

**Proteomics Informatics -
Protein characterization: post-translational
modifications and protein-protein
interactions (Week 10)**

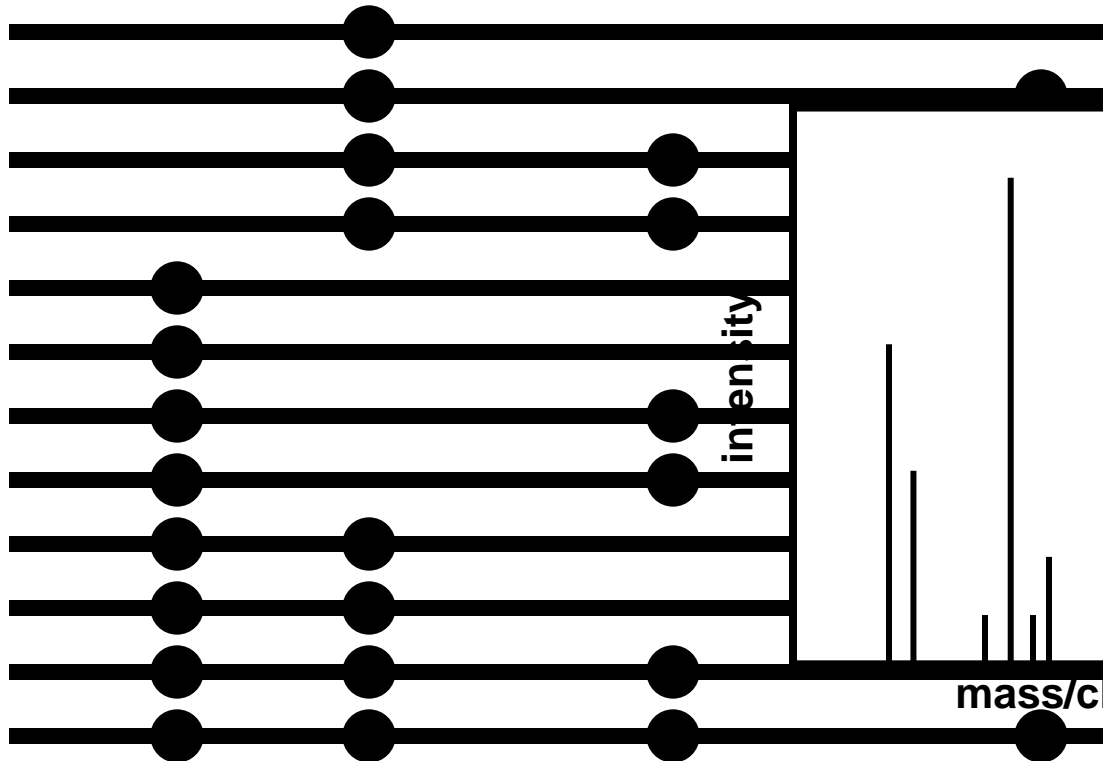
Top down / bottom up



Top down

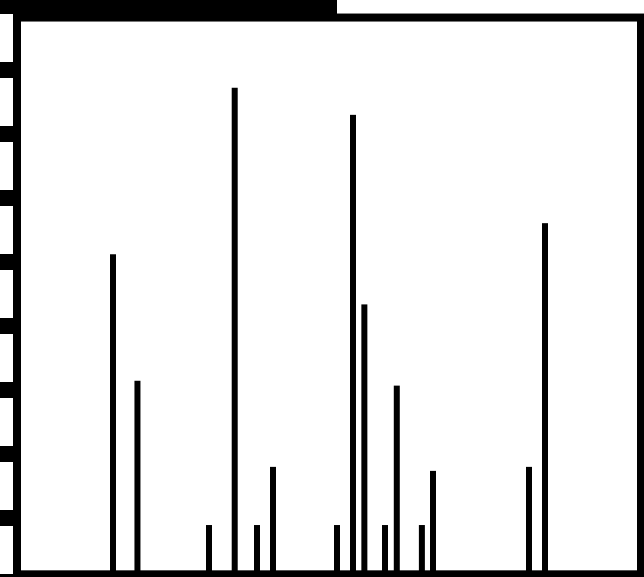


Bottom up



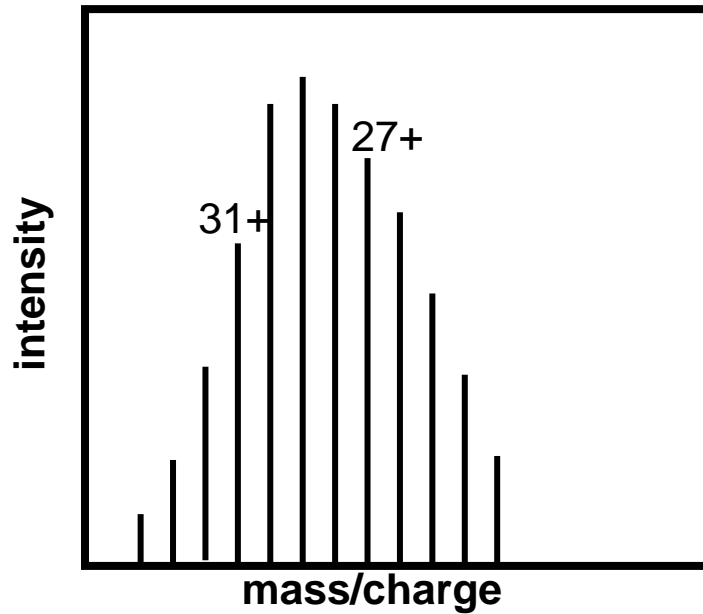
intensity

mass/charge

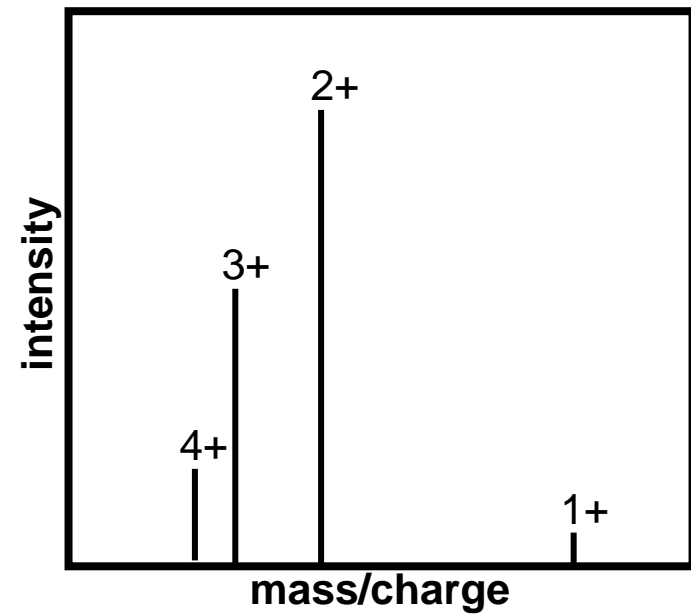


Charge distribution

Top down

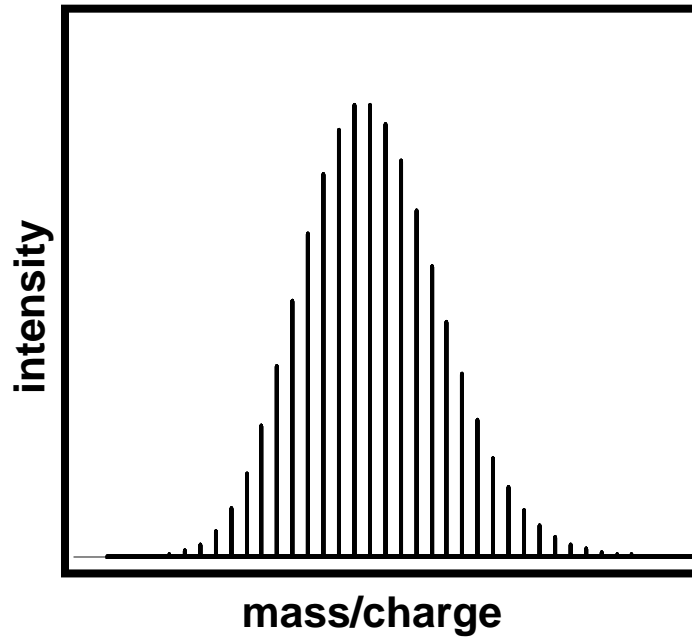


Bottom up

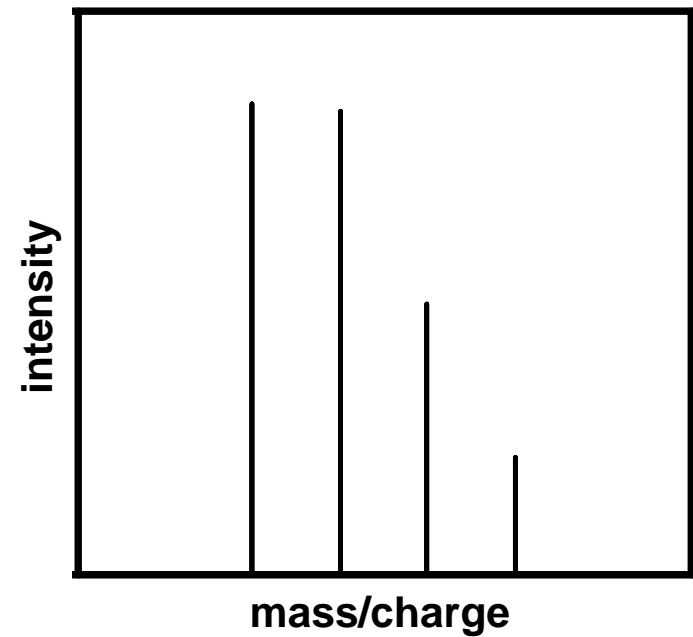


Isotope distribution

Top down



Bottom up

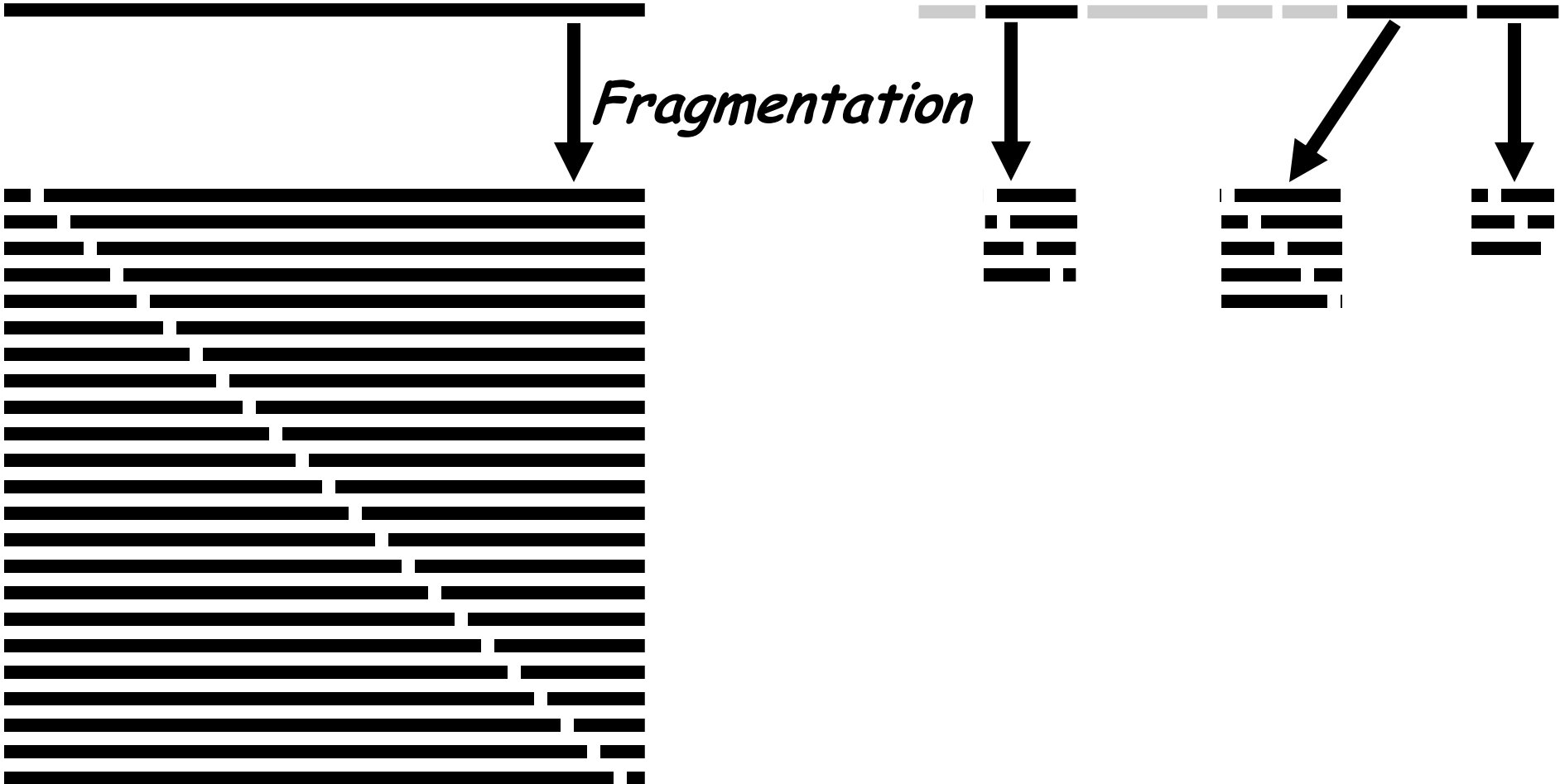


Fragmentation

Top down

Bottom up

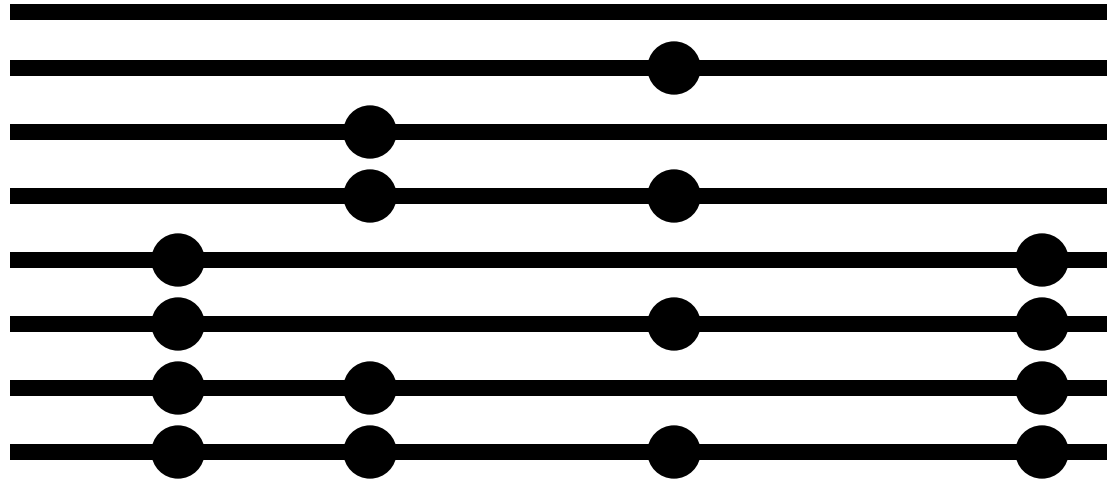
Fragmentation



Correlations between modifications



Top down

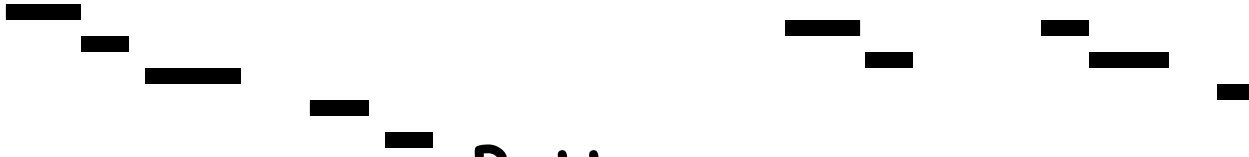
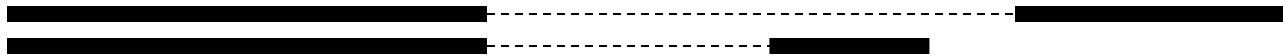


Bottom up



Alternative Splicing

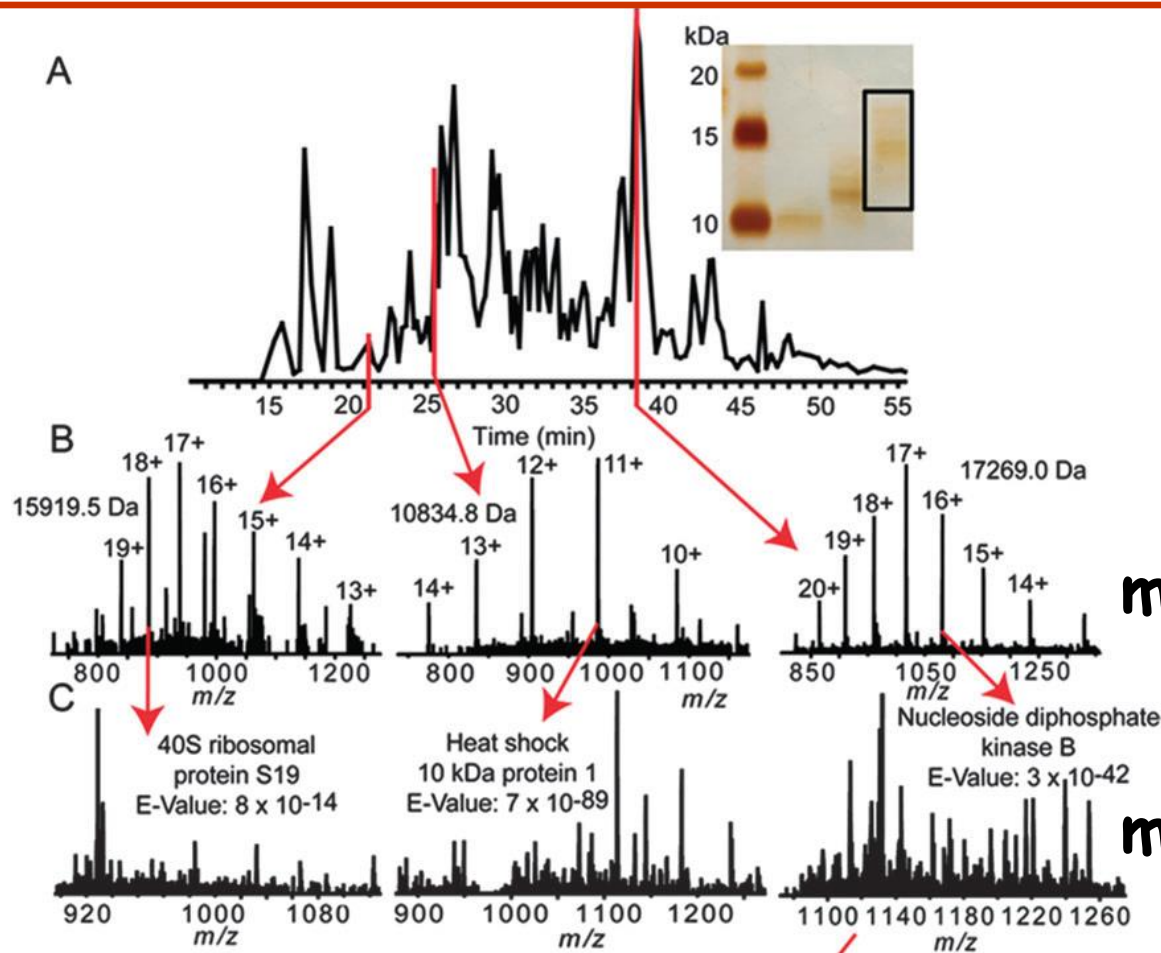
Top down



Bottom up



Top down



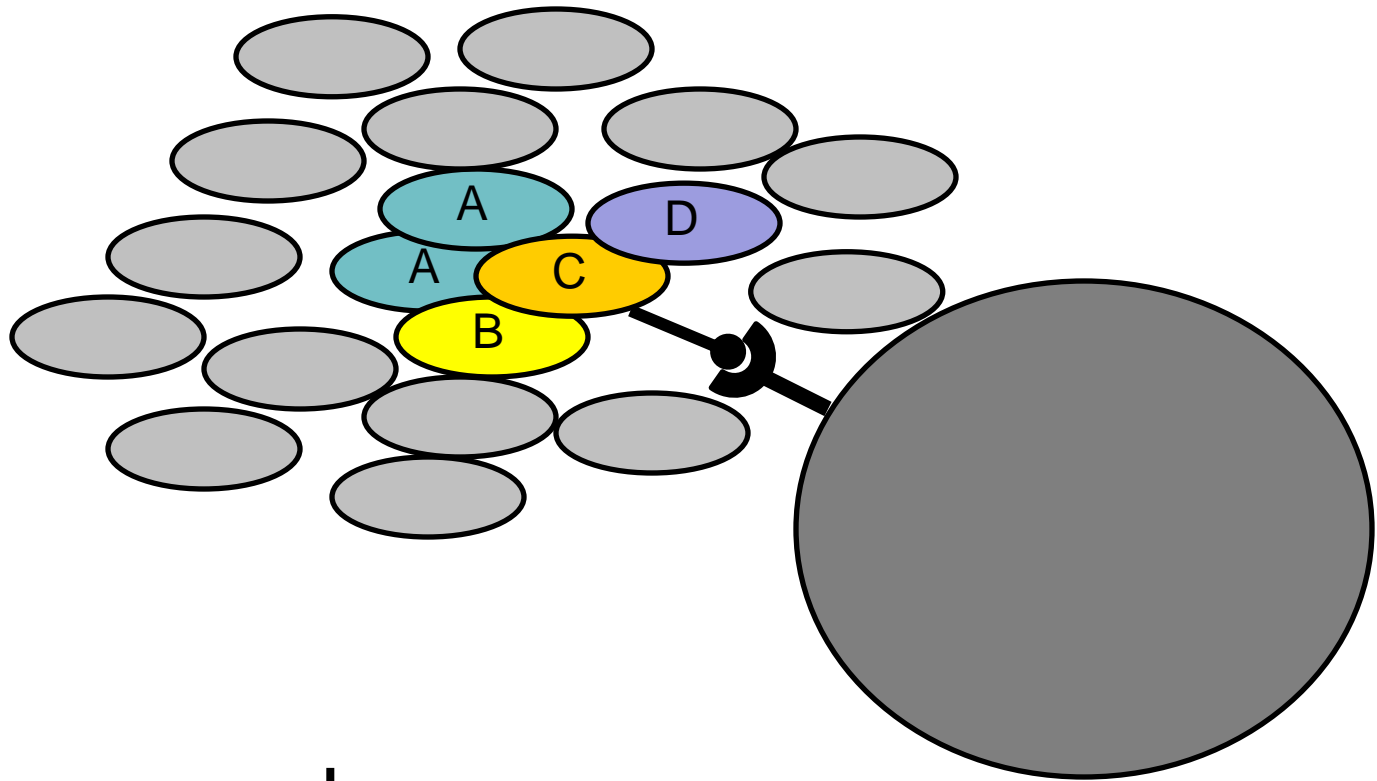
Protein
mass spectra

Fragment
mass spectra

D

·A·N·L·E·R·T·F·I·A·I·K·P·D·I·G·V·Q·R·G·L·V·G·E·I·I·I·K·
 ·R·F·E·Q·K·G·F·R·L·V·A·M·K·F·L·R·A·S·E·E·H·L·K·I·Q·H·
 ·Y·I·D·I·L·K·D·R·P·F·F·P·G·L·V·K·Y·M·N·I·S·I·G·I·P·V·V·A·M·I·
 ·V·I·W·E·I·G·L·N·V·I·V·K·T·G·R·V·M·L·G·E·T·N·I·P·A·D·I·S·K·I·P·
 ·G·T·I·R·G·D·F·C·I·Q·V·G·R·N·I·I·H·G·S·D·I·S·V·K·S·A·
 ·E·K·E·I·S·L·W·F·K·P·E·E·L·V·D·I·Y·K·S·C·A·H·D·I·W·I·V·Y·
 ·E·

Protein Complexes



↓
Digestion

↓
Mass spectrometry

Protein Complexes - specific/non-specific binding

E

Stats Table

	Bait 1	Bait 2	Bait 3	Bait 4	Bait k	
Interactor 1	$X_{1,1}$	$X_{2,1}$	$X_{3,1}$	$X_{4,1}$	$X_{k,1}$	\bar{X}_1
Interactor 2	$X_{1,2}$	$X_{2,2}$	$X_{3,2}$	$X_{4,2}$	$X_{k,2}$	\bar{X}_2
Interactor 3	$X_{1,3}$	$X_{2,3}$	$X_{3,3}$	$X_{4,3}$	$X_{k,3}$	\vdots
Interactor 4	$X_{1,4}$	$X_{2,4}$	$X_{3,4}$	$X_{4,4}$	$X_{k,4}$	\vdots
Interactor m	$X_{1,m}$	$X_{2,m}$	$X_{3,m}$	$X_{4,m}$	$X_{k,m}$	\bar{X}_m

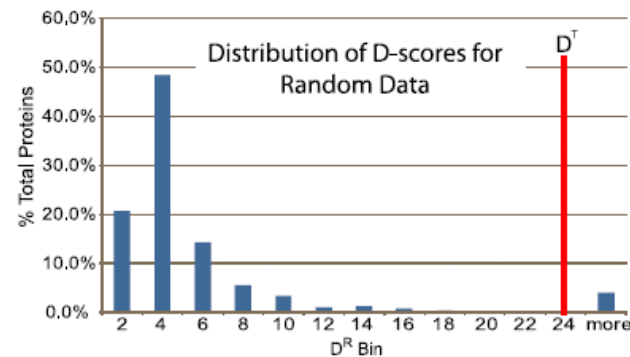
$X_{i,j}$ = total spectral counts for interactor j from bait i

$$\bar{X}_j = \frac{\sum_{i=1, j=n}^{i=k} X_{i,j}}{k} ; n = 1, 2, \dots, m \quad (\text{Eq. 1})$$

$$z_{i,j} = \frac{X_{i,j} - \bar{X}_j}{\sigma_j} \quad (\text{Eq. 2})$$

$$f_{i,j} = \begin{cases} 1 ; X_{i,j} > 0 \\ X_{i,j} \end{cases} \quad p = \begin{cases} \text{number of replicates} \\ \text{runs in which} \\ \text{the interactor is present} \end{cases}$$

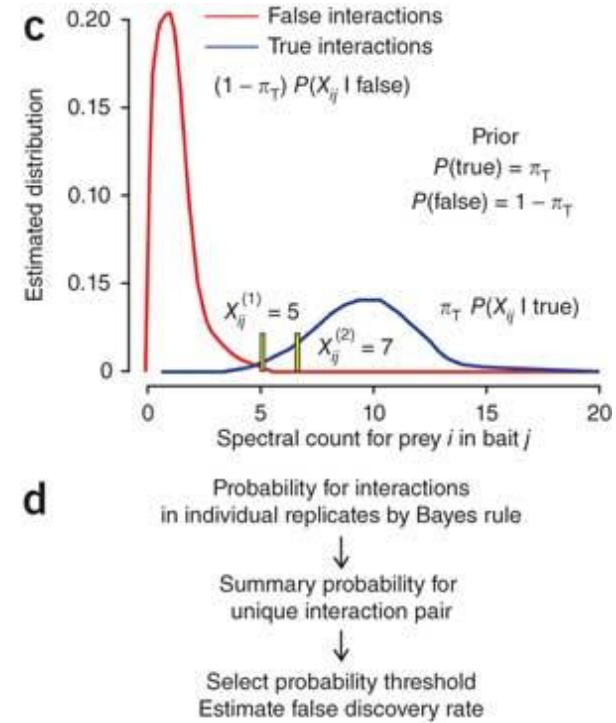
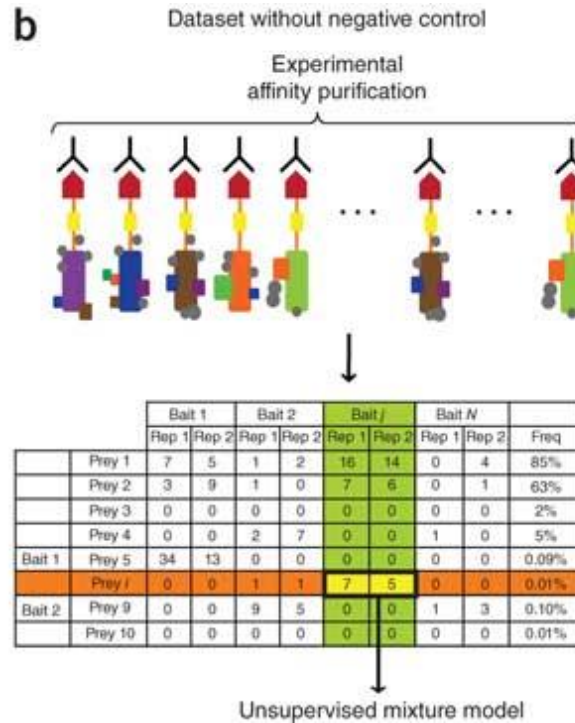
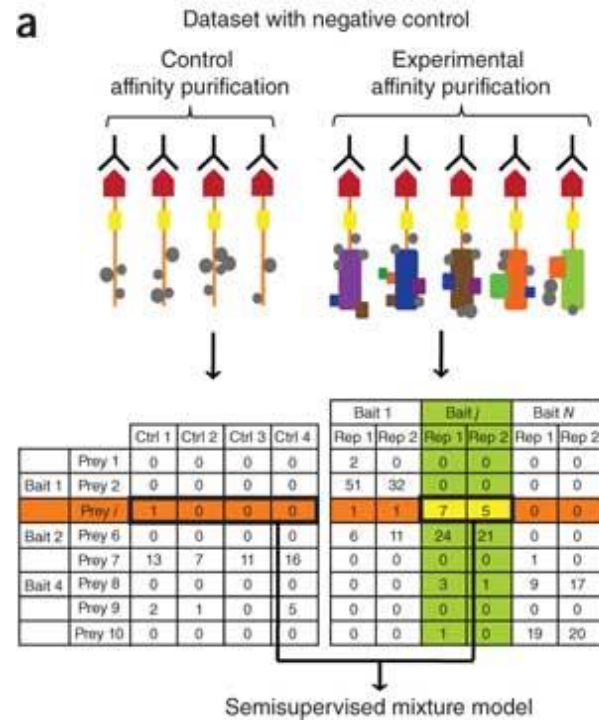
$$D_{i,j}^R = \sqrt{\left(\frac{k}{\sum_{i=1}^k f_{i,j}}\right)^p X_{i,j}} \quad (\text{Eq. 3})$$



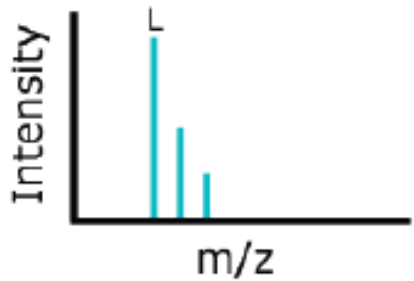
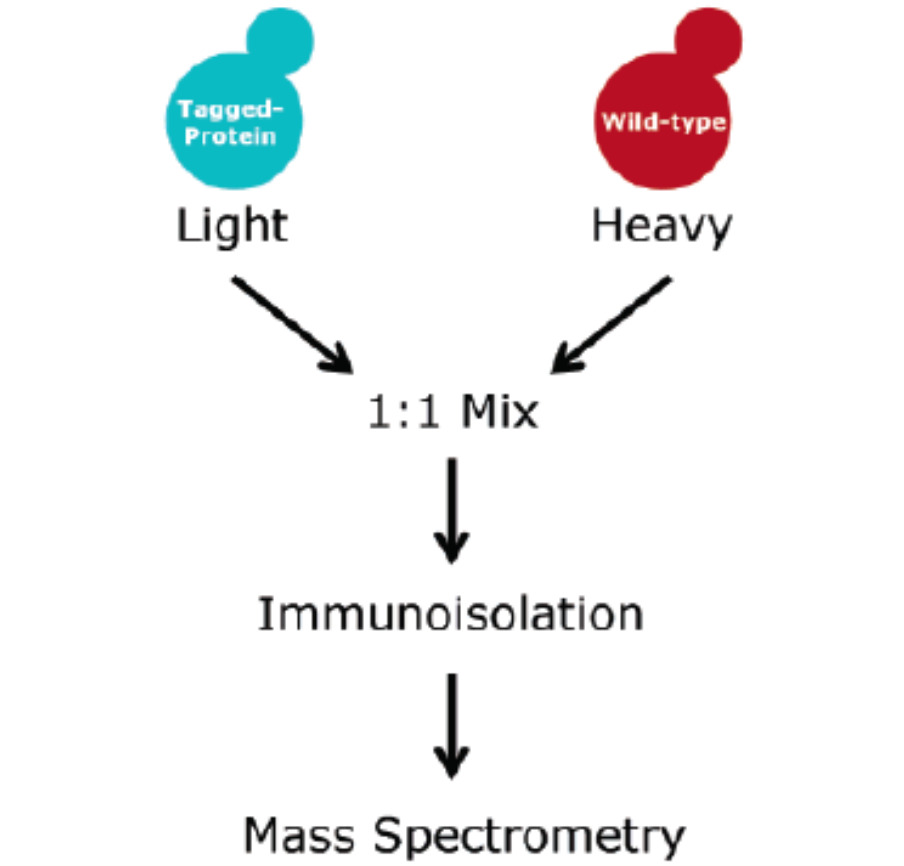
D^T = D-score threshold below which 95% of simulated data falls

$$D_{i,j}^N = D_{i,j}^R / D^T \quad (\text{Eq. 4})$$

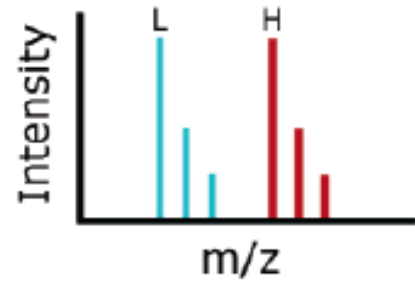
Protein Complexes - specific/non-specific binding



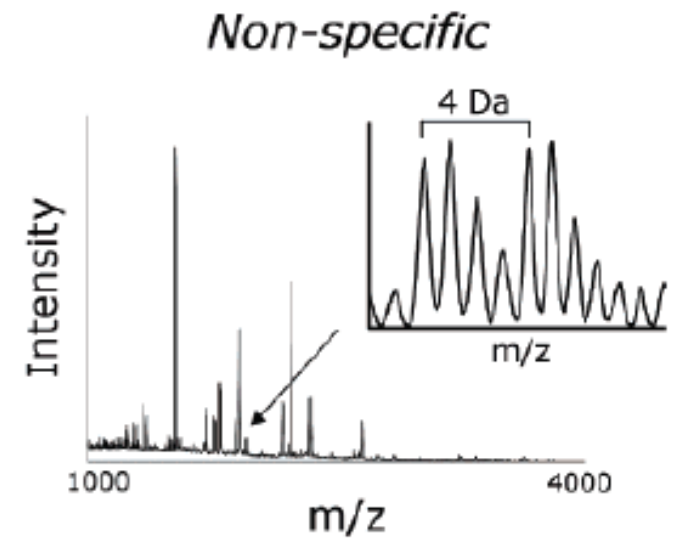
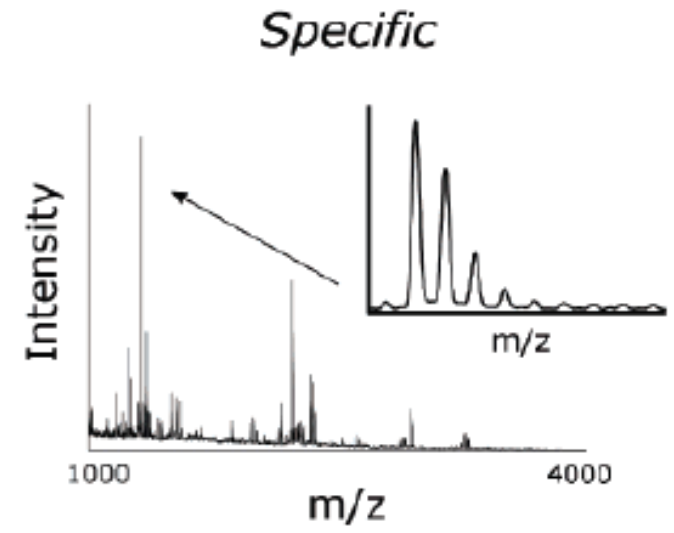
Protein Complexes - specific/non-specific binding



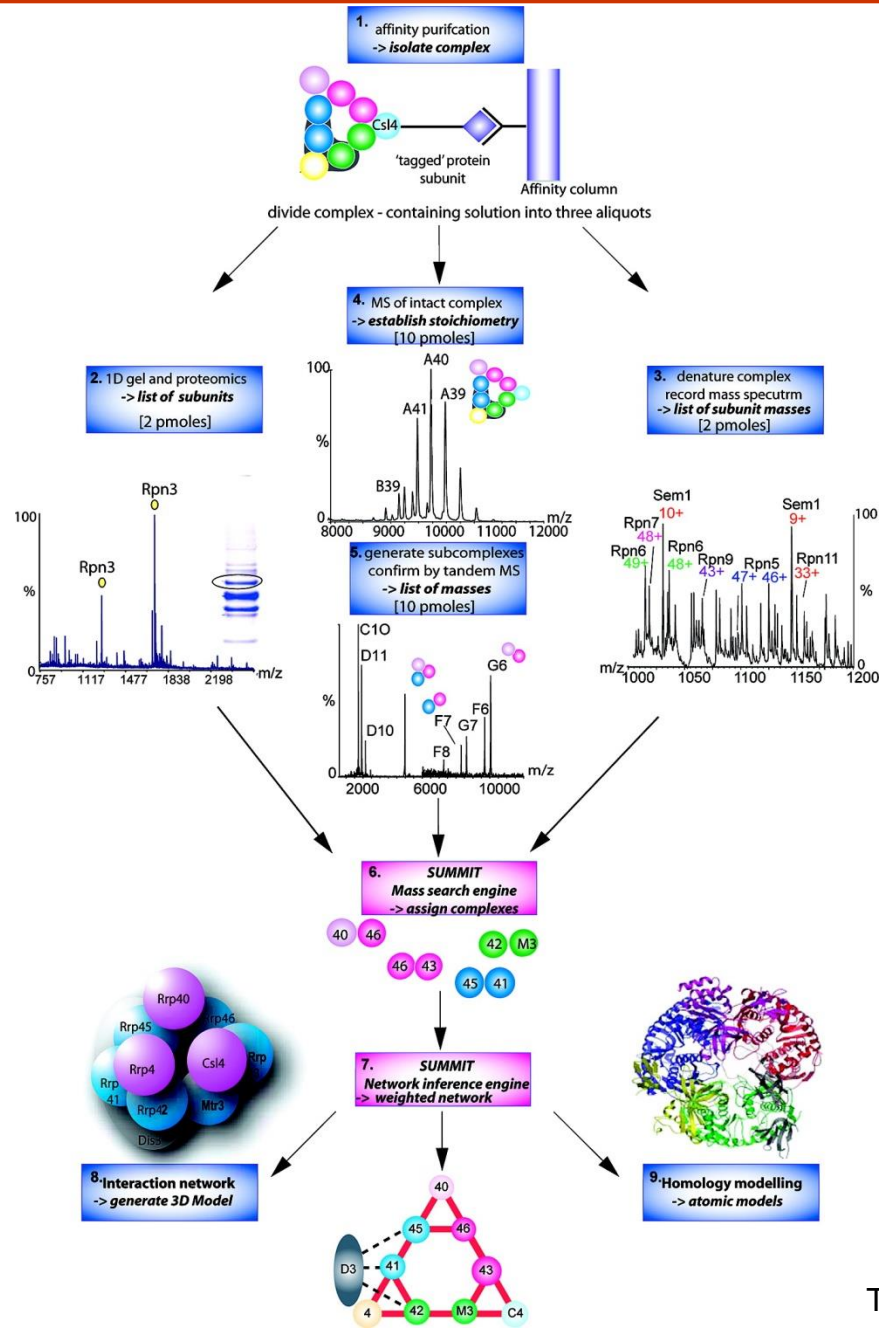
Specific



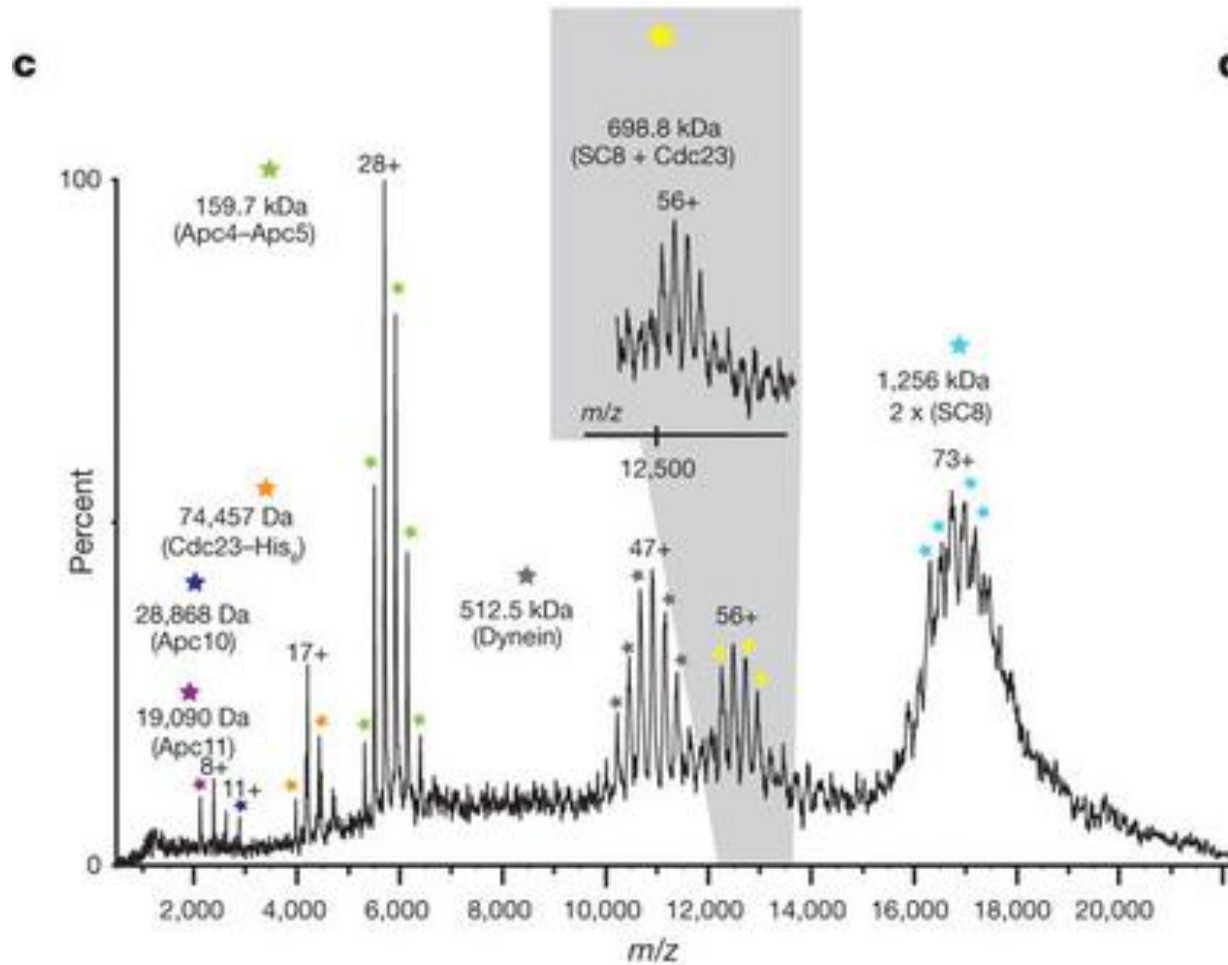
Non-specific



Analysis of Non-Covalent Protein Complexes



Non-Covalent Protein Complexes

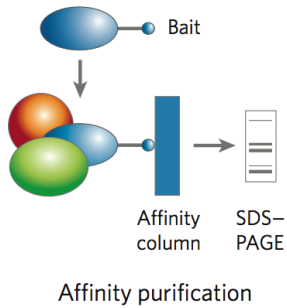


d

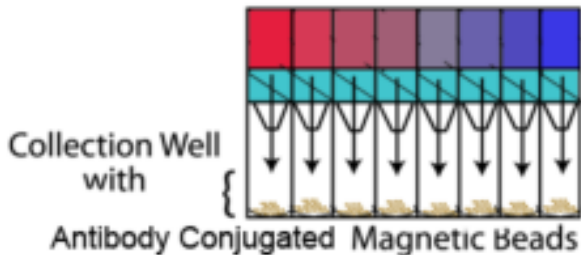
<i>S. cerevisiae</i> APC/C subunit	Stoichiometry
Apc1	1
Apc2	1
Apc4	1
Apc5	1
Apc10	1
Apc11	1
Cdc23	2
Cdc16	2
Cdc27	2
Apc9	ND
Mnd2	1
Apc13	1
Cdc26	2
Total MW (kDa)	1,127–1,158

Affinity Capture Optimization Screen

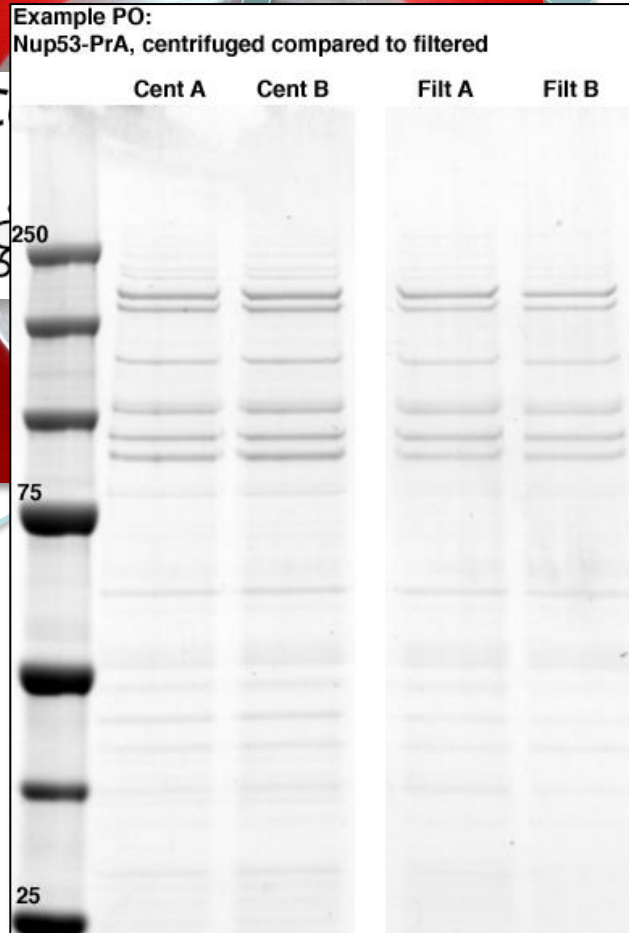
More / better quality interactions



Filtration



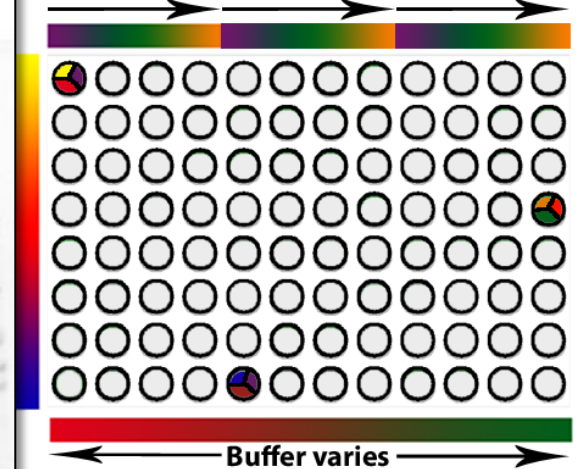
Lysate clearance/
Batch Binding



Binding/Washing/Eluting

Cell extraction

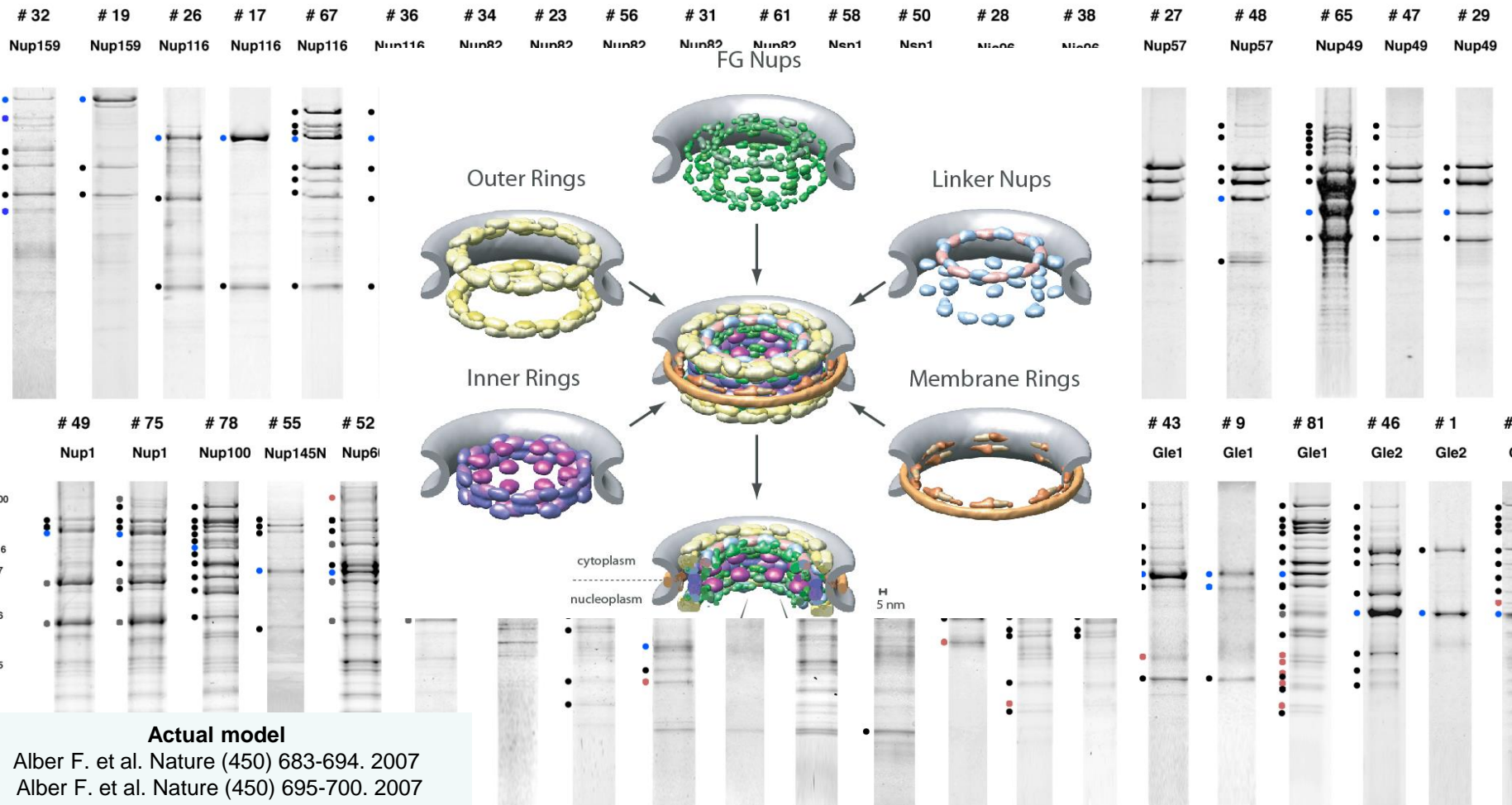
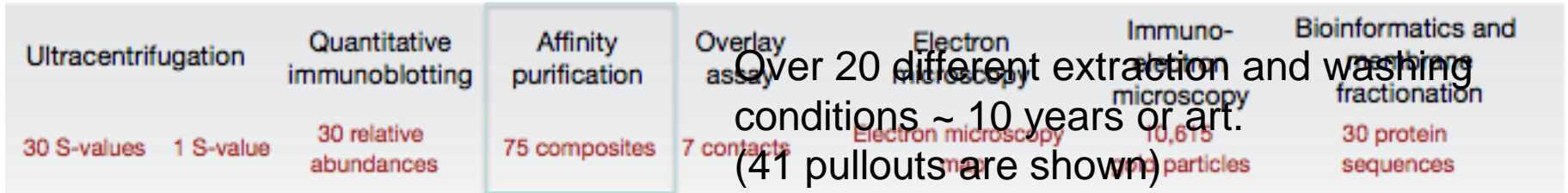
Sample extraction buffer matrix
Salt concentration changes



SDS-PAGE

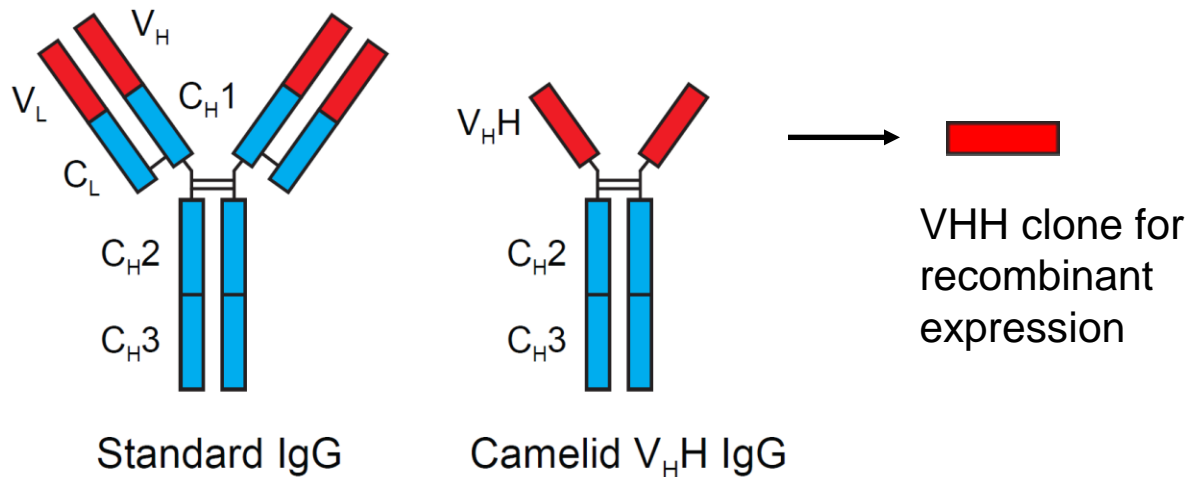


Molecular Architecture of the NPC



Cloning nanobodies for GFP pullouts

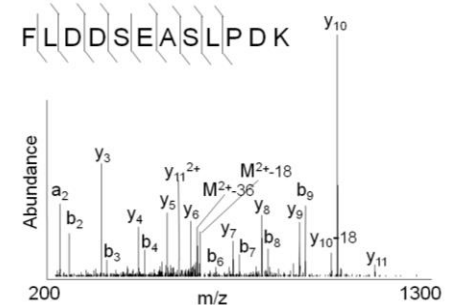
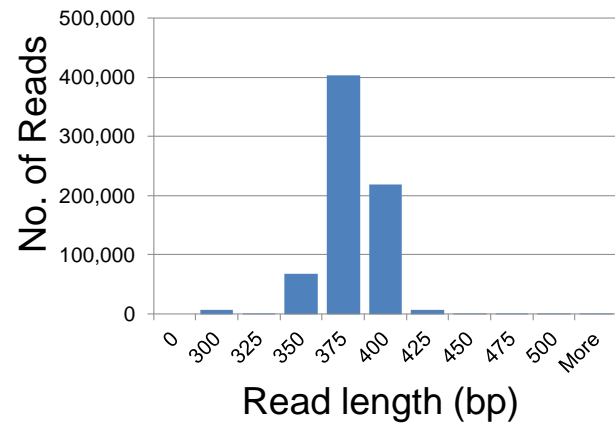
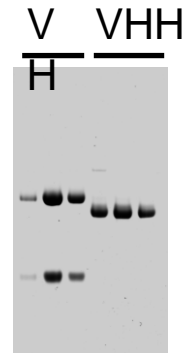
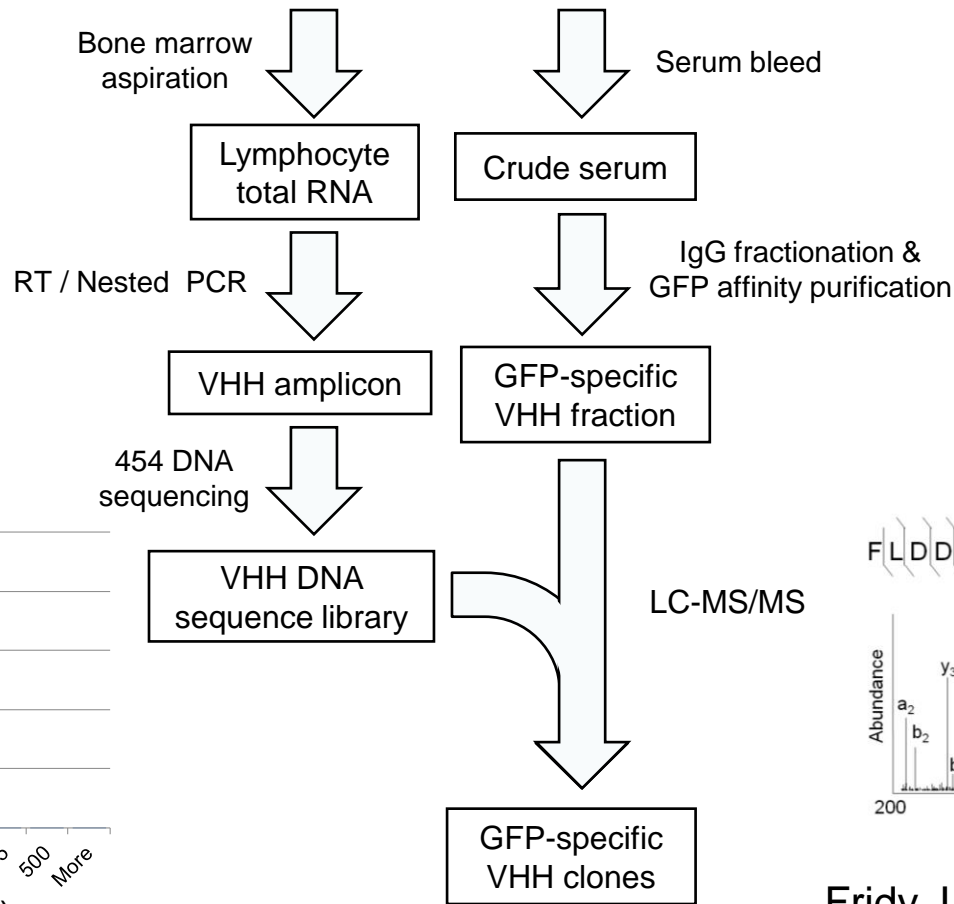
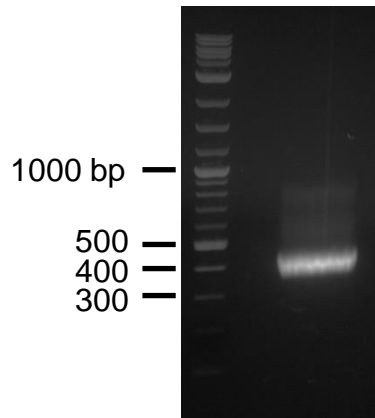
- Atypical heavy chain-only IgG antibody produced in camelid family – retain high affinity for antigen without light chain
- Aimed to clone individual single-domain VHH antibodies against GFP – only ~15 kDa, can be recombinantly expressed, used as bait for pullouts, etc.
- To identify full repertoire, will identify GFP binders through combination of high-throughput DNA sequencing and mass spectrometry



Cloning llama antibodies for GFP pullouts



Llama GFP immunization



Identifying full-length sequences from peptides

Underlined regions are covered by MS

CDR1

CDR2

CDR3

CDR3: 100.0% (14/14); combined CDR: 100.0% (33/33); DNA count: 10

MAQVQLVESGGGLVQAGGSLRLSCVASGRTFSGYAMGWERQTPGREREAVAAITWSAHSTYYSDSVKDRFTISIDNTRNTGYLQMNSLKPEDTAVYYCTVRHGTWFTTSRYWTDWGOGTQVTVS



CDR3: 100.0% (14/14); combined CDR: 72.7% (24/33); DNA count: 1

MAQVQLVESGGALVQAGASLSVSCAASGGTISKYINMAWFRRAPGREREAVAAITWSAHSTYYSDSVKDRFTISIDNTRNTGYLQMNSLKPEDTAVYYCTVRHGTWFTTSRYWTDWGOGTQVTVS

CDR3: 100.0% (14/14); combined CDR: 72.7% (24/33); DNA count: 1

MADVQLVESGGGLVQSGGSRTLSCAASGRVLATYHLGWFRQSPGREREAVAAITWSAHSTYYSDSVKGRFTISIDNARNRTGYLQMNSLKPEDTAVYYCTVRHGTWFTVSRYWTDWGOGTQVTVS

CDR3: 100.0% (14/14); combined CDR: 42.4% (14/33); DNA count: 1

MAQVQLEESGGGLVQAGDSLTLSCSASGRFTFTNYAMAWSRQAPGKERELLAIDAAGGATYYSDSVKGRFTISIDNTRNTGYLQMNSLKPEDTAVYYCTVRHGTWFTTSRYWTDWGOGTQVTVS

CDR3: 100.0% (14/14); combined CDR: 42.4% (14/33); DNA count: 1

MAQVQLVESGGGRVQAGGSLTLSCVSGSEGI FWNHVMGWFRQSPGKDREFVARISKIGTTNYADSVKGRFTISIDNTRNTGYLQMNSLKPEDTAVYYCTVRHGTWFTTSRYWTDWGOGTQVTVS

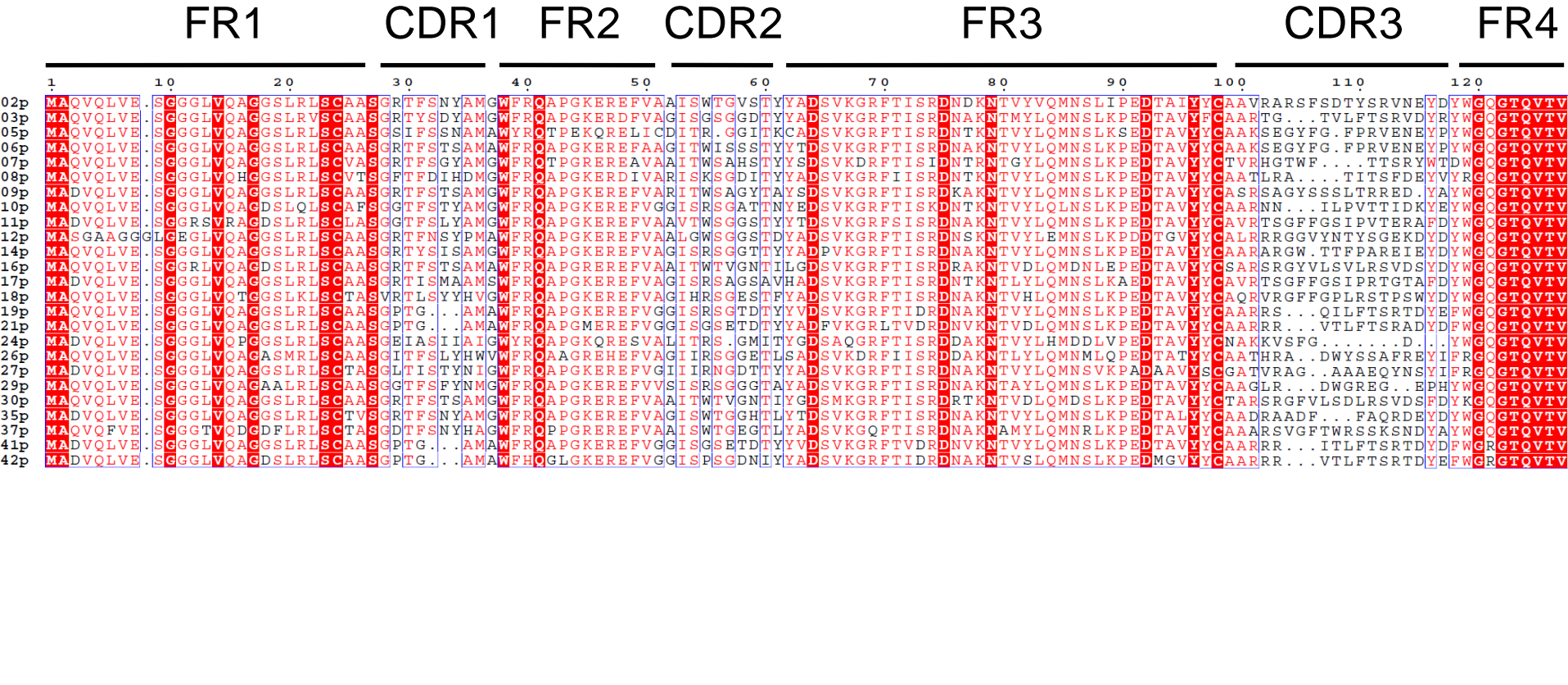
Rank sequences according to:

CDR3 coverage; Overall coverage;

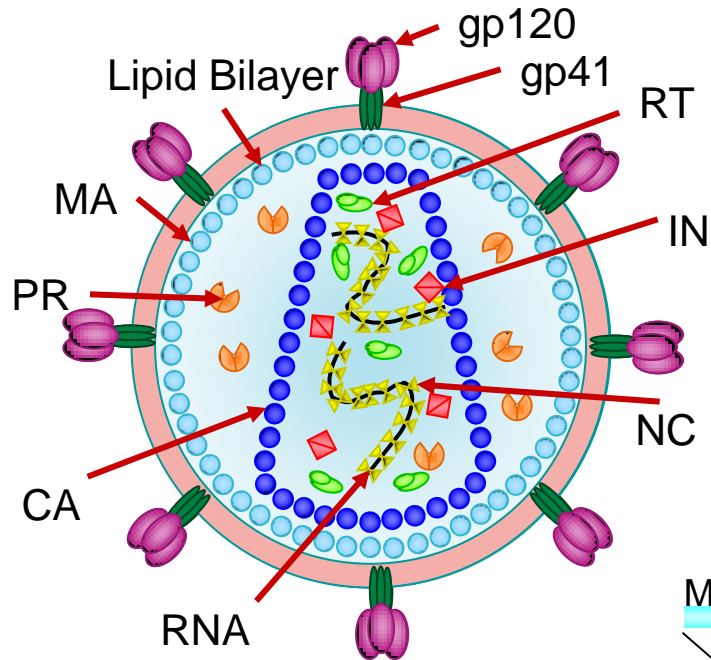
Combined CDR coverage; DNA counts;

Sequence diversity of 26 verified anti-GFP nanobodies

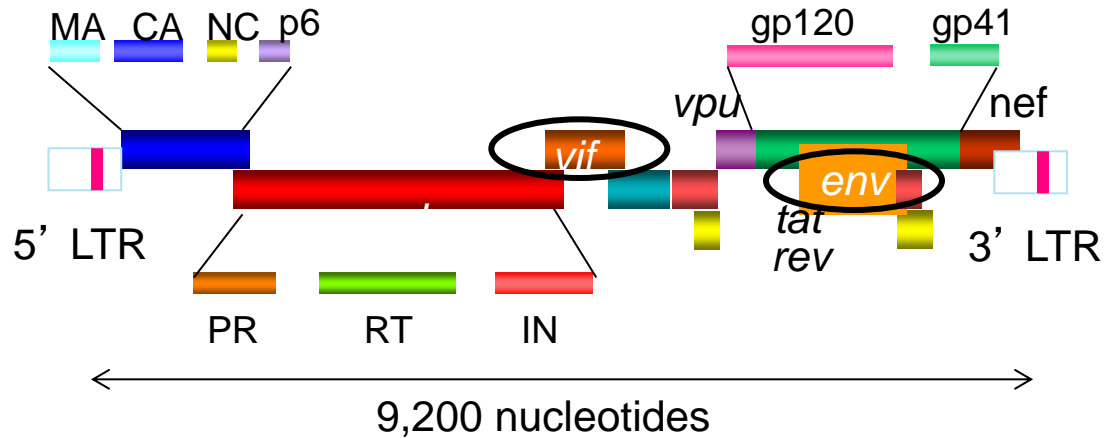
- Of ~200 positive sequence hits, 44 high confidence clones were synthesized and tested for expression and GFP binding: 26 were confirmed GFP binders.
- Sequences have characteristic conserved VHH residues, but significant diversity in CDR regions.



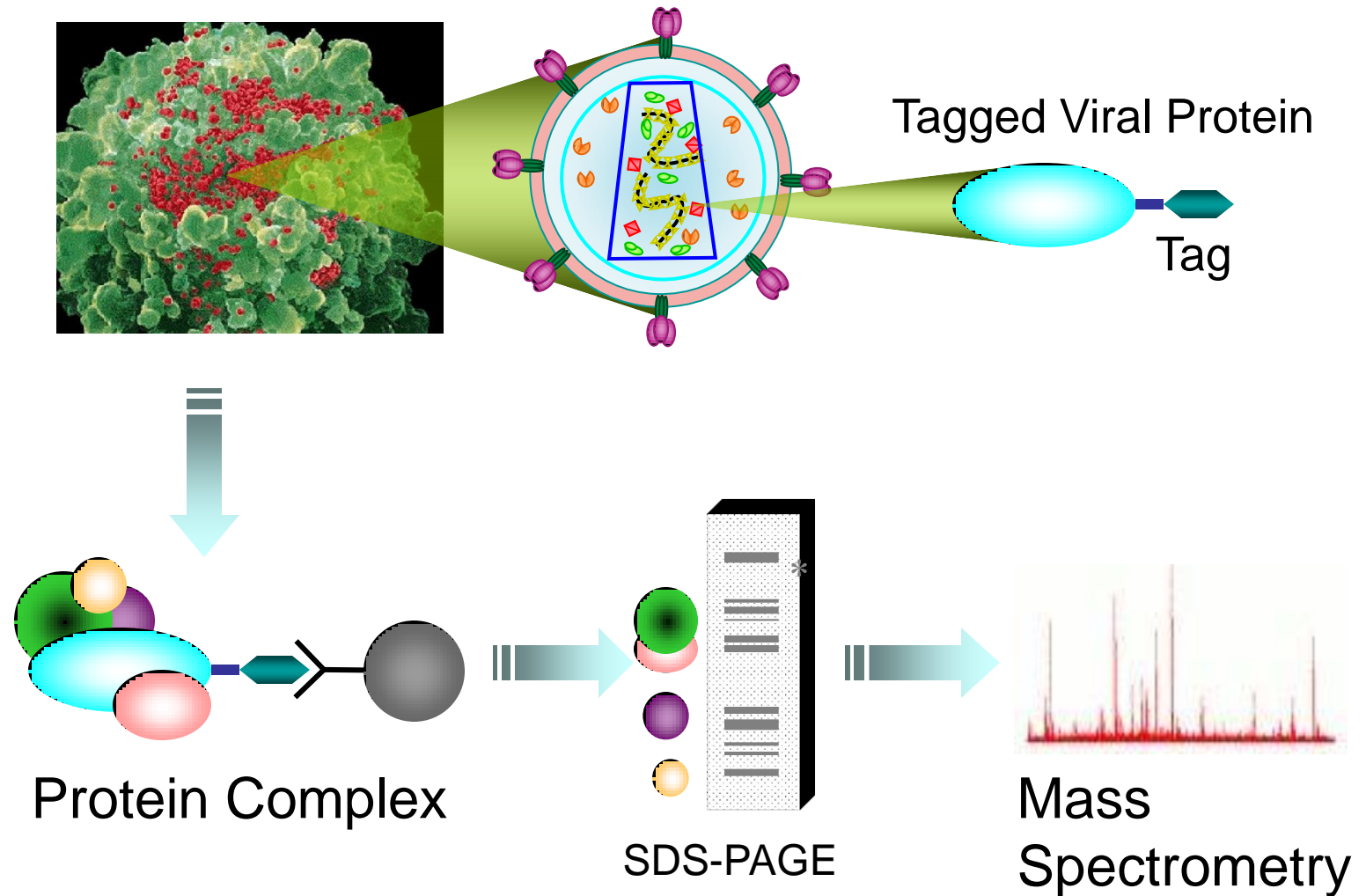
HIV-1



Genome

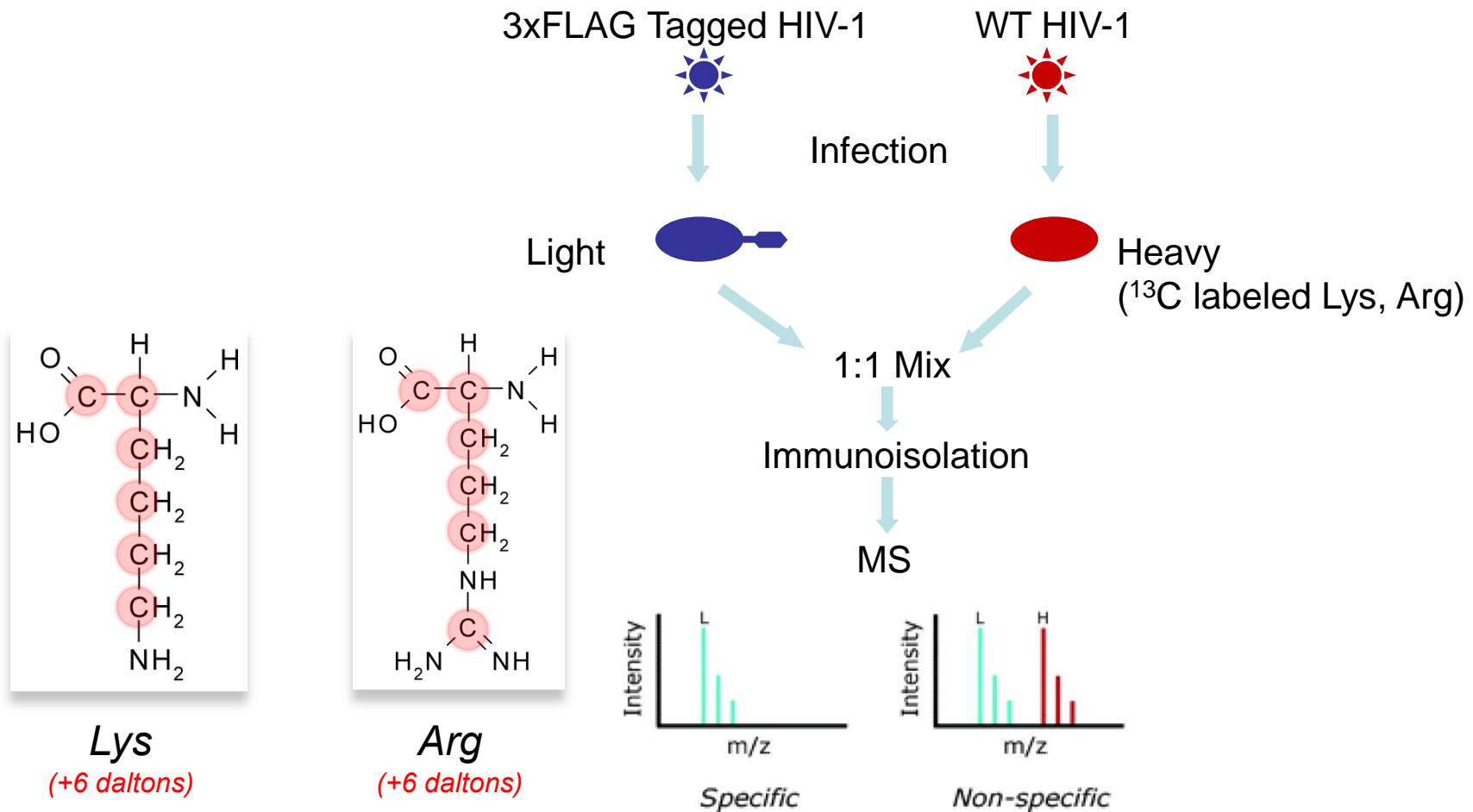


Genetic-Proteomic Approach



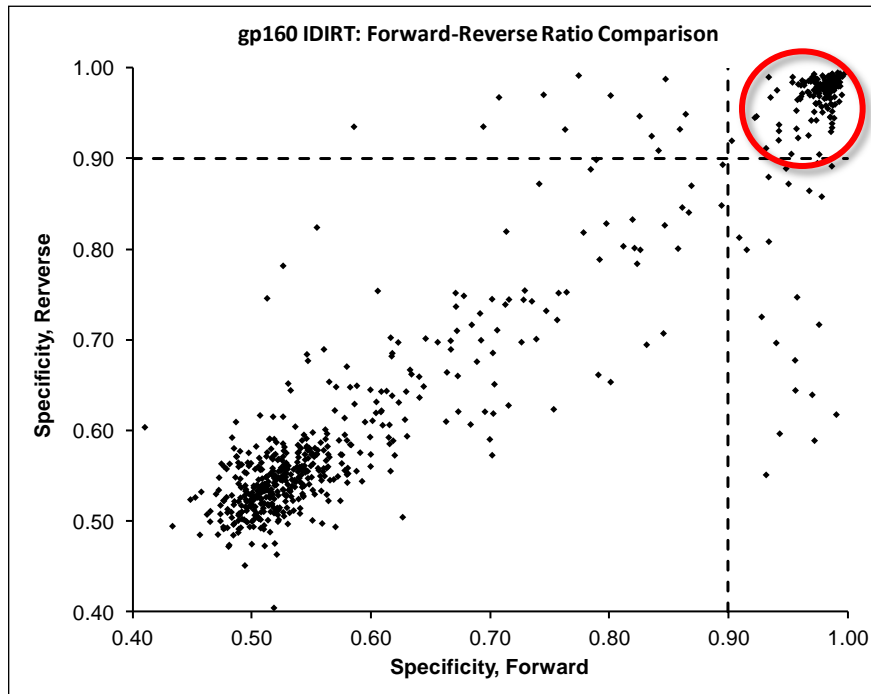
I-DIRT for Specific Interaction

I-DIRT = Isotopic **D**ifferentiation of **I**nteractions as **R**andom or **T**argeted

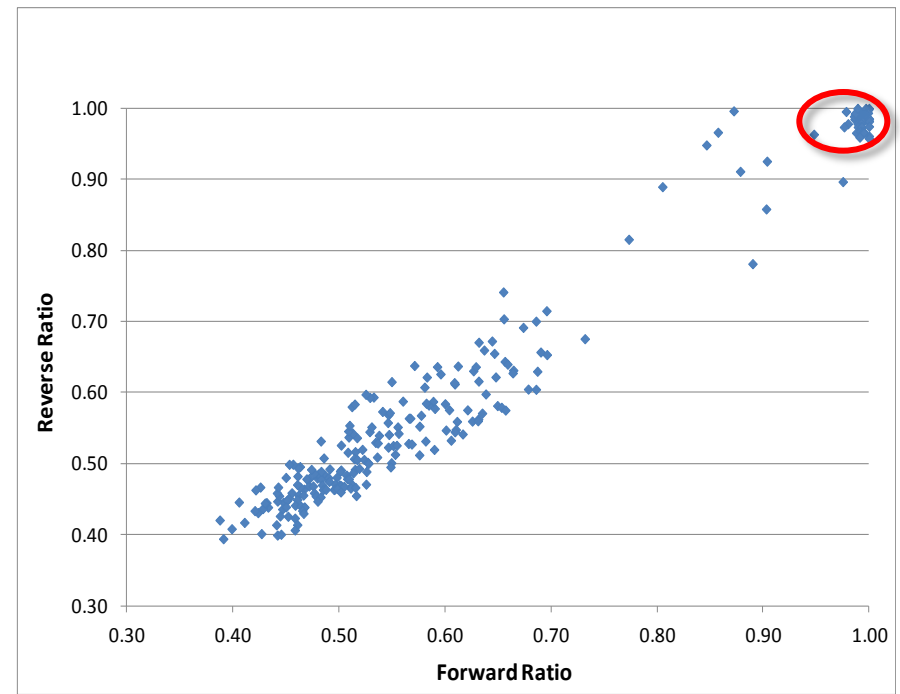


IDIRT and Reverse IDIRT

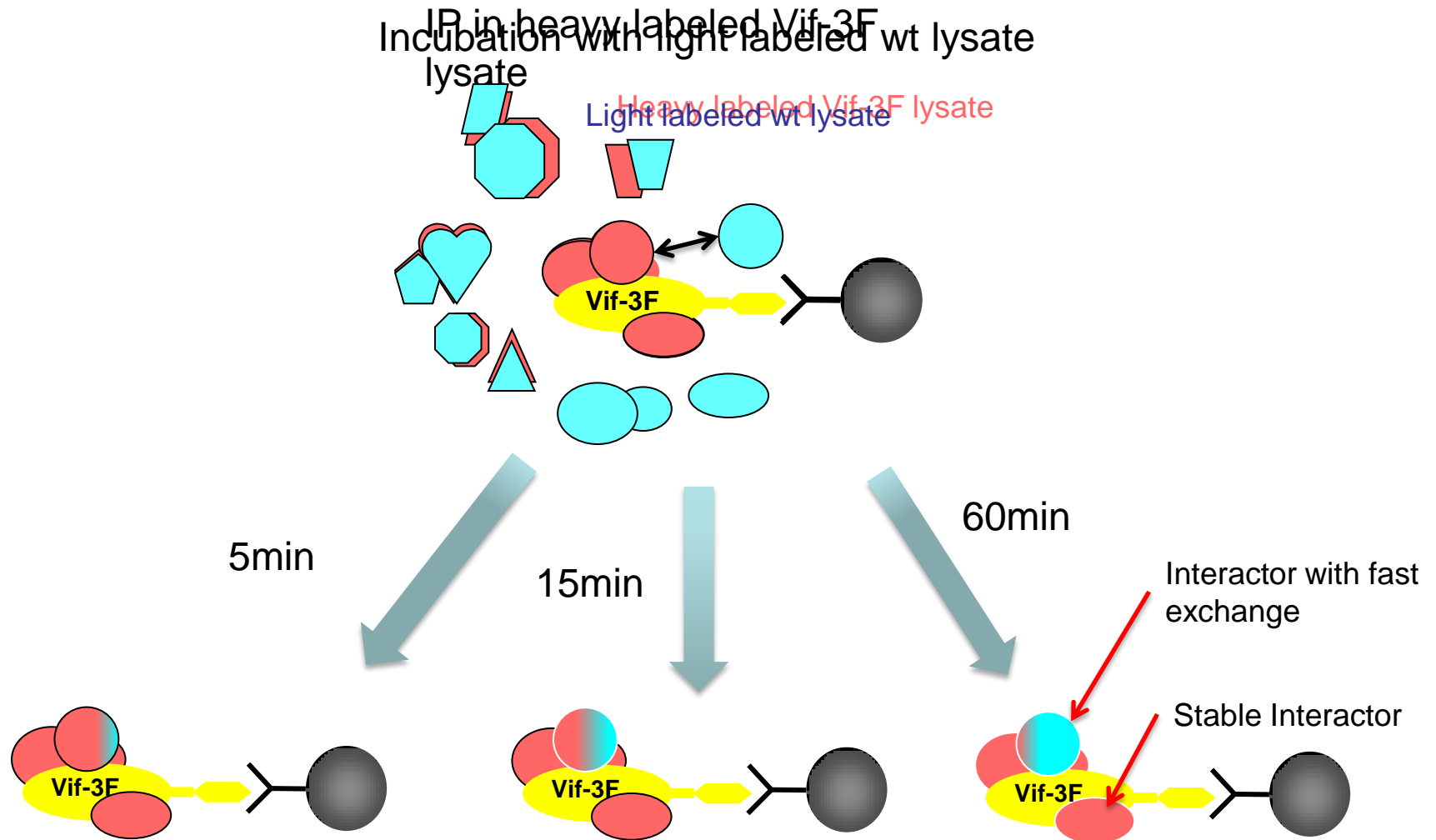
Env-3xFLAG



Vif-3xFLAG

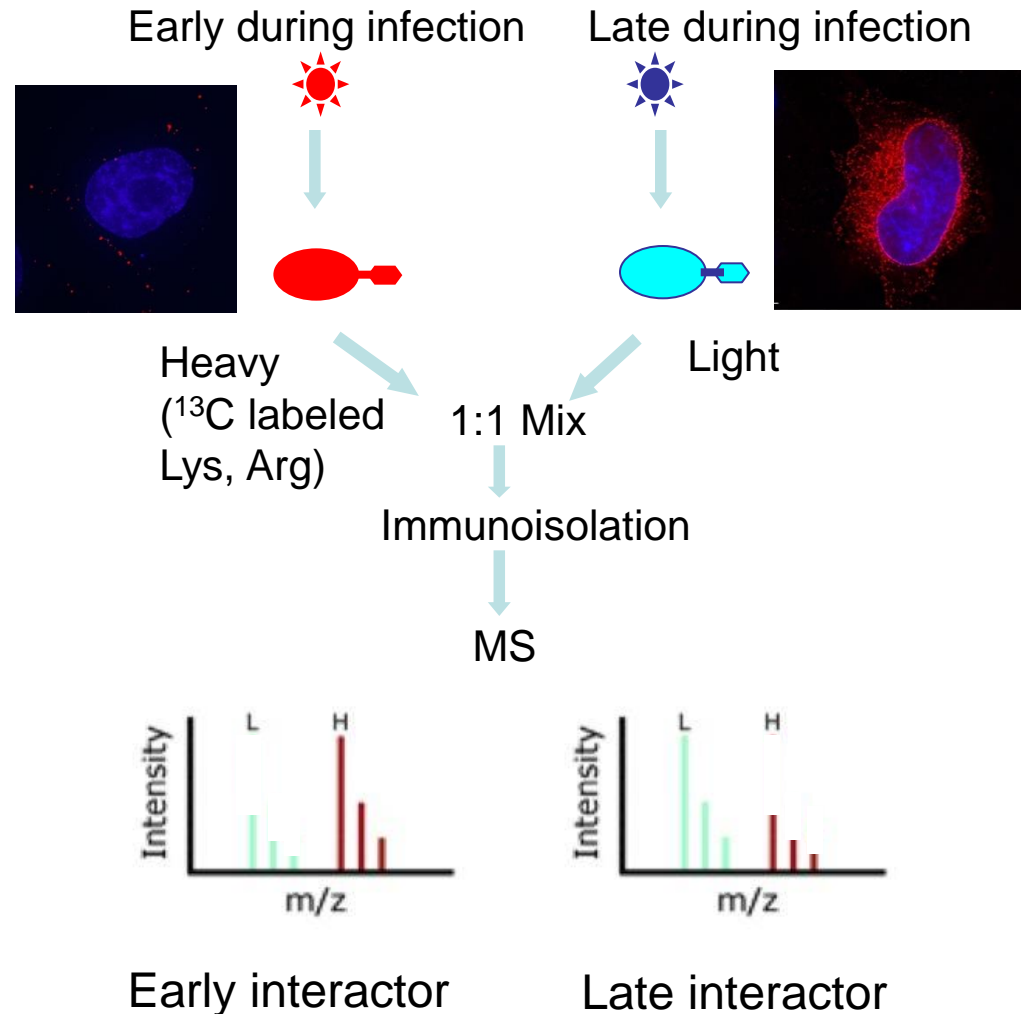


Protein Exchange



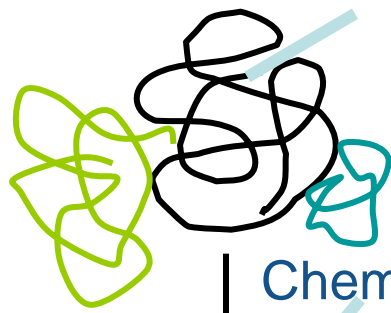
Env Time Course SILAC

- Differentially labeled infection harvested at early or late stage of infection
- Distinguish proteins that interact with Env at early or late stage during infection



Interaction Partners by Chemical Cross-Linking

Protein Complex



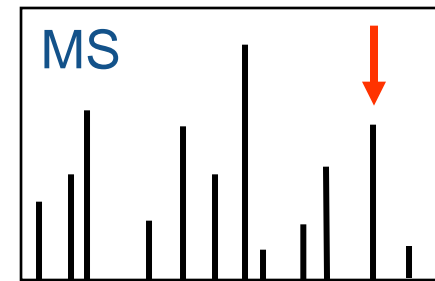
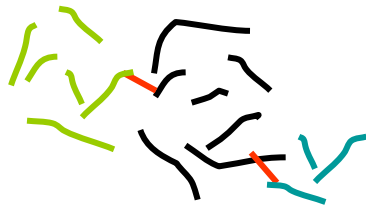
Chemical Cross-Linking

Cross-Linked Protein Complex



Enzymatic Digestion

Proteolytic Peptides

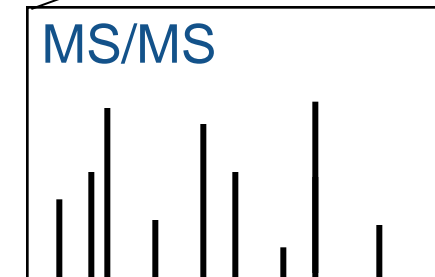


Isolation



Fragmentation

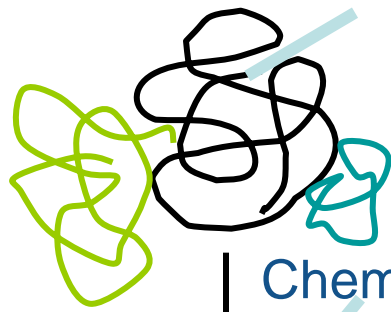
Peptides Fragments



M/Z

Interaction Sites by Chemical Cross-Linking

Protein Complex



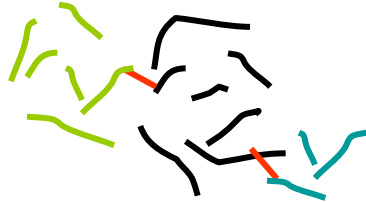
Chemical Cross-Linking

Cross-Linked Protein Complex

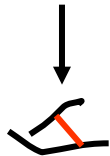


Enzymatic Digestion

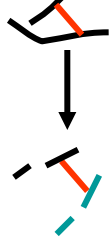
Proteolytic Peptides



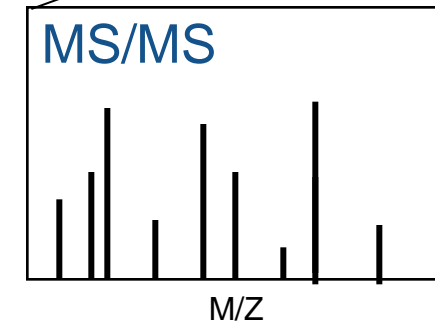
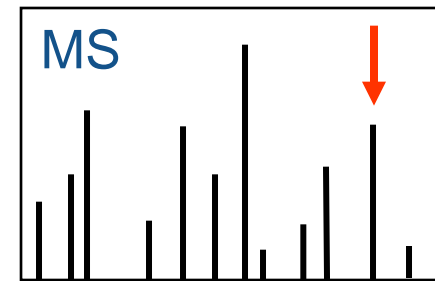
Isolation



Fragmentation



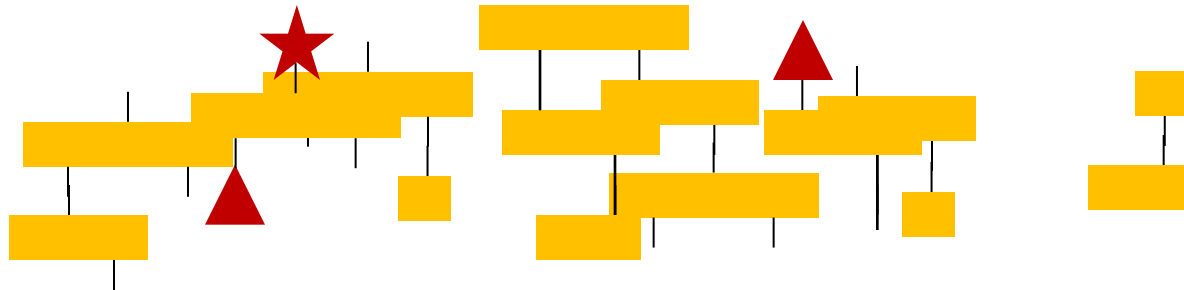
Peptides Fragments



M/Z

Cross-linking

protein 



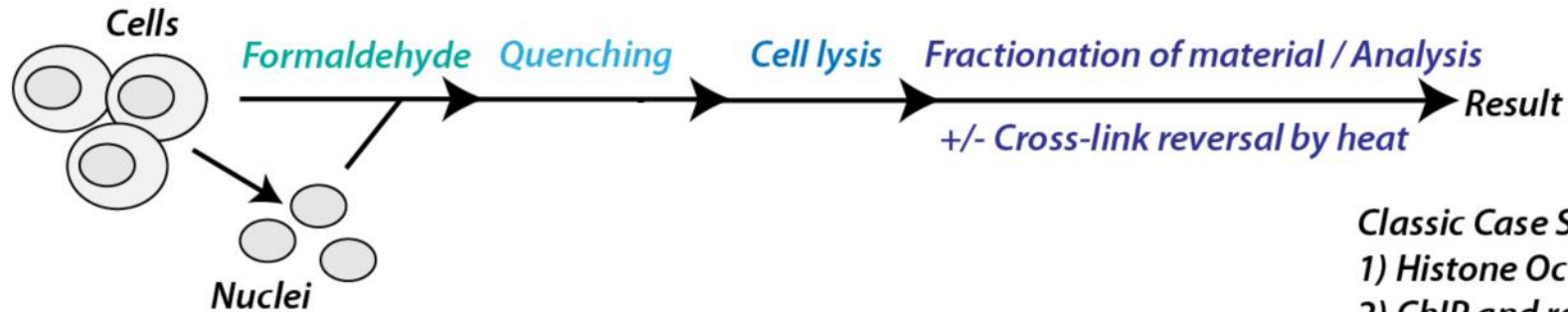
n peptides with reactive groups



$(n-1)n/2$ potential ways to cross-link peptides pairwise
+ many additional uninformative forms

Protein Crosslinking by Formaldehyde

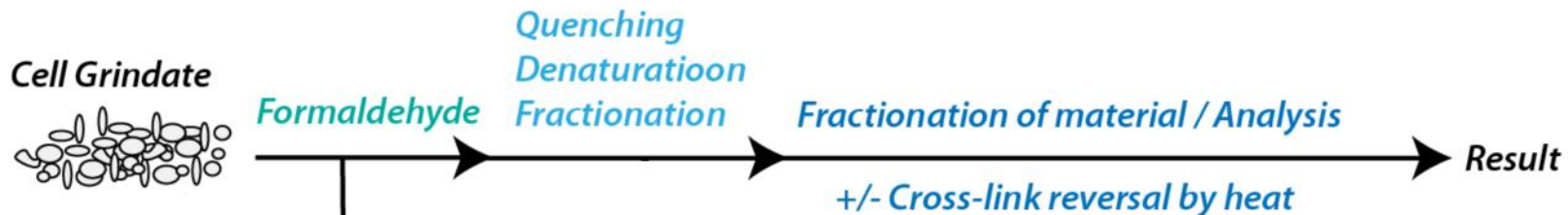
Example Traditional Formaldehyde Cross-linking:



Classic Case Studies:
1) Histone Octamer
2) ChIP and related methods

~1% w/v Fal
20 – 60 min

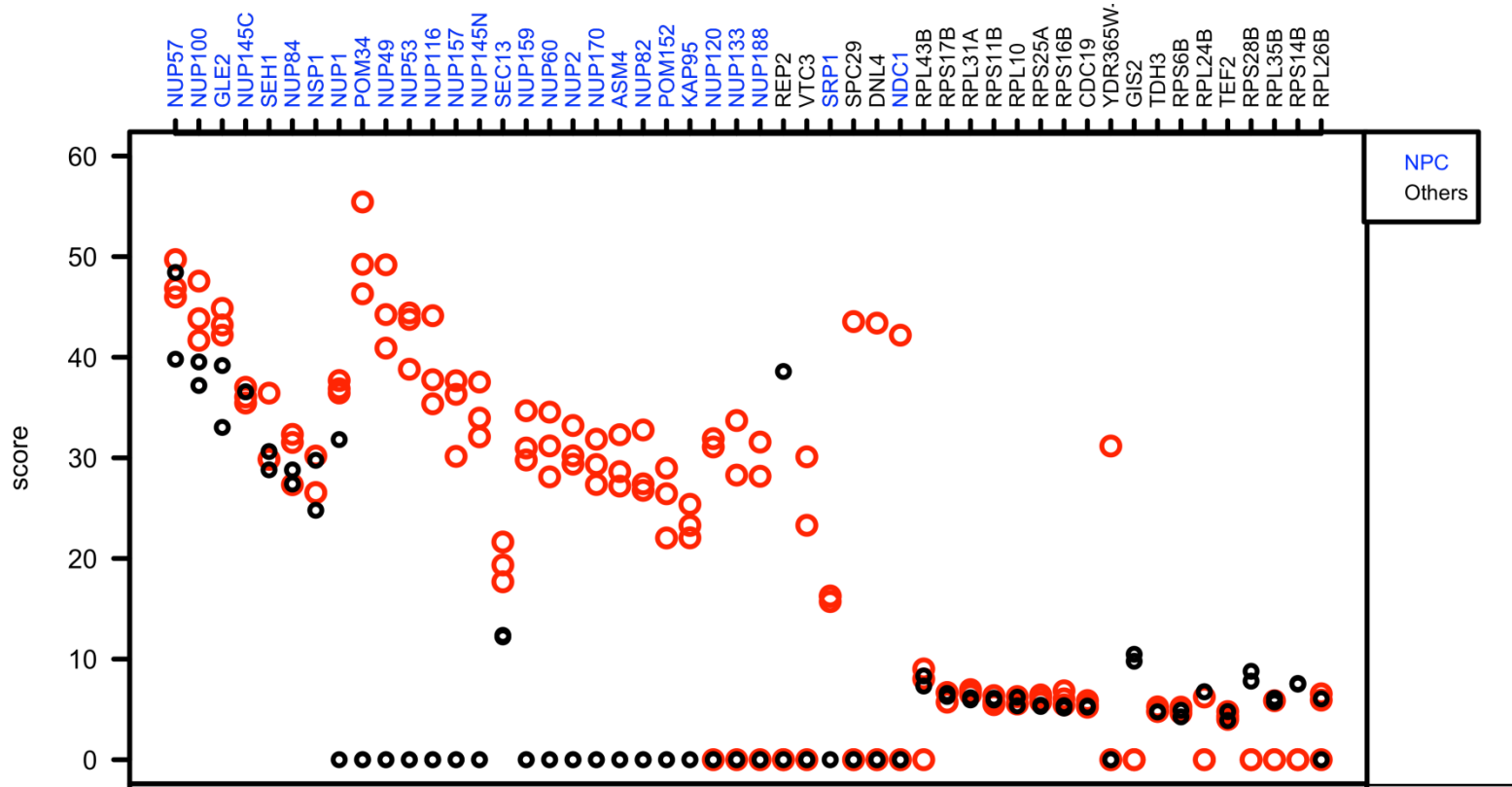
Example Modified Approach:



*Extensive control
over conditions of
XL*

~0.3% w/v Fal
5 – 20 min
1/100 the volume

Protein Crosslinking by Formaldehyde



RED: triplicate experiments, FAI treated grindate

BLACK: duplicated experiments, FAI treated cells (then ground)

SCORE: $\text{Log Ion Current} / \text{Log protein abundance}$

Cross-linking

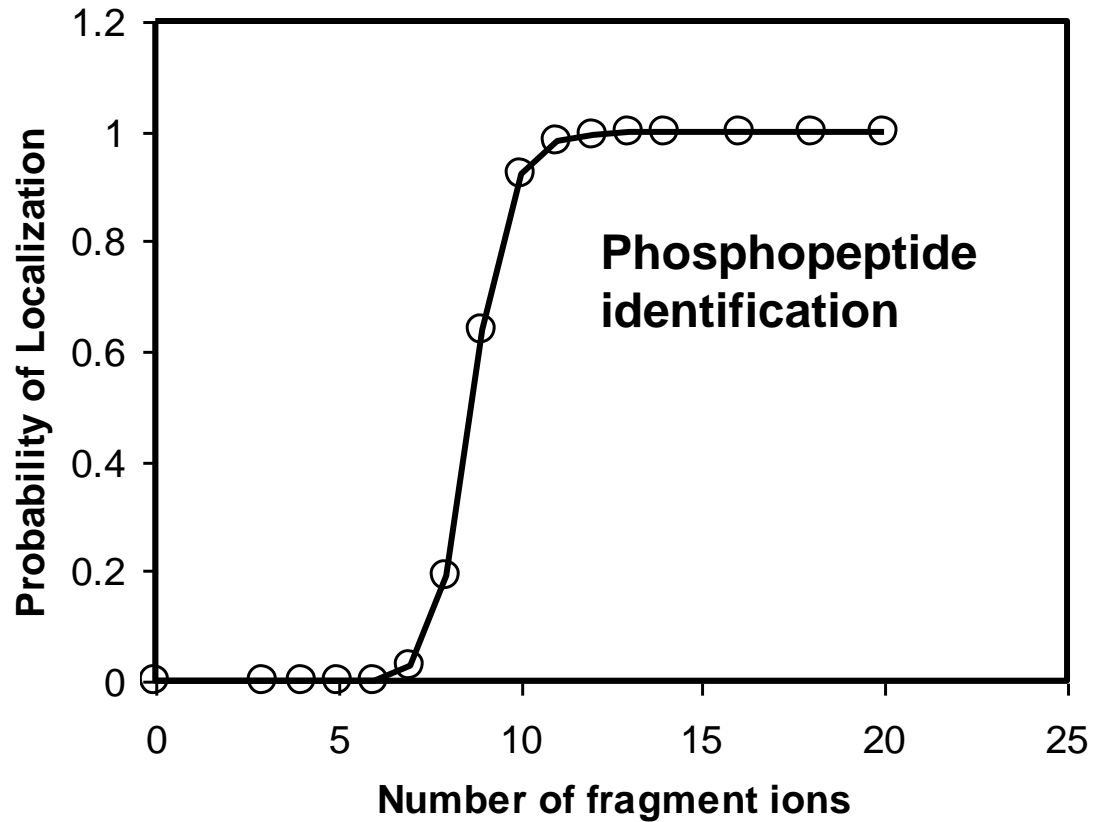
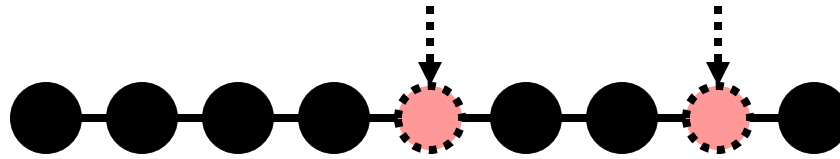
Mass spectrometers have a limited dynamic range and it therefore important to limit the number of possible reactions not to dilute the cross-linked peptides.

For identification of a cross-linked peptide pair, both peptides have to be sufficiently long and required to give informative fragmentation.

High mass accuracy MS/MS is recommended because the spectrum will be a mixture of fragment ions from two peptides.

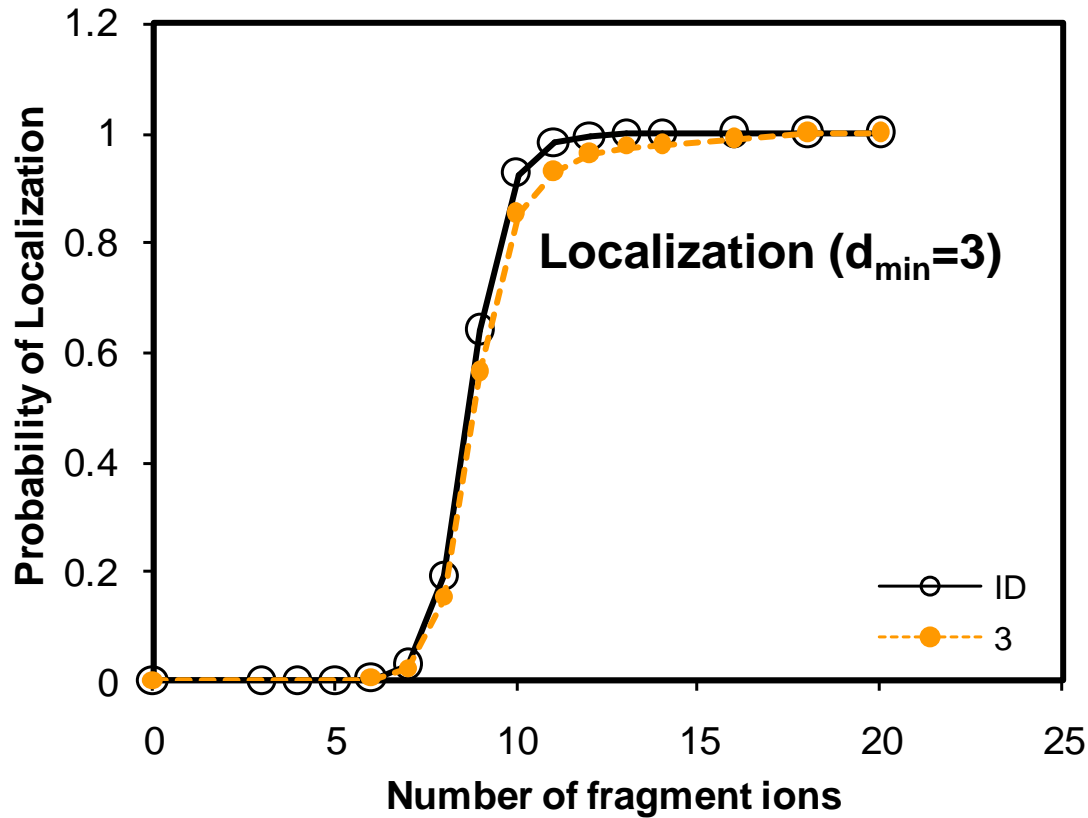
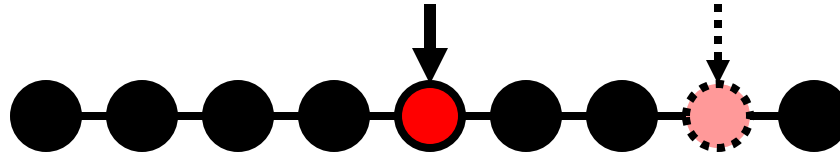
Because the cross-linked peptides are often large, CAD is not ideal, but instead ETD is recommended.

Localization of modifications



$m_{\text{precursor}} = 2000 \text{ Da}$
 $\Delta m_{\text{precursor}} = 1 \text{ Da}$
 $\Delta m_{\text{fragment}} = 0.5 \text{ Da}$
Phosphorylation

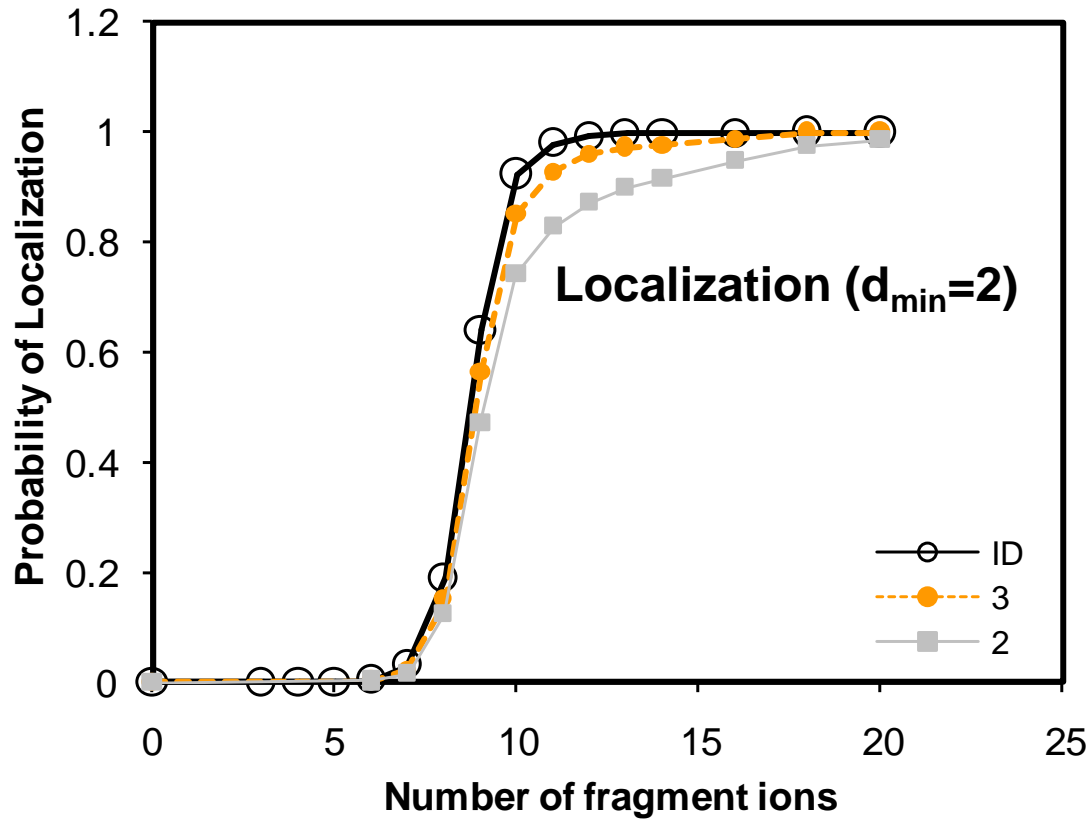
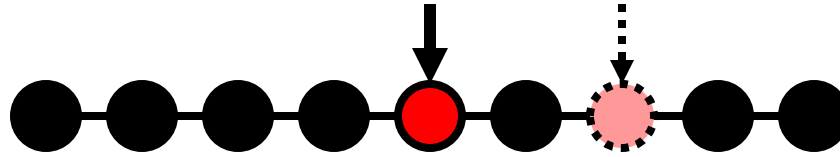
Localization of modifications



$d_{\min} \geq 3$ for 47%
of human tryptic
peptides

$m_{\text{precursor}} = 2000 \text{ Da}$
 $\Delta m_{\text{precursor}} = 1 \text{ Da}$
 $\Delta m_{\text{fragment}} = 0.5 \text{ Da}$
Phosphorylation

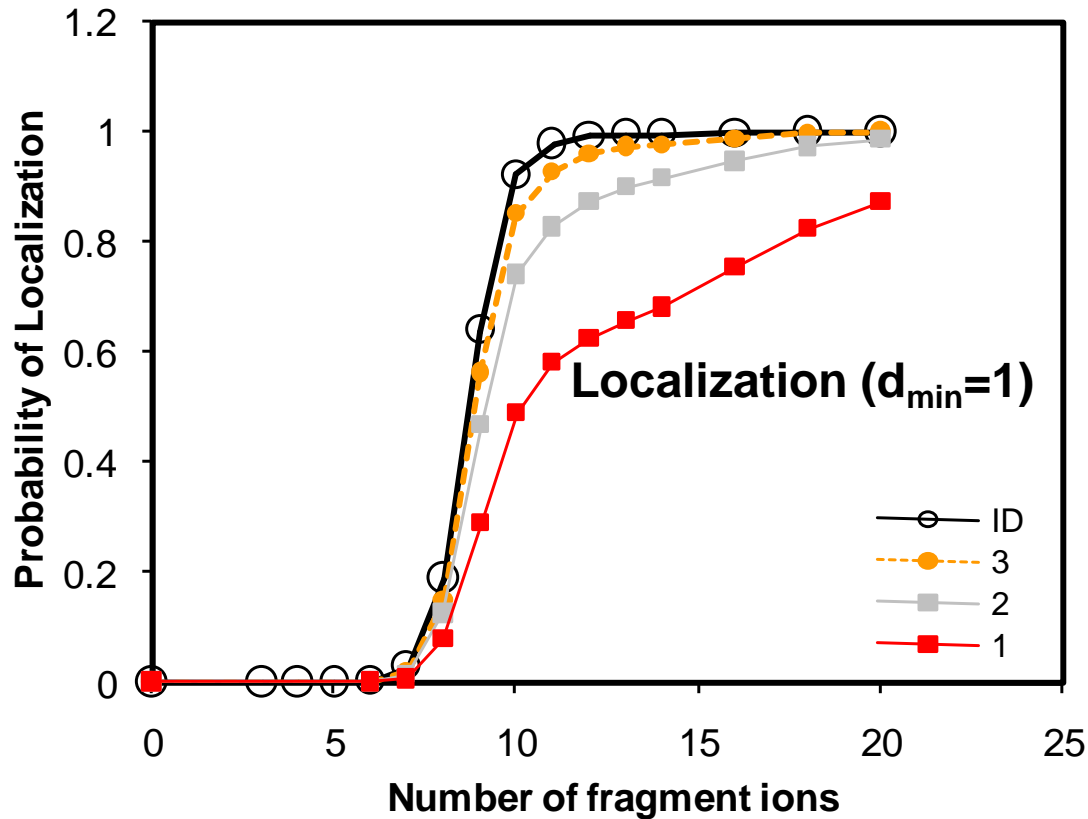
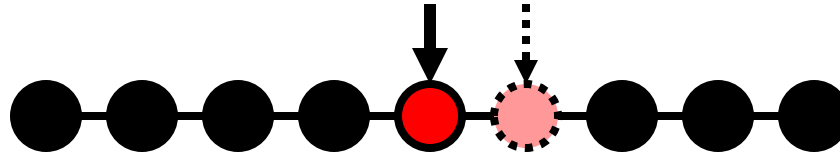
Localization of modifications



$d_{\min}=2$ for 33% of human tryptic peptides

$m_{\text{precursor}} = 2000 \text{ Da}$
 $\Delta m_{\text{precursor}} = 1 \text{ Da}$
 $\Delta m_{\text{fragment}} = 0.5 \text{ Da}$
Phosphorylation

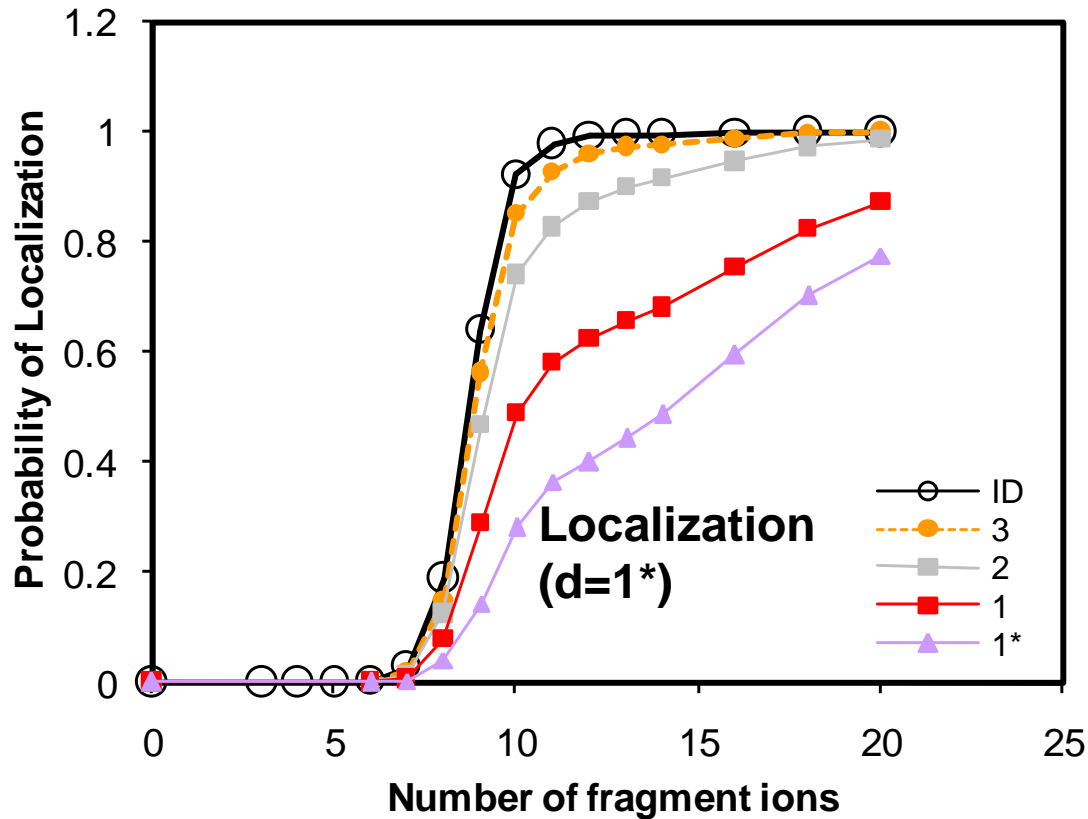
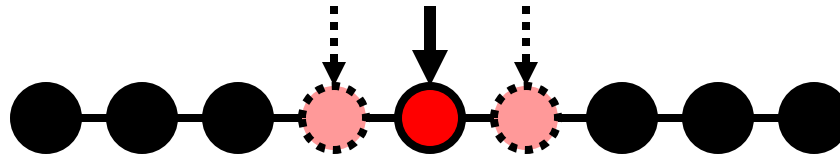
Localization of modifications



$d_{\min}=1$ for 20% of human tryptic peptides

$m_{\text{precursor}} = 2000 \text{ Da}$
 $\Delta m_{\text{precursor}} = 1 \text{ Da}$
 $\Delta m_{\text{fragment}} = 0.5 \text{ Da}$
Phosphorylation

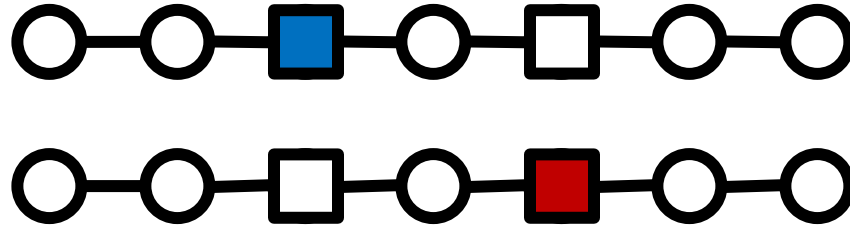
Localization of modifications



$m_{\text{precursor}} = 2000 \text{ Da}$
 $\Delta m_{\text{precursor}} = 1 \text{ Da}$
 $\Delta m_{\text{fragment}} = 0.5 \text{ Da}$
Phosphorylation

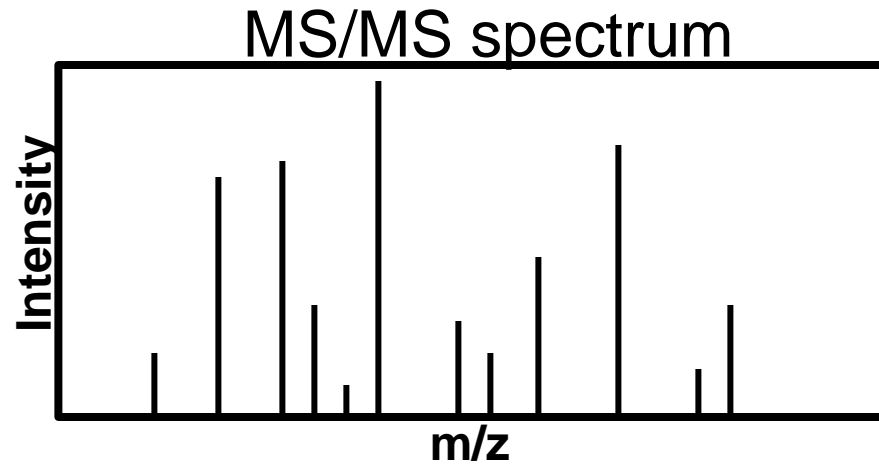
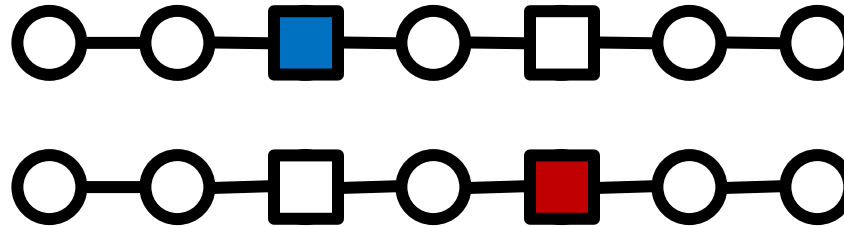
Localization of modifications

Peptide with two possible modification sites



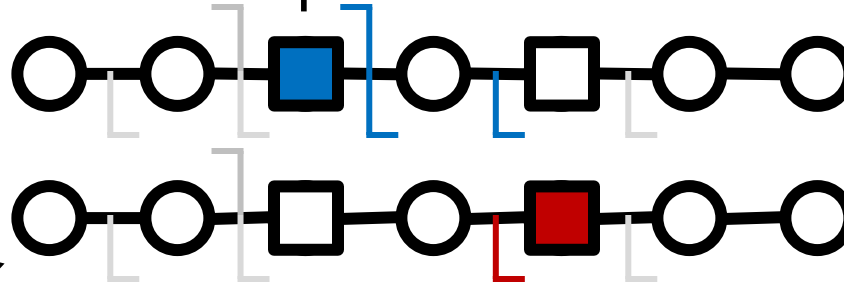
Localization of modifications

Peptide with two possible modification sites

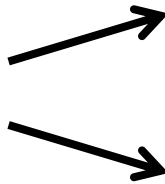


Localization of modifications

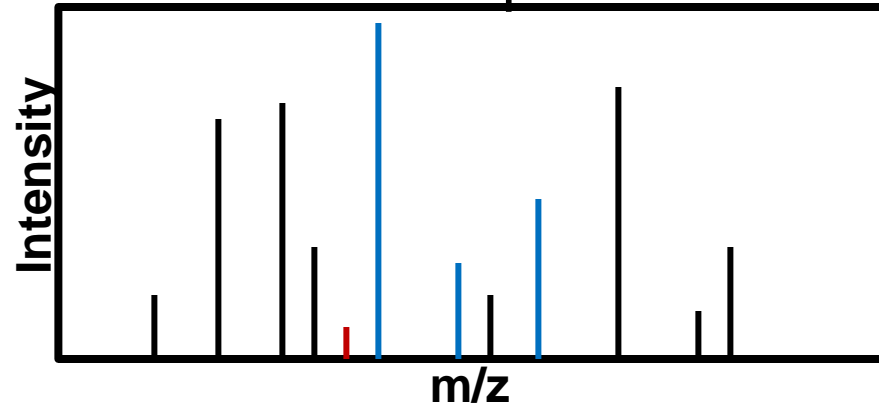
Peptide with two possible modification sites



Matching

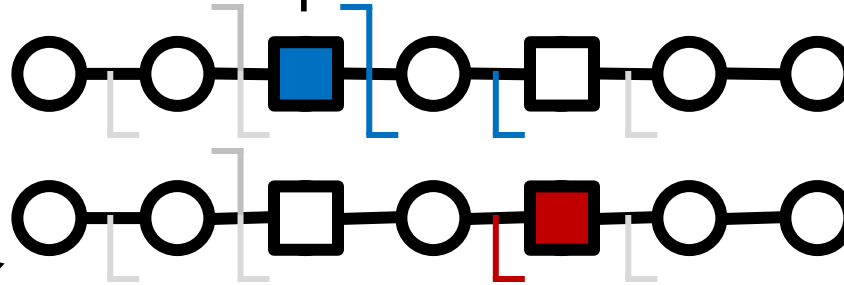


MS/MS spectrum

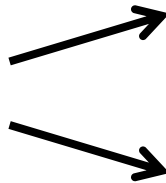


Localization of modifications

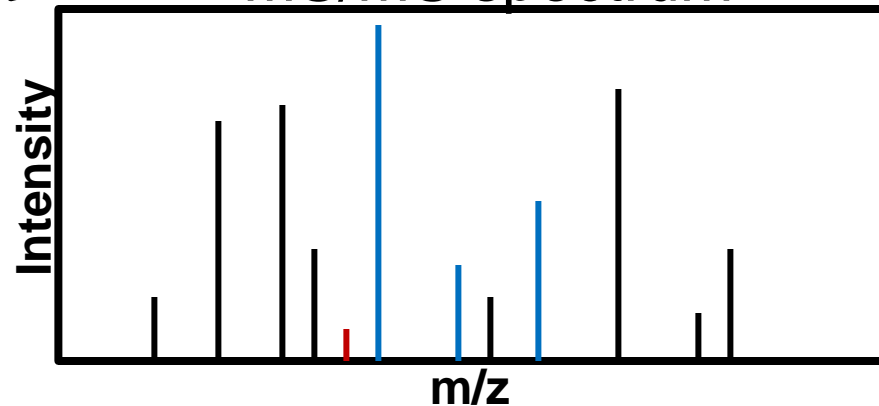
Peptide with two possible modification sites



Matching



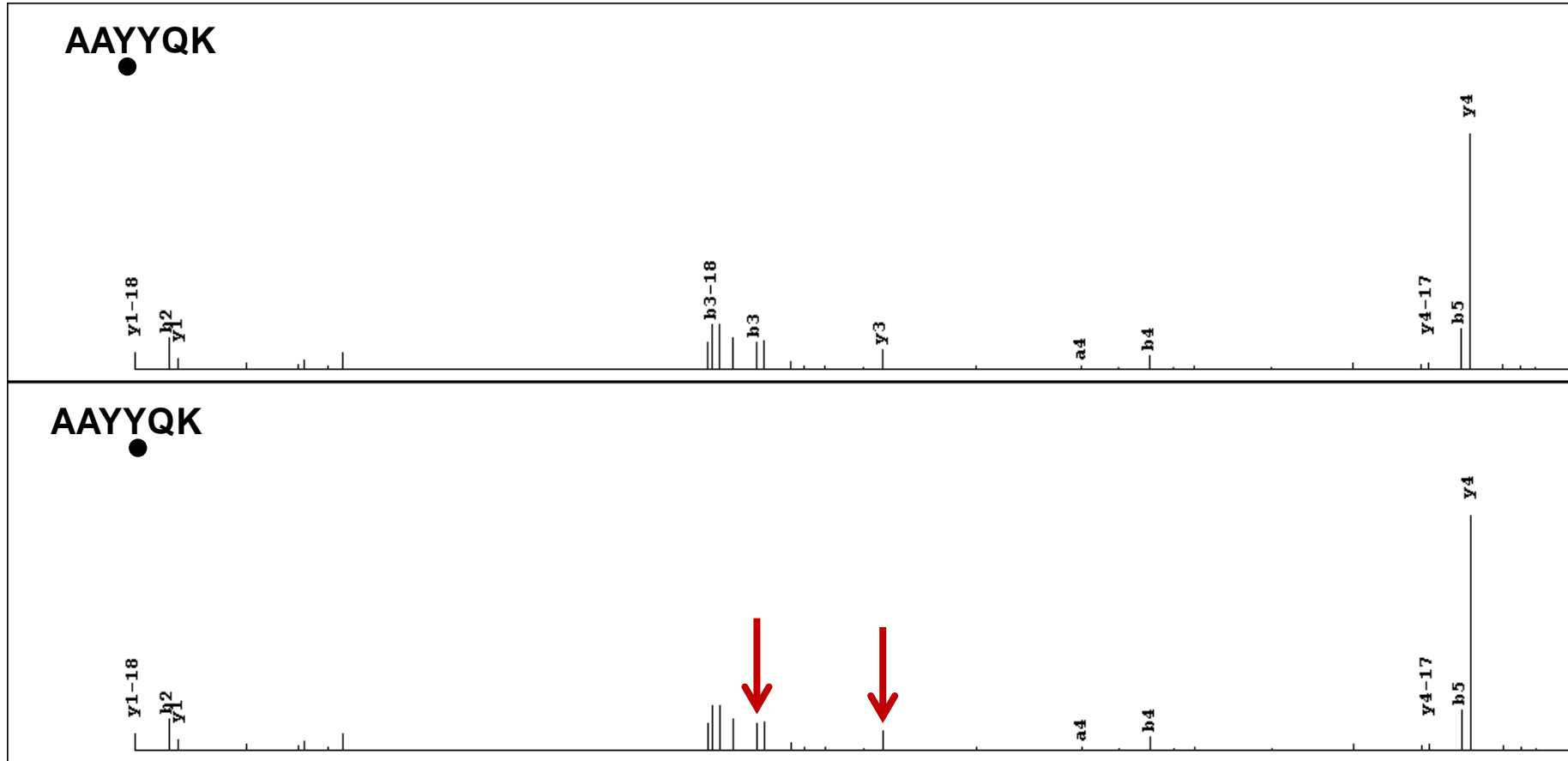
MS/MS spectrum



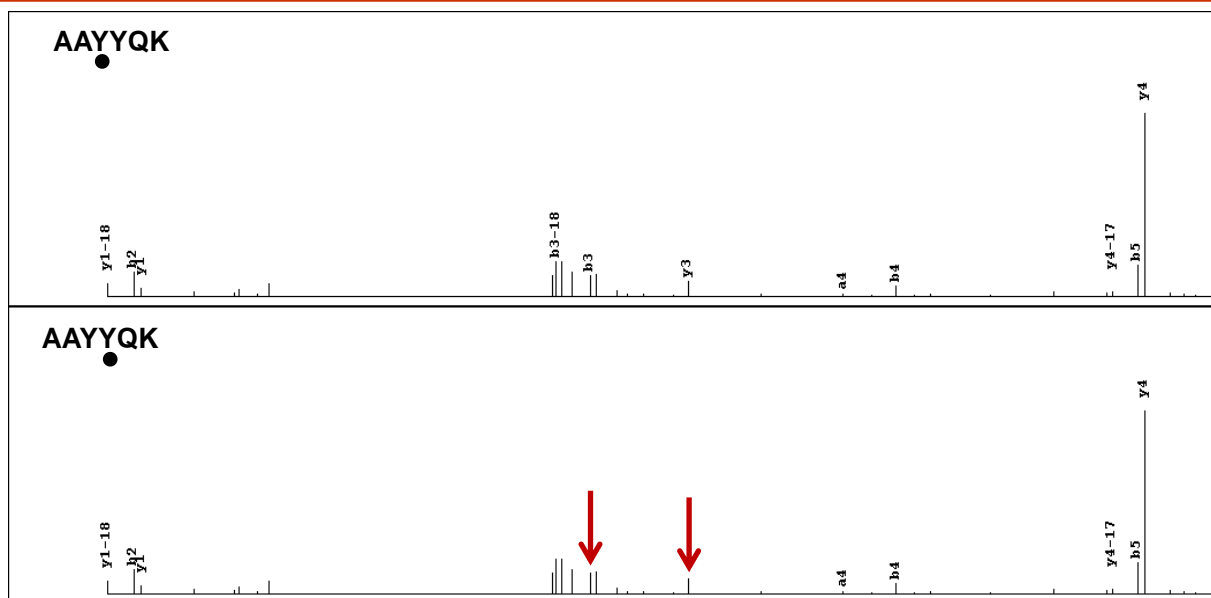
Which assignment does the data support?





1, 1 or 2, or 1 and 2?






Visualization of evidence for localization








Visualization of evidence for localization



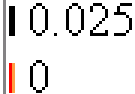




Rank	AAYYQK	Total (matched)	Difference (matched)
1	AAYYQK	 8	-
2	AA Y YQK	  6	 2 0

Rank	AAYYQK	Total (intensity)	Difference (intensity)
1	AA Y YQK	  0.809	-
2	AA Y YQK	  0.735	 0.074 0

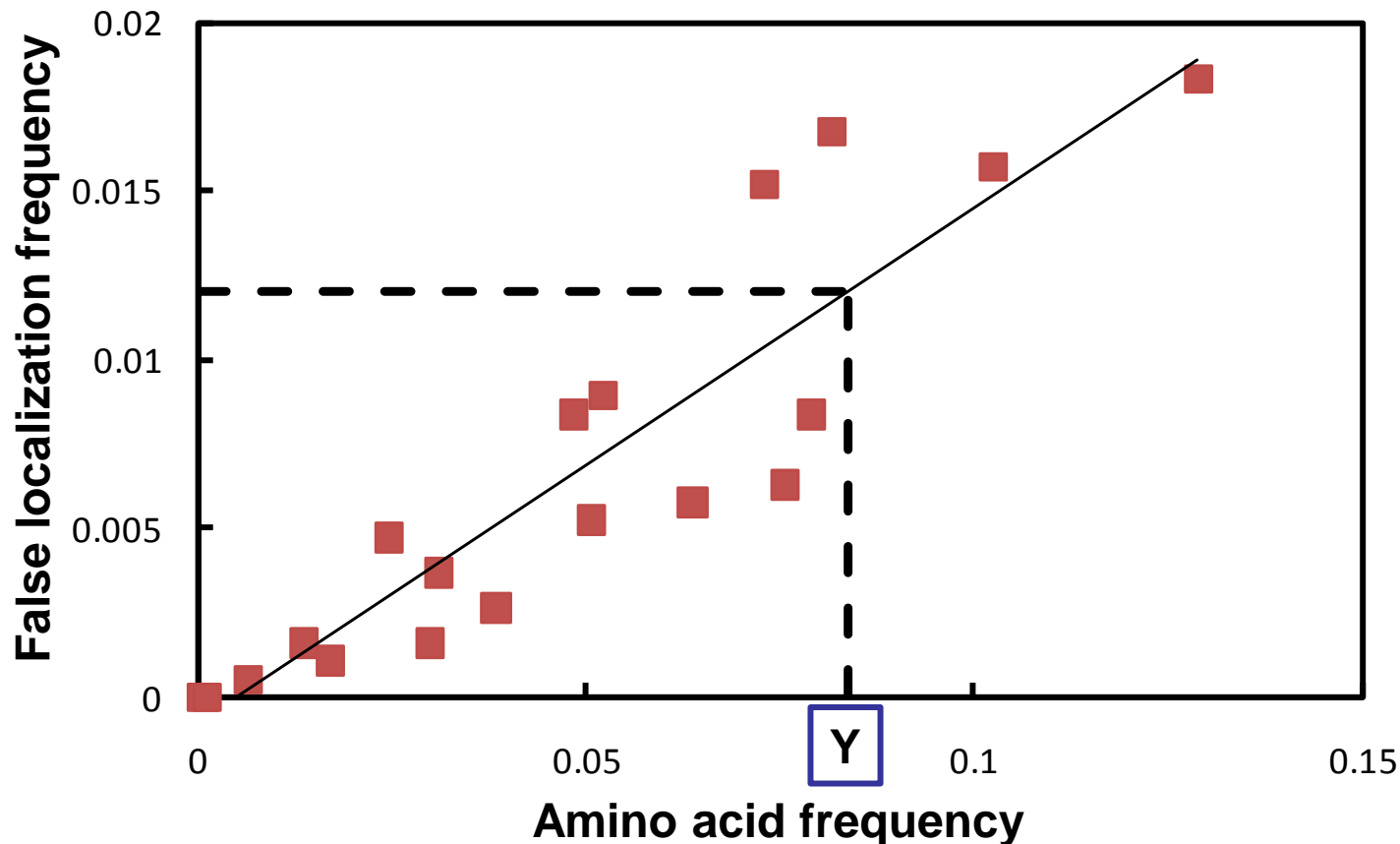
Visualization of evidence for localization

Rank	AAVPSGASTGIYEAL Y LELR	Total (matched)	Difference (matched)
1	AAVPSGASTGIYEAL Y LELR	 15	-
2	AAVPSGASTGIYEAL T LELR	 14	
3	AAVPSGASTGIYEAL S LELR	 13	

Rank	AAVPSGASTGIYEAL T LELR	Total (intensity)	Difference (intensity)
1	AAVPSGASTGIYEAL T LELR	 0.315	-
2	AAVPSGASTGIYEAL S LELR	 0.29	
3	AAVPSGASTGIYEAL Y LELR	 0.187	

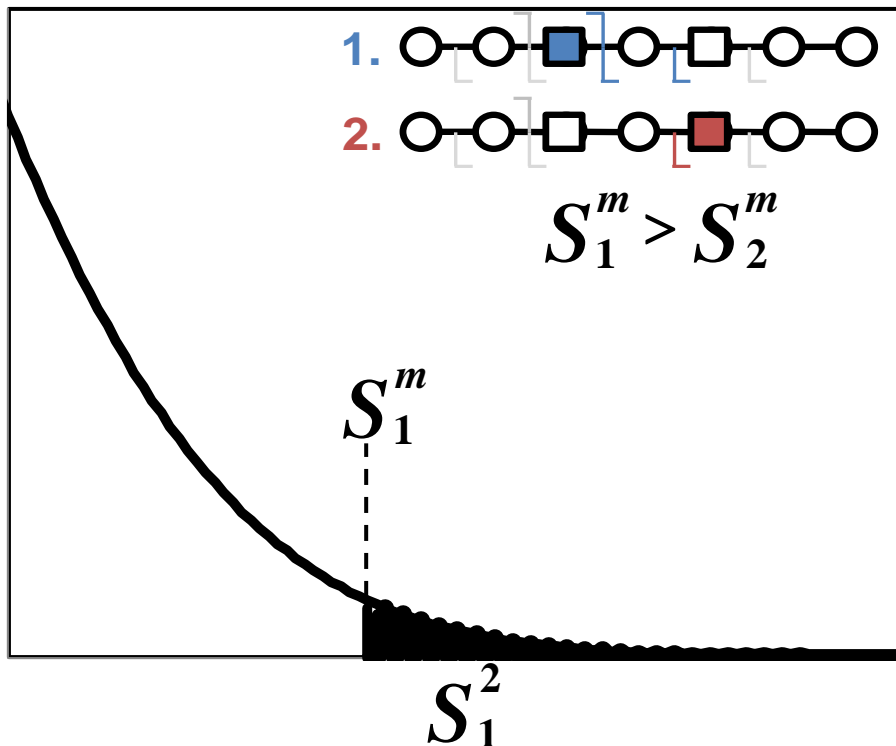
Estimation of global false localization rate using decoy sites

By counting how many times the phosphorylation is localized to amino acids that can not be phosphorylated we can estimate the false localization rate as a function of amino acid frequency.



How much can we trust a single localization assignment?

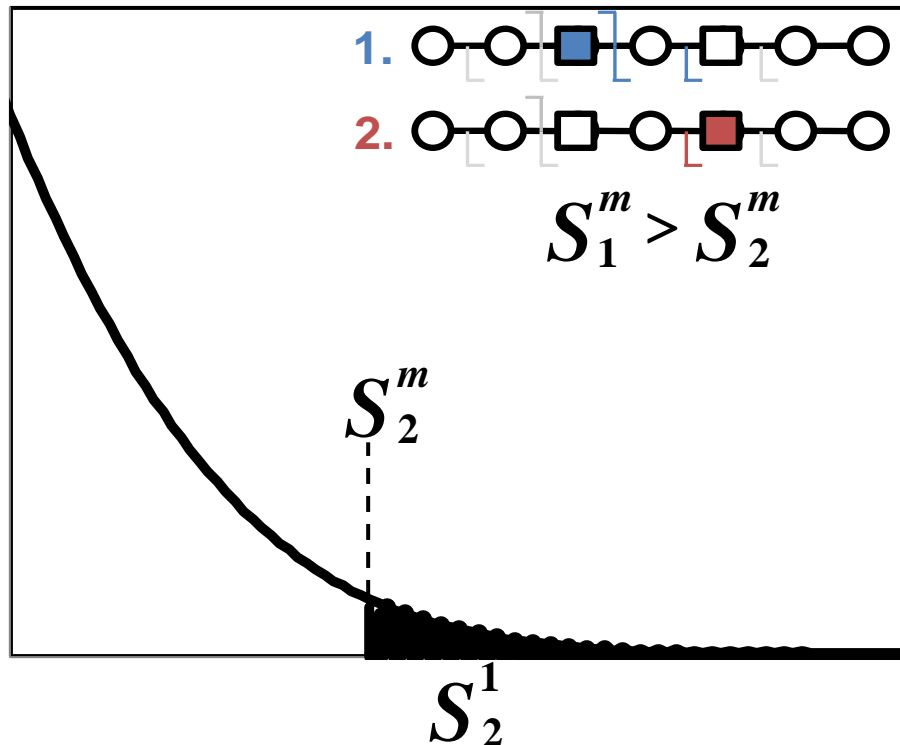
If we can generate the distribution of scores for assignment 1 when 2 is the correct assignment, it is possible to estimate the probability of obtaining a certain score by chance for a given peptide sequence and MS/MS spectrum assignment.



$$p_1^2 = \frac{\int_0^{S_1^m} F^2(S_1^2) dS_1^2}{\int_0^{\infty} F^2(S_1^2) dS_1^2}$$

Is it a mixture or not?

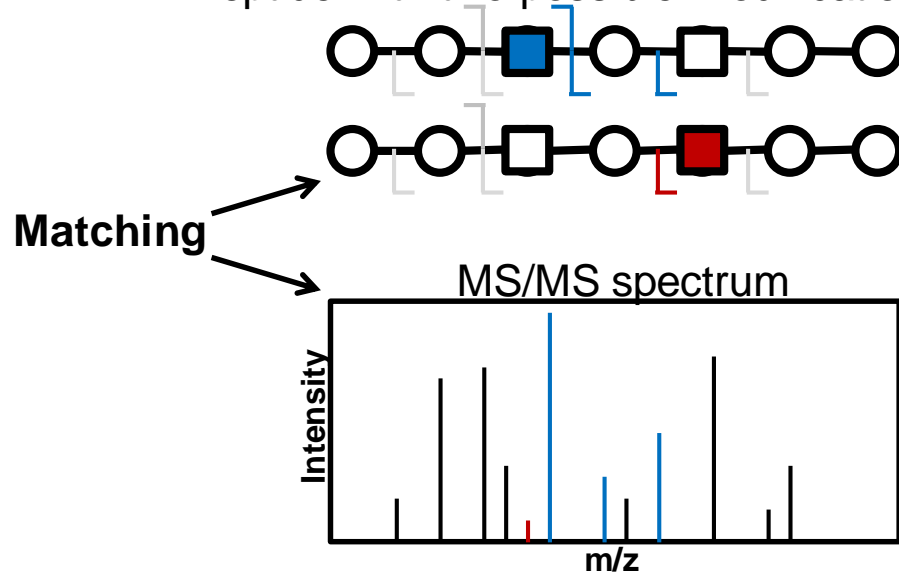
If we can generate the distribution of scores for assignment 2 when 1 is the correct assignment, it is possible to estimate the probability of obtaining a certain score by chance for a given peptide sequence and MS/MS spectrum assignment.



$$p_2^1 = \frac{\int_0^{S_2^m} F^1(S_2^1) dS_2^1}{\int_0^\infty F^1(S_2^1) dS_2^1}$$

Localization of modifications

Peptide with two possible modification sites



Which assignment does the data support?

1, 1 or 2, or 1 and 2?

$$p_1^2 \leq p_{th} \text{ and } p_2^1 \leq p_{th} \Rightarrow \mathbf{1 \text{ and } 2}$$

$$p_1^2 \leq p_{th} \text{ and } p_2^1 > p_{th} \Rightarrow \mathbf{1}$$

$$p_1^2 > p_{th} \text{ and } p_2^1 \leq p_{th} \Rightarrow \emptyset \quad (S_1^m \geq S_2^m \Rightarrow p_1^2 \leq p_2^1)$$

$$p_1^2 > p_{th} \text{ and } p_2^1 > p_{th} \Rightarrow \mathbf{1 \text{ or } 2}$$

**Proteomics Informatics -
Protein characterization: post-translational
modifications and protein-protein
interactions (Week 10)**
