

# DiNovo enables high-coverage and high-confidence de novo peptide sequencing via mirror proteases and deep learning

Zixuan Cao<sup>1</sup>, Xueli Peng<sup>1</sup>, Di Zhang<sup>2</sup>, Piyu Zhou<sup>1</sup>, Li Kang<sup>3</sup>, Hao Chi<sup>4</sup>, Ruitao Wu<sup>2</sup>, Zhiyuan Cheng<sup>1</sup>, Yao Zhang<sup>3</sup>, Jiaying Dai<sup>3</sup>, Yanchang Li<sup>3</sup>, Lijin Yao<sup>2</sup>, Xinming Li<sup>2</sup>, Yaoyu He<sup>1</sup>, Jinghan Yang<sup>1</sup>, Haipeng Wang<sup>2</sup>, Ping Xu<sup>3</sup>, Yan Fu<sup>1\*</sup>

<sup>1</sup> Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing, China

<sup>2</sup> Shandong University of Technology, Zibo, China

<sup>3</sup> Beijing Institute of Lifeomics, Beijing, China

<sup>4</sup> Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

\*Presenting and corresponding author

Despite the recent advancements driven by deep learning, de novo peptide sequencing is still constrained by incomplete peptide fragmentation and insufficient protein digestion in current single protease-based proteomic experiments. Furthermore, the evaluation of de novo sequencing algorithms heavily depends on constructing restricted benchmark datasets from peptide-spectrum matches obtained by database search. We present a software system, named DiNovo, for high-coverage and high-confidence de novo peptide sequencing by leveraging the complementarity of mirror proteases. DiNovo is empowered by several innovative algorithms, including a mirror-spectra recognition algorithm independent of pre-sequencing, two sequencing algorithms based on deep learning and graph theory, respectively, and target-decoy mapping, a method for sequencing result evaluation free of prior peptide identification. We evaluate the performance of DiNovo using two pairs of mirror proteases including trypsin/LysargiNase, and Lys-C/Lys-N, which were used to digest proteins from *E. coli* and yeast proteomes. By taking full advantage of the complementarity of mirror spectra, DiNovo achieves much higher sequence coverage and confidence than state-of-the-art single-protease de novo sequencing algorithms. Compared with the trypsin protease used alone, DiNovo using two pairs of mirror proteases leads to two to three times high-confidence amino acids sequenced. Furthermore, DiNovo identifies a comparable number of proteins to database search at the same false discovery rate level, showing great potential as a powerful complement or even an alternative to database search for protein identification. DiNovo is designed to dramatically increase the coverage and confidence of de novo peptide sequencing and to facilitate mirror-protease proteomics.