

Identifying intact crosslinks with Mascot Server

Ville Koskinen
Matrix Science

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



A new feature in Mascot Server 2.7 is the ability to identify crosslinked peptides from MS/MS data. I'll show you how to use the new feature and discuss validating crosslinked search results.

Crosslinking experiments

- **Naturally occurring**
 - E.g. disulfide bond
- **Artificially created for elucidating:**
 - Protein tertiary structure
 - Protein complex topology
 - Conformational changes
 - Protein-protein interactions

Crosslinks can be roughly divided into naturally occurring and artificially created. An example of a naturally occurring crosslink is the disulfide bond. In most experiments, disulfide bonds are of no interest, so they are cleaved with a reducing agent prior to analysis. But, if the goal is to study protein folding, identifying the sites linked by a disulfide bond is crucial.

Crosslinks can also be created artificially by a chemical reaction. A crosslink can only be formed between sites that are adjacent in three-dimensional space. The maximum distance depends on the spacer arm length of the linker. So, the locations of the identified links can be used for constraining the three-dimensional shape of the molecule, which has obvious applications in elucidating protein tertiary structure, protein complex topology, conformational changes and protein-protein interactions.

Crosslinking experiments

- **Enzymatic digestion produces:**
 - Linear peptides
 - Crosslinked peptides: *alpha* linked with *beta*
- **Cleavable crosslink**
 - Cleaved prior to analysis, e.g. irradiation or CID
 - Leaves behind a fragment at linked site
- **Intact crosslink**
 - Survives MS/MS analysis
 - Spectrum has fragments from both peptides

MASCOT : Identifying intact crosslinks

© 2020 Matrix Science



Whether natural or artificial, a typical bottom-up crosslinking experiment starts by digesting the crosslinked proteins. This produces two types of precursor: linear peptides and crosslinked peptides. A linear peptide either doesn't span a linkable site or if it does, the linker did not form a crosslink. Otherwise, the precursor is a crosslinked peptide, where an alpha peptide is linked with a beta peptide. The alpha peptide is conventionally the longer or more massive of the two.

How crosslinked peptides appear to a search engine like Mascot depends on the type of the linker.

Cleavable crosslinks are designed to be cleaved prior to analysis by various means like irradiation or CID. The link cleaves asymmetrically, so that the alpha and beta peptides carry linker fragments with different mass. The fragment is modelled exactly like a variable modification on a linear peptide. The linear peptide matches are then paired up by a postprocessing algorithm, which figures out which peptides come from the same crosslinked precursor.

Conversely, intact crosslinks survive MS/MS analysis. There is additional complexity to the crosslinked spectrum: alpha fragments that contain the linked site have a mass offset equal to the beta peptide mass plus the linker molecule. Vice versa for beta

fragments. The search engine must be able to fragment crosslinked peptides *in silico*. However, the advantage is, no postprocessing step is needed.

The boundary between cleavable and intact is a bit fuzzy. If the cleaving efficiency isn't 100%, there will be some precursors with an intact link. Similarly, an intact linker can fail to form a link, which leaves behind a linker fragment.

Mascot Server 2.7

- **Model the experiment with a *crosslinking method***
 - Linker, linkable sites, directionality, crosslinked products (intralink, interlink, looplink, monolink), protein accessions, filters
- **Plus standard search parameters**
 - Applies to both linear and crosslinked peptides
- **Mascot ships with 13 linker definitions (Unimod), or create your own**

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



Because of the great variety and complexity in crosslinking experiments, we've chosen an approach similar to quantitation. We want to keep the software universal and configurable, and not tied to a specific experimental workflow.

The crosslinking experiment is modeled with a crosslinking method. The method defines which linker is used, which linkable sites are considered, linker directionality, what are the expected crosslinked products and which protein or proteins are being linked. There are also a few filters available.

There are four possible crosslinked products. An intralink is an intact crosslinked peptide within the same protein, while an interlink is an intact link between two different proteins. A looplink is an intact link within a peptide. Finally, monolinks model cleavable linker fragments as well as dead-end or quenched links where only one end is attached to a peptide. I'll show examples of crosslinked products later.

The standard search parameters like enzyme, fixed mods and variable mods are applied to both linear peptides and crosslinked peptides.

The linker definitions are standard Unimod entries. Mascot 2.7 ships with 13

definitions: 7 for intact crosslinks and 6 for cleavable crosslinks. You can easily create your own definition for any linker chemistry.

Interlinked HOP2-MND1 complex

- Rampler et al: “Comprehensive Cross-Linking Mass Spectrometry Reveals Parallel Orientation and Flexible Conformations of Plant HOP2-MND1”, J Prot Res 14(12):5048-5062, 2015
- PRIDE: PXD001538
- Purified *A. thaliana* HOP2-MND1 complex
- DSS/BS3 linker, biological replicate 1

MASCOT : Identifying intact crosslinks

© 2020 Matrix Science



I'll illustrate intact crosslinking with this data set from the Mechtler lab. The data is available as the PRIDE project PXD001538. The study uses a purified HOP2-MND1 complex from *Arabidopsis thaliana*, which is crosslinked using three different linkers. I've chosen the biological replicate 1 linked with disuccinimidyl suberate (DSS). DSS has two identical reactive groups, so it's homobifunctional. It's amine reactive and creates intact crosslinks between lysines and protein N-term. The study was designed so that there are only lysine-lysine links.

Interlinked HOP2-MND1 complex

- **Peak pick in Distiller**
 - Crosslinked precursor charge often 3+ or more, so output fragment masses as MH+
- **Set up a crosslinking method**
 - Linker specificity: Xlink:DSS (K)
 - Monolinks: enable W (156 Da)
 - Accessions: HOP2_ARATH, MND1_ARATH
 - Scope: enable InterLink, IntraLink, LoopLink
 - Filters: MinLen=4

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



The raw data was peak picked in Mascot Distiller. Crosslinked precursors often have charge of 3+ or more, which means the product ions could also be highly charged. It's best to decharge the fragments by choosing to output fragment masses as MH+.

Setting up the crosslinking method is straightforward. DSS has a Unimod definition titled Xlink:DSS, and we only want the (K) specificity for lysine-lysine links. The authors determined only one of the three monolinks is possible, in this case the water quenched 156 Da monolink. This has the code W in the Unimod definition.

The accessions in the crosslinking method act as selectors and instruct Mascot to choose only the listed proteins for crosslinking. The sequences are in SwissProt, HOP2_ARATH and MND1_ARATH, so we can just list them here and select SwissProt as the search database when the search is submitted.

The crosslinked products to consider are protein interlinking and intralinking, and peptide looplinking. Finally, allow the beta peptide to be as short as 4 residues. If no MinLen is specified, Mascot filters out alpha and beta peptides shorter than MinPepLenInSearch, which is 7 by default.

Interlinked HOP2-MND1 complex

```
<mxm:method description="" name="PDX001538 DSS" strategy="Brute-force">
  <mxm:linkers>
    <mxm:linker ModFileName="Xlink:DSS (K)">
      <mxm:monolink>W</mxm:monolink>
    </mxm:linker>
  </mxm:linkers>
  <mxm:accessions>
    <mxm:accession>MND1_ARATH</mxm:accession>
    <mxm:accession>HOP2_ARATH</mxm:accession>
  </mxm:accessions>
  <mxm:scope>
    <mxm:parameter name="InterLink">True</mxm:parameter>
    <mxm:parameter name="IntraLink">True</mxm:parameter>
    <mxm:parameter name="LoopLink">True</mxm:parameter>
  </mxm:scope>
  <mxm:filters>
    <mxm:parameter name="MinLen">4</mxm:parameter>
  </mxm:filters>
</mxm:method>
```

Linker, linkable sites,
monolinks

Protein selectors

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



The crosslinking method is a simple XML file, shown here. Monolinks are toggled in the mxm:linker section. The mxm:accession elements can optionally specify a database name. If there is no database name, Mascot chooses proteins for crosslinking by accession, regardless of database. The scope and filters sections contain simple key-value parameters.

Interlinked HOP2-MND1 complex

- **Peak pick in Distiller (done)**
- **Set up a crosslinking method (done)**
- **Submit from Distiller**
 - Database: contaminants, SwissProt, *A. thaliana* taxonomy
 - Choose crosslinking method in search form
- **Open full report in browser**
 - Distiller support for displaying crosslinked matches will be added in a future release

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



So, the raw data has been peak picked and we have a crosslinking method. The search is submitted as normal from Distiller. Since the two proteins are in SwissProt, I've chosen to search SwissProt with *A. thaliana* taxonomy for convenience. Also, remember to include a contaminants database. Choose the right crosslinking method in the search form.

Once the search finishes, choose to open the full report in browser. At the moment, Distiller can only display linear peptides. Support for crosslinked matches will be added in a future release.

MASCOT Search Results

User :
E-mail :
Search title : 20140603_QEx3_RSLC4_Ramplер_Mechtler_IMP_shotgun_DSS12.raw
MS data file : 20140603_QEx3_RSLC4_Ramplер_Mechtler_IMP_shotgun_DSS12.temp.mgf
Databases : **1:** contaminants 20090624 (262 sequences; 133,770 residues)
2: SwissProt 2020_01 (561,911 sequences; 202,173,710 residues)
Taxonomy : **1:** (none)
2: Arabidopsis thaliana (thale cress) (15,940 sequences)
Timestamp : 27 Apr 2020 at 15:21:14 GMT

Re-search All Non-significant Unassigned [\[help\]](#) Export As FASTA

► Search parameters
 ► Score distribution
 ▼ **Modification statistics for all protein families**

Modification	Delta	Type	Site	Total matches
Xlink:DSS[W]	156.078645	variable	K	465
Oxidation	15.994915	variable	M	340
Xlink:DSS	138.06808	crosslink	K - K	218
Xlink:DSS	138.06808	looplink	K - K	140
Carboxymethyl	58.005479	fixed	C	21

► Legend
Protein Family Summary

MASCOT : Identifying intact crosslinks © 2020 Matrix Science

The Protein Family Summary report has been extended to display crosslinked search results. The modification statistics table gives counts of discovered intact links, monolinks and looplinks. Xlink:DSS[W] on the first row is the 156Da monolink enabled in the crosslinking method. The crosslink and looplink rows display the intact linker mass as well as the linked sites.

The screenshot shows the Mascot search results interface. At the top, the 'Sensitivity' section is visible, with 'Target' set to '1180'. Below this, there are buttons for 'Proteins (35)', 'Report Builder', and 'Unassigned (9195)'. The main section is titled 'Protein families 1-10 (out of 35)'. A search bar contains 'Accession' and 'contains'. The results are listed in a table with columns for family ID, protein name, and description. The entry '2::HOP2_ARATH' is circled in red. The footer includes the Mascot logo, the tagline 'Identifying intact crosslinks', the copyright '© 2020 Matrix Science', and the Matrix Science logo.

Family ID	Protein Name	Description
1	2::HOP2_ARATH	17682 Homologous-pairing protein 2 homolog OS=Arabidopsis thaliana OX=3702 GN=HO...
2	2::MND1_ARATH	17505 Meiotic nuclear division protein 1 homolog OS=Arabidopsis thaliana OX=3702 GN=...
3	1::IPI:CON_Trypsin ...	202 Trypsin - Sus scrofa (Pig).
4	1::IPI:CON_0000986...	75 Tax_Id=9606 Gene_Symbol=KRT10 Keratin, type I cytoskeletal 10
5	1::IPI:CON_0013136...	62 Tax_Id=10090 Gene_Symbol=Krt6a Keratin, type II cytoskeletal 6A
6	2::TBA2_ARATH	53 Tubulin alpha-2 chain OS=Arabidopsis thaliana OX=3702 GN=TUBA2 PE=2 SV=2
7	1::IPI:CON_0069890...	42 (Bos taurus) Actin, cytoplasmic 1
8	2::EDL18_ARATH	41 Sugar transporter ERD6-like 18 OS=Arabidopsis thaliana OX=3702 GN=SFP2 PE=2...
9	1::IPI:CON_IPI00121...	37 Tax_Id=10090 Gene_Symbol=Krt23 Keratin, type I cytoskeletal 23
10	1::IPI:CON_IPI00311...	36 Tax_Id=10090 Gene_Symbol=Krt18 Keratin, type I cytoskeletal 18

Protein inference is not affected by interlinking. Because HOP2_ARATH and MND1_ARATH don't share any significant peptide matches, protein inference puts them in different families, regardless of whether there are any interlinked peptide matches.

Notice also the small amounts of evidence for other Arabidopsis proteins in the sample, despite the high levels of purification.

88	6.1e-09	▶1	U K.SLQSNLTLEEIQERDAKLRK.E + K17<-Xlink:DSS[I]->K20
106	1.2e-10	▶1	U K.SLQSNLTLEEIQERDAKLRK.E + Xlink:DSS[W]
123	3e-12	▶1	U K.SLQSNLTLEEIQERDAKLRK.E + Xlink:DSS[W]
49	7.5e-05	▶1	U K.EMEEKLVKLRREGITLVRPEDK.K
59	3.2e-06	▶1	U K.TAVQKALDSLADAGK.I K5<-Xlink:DSS->K2 K.QKIYIAR.Q
34	0.00062	▶1	U K.TAVQKALDSLADAGK.I K5<-Xlink:DSS->K2 K.QKIYIAR.Q
46	5.4e-05	▶1*	[2:;MND1_ARATH]LESDLQGSNK K10<-Xlink:DSS->K8 K.LQEQLQEKK.K

Annotations in the image:

- Looplinded peptide**: points to the first row (88).
- Monolinks**: points to the second and third rows (106 and 123).
- Intralinks in HOP2_ARATH**: points to the fourth, fifth, and sixth rows (49, 59, and 34).
- Interlink to MND1_ARATH**: points to the seventh row (46).

MASCOT : Identifying intact crosslinks © 2020 Matrix Science

Expand the first protein hit, HOP2_ARATH. The peptide match table displays crosslinking data in a compact form, and I've selected a few example rows here.

On the first line is a looplinked peptide. A looplink is an intact link between two residues in the same peptide. Here the intact link is between lysines K17 and K20. Looplinks use the symbol [I] in brackets to indicate it's an intact link.

The next two rows show matches to the same peptide sequence, this time with the W monolink. Monolinks are variable modifications generated from the linker definition. With cleavable linkers, this is usually the only crosslinked product of interest, since monolinks represent the linker fragments. Monolinks can also be used for modelling dead-end or quenched links, such as the water quenched monolink here.

Further down, there are three matches with an intact crosslink. The first two are protein intralinks, which means it's an intact link between two non-overlapping peptides within a single protein.

The last line shows a match to a protein interlink. An interlink is an intact link between two peptides from different proteins. The alpha peptide is in MND1_ARATH

and the beta peptide in HOP2_ARATH.

The scores for all the matches are standard Mascot scores, and the identity and homology thresholds are calculated using the same statistical model as for linear peptide matches.

Validating intact crosslinked matches

- **Target-decoy searching?**
 - Works with cleavable linker, because linker fragments are just variable mods
- **Intact linker**
 - Sample has just two purified proteins
 - Vast majority of peptide matches will cluster in the two proteins
 - Too few sequences for meaningful FDR estimate

MASCOT : *Identifying intact crosslinks*

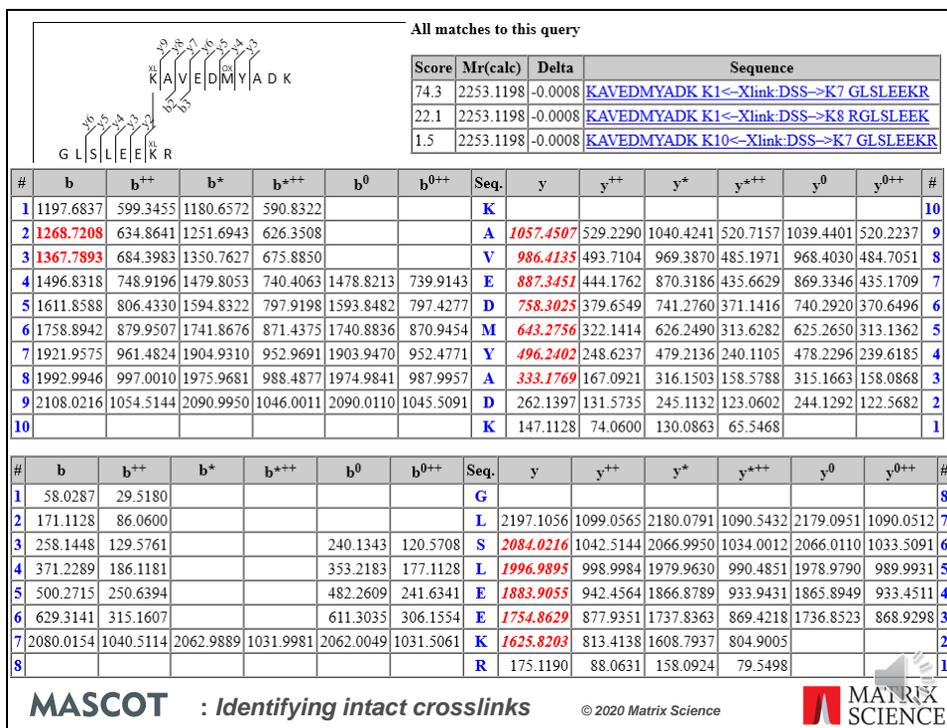
© 2020 Matrix Science



Which of the matches on the previous slide are correct? The usual approach is to do a target-decoy search, which allows you to estimate the false discovery rate. This approach works if you're using a cleavable linker. The linker fragment is just a variable modification and all the peptide matches are ordinary linear peptides., so standard FDR estimation techniques work well in this case.

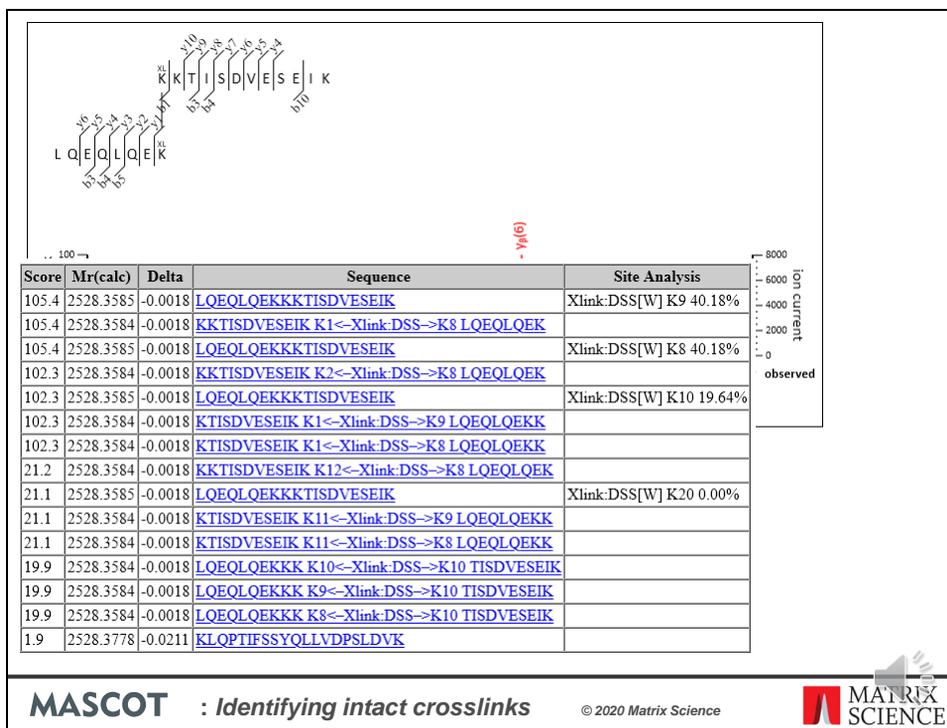
However, when using intact linkers, the search space is often too small for a meaningful estimate of the FDR. In the present case, the sample is expected to contain only the two purified proteins. If you set up a two-sequence database and search ten thousand MS/MS spectra against it, no statistical model can give a meaningful estimate of the FDR. In this case, I've used as the target database the fifteen thousand *A. thaliana* proteins in SwissProt, plus a couple hundred contaminant sequences, which allows incorrect matches to distribute a bit more evenly. Even so, the peptide matches completely saturate the available peptide sequences from the two target proteins, and there are few matches in other proteins.

The Mascot statistical model is robust and will give you reasonable separation between correct and incorrect matches, even in these adverse conditions. Nonetheless, it's prudent to keep a sceptical mind. Let's look at a few example spectra.



Here's a Peptide View report of a crosslinked match with excellent fragmentation from both alpha and beta, with score 74. Further down the page, there is a fragmentation table for the alpha peptide and another for the beta peptide. The fragment masses have no overlap between the two, which shows both peptides have independent evidence.

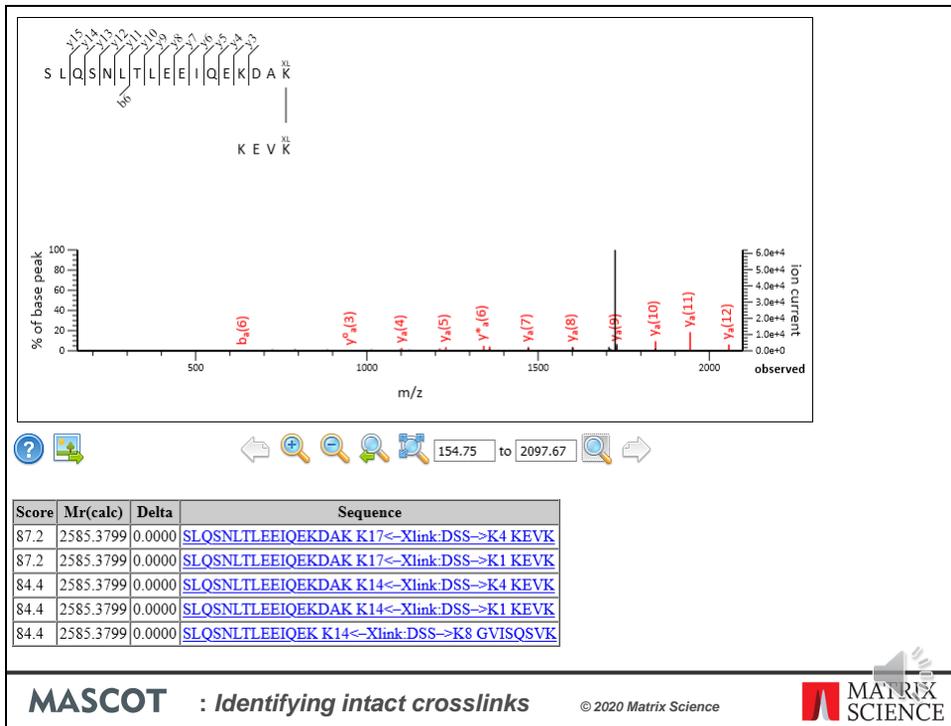
At the bottom of the page, you can see the query also matches a slightly different beta peptide. However, the rank 2 score is much lower, only 22, and there's no ambiguity in the alpha peptide, so we can be confident the top match is the correct match.



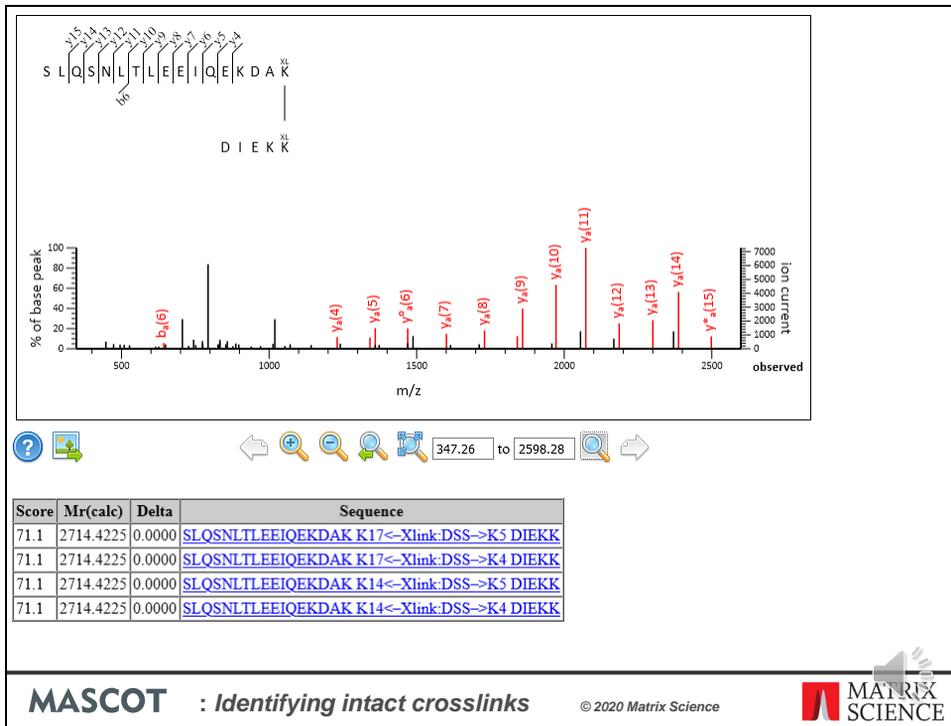
Here is another crosslinked peptide match with a high score and pretty nice fragmentation. At such a high score, 105, it looks very promising, but I have to dash your hopes.

Here are the matches at the other ranks. It turns out, beta and alpha are sequential in the protein sequence. The concatenated beta-alpha sequence with a water quenched monolink has exactly the same score and equally convincing fragment matches. The location of the potential monolink or crosslink is also in question. The MS/MS spectrum has no evidence to differentiate between the alternatives, so choosing the correct match requires consulting structural data.

This also shows the importance of presenting an integrated report. If the search engine matched only crosslinked or only linear peptides, the ambiguity could be easily missed.



Here's a crosslinked match where the identity of the beta peptide should be questioned. The alpha peptide has a nice sequence ladder, but there's no fragmentation from the beta. There's an almost as good a match at rank 5, where the beta peptide is completely different. The alternatives score almost the same because there are no alpha fragments that span the linker site. The two alpha peptides match exactly the same observed fragments.



Finally, this crosslinked match is very similar to the previous one with no fragmentation from the beta peptide. But in this case, the search space has no other alternative for beta. If we assume the search space is complete, there's no reason to doubt the identity of the beta peptide, even if there are no beta fragments in the spectrum.

Visualising in xiVIEW

- xiVIEW is a search engine independent, web-based visualisation tool
- Maintained by the Rappsilber laboratory
- <https://xiview.org>

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



Other common ways to validate intact crosslinked matches is by comparing with existing annotations, and comparing with structural data. If the proposed crosslink would have a length longer than the linker's spacer arms, it's either a false positive or the structural data is incorrect. Mascot doesn't currently make use of structural data, so the validation must be done with external tools.

One way to visualise the results is by viewing them in xiVIEW. This is a search engine independent tool maintained by the Rappsilber laboratory.

Visualising in xiVIEW

- **Choose xiVIEW CSV format in Mascot export form**
- **Also export peak lists as MGF**
- **Export protein sequences as FASTA**
 - Not needed if accessions are in sp|Q9FX64|HOP2_ARATH format
- **Upload the data to <https://xiview.org>**

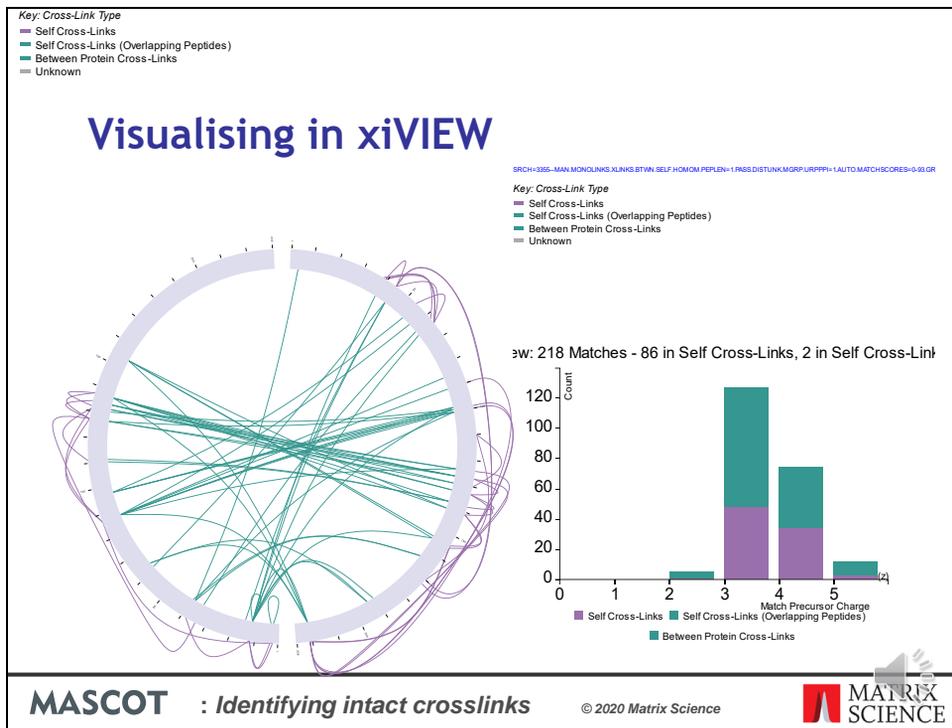
MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



Export the results in xiVIEW CSV format from the Mascot export form. If you want to view the annotated spectra in xiVIEW, also export the peak lists as MGF. Finally, export the sequences as FASTA. xiVIEW can fetch the protein sequences and structural data automatically. Unfortunately, the default accession format for SwissProt in Mascot is the AC format, and xiVIEW expects the raw accession.

Upload the files to xiview.org.



Here's an example circular plot of the crosslinks between the two proteins, based on the Mascot search data. xiVIEW can also plot histograms like this one showing the crosslinked match charge states. As expected, nearly all of the crosslinked matches have charge at least 3+ or higher.

Unfortunately, there is currently no crystallography data for the HOP2-MND1 complex from *Arabidopsis thaliana*, so I can't show you the 3D view.

Visualising in xiVIEW

	A	B	C	D	E	F	G	H	I	J	K
1	Id	Protein1	PepPos1	PepSeq1	LinkPos1	Protein2	PepPos2	PepSeq2	LinkPos2	Score	CrossLinkerModMass
2	q1645_p1_1	MND1_ARATH	6	GLSLEEKR	7	HOP2_ARATH	189	DVKELK	6	16.18	138.06808
3	q2520_p1_1	MND1_ARATH	6	GLSLEEK	7	MND1_ARATH	35	M15.994915GPK	4	15.98	138.06808
4	q3722_p1_1	MND1_ARATH	6	GLSLEEK	7	MND1_ARATH	2	SKKR	3	16.33	138.06808
5	q3723_p1_1	MND1_ARATH	6	GLSLEEK	7	MND1_ARATH	2	SKKR	3	16.07	138.06808
6	q4017_p1_1	MND1_ARATH	4	KRGLSLEEK156.078645R	1	MND1_ARATH	1	M15.994915SKK156.078645	3	24.69	138.06808
7	q4974_p1_1	MND1_ARATH	120	TEALTQLK	8	MND1_ARATH	31	ELEK	4	33.83	138.06808
8	q5157_p1_1	MND1_ARATH	6	GLSLEEK	7	HOP2_ARATH	39	KTAVQK	6	20.99	138.06808
9	q5158_p1_1	MND1_ARATH	6	GLSLEEK	7	HOP2_ARATH	39	KTAVQK	6	28.67	138.06808
10	q5284_p1_1	MND1_ARATH	6	GLSLEEK	7	HOP2_ARATH	221	KRPR	1	24.18	138.06808
11	q5371_p1_1	MND1_ARATH	6	GLSLEEK	7	HOP2_ARATH	189	DVKELK	6	34.83	138.06808
12	q6018_p1_1	MND1_ARATH	6	GLSLEEK	7	HOP2_ARATH	39	KTAVQK	6	33.64	138.06808

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



The xiVIEW CSV format is a simple tabular format, which you can load in Excel or LibreOffice. Modifications like oxidation are encoded as a delta mass within the peptide sequence. The CSV file can be edited before uploading to xiVIEW, for example if you want to remove a putative crosslink because fails a quality check.

xiVIEW can also read crosslink data in mzIdentML 1.2 format, which we will support in a future release of Mascot.

More examples on our website

- **Reference documentation:**
<http://www.matrixscience.com/help/crosslink.html>
- **February 2020 blog article “Disulfide bond characterisation”**
- **Heterofunctional linker example:**
http://www.matrixscience.com/help/xl_examples.html

MASCOT : *Identifying intact crosslinks*

© 2020 Matrix Science



If you'd like to see more examples of crosslinked searches, we have some on our website.