK
A3	3	1463.7665	1464.7738	VYKPHYASHPSPVYK
J1	4	1504.7316	1505.7389	GVPGADIIFYEANPR
A1	5	2327.1340	2328.1413	FPHANSGEWPDLYVYVNER

Hyp, hydroxyproline; masses listed are monoisotopic

Table 2. Summary of Instrument Configuration and Ionization Mode Utilization

Mass analyzer ¹	Number of laboratories ²	Average score ³
<i>Single instrument used</i>		
qTOF	13	38.5
Ion trap	10	24.1
TIOF	9	43.2
<i>One or more instruments used⁴</i>		
qTOF +	22	41.7
Ion trap +	18	31.1
TIOF +	13	46.4
Ionization mode		
ES	24	34.9
MALDI	13	43.8
ES and MALDI	6	35.0

¹qTOF, quadrupole time-of-flight; Ion trap, 3-D or linear trap; TIOF, tandem time-of-flight

²Four laboratories did not report the ionization method used.

³Average total score for all peptides analyzed by the indicated instrument or ionization mode

⁴This category represents each instance of the use of the indicated instrument. A number of laboratories reported use of more than one mass spectrometer to generate sequence information; however, details were not always provided about which specific instruments were used for each sequence analysis. For this table, if a specific instrument was listed, it was included in the appropriate category.

Table 3. Results compilation

Identifier	Total score	Peptide sequence (first choice) and score					Ionization method	Instrument type						
		Score T50	Score A2	Score J1	Score A3	Score A1								
13579A	66.0	7.0	11.0	LGAIILKLIKLPK	12.0	AYTFNMGQHSLSK	13.0	VYKPHYASHPSPVYK	14.0	GVPGADIIFYEANPR	20.0	FPHANSGEWPDLYVYVNER	MALDI	TIOF
72079	84.0	8.5	11.0	(L)GAI(L)I(L)K(K)I(L)I(L)PK	12.0	AYTFN MoxGQHSLSK	13.0	VYKPHYASHPSPVYK	14.0	GVPGADIIFYEANPR	17.5	FPHANSWWPDLYVYVNER	ES	qTOF
715	84.0	9.0	11.0	LGAIILKLIKLPK	12.0	AYTFNMGQHSLSK	13.0	VYKPHYASHPSPVYK	14.0	GVPGADIIFYEANPR	18.0	FPHANSWPDLYVYVNIQDQR	ES, E	qTOF
20107	7.0	63.0	6.5	(L)GAI(L)I(L)K(K)I(L)I(L)PK	11.5	AYTFNMGQHSLSK	13.0	VYKPHYASHPSPVYK	14.0	GVPGADIIFYEANPR	16.5	(M)I(L)VANSGEWPDLYVYVNER	MALDI	qTOF
26019	62.3	7.5	(7)GAI(L)I(L)K(K)I(L)I(L)PK	12.3	VYKPHYASHPSPVYK	12.3	VYKPHYASHPSPVYK	12.0	GVPGADIIFYEAGGPR	19.0	FPHANSGEWPDLYVYVNER	ES	qTOF	
65214	61.5	6.0	KHPY(L)I(L)K(K)I(L)I(L)PK	11.5	AYTFNMGQHSLSK	11.5	VYKPHYASHPSPVYK	11.5	IRPGAD(L)IFYEANPR	19.5	FPHANSGEWPDLYVYVNER	MALDI	qTOF, TIOF	
10286	59.0	8.5	(L)GAI(L)I(L)K(K)I(L)I(L)PK	11.5	(Y)TFNMGQHSLSK	11.0	VYKPHYASHPSPVYK	8.5	rsQDL(L)IFYEANPR	19.5	FPHANSGEWPDLYVYVNER	ES	LIT, LIT-FT	
48011	58.0	7.0	LGAIILKLIKLPK	12.0	AYTFNMGQHSLSK	5.0	VYKPHYASHPSPVYK	5.0	VYKPHYASHPSPVYK	20.0	FPHANSGEWPDLYVYVNER	MALDI	TIOF, PSD	
dsu	V	56.0	4.0	kllkllkllkllk	9.0	YATFNMGQHSLSK	10.0	VYKPHYASHPSPVYK	12.0	GVPGADIIFYEAGGPR	18.0	FPHANSGEWPDLYVYVNER	MALDI	TIOF
12800	52.5	4.5	wr(L)IKKHhyhyPK	8.0	(A)YTFEGMLHSLSK	11.0	vykpslsppvyk	13.0	GVPGAD(L)IFYEANPR	16.0	FPHANSWPDLYVYVNER	MALDI	TIOF	
78384	52.0	7.5	(L)GAI(L)I(L)K(K)I(L)I(L)PK	9.5	(A)YTFMoxGQHSLSK	7.0	FDIK(Q)I(L)AS(L)I(L)SPVYK	10.5	IRPGAD(L)IFMoxYEAANPR	17.5	(L)I(L)MVANSGEWPDLYVYVNER	ES	LIT	
13579B	52.0	8.0	(L)GAI(L)I(L)K(K)I(L)I(L)PK	11.5	AYTFNMGQHSLSK	12.0	VYKPHYASHPSPVYK	13.5	GVPGAD(L)IFYEANPR	17.5	FPHANSWPDLYVYVNER	ES, MALDI	3DIT, qTOF	
8500	51.0	6.0	rsQDL(L)I(L)K(K)I(L)I(L)PK	11.0	AYTFNMGQHSLSK	13.0	VYKPHYASHPSPVYK	13.5	GVPGAD(L)IFYEANPR	16.5	FPHANSGEWPDLYVYVNER	ES, MALDI	LIT, qTOF	
51565	51.0	4.0	KLLKLLKLIKLPK	9.0	HPTFNMGQHSLSK	9.0	VYQPLAS(L)SPVYK	11.0	IRPGAD(L)IFYEANPR	18.0	FPHANSWPDLYVYVNER	MALDI	TIOF, PSD	
30109	48.8	5.0	QHL(L)hy(L)hy(L)K(K)I(L)hy(L)hy(L)PK	9.5	PHTFNMGQHSLSK	12.0	VYKPHYASHPSPVYK	12.0	RPQADIIFYEANPR	10.3	[568.3]SPWPDLYVYVNER	MALDI	TIOF, qTOF	
11010	48.3	7.0	L7IAILKLIKDL	10.0	AYTFNMGQHSLSK	13.0	VYKPHYASHPSPVYK	12.0	GVPGADIIFYEANPR	5.3	[235.19](214.05)(m[212.15])Lhy(L)hy(VVV[243.15])ES	MALDI, E3DIT, qTOF		
47223	44.5	6.0	rsQDL(L)I(L)K(K)I(L)I(L)PK	1.0	ag(L)SPPAG(L)I(L)SPVYK	7.0	vykpslsppvyk	14.0	GVPGADIIFYEANPR	16.5	FPHANSWPDLYVYVNER	MALDI	TIOF, PSD	
4318	V	42.0	2.0	K(L)I(L)I(L)K(K)I(L)I(L)PK	10.0	AYTFNMGQHSLSK	8.0	VYKPHYASHPSPVYK	4.0	GVDSAEYLAPGPR	18.0	FPHANSGEWPDLYVYVNER	MALDI	TIOF
51952	41.0	2.0	K(Q)I(L)I(L)K(Q)I(L)I(L)PK	11.0	AYTFNMGQHSLSK	13.0	VYQ(K)PhyASHPSPVYK	7.0	rsQDL(L)IFYEANPR	8.0	[938.57]W(L)VVYV[243.14]	ES	qTOF	
99999	41.0	5.0	(L)I(L)Q(I)I(L)I(L)K(Q)I(L)I(L)PK	11.0	AYTFNMGQHSLSK	11.0	VYKPHYASHPSPVYK	11.0	GVPGAD(L)IFYEAGGPR	11.0	GVPGAD(L)IFYEAGGPR	ES	qTOF	
73108	40.0	2.0	K(L)I(L)K(L)I(L)I(L)PK	11.0	AYTFNMGQHSLSK	7.0	GVPSLIFYEAGGPR	7.0	GVPSLIFYEAGGPR	0.0	PDL(V)FGWPDLYVYVNER	ES	qTOF	
17989	40.0	6.0	(L)GAI(L)I(L)K(Q)I(L)I(L)AG(L)hyPK	9.5	YATFNMGQHSLSK	13.5	GVPGAD(L)IFYEANPR	13.5	GVPGAD(L)IFYEANPR	11.0	GVPGAD(L)IFYEANPR	ES	qTOF	
98186	38.0	7.5	(L)GAI(L)I(L)K(Q)I(L)I(L)PK	10.5	AYTFNFGQI(L)SHyPK	9.5	VYI(Q)I(L)I(L)AS(I)I(L)PK	10.5	IRPGAD(L)IFYEAGGPR	0.0	TFNFI(L)IKHShyK	ES	3DIT, qTOF	
27408	38.0	7.0	(L)GAI(L)I(L)K(Q)I(L)I(L)PK	11.0	AYTFNMGQHSLSK	8.5	(V)I(Q)I(L)I(L)AS(L)I(L)I(L)PK	11.5	GVPGAD(L)IFYEANPR	0.0	AYTFNMGQHSLSK	ES	qTOF	
91741	34.5	3.5	(Q)I(L)I(L)I(L)K(Q)I(L)I(L)PK	11.5	HPTFNMGQHSLSK	11.5	PHYPIASPPVYK	14.0	GVPGADIIFYEANPR	5.0	FNFASEGW(L)VVYVYVDRK	MALDI	TIOF	
91573	34.0	1.0	TFNMGQHSLSK	11.5	AYTFNMGQHSLSK	11.5	VYKPHYASHPSPVYK	11.5	[558]GVPGAD(L)IFYEANPR	0.0	T(L)I(L)WANSGEWPDLYVYVNER	ES	qTOF, 3DIT	
51583	V	32.0	2.0	KLLKLLKLIKLPK	1.0	HPTFNMGQHSLSK	12.0	RPQADIIFYEANPR	12.0	EDHWANSVShyPhyHyVYVNER	MALDI	TIOF		
70091	31.0	6.0	LGAIK467.2PK	1.0	[221.12](427.0)NFSTK	5.0	VYKPHYASHPSPVYK	14.0	GVPGADIIFYEANPR	0.0	[878.2](424.1)wesi(L)Hy(381.3)	ES, E	LIT	
19351	30.0	6.0	AYTFN	6.0	AYTFN	12.0	VYKPHYASHPSPVYK	12.0	GVPGADIIFYEANPR	0.0	GVPGAD(L)IFYEANPR	ES, E	qTOF, LIT-FT	
27974	29.5	0.0	[243.8](L)I(L)K(Q)I(L)I(L)K(Q)I(L)I(L)PK	12.0	AYTFNMGQHSLSK	6.0	VYKPHYASHPSPVYK	8.5	RPQADIIFYEANPR	2.0	(7)VVYV(L)DPW(7)	ES	qTOF	
17017	26.0	4.5	(Q)I(L)I(L)K(Q)I(L)I(L)PK-NH ₂	10.5	VYKPHYASHPSPVYK	10.0	VYKPHYASHPSPVYK	8.5	RPQADIIFYEANPR	0.0	FPHANSWPDLYVYVNER	ES, MALDI, E3DIT, PSD		
12144	25.0	0.0	ayghpsvvalpr	2.0	agplssppvyk	11.0	VYKPLAS(L)SPVYK	12.0	gvpsdlyyeaggr	ES	qTOF			
32486	22.0	1.0	K(D)I(L)I(L)K(Q)I(L)I(L)PK	11.5	AYTFNMGQHSLSK	1.0	cmTFNkthLSK	8.5	RPQADIIFYEANPR	ES	MALDI 3DIT			
78544	21.5	0.0	(K)I(L)I(L)K(Q)I(L)I(L)ASHPYK	0.0	[235](MoxF)(114)(MoxF)(185)	9.8	VYI(K)P(L)I(L)hyAS200PVYK	11.5	RPQADIIFYEANPR	ES	LIT			
1487	19.5	5.0	rsQDL(L)I(L)PK	1.0	(L)I(L)K(Q)I(S)HPTFNSTK	1.0	[849]ssppvk	4.0	RPQADIIFYEANPR	ES	3DIT			
52104	19.5	1.0	VYKPHYASHPSPVYK	12.5	VYI(Q)I(L)I(L)ASHPSPVYK	12.5	VYI(Q)I(L)I(L)ASHPSPVYK	6.0	PDQLIFYEAGGPR + neutral loss of 157	8.5	RPQADIIFYEANPR	ES, MALDI	qTOF	
80053	19.5									19.5	FPHANSGEWPDLYVYVNER	ES	qTOF	
54321	18.5													
85035	V	18.0	5.0	LGAF(V)IQ(I)(7)PK	1.0	FST(F)FVHTHYK	8.0	VYKPLASHPSPVYK	0.0	VYKPHYASHPSPVYK	10.5	IRPGAD(L)IFYEAGGPR	ES	qTOF
10523	V	12.0	0.0	[243.0](L)I(L)I(L)2[363]3SK	1.0	[145.095]DS(L)I(L)I(L)GQ[382.44]	3.0	VYKPHYASHPSPVYK	4.0	RPQADIIFYEANPR	5.0	[1449.608](V)VVYVGR	ES	qTOF
87458	10.5</													