

# Comprehensive Evaluation of CAZyme Prediction Tools in Fungal and Bacterial Species

## Introduction

Carbohydrate Active enZymes (CAZymes) are pivotal in pathogen recognition, signalling, structure and energy metabolism. CAZY is the most comprehensive CAZyme database, cataloguing CAZymes into sequence-based CAZY families [1]. The CAZyme prediction tools **dbCAN** [2], **CUPP** [3] and **eCAMI** [4] annotate CAZymes with CAZY families. However, these tools have not been independently evaluated on a common high-quality dataset. Additionally, previous evaluations did not evaluate the **binary classification** of CAZymes/non-CAZymes, and the **multilabel classification** of CAZymes to multiple CAZY families.

## Method

The bioinformatic pipeline **pyrewton** was developed for this independent evaluation (Fig.1).

**GitHub:** <https://github.com/HobnobMancer/pyrewton>

The ground truths were retrieved using **cazy\_webscraper**.

**GitHub:** [https://github.com/HobnobMancer/cazy\\_webscraper](https://github.com/HobnobMancer/cazy_webscraper)

70 genomic assemblies:

40 Bacteria  
16 Fungal  
11 Yeast  
3 Eukaryotes

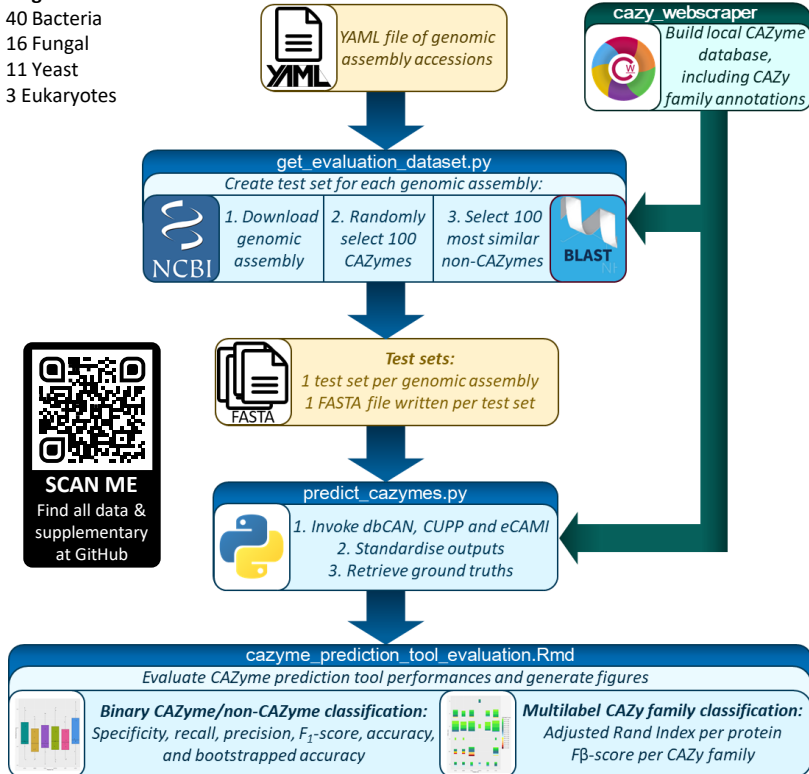


Fig.1 Schematic of the bioinformatic pipeline **pyrewton** for evaluating CAZyme prediction tools

The University of St Andrews is a charity registered in Scotland, No: SC013532

## Results

### Binary CAZyme/non-CAZymes classification evaluation

dbCAN invokes the function prediction tools HMMER, Hotpep and DIAMOND. All prediction tools showed a low probability of misidentifying non-CAZymes as CAZymes, but also showed a tendency to miss identify a small proportion of CAZymes as non-CAZymes (Fig.2).

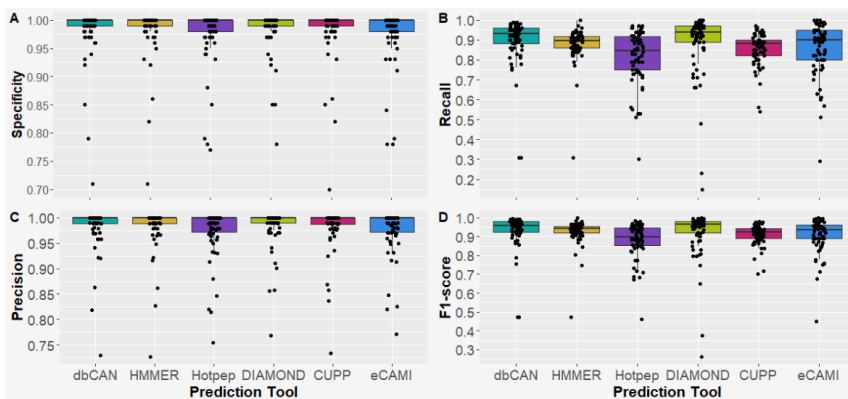


Fig.2 Evaluation of CAZyme/non-CAZyme differentiation performance. One-dimensional scatterplots overlaying boxplots for [A] specificity, [B] recall, [C] precision and [D] F1-score.

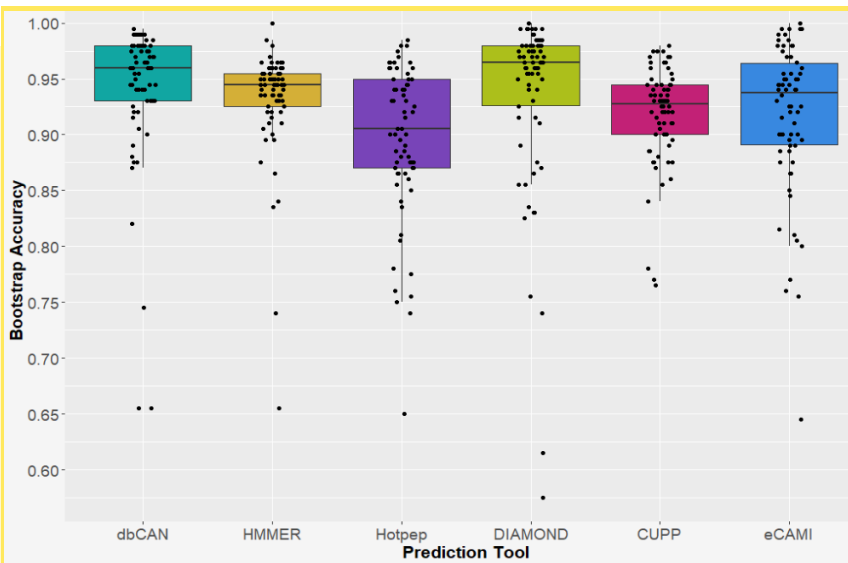


Fig.3 Expected range of performance of CAZyme prediction tools **dbCAN**, **CUPP** and **eCAMI**. Bootstrapping the CAZyme/non-CAZyme prediction accuracy was performed 10,000 times per test set. The median bootstrapped accuracy of each test set was plotted.

### Multilabel CAZY family classification evaluation

Multilabel classification arises from the ability of a CAZyme to be assigned multiple CAZY families. The Adjusted Rand Index (ARI) was calculated per protein (Fig.4[A]) and the F $\beta$ -score ( $\beta=1$ ) calculated for each CAZY family, true negative non-CAZyme predictions were excluded (Fig.4[B]).

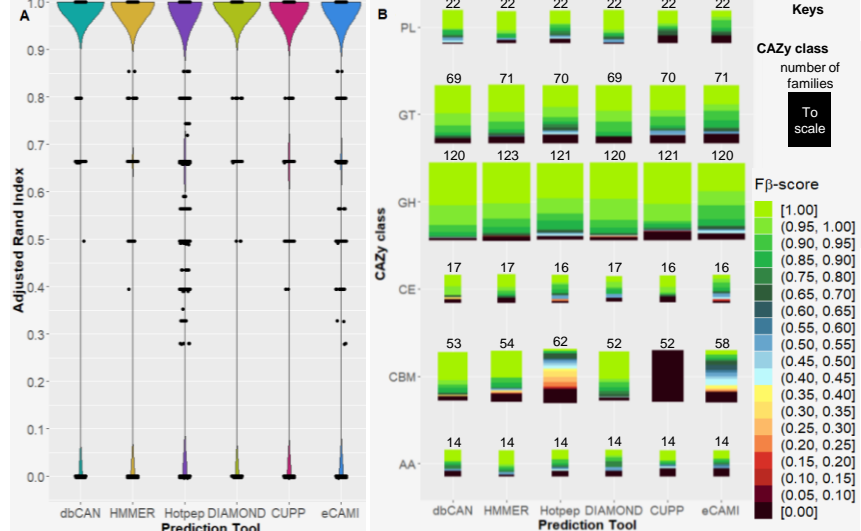


Fig.4 Evaluation of CAZY family multilabel classification

[A] Adjusted Rand Index per protein sequence. [B] Proportional area plot of CAZY classes sized by the number of families analysed, and coloured by the proportion of CAZY family F $\beta$ -scores within each range of the scale, ( $\beta=1$ ).

## Conclusions

- Created a bioinformatic pipeline for the reproducible evaluation of CAZyme predictions tools, and benchmarked dbCAN, CUPP and eCAMI against a high quality test set
- Evaluated the binary and multilabel classification of CAZymes for the first time
- Statistically evaluated the expected range of performance for the first time
- dbCAN was best overall but the weakest was Hotpep, which is incorporated into dbCAN
- Best performance may be achieved by replacing Hotpep with CUPP and/or eCAMI
- Next steps are to expand the dataset and evaluate substituting Hotpep with CUPP and/or eCAMI

## References

- Lombard, V. et al. (2014) 'The carbohydrate-active enzymes database (CAZY) in 2013', *Nucleic Acids Research*, 42, pp.D490-D495
- Zhang et al. (2018) 'dbCAN2: a meta server for automated carbohydrate-active enzyme annotation', *Nucleic Acids Research*, 46, W1, pp. W95-W101
- Barrett, K., Lange, L. (2019) 'Peptide-based functional annotation of carbohydrate active enzymes by conserved unique peptide patterns (CUPP)', *Biotechnology for biofuels*, 12, 102
- Xu et al. (2020) 'eCAMI: simultaneous classification and motif identification for enzyme annotation', *Bioinformatics*, 36, 7, pp.2068-2075

## Acknowledgements

We would like to thank the EASTBIO Doctoral Training Partnership (BBSRC) for funding our research.