# Evaluating a new Workflow for Top-Down Protein Sequence Confirmation and *de novo* Sequencing

**BRUKER**

Mariangela Kosmopoulou[1]; George Alevizos[1]; Georgia Orfanoudaki[1]; Athanasios Smyrnakis[1]; Michael Greig[2]; Detlev Suckau[3], [1]Fasmatech, Athens, Greece; [2]Bruker Scientific, Billerica, MA; [3]Bruker Daltonics, Bremen, Germany

## Introduction

- **Top-Down MS sequencing** efficiently characterizes proteins N- and C-terminal sequences and assists with PTM localization.

- The complexity of high-resolution MS data limited ESI-Top-Down (TD) analyses to terminal sequence confirmation.

- Thus, the *de novo determination of protein sequences* requires new data analysis software to reliably extract and timely process such high-density information.

- Here, we describe a new software – OmniScape™ - with workflows for protein *de novo sequencing* and sequence confirmation using complex ESI and MALDI TD MS/MS data Fig 1.

- Maximized sequence coverage can be acquired from multiple experiments, while the identification of the best matching proteoforms is another crucial feature.



Fig. 1 Scheme of the workflow for seamless proteoform discovery and detailed characterization

## Methods

Top-Down ESI-QTOF-MS spectra were acquired after ETD, ECD and EID fragmentation of protein precursor ions.

A new analysis software - OmniScape - was developed 1) for the confirmation and spectral annotation of known protein lead sequences, and 2) for the *de novo* sequencing and homology searching of unknown proteins. First, the core algorithm - OmniWave™ - exhaustively extracts possible fragment ion masses, including all charge states and isotopes. In a second step, it determines all possible monoisotopic masses of fragment ions and matches all plausible sequences.
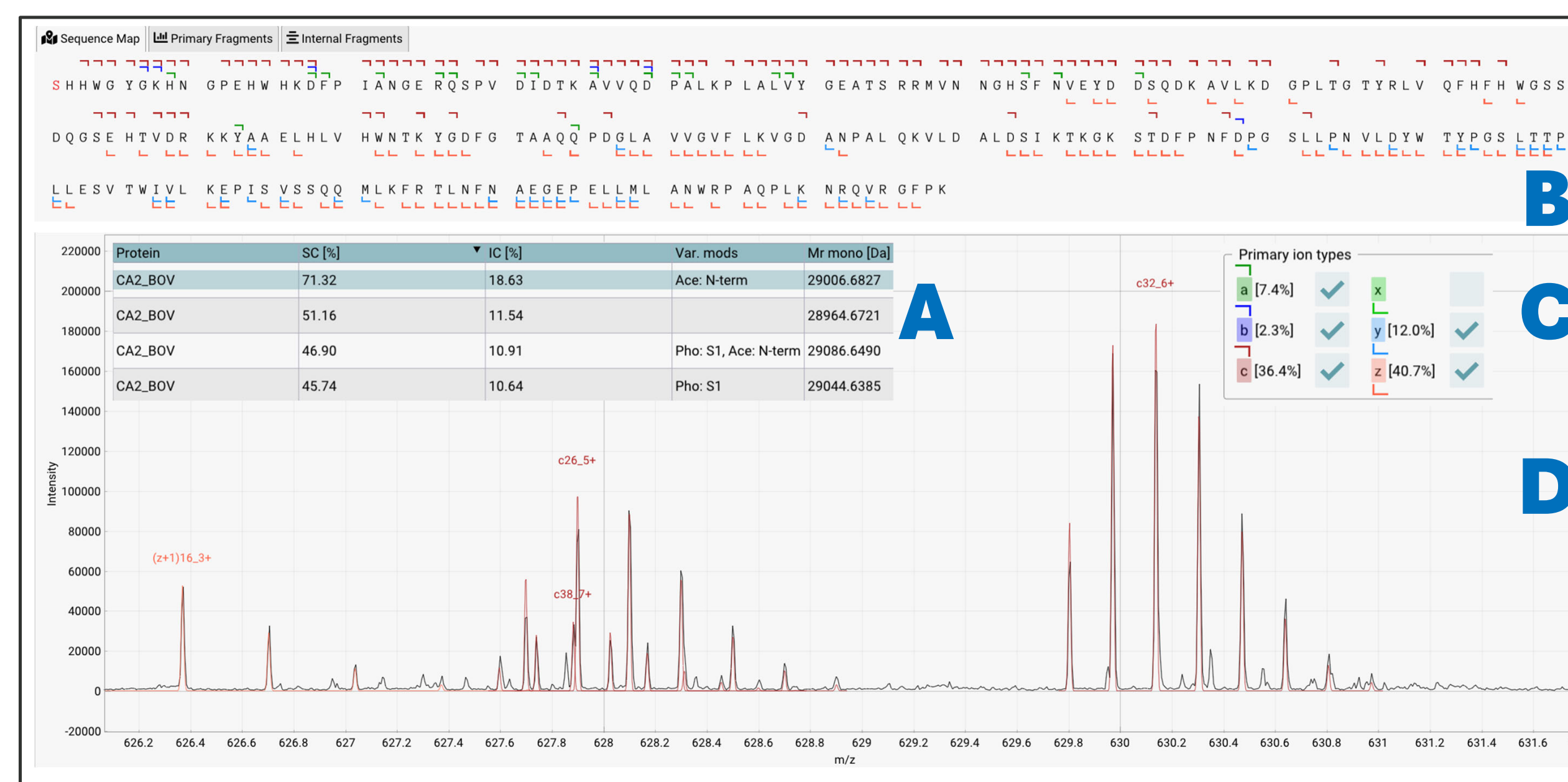
## Results



Fig. 2 User interface essentials: A) scored list of proteoforms, highest Sequence Coverage (SC) for N-acetyl CAH2, B) Sequence Map with matching fragments which are described in legend C), D) zoomed in view of the ETD spectrum (maXis II ETD) with annotated fragment isotope patterns.
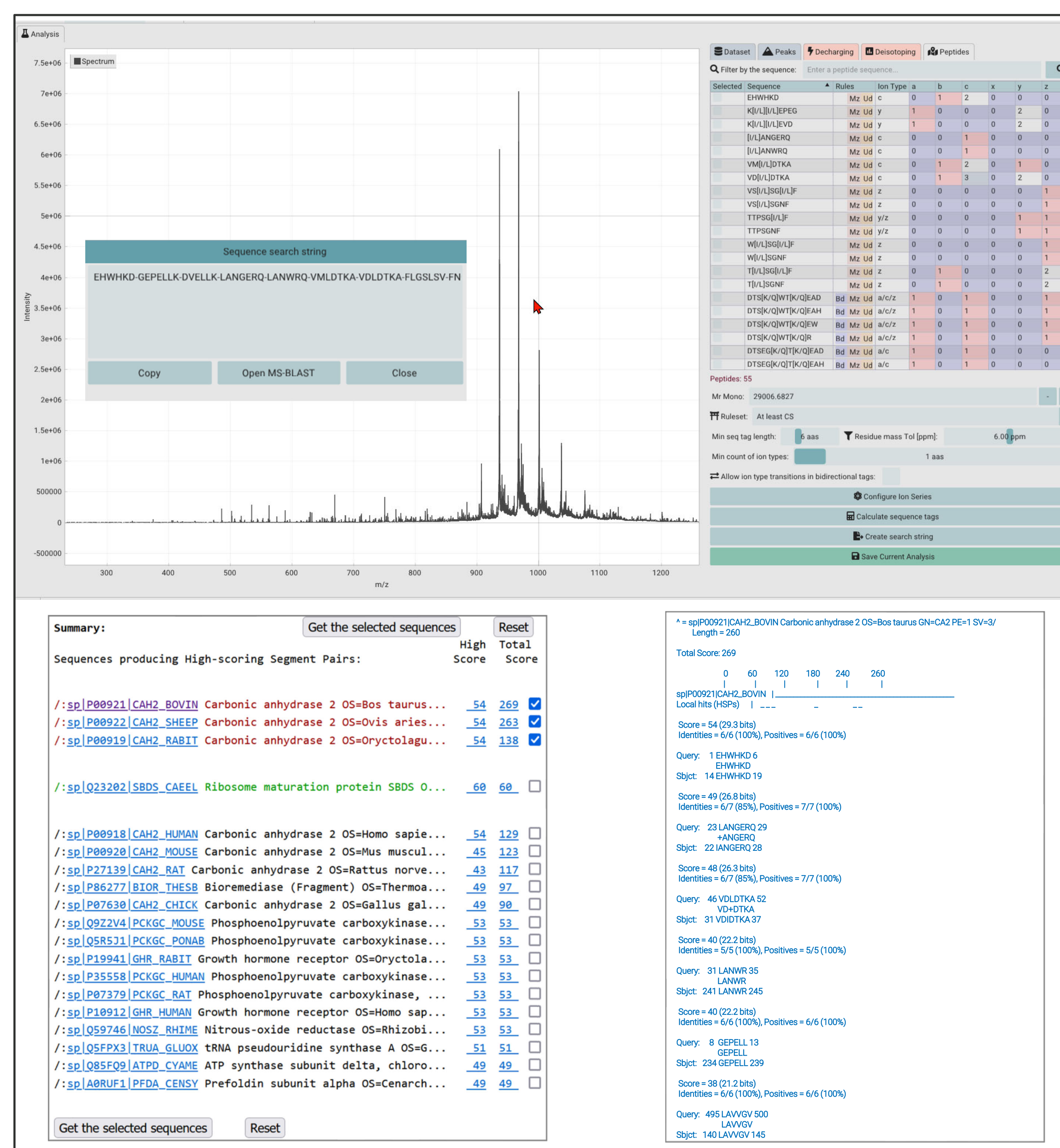


Fig. 3 *De novo* sequence tag generation (top) and MS-BLAST (Harvard, Sunaev lab) homology search result with sequence tags matching to the CAH2 sequence (bottom).

*De novo* sequence analysis results in several sequence tags that can be used for homology searches in protein sequence databases such as SwissProt. CAH2 was identified from an ETD spectrum in Fig 3.
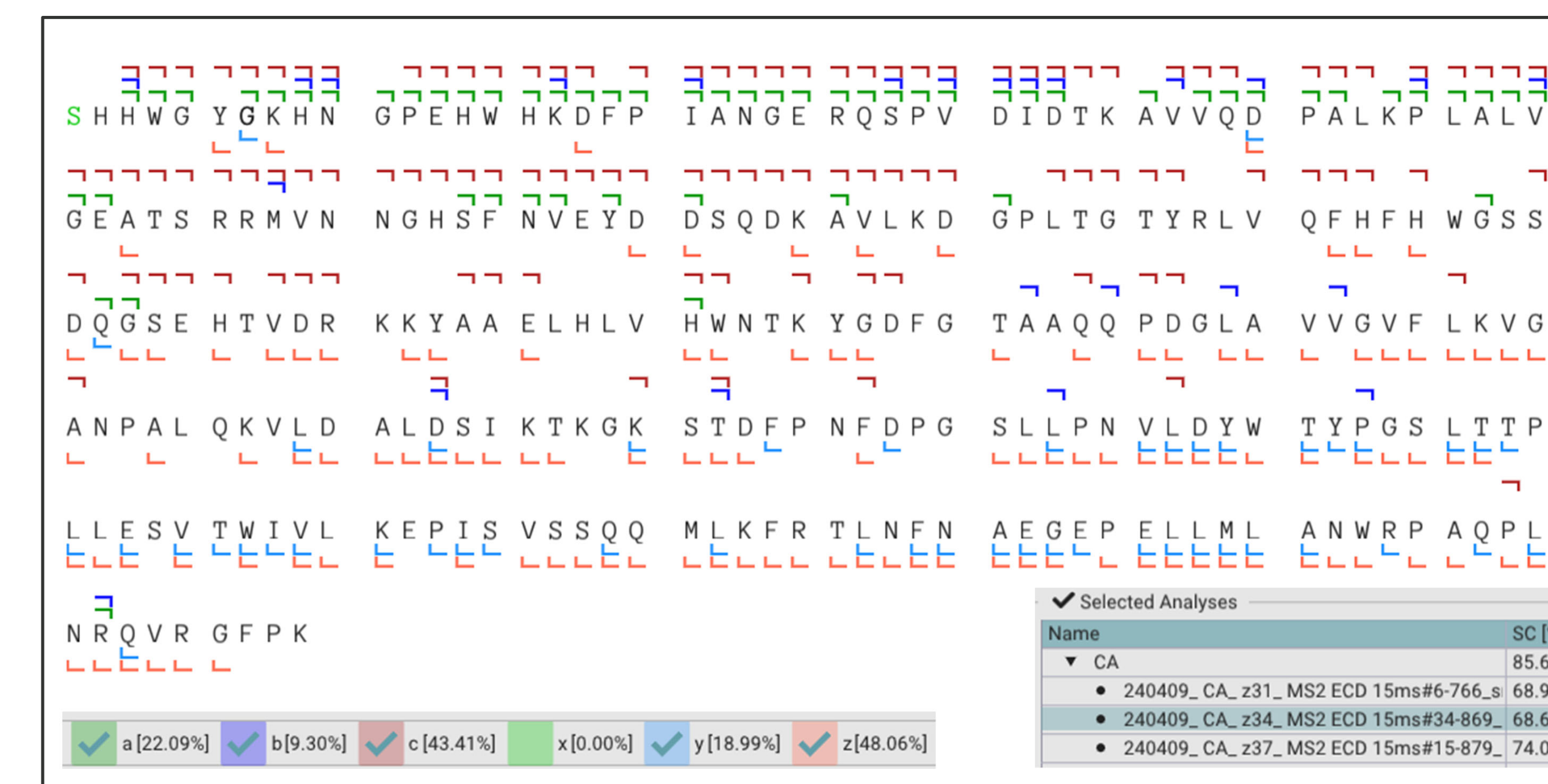


Fig. 4 Sequence map of CAH2 resulting from the combined maps obtained by ECD analysis of the 31+, 34+ and 37+ precursor ions (timsTOF Omnitrap prototype)

The match between experimental and calculated isotope patterns in the Top-Down ETD spectrum of bovine Carbonic Anhydrase 2 (CAH2) was used for interactive peak list curation. Of the 4 possible proteoforms of CA2 the N-terminally acetylated one scored highest – in agreement with the intact mass – yielding a sequence coverage of 71.3% Fig 2.

In OmniScape, multiple analysis results can be combined into a single sequence map, which increased the sequence coverage of CAH2 from 69% of the individual charge states to 86% for the combined map Fig 4.
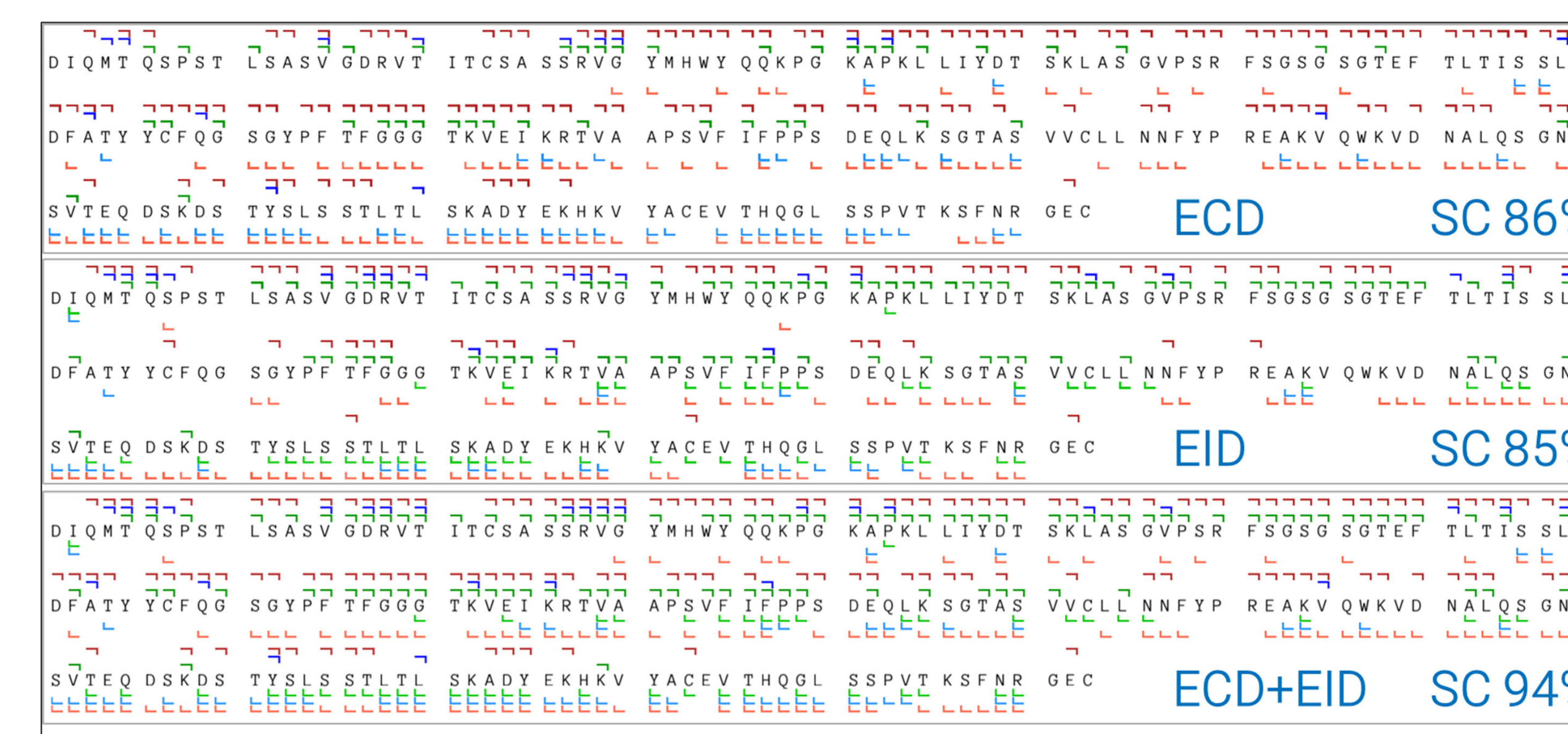


Fig. 5 NISTmAb LC subunit analysis with ECD and EID on a timsTOF Omnitrap prototype. Combination of both sequence maps into a combined map increased the sequence coverage (SC) from ~85% each to 94%.

In an analysis of NISTmAb Light Chain different fragment-ation methods were combined. Due to the complementarity of EID and ECD fragmentation the sequence coverage could be improved to 94% for the combined map Fig 5.

## Summary

The OmniWave algorithm embedded in OmniScape is responsible for the determination of isotope clusters and charge state assignment. In addition, it determines amino acid distances between isotope patterns in the neutral domain to establish sequence tags for the *de novo* sequence determination of a lead sequence Fig 4.

Such a lead sequence can be further analyzed in OmniScape. Proteoforms with various modifications can be scored for the presence of individual structures Fig 2. Even multiple datasets obtained under different conditions or using different fragmentation methods can be combined Fig 5.

Together with the intact neutral mass of a protein of interest, these analyses allow to reliably determine the sequence and its modifications if high resolution/high mass accuracy data are used (here, all analyses used an MS/MS tolerance < 6 ppm) and high-fidelity isotopic clusters intensity is measured by the selected instrument.

The fragment ion representation in OmniScape provides a high information density, enabling the representation of multiple experiment results on limited space Fig 5.

Although high resolution data from ESI-QTOF instruments were discussed here, axial MALDI-ISD data can be reliably processed in OmniScape, as well.

## Conclusion

- OmniScape provides full analysis support from Top-Down protein identification through homology searches to proteoform elucidation using a lead sequence
- Multiple datasets across instrument platforms or parameter variation can be combined into a higher sequence coverage representation of the protein
- Powerful interactive data curation in OmniScape safeguards meaningful sequence assessment

**Top-Down Protein Sequencing**

*The authors declare no competing financial interest*