

# There be Dragons: the dangers associated with assigning internal fragment ions of proteins

David PA Kilgour<sup>1</sup>, Harry Taylor<sup>1</sup>, Yury Tsybin<sup>2</sup>, David Clarke<sup>3</sup>, Logan Mackay<sup>3</sup>, Luca Fornelli<sup>4</sup>

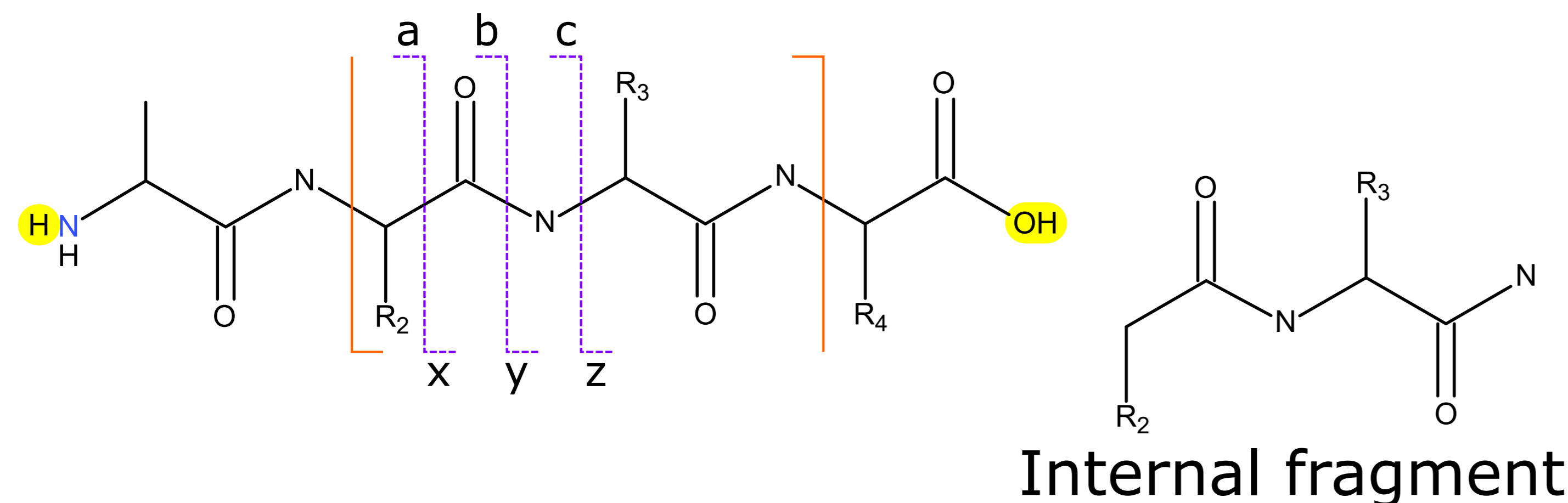
1 Nottingham Trent University; 2 Spectroswiss Sàrl; 3 University of Edinburgh; 4 University of Oklahoma

## Introduction

There is a great deal of current interest in the potential utility of internal fragments generated during top-down tandem mass spectrometry of intact proteins. But, there are a wide variety of risks associated with attempting to make use of these internal fragment assignments for samples that are either potentially impure or are not already very well characterised, and, as there are so many more possible internal fragments, the value in the assignments appears to be orders of magnitude lower than for terminal fragments. Finally, there are differences between the various nomenclatures used to describe internal fragment ions, and this lack of consistency makes it harder for us, as a community, to discuss these issues as this potentially exciting area develops.

## Internal fragments

Conventionally, terminal fragments are used in top-down protein analyses. Internal fragments are those created by breaking the peptide backbone in at least two places, to produce a fragment from the interior of the sequence.



## A numbers game

The main cause of the difficulties involved in making use of internal fragments is down to huge numbers of potential fragments. Take Bovine Carbonic Anhydrase as an example. For this 263 residue protein, there are 525 possible terminal fragments, but over 33k possible internal fragments.

As internal fragments are presumed to be caused by a terminal fragment being fragmented a second time, we assume that the ends of the internal fragment can mirror the possibilities for the terminal fragments in the same data. The total variety of internal fragments is the product of the numbers of N and C terminal fragment types. Therefore, if your fragmentation method produces only one major type of N and C terminal fragments, then you have only one type of internal fragment to find. But, if you have 3x N terminal fragment types and 3x C terminal fragment types then you have 9 times as many possible internal fragment classes to search - which may each appear in many different charge states.

The libraries become so dense that spurious hits are almost inevitable. It may require mass accuracy levels that are beyond the routine performance levels of current technology to permit their routine use.

This can be easily demonstrated by assigning the same spectrum against internal fragment libraries that match the true sequence and against libraries built from false sequences. Commonly we see very similar internal fragment assignment statistics against false or scrambled sequences - which greatly reduces confidence in assignments. Even for FT-ICR MS data.

True sequence	Scrambled Sequence
1 G A I M A S I H H W G L Y G K H I N G P I E H W H 244	1 G A I M A S I H H W G Y E P L L G P I S G L L I A A 244
21 L K D F L P I T A N G E I R Q S P I V D I D T K I A 224	21 N F F V V V E G H Y E L L E K K N P D L L L H 224
41 V V Q D P I A L L K P L A L V I Y G E A T S R 204	41 F G S S E D V D D P I V Q S D I K P R P A Q 204
61 R T M V I N N G H S F N V E Y D D I S Q D I K A T 184	61 N W W A F K Y A I T F L V T T V P Q L V I T 184
81 V L L K D G I P L L T G T Y R L V Q F H F H I W 164	81 A I I K V N I D L Q P N G K N I T W K L D E I 164
101 G S S D D Q I G S E L H T V I D R K I K I Y A A E I 144	101 T R G G G P S Q P L P N Y W L L K I T F A 144
121 L H L V H W N T K I Y G D F I G T A A Q Q P I 124	121 G D I P R D Y G W P I V T Q A T V I R P H E K 124
141 D G L A V V G V F L K I V G D I A N P I A L Q I 104	141 R Q H L L G Y S H I A L P S I A T E A S L V G I 104
161 K I V L D I A L D S I K I T K G K S T I D I F P N 84	161 N H G M W K R L L N V Q N K L H K D H L L I 84
181 F D I P G S L I L P I N V L D Y W T Y P G S L I 64	181 S S V G V F G I T G K I D N Y N E I E D I F E I 64
201 T T P P L L L E S I V T W I V L K E P I S V I 44	201 K A A I I F P D D L S K A V S Q L V I D I Q V I 44
221 S I S Q Q M L K F R T L I N F N I A T E G E P E 24	221 R A D F M I Q D L L S L M L I T D K I L S T V I A 24
241 L L M L A N W I R P A Q P L I K N I R Q V I R G I 4	241 V P P L I T G V I I L T S G I A K I N I R Q V I R G I 4
261 F P K I 1	261 F P K I 1

Showing internal fragments assigned for ECD fragments of Bovine carbonic anhydrase (22+ charge state) against the true sequence (left) and a randomly scrambled sequence (right) - where the 10 terminal residues are left in their normal sequence.

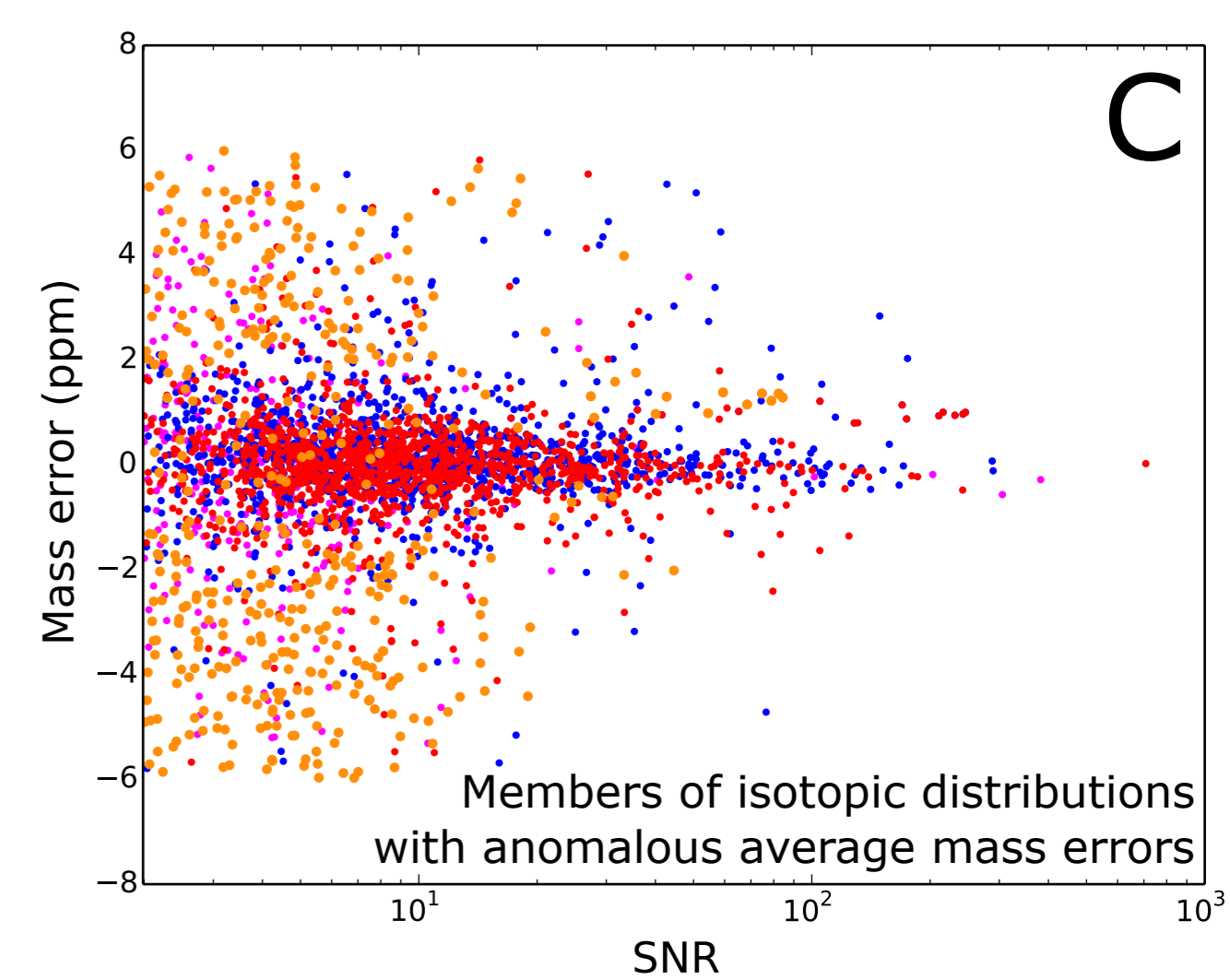
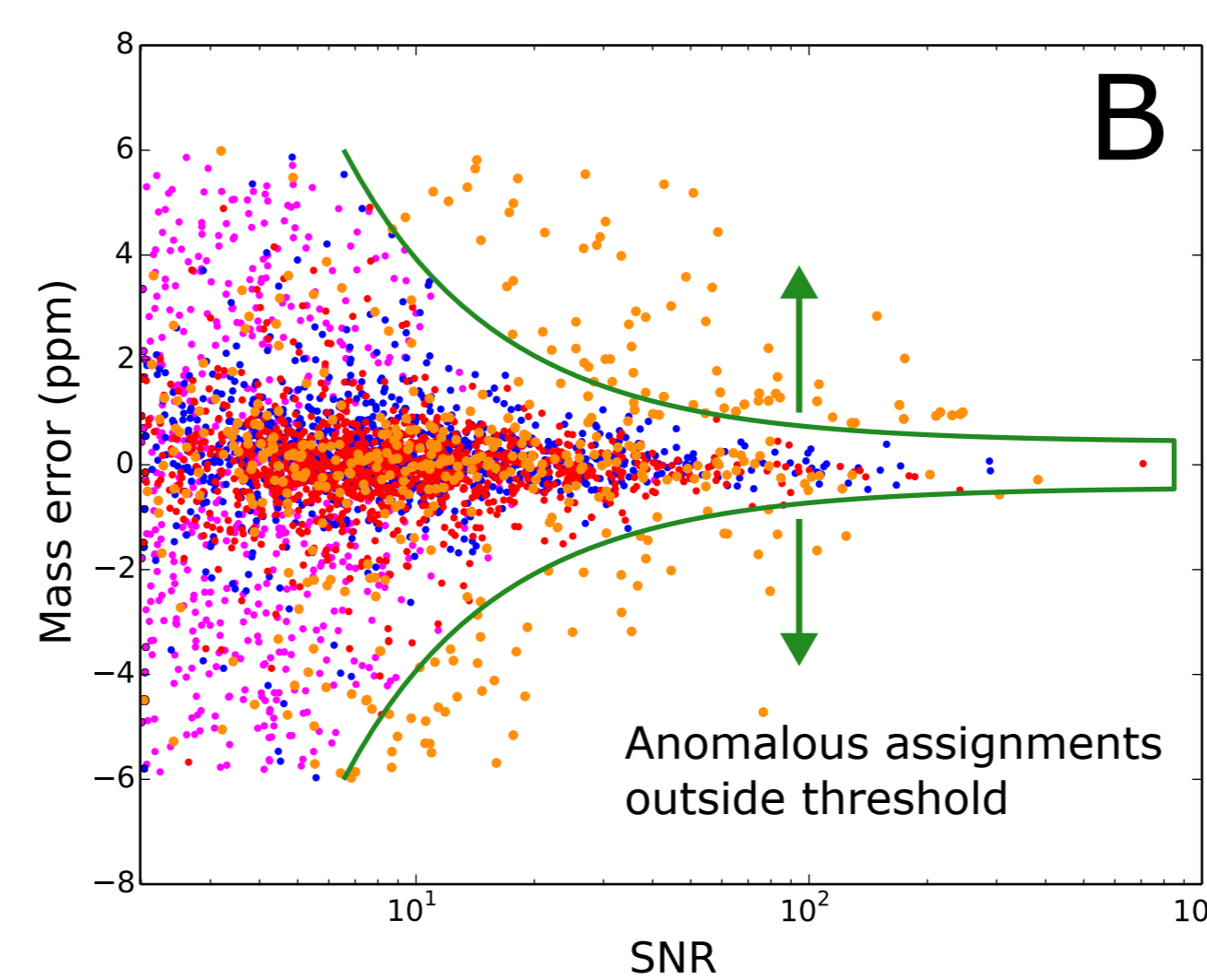
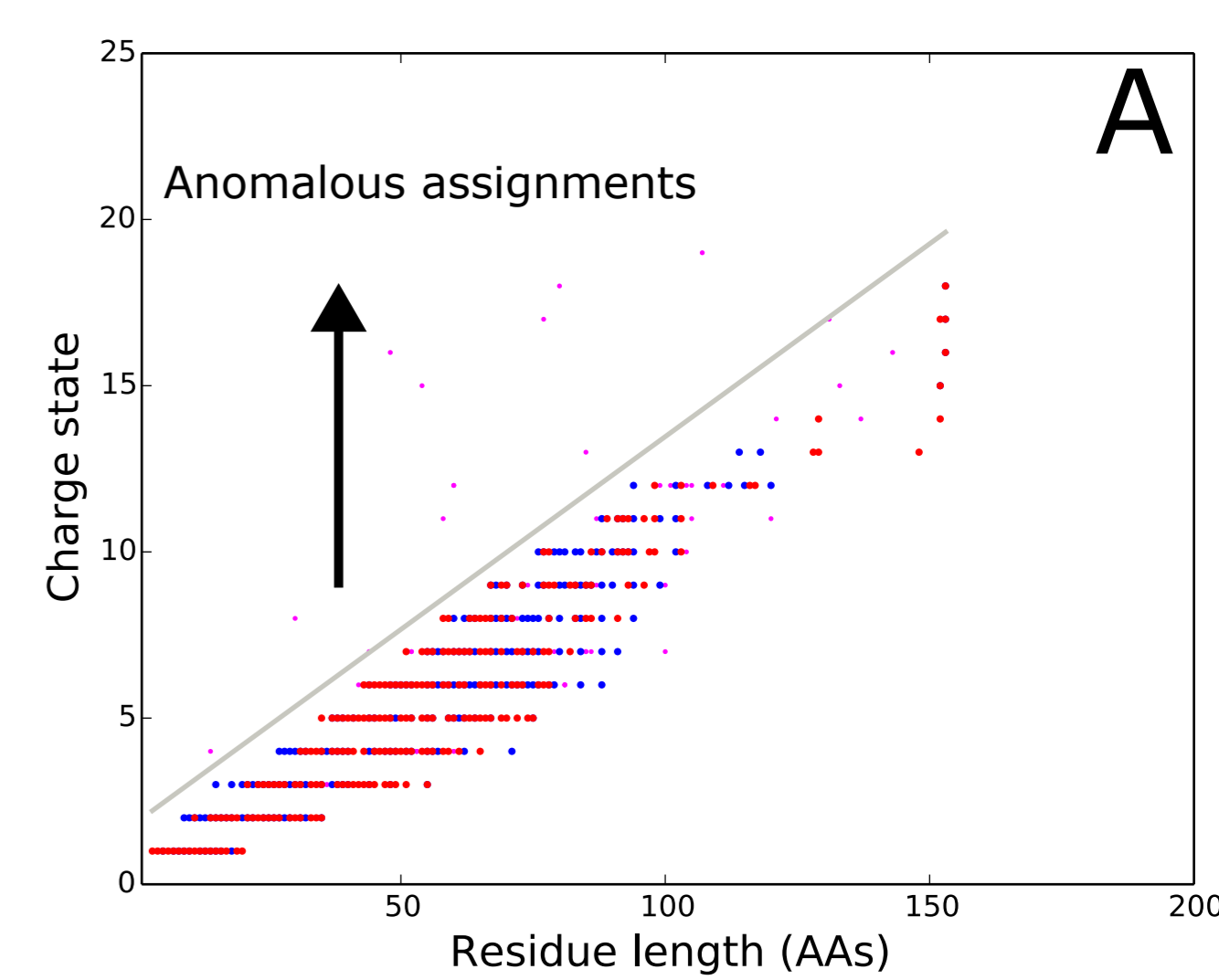
## Improving confidence

There are several approaches that can be used to improve confidence in the assignment of fragments, including internal fragments.

- Basic methods
  - no charge state deconvolution - mass accuracy is paramount
  - ppm error cut-off - tighter the better & 5ppm is not good enough
  - isotopic fit
- More advanced filters
  - charges per amino acid (length of fragment)
  - statistical tests - S/N vs mass error ("Prosaic") and average isotopic mass error ("Cookie Cutter")

All of these methods have been implemented in the AutoVectis Suite top-down assignment tool "AutoSeequer".

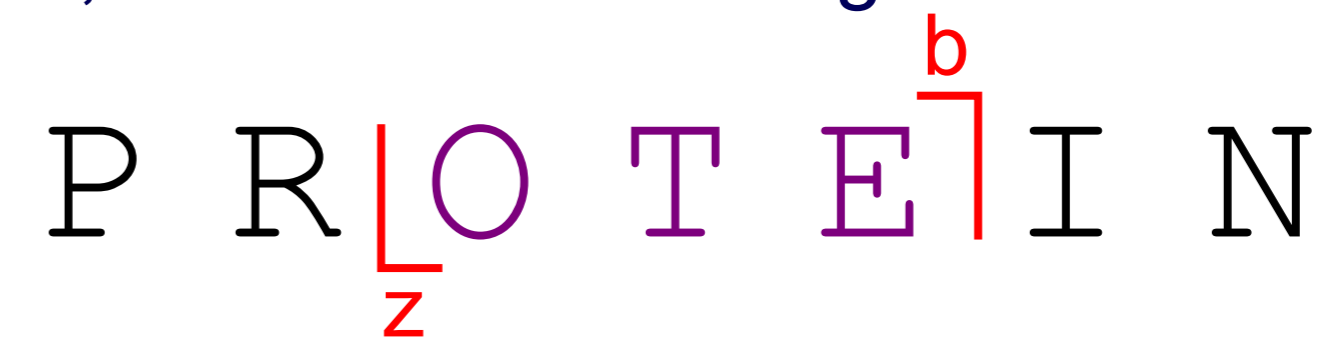
Demonstrating filtering of assignments, including internal fragments, based on looking for anomalous combinations of fragment charge state versus length [A]; S/N versus mass error [B]; and average mass error across an isotopic distribution [C] provides useful increased confidence. But, while these are very successful for terminal fragments, the confidence of internal fragments remains problematic.



## Nomenclature

There is also an issue with the lack of standardization for the nomenclature used to define internal fragments of proteins. All tools we have reviewed agree on the numbering of internal fragments. However, there is less consistency when it comes to the classification - which affects the ionic formulae. Some software, e.g. ClipsMS<sup>(1)</sup>, defines the fragments according to what has been lost from each end of the internal fragment whereas other software, including AutoVectis, defines the fragments by what is present in the internal fragment. This adds an additional level of complexity when comparing results produced by different processing pipelines and is something that the community may wish to address.

Consider this example, and the internal fragment defined.



Is this a b&z fragment (describing the ends of the fragment detected or c&y type fragment (describing the lost terminal fragments)?

## Conclusions

Internal fragments offer some potential to assist in top-down protein characterization - but there are many problems to overcome before they may become robustly applicable.

## Acknowledgments

This research project has received funding from the European Horizon 2020 research and innovation program under grant agreement No 829157.

## References

1. Lantz C et al. ClipsMS: An Algorithm for Analyzing Internal Fragments JPR. 2021 Mar 2;20(4):1928-35.

