

Bioinformatics for Top-Down Proteomics

DAVID L. TABB, PH.D.

Overview

- Motivation: the *proteoform*
- Separations and fragmentation
- Identification in ProSight PTM
- Identification by spectral alignment
- Consortium and standardization

Humans have how many genes, transcripts, and proteins?

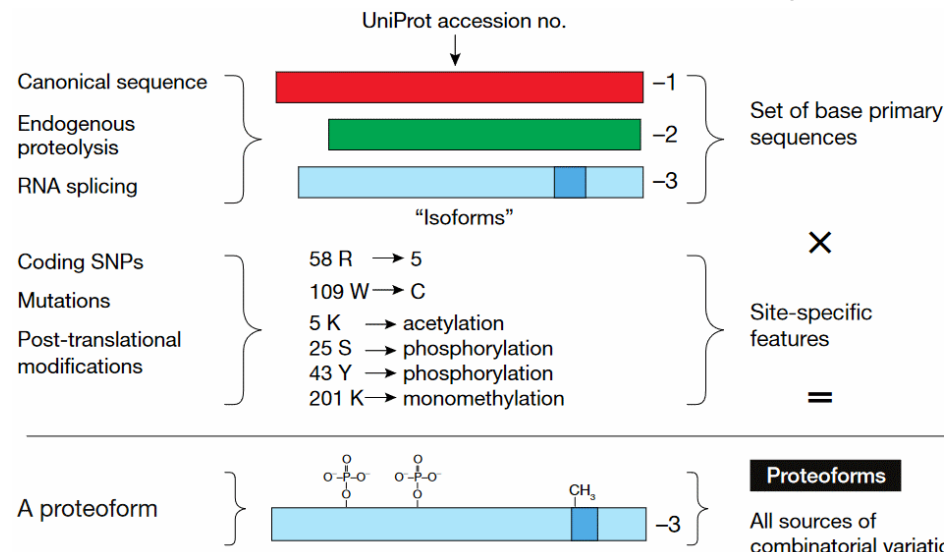
GENCODE 34, released April 2020:

- Total number of genes: 60,669
- Protein-coding genes: 19,959
(See also pseudogenes and lncRNAs)
- Protein-coding transcripts: 84,068
- Total # of distinct translations: 62,269
- Genes with multiple distinct translations: 13,717
(69% of all protein-coding genes)

A gene yields mRNA isoforms that yield multiple proteoforms

Proteoforms differ through three primary factors:

- Genetic variation and RNA splicing give different mRNA sequences.
- Proteolysis and nonsense mutations truncate mature proteins.
- Post-translational modification may dramatically alter activity.



Phenotype may be specific to a particular proteoform.

DNA sequencing is blind to post-translational modifications.

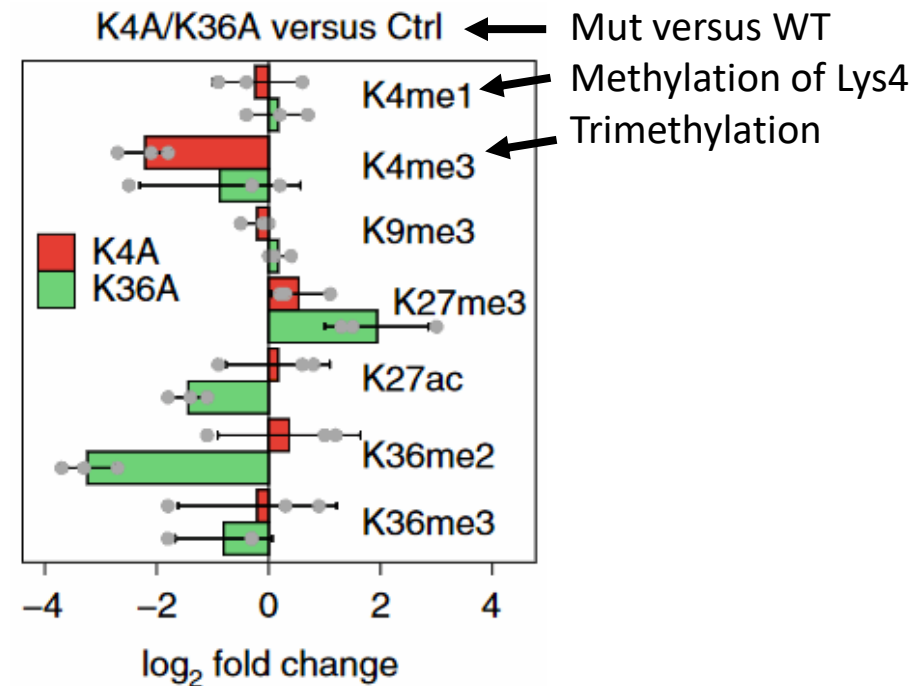
Known blind spots for shotgun proteomics

WHICH VEGF ISOFORM RISES IN CANCER COHORT?

```
>sp|P15692-2|VEGFA_HUMAN Isoform VEGF189
MNFLLSWVHWSLALLLYLHHAKWSQAAPMAEAGGGQNHHE
VVKFMDVYQRSYCHPIETLVDIFQEYPDEIEYIFKPSCVP
LMRCGGCCNDEGLECVPTESNITMQIMRIKPHQGQHIGE
MSFLQHNTCECRPKKDRARQEKKSVRGKGKGQKRKRKKS
YKSWSVPCGPCSERRKHLFVQDPQTCKCCKNTDSRCKAR
QLELNERTCRCDKPRR
```

```
>sp|P15692|VEGFA_HUMAN Isoform VEGF206
AMNFLLSWVHWSLALLLYLHHAKWSQAAPMAEAGGGQNHHE
VVKFMDVYQRSYCHPIETLVDIFQEYPDEIEYIFKPSCVP
LMRCGGCCNDEGLECVPTESNITMQIMRIKPHQGQHIGE
MSFLQHNTCECRPKKDRARQEKKSVRGKGKGQKRKRKKS
YKSWSVYVGARCCCLMPWSLPGPHPCGPCSERRKHLFVQDP
QTCKCCKNTDSRCKARQLELNERTCRCDKPRR
```

WHICH OF THESE PTMS CO-OCCUR ON HISTONE H3.3?



Analytical chemistry for intact proteins is more challenging.

“Numerically, the intact proteome appears to be a much simpler mixture than its corresponding peptide digests. In practice, however, protein-level fractionation and separation are daunting tasks due to the diverse physicochemical properties (e.g., size, charge, and hydrophobicity) and the wide dynamic range of the proteome.”

>sp|P01308|INS_HUMAN Insulin (11,981 Da)
MALWMRLPLALLALWGPDPAAAFVQHLGSHLVEALYLVCGERGFFYTPKTRREAED
LQGVQVELGGGPGAGSLQLPLALEGSLQKRGIVEQCCTSIICSLYQLENYCN

>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha (15,258 Da)
MVLSPADKTNVKAAWGVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
KKVADALNAVAHVDDMPNALSALSDLHAHLKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR

>sp|P07477|TRY1_HUMAN Trypsin-1 (26,558 Da)
MNPLLILTFVAAAALAPFDDDDKIVGGYNCSENSVPYQVLSNGYHFCGGSLINEQWVVS
AGHCYKSRIQVRLGEHNI EVL EGNQF INAAKIIRHPQYDRKTLNNDIMLIKLSRAVIN
ARVSTISLPTAPPATGKCLISGWGNTASSGADYPDELQCLDAPVLSQAKCEASYPGKIT
SNMFCVGFLEGGKDSQGDGSGGPPVVCNGQLQGVVSWGDGCAQKNKPGVYTKVYNYVKWIK
NTIAANS

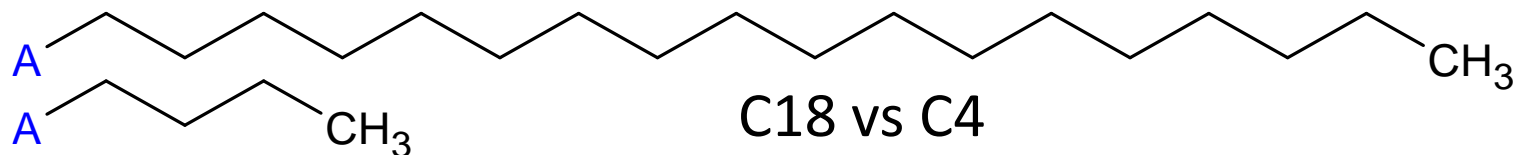
>sp|P60709|ACTB_HUMAN Actin, cytoplasmic 1 (41,737 Da)
MDDIAALVVDNGSGMCKAGFAGDDAPRAVFPISIVGRPRHQGVVMGMGQKDSYVGDEAQS
KRGILTLYPIEHGIVTNWDDMEKIWHHTFYNELRVAPEEHPVLLTEAPLNPKANREKMT
QIMFETFNTPAMYVAIQAVLSLYASGRITTGIVMDSGDGVTHTVPIYEGYALPHAAILRLDL
AGRDLDYLMKILTERGYSFTTTAEREIVRDIKEKLCYVALDFEQEMATAASSSSLEKSY
ELPDGQVITIGNERFRCPALFQPSFLGMESCGIHEHTTNSIMKCDVDIRKDLANTVLS
GGTTMYPGIADRMQKEITALAPSTMKIKIIAPPERKYSVWIGGSILASLSTFQQMWISKQ
EYDESGPSIVHRKCF

>sp|Q14533|KRT81_HUMAN Keratin, type II cuticular Hb1 (54,928 Da)
MTCGSGFGGRAFCISACGPRPRGCCITAAPYRGISCYRGLTGGFGSHSVCGGFRAGSCG
RSFGYRSGGVCSPSPCITTVSVNESLLTPLNLEIDPNAQCVKQEEKEQIKSLNSRFAAF
IDKVRFLQKQNKLLLETKLQFYQNRECCQSNLEPLFEGYIETLRREAECVEADSGRLASEL
NHVQEVLEGYKKKYEEVSLRATAENEFVALKKDVCAYLRKSDLEANVEALIQEIDFLR
RLYEEIEILQLSHISDTSVVVKLDNSRDLNMDCIIAEIKAQYDDIVTRSRAEASWYRSK
CEEMKATVIRHGETLRRTKEEINELNRMIQRLTAEVENAKQNSKLEAAVAQSEQQGEAA
LSDARCKLAELEGALQKAKQDMACLIREYQEVMMNSKGLDIEIATYRRLLEGEQRLCEG
IGAVNVVCSRSSRGVVCGLCVSGSRPVTGSVCSAPCNGNVAVSTGLCAPCGQLNTTCGG
GSCGVGSCGISLGVGSCGSSCRKC

>sp|P02768|ALBU_HUMAN Serum albumin (69,367)
MKWVTFISLLFLFSSAYSRGVFRDAHKSEVAHRFKDLGEENFKALVLIIFAQYLQQCPF
EDHVKLVNEVTEFAKTCVADESAENCDSLHTLFGDKLCTVATLRETYGEMADCCAKQEP
ERNECFLLQHKDDNPNLPRLVRPEVDVMCTAFHDNEETFLLKKYLYEIAARRHPYFYAPELLF
FAKRYKAAFTCECQAADKAACLLPKLDELDEGKASSAKQRLKCSLQKGFGERAFKAWAV
ARLSQRFPKAEFAEVSKLVTDLTKVHTECCCHGDLLECADRADLAKYICENQDSISSKLLK
ECCEKPLLEKSHCIAEVENDEMPADLPSLAADFVESKDVCKNYAEAKDVLFGMLFYEAR
RHPDYSVVLRLRAKTYETTTLEKCCAAADPHECYAKVFDEFKPLVEEPQNLIKQNCLEFE
QLGEYKFQNALLVRYTKKVPQVSTPTLVEVSRNLGKVGSKCKHPEAKRMPCAEDYLSVV
LNQLCVLHEKTPVSDRVTCKCTESLVNRRPCFSALEVDETYVPKEFAETFTFHADICTL
SEKERQIKKQ TALVELVKHKPKATKEQLKAVMDDFAAFVEKCKKADDKETCF AE EGKLLV
AASQAALGL

>sp|P02452|CO1A1_HUMAN Collagen alpha-1(I) chain (138,941 Da)
MFSFVDLRLLLLAATALLTHGQEEQVEGQDEIPITCVQNGRLRYHDRVWVKEPCRI
CVCNDGKVLCDVICDETKNCPGAEVPEGECCPVC PDGSESPDTQETTGEVGP KGTGPR
GPRGPAGPPGRDGI PGQPLPGPPGPPGPPGPPGLGNGFAPQLSYGYDEKSTGGISVPGP
MGPSGPRGLPGPPGAPGPQGFQGPPEGEPEGASGPMGRPPGPPGKNGDDGEAGKPRG
PGERGPPGPQAGRLPGTAGLPGMKGHRGFSGLDGAKG DAGAPGKPEGSPGENGAPGQ
MGPRGLPGERGRPGAPGAPAGARGNDGATGAAGPPGPTGPAGPPGFPAGVAKGEAGPQGP
RGSEGPQGVRGEPGPPGAGAAGPAGNPGADGQPGAKGANGAPGAGAPGFPAGRPSGP
QGP GPPGPPGKNSGEPGAPGSKGDTGAKGEPGPVGVQGP PGAGEEGKRGARGEPGPTGL
PGPPGERGGPGSRGFPAGDVAGPKGPAGERGSPGAPGKSGPEAGRPEAGLPGAKGL
TGSPGSPGPDGKTGPPGAPAGQDGRPPGPPGARGAQAGVMGFPGPKAAGEPGKAGERGV
PGPPGAVGAPAGKDGEAGAQQPPGAPGAGERGEQGPAGSGPFGQLPGPAGPPGEAGKPE
QGVPGDLGAPGSPGARGERGFPERGVPVQPPGAPRGANGAPGNDGAKGDAGAPGAPGS
QGAPGLQGMPPERGAAGLPGPKGDRGDAGPKGADGSPGKDGVRGLTGP IGP PGAPAGPD
KGESGSPGAPGTGARGAPGDRGEPGPPGAPGAGPPGADGQPAKGE PGDAGAKGDAGP
PGPAGPAGPPGIGNVGAPGAKGARGSAGPPGATGFPGAAGRVGPPGSGNAGPPGPPGP
AGKEGKGRPGRTGPAGRPGEVGP GPPGAGEKSGPGADGAPAGPTGPGQGIAGQRGV
VGLPGQRGERGFPLPGPSGEPGKQGPSGASGERGPPGPMGPPGLAGPPGESGREGAPGA
EGSPGRDGS PGAKGDRGETGPAGPPGAPGAPGAPGVGAPGKSGDRGETGPAGPTGPVGP
VGARGPAGPQGPGRDKGETGEQGDRIKGHRGFSGLQGP GPPGPPGSGPEQGPSGASGAPG
RGPPGSAGAPGKDG LNLPGPIGPPGPRGRGTGADGVPVPPGPPGPPGPPGPPSAGDFSF
LPQPPQEKAHDGGRYRADDANVVRDRDLEVDTTLSLSQQIENIRSPESRKNPARTCR
DLKMHSDWKSGEYWDPNQGCNLDAIKVFCNMETGETCVYPTQPSVAQKNWYISKNP KD
KRHWVFGESMTDGFQFEYGGQGS DPADVAIQLTFLRLMSTEASQNITYHCKNSVAYMDQQ
TGNLKKALLLQGSNEIEIRAEGNSRFTYSVTVDGCTSHTGAWGKTIVIEYKTKTSRLPII
DVAPLDVGAPDQEF GFDVGPVCFL

Protein recovery from LC columns requires shorter chains.



- Peptides / proteins are lured away from beads through increasing hydrophobicity.
- Column length, pump pressure, and pore size are substantial factors for separation.

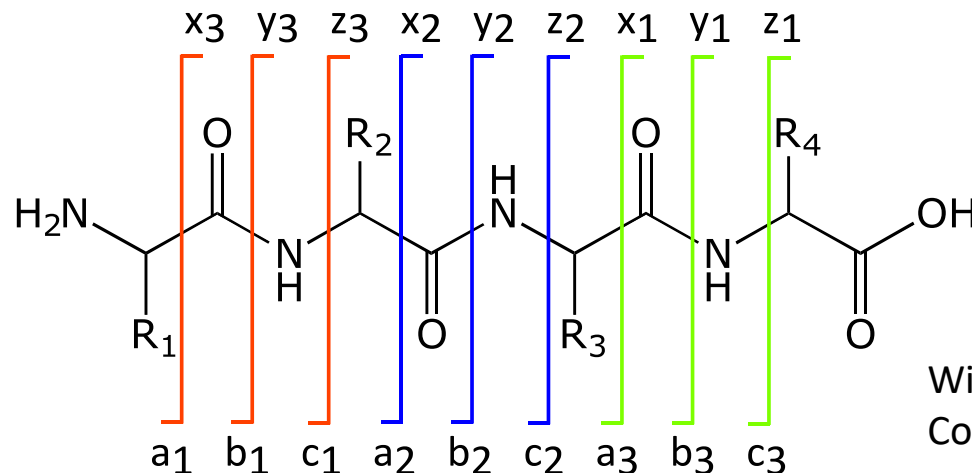
$$\text{pressure drop} \rightarrow P = \frac{\eta v L}{d_p^2}$$

← viscosity, velocity, and length

← particle diameter

Fragmenting peptides and proteins

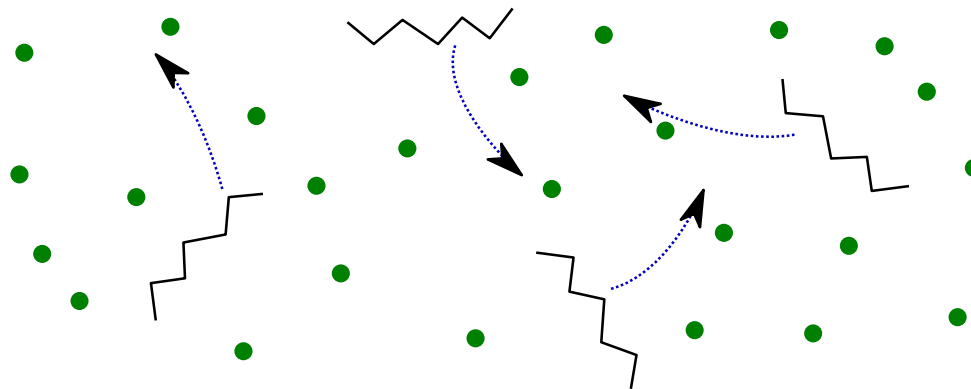
- Collision-induced dissociation
 - Standard quadrupole technique
- Electron transfer dissociation
 - Ion-ion reaction for gentle bond cleavage



Wikimedia
Commons

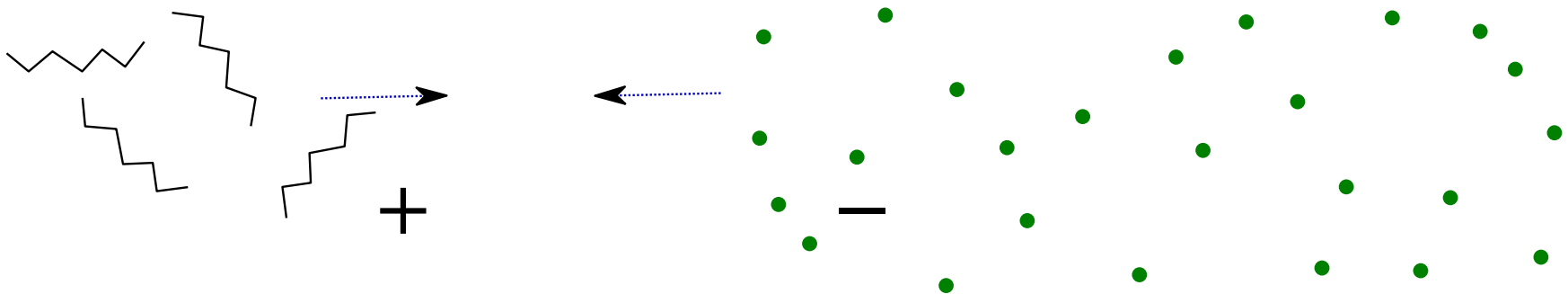
Collision-Induced Dissociation (CID)

- When the quadrupole adds energy to ions, they collide more frequently with gas molecules, gaining energy.
- Protons become mobile, destabilizing peptide bonds (creating *b-y* fragments).



Electron Transfer Dissociation (ETD)

- Charge draws positively-charged proteins to accept electrons from radical anions.
- An amino acid backbone cleaves between nitrogen and alpha carbon (c-z fragments).



Radical chemistry in peptides

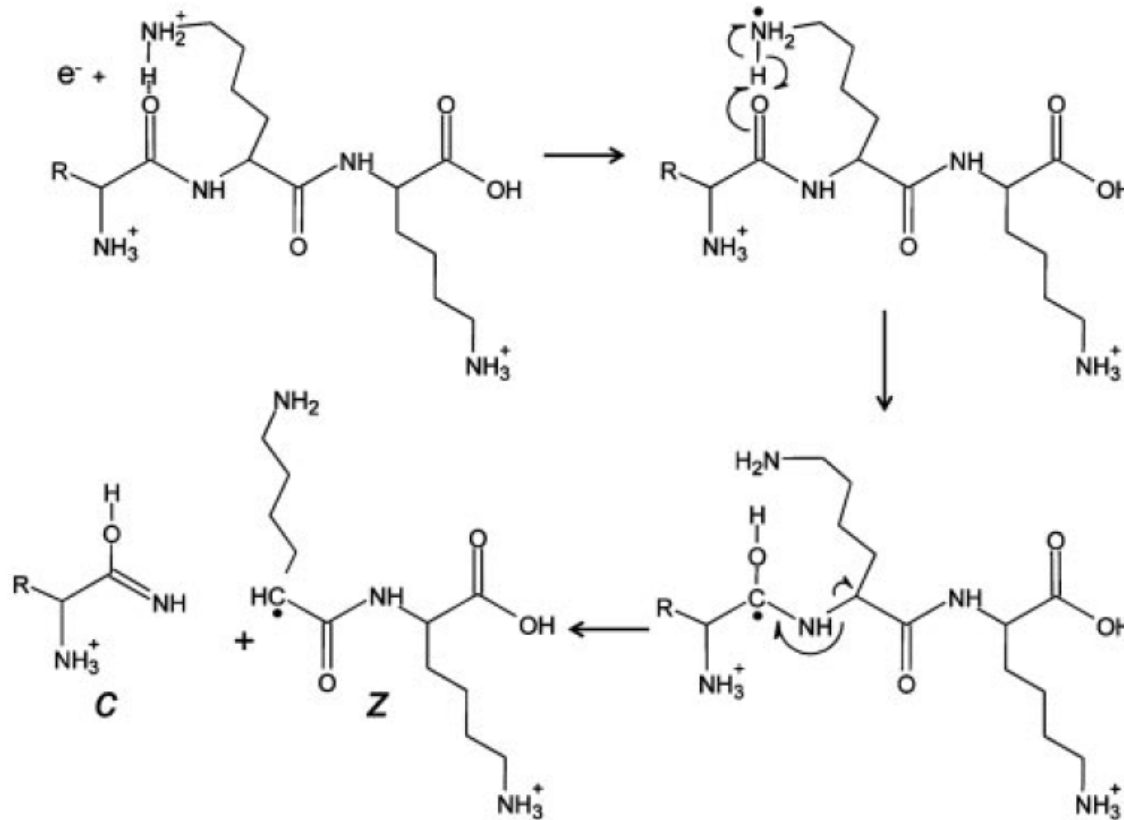
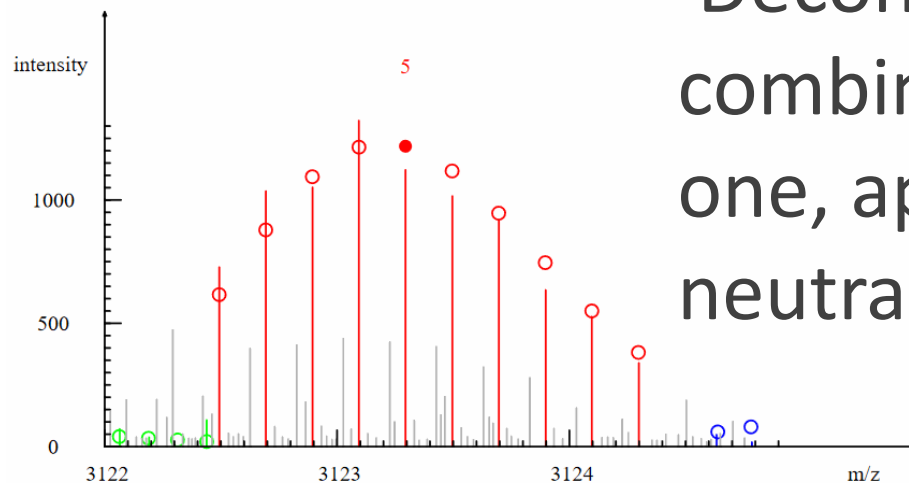


Fig. 1. Fragmentation scheme for production of c- and z-type ions after reaction of a low-energy electron with a multiply protonated peptide.

Intermission

Deconvolution goals

- A molecule appears in many isotopes and at many charges to produce isotopic *envelopes* in both MS and MS/MS scans.



- Deconvolution attempts to combine these peaks to only one, appearing at +1 or neutral monoisotopic mass.

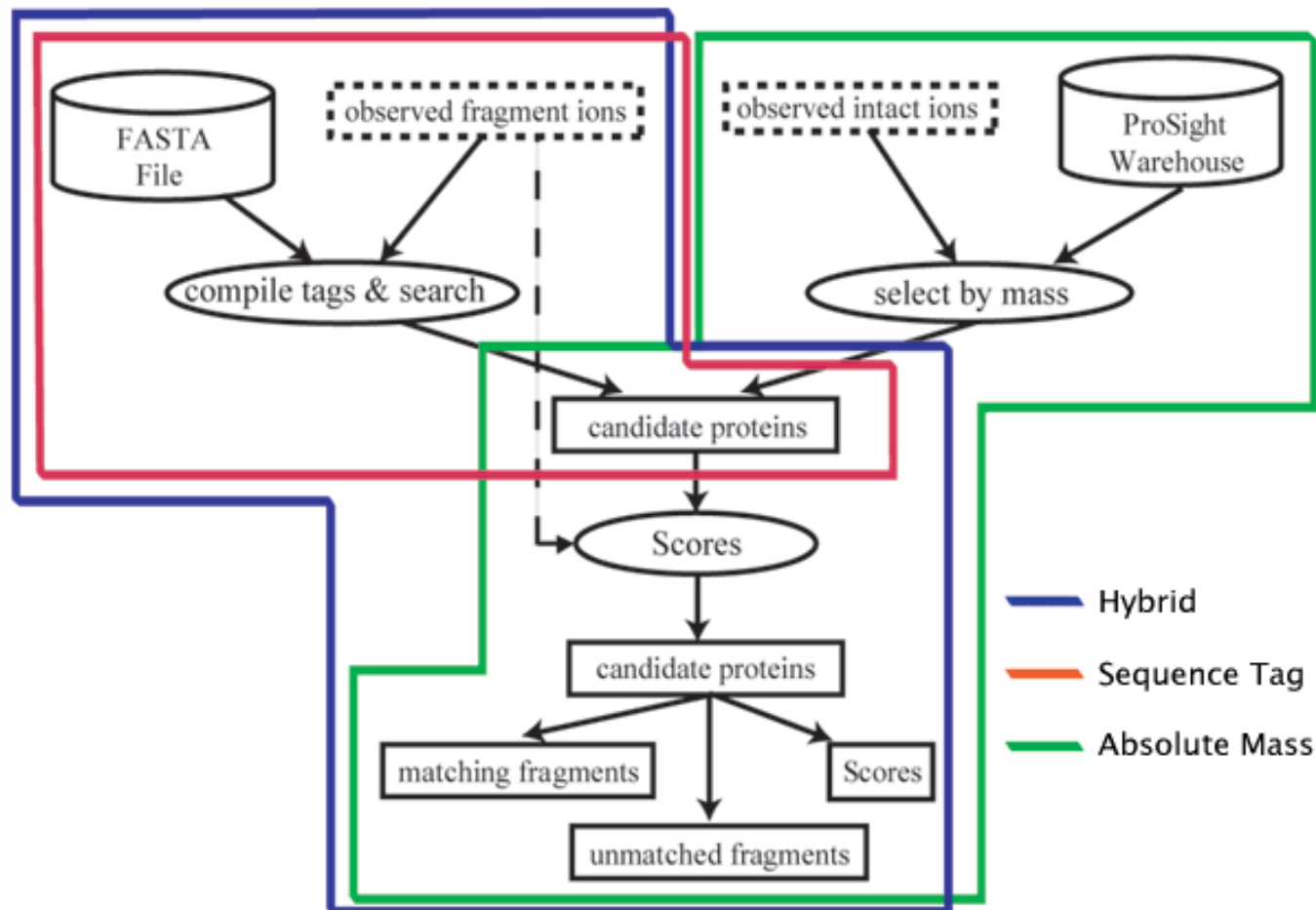
Limiting the sequence expansion for *PrSMs*

Proteoform
Spectrum Matches



- In shotgun PTM ID, we decorate peptide sequences with allowable mass shifts. Top-down benefits from known PTM annotation.
- Partial inferred sequence tags can narrow protein candidate list considerably.
- Signal peptides and other backbone cleavages take on special importance.

ProSight PTM schematic



Poisson scoring model

$$\blacksquare P = \frac{(xf)^{n_*} e^{-xf}}{n!}$$

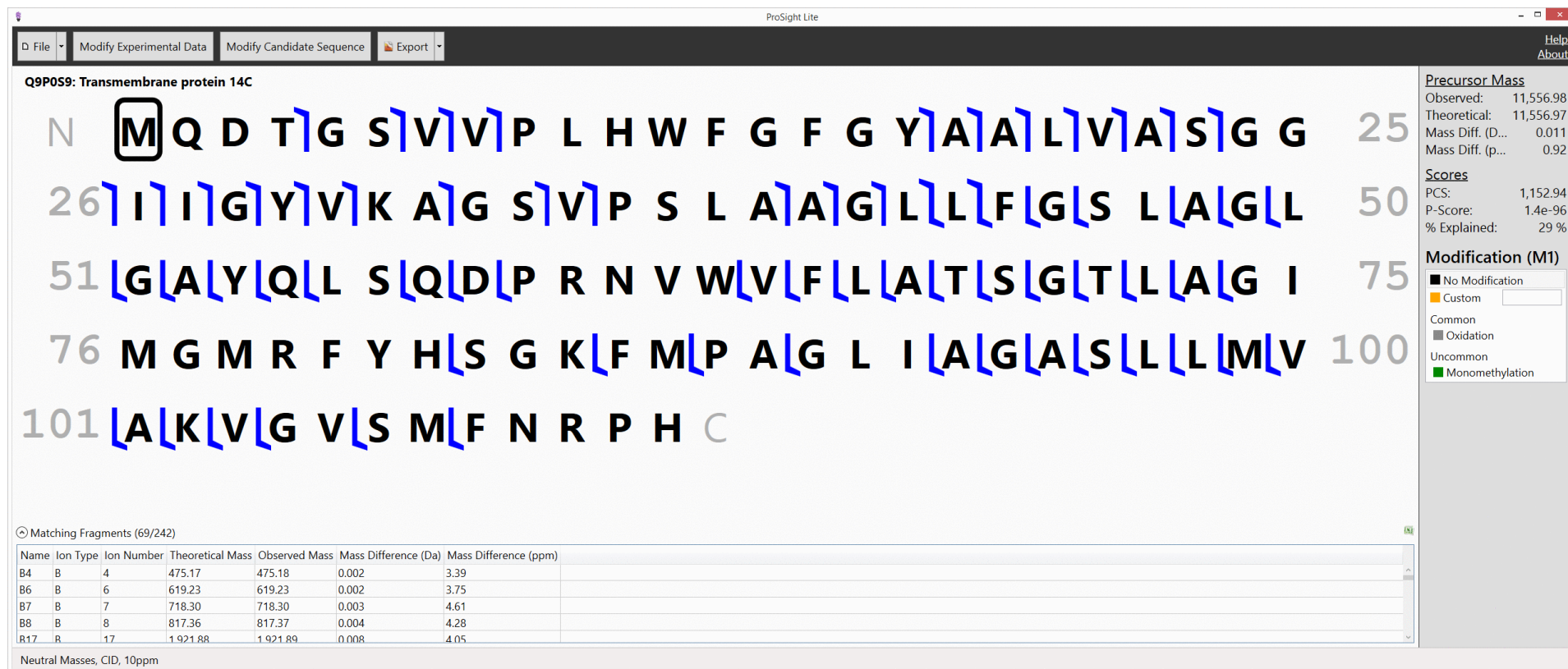
↑ Probability of random match
 ↑ Number of matched ions
 ↑ Number of predicted ions

See also OMSSA:
Geer J. *Proteome Res.*
(2004) 3: 958-964.

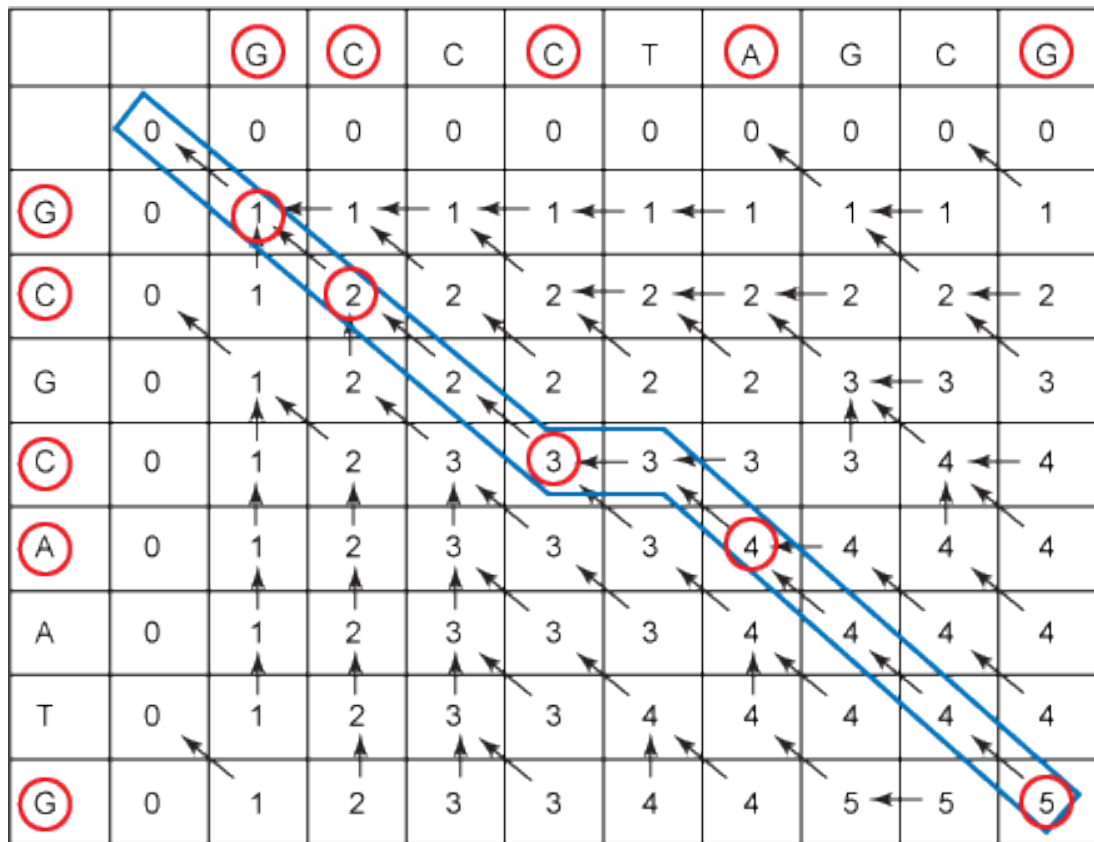
$$\blacksquare x = \frac{1}{1111.1} * 2 * (M_a * 2)$$

↑ Mean prob. of random match
 ↑ Freq-weighted avg AA mass
 ↑ Fragments per bond
 ↑ Mass tolerance

Visualizing supporting fragment ions



Dynamic programming is for more than sequence alignment.



- In Smith-Waterman, we use gaps to represent INDEL differences between sequences.
- Approach can be adapted to many additive optimization problems in proteomics!

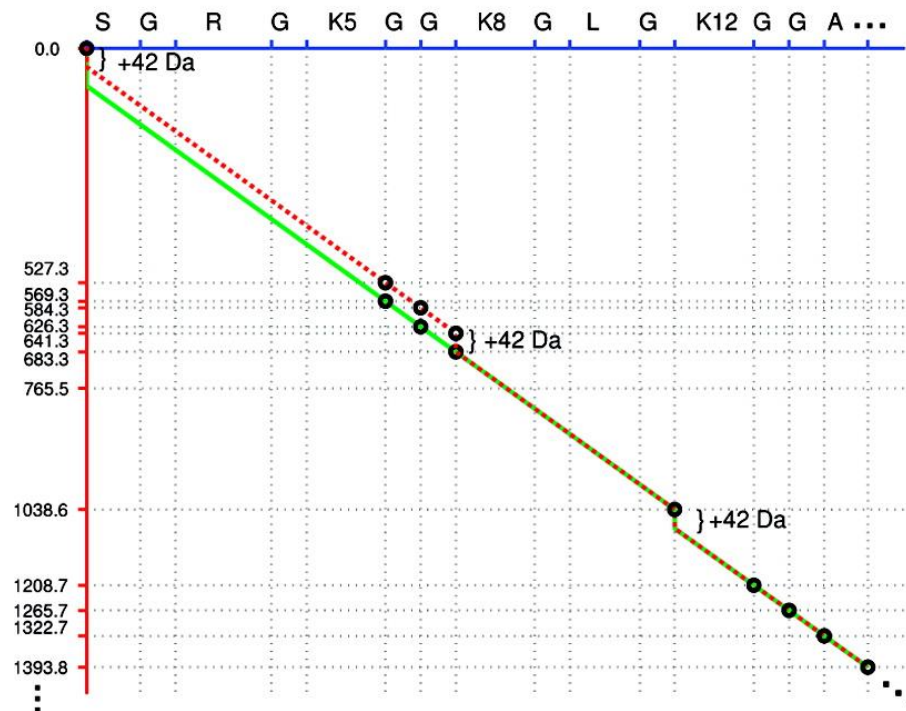
Dynamic programming in shotgun proteomics

- Infer sequences from MS/MS *de novo*
 - V Dancik et al. *J. Computat. Bio.* (1999) 6: 327.
- Align LC retention times of features
 - M Ono et al. *Mol. Cell. Proteomics* (2006) 5: 1338.
- Localize phosphorylations within peptide
 - F Saeed et al. *IEEE* (2012) 10.1109/BIBMW.2012.6470210
- Compute exact p-values for XCorrs of PSMs
 - JJ Howbert et al. *Mol. Cell. Proteomics* (2014) 13: 2467.

Problems leveraging sequence and PTM composition

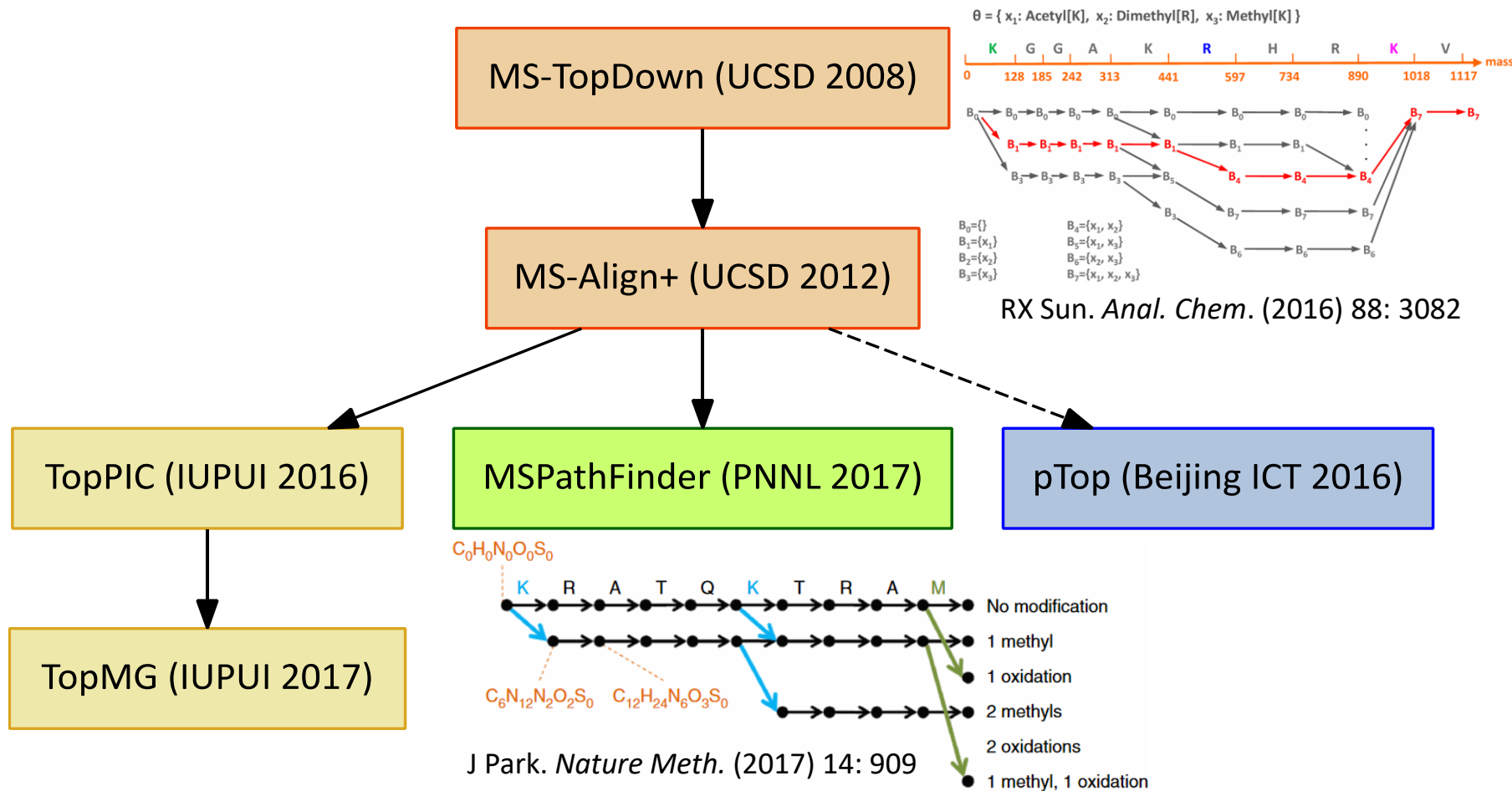
- “The candidate expansion method... leads to an exponential growth in the number of candidate protein forms that need to be considered.”
- “top-down spectral alignment may deal with as many as 10-20 PTMs to a protein”
- “one often deals with multiple isobaric protein forms in the same spectrum”

MS-TopDown and histone H4



- Acetylation adds 42 Da to N-terminus and two other sites in first 15 amino acids.
- Lys5 and Lys8 ambiguity may result from co-fragmenting variants in one MS/MS.

Inheritors of MS-Align+



How do we communicate proteoforms?

ProForma Proteoform Notation Rules

The Basics

1. The amino acid sequence is written. Ambiguous amino acids can be specified.

SEQUENCE SEQXXNCE

2. Modifications and are written inside square brackets.

SEQVK[Unimod:Label:13C(3)][Acetyl]ENCE

3. Tags contain descriptors in key : value pairs.

SEQVEN[mass:+14.02]CE

4. Multiple descriptors are separated by pipes.

SEQVEN[mod:Methyl|mass:+14.02]CE

Advanced Usage

6. Prefix tags define the key for all subsequent tags.

[RESID]+S[AA0037]EQVE[AA0234]NCE

[mass]+S[80]EQVE[14]NCE

[formula]+S[HPO(3)]EQVE[CH(2)]NCE

7. Terminal modifications are separated from the sequence by a dash.

[mass:-17.027]-QVENCE-[Amidation]

The Specifics

5a. Modification Name

PRT[Phospho]EFRM

PRT[Phosphothreonine(UniProt)]EFRM

PRT[O-phospho-L-threonine(RESID)]EFRM

PRT[O-phospho-L-threonine(PSIMod)]EFRM

5b. Database Accession

PRT[Unimod:21]EFRM

PRT[UniProt:PTM-0254]EFRM

PRT[RESID:AA0038]EFRM

PRT[PSI-MOD:MOD:00047]EFRM

5c. Mass

SEQ[mass:+15.995]VENCE

SEQ[mass:+16]VENCE

SEQ[mass:16]VENCE

5d. Chemical Formula

SEQVEN[Methyl|formula:H(2)C]CE

5e. Additional Information

SEQ[info: unstructured text]VENCE

ProForma in practice

Examples of Best Practices

- i. Histone H4 with several modifications. This example is human-readable and conforms to best practices.

[Acetyl]-S[Phospho|mass:79.966331]GRGK[Acetyl|Unimod:1|mass:42.010565]QGGKARAKATRSSRAGLQFPVGRVHLLRKGNYAERVGAGAPVYLAADVLEYLTAEILELAGNAARDNKKTRIIPRHLQLAIRNDEELNKLLGKVTIAQGGVLPNIQAVLLPKKT[Unimod:21]ESHHKAKGK

- ii. This is a valid and compact way of specifying Unimod accessions in multiple locations in the sequence.

[Unimod]+[1]-S[21]GRGK[1]QGGKARAKATRSSRAGKVTIAQGGVLPNIQAVLLPKKT[21]ESHHKAKGK

- iii. Extensive description of a modification using descriptors and IDs from different databases.

MTLFQLREHWFVYKDDEKLTAFRNK[p-adenosine| N6-(phospho-5'-adenosine)-L-lysine (RESID)| RESID:AA0227| PSI-MOD:00232|N6AMPLys(PSI-MOD)]SMLFQREL RPNEEVTWK

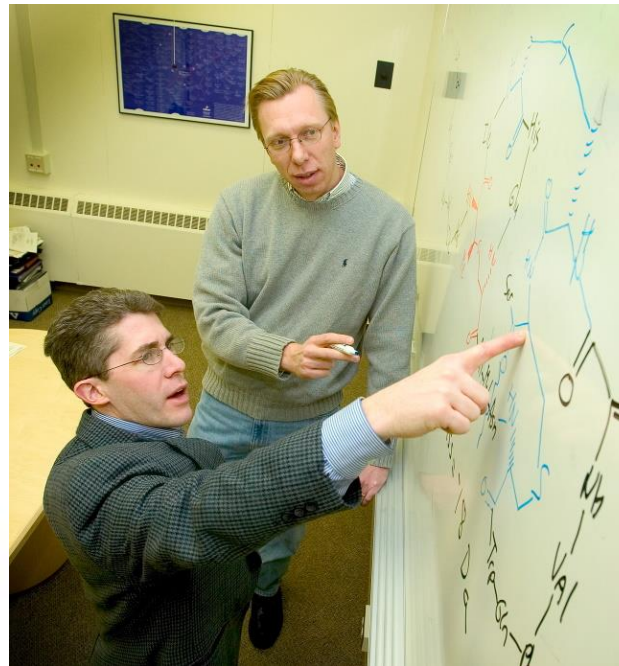
- iv. Unknown modifications are best described by their mass shift and marked as unknown.

MTLFQLDEKLTA[mass:-37.995001|info:unknown modification]FRNKSMLFQREL RPNEEVTWK

People of interest

- Neil Kelleher
- Lloyd Smith
- Ying Ge
- Jeff Agar
- Ljiljana Pasa-Tolic

- Pavel Pevzner
- Xiaowen Liu
- Rui-Xiang Sun



Takeaway messages

- If we digest proteins to measure them, we lose peptide relationships.
- Top-down proteomics relies upon high-resolution MS/MS and good separations.
- The dominant ProSight PTM framework has competition from alignment-based software.
- The CTDP seeks to broaden the use of top-down tech throughout biomedical research.