

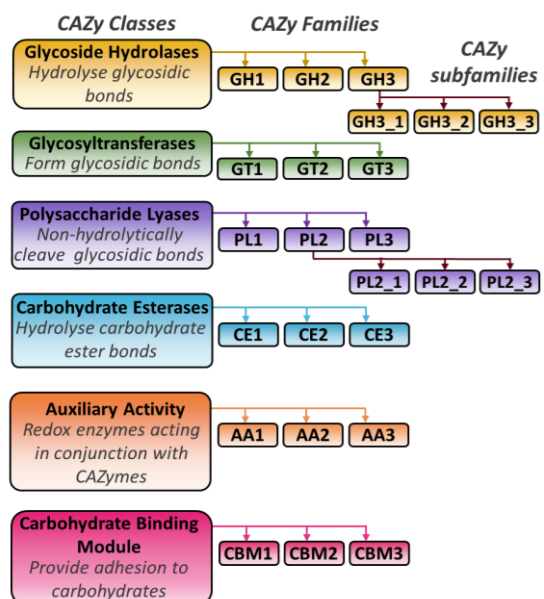
# cazy\_webscraper

## For creating a local CAZy database

### Introduction

Carbohydrate Active enZymes (CAZymes) are pivotal in pathogen recognition, signalling, structure and energy metabolism. CAZY ([www.cazy.org](http://www.cazy.org)) is the most comprehensive CAZyme database [1], but it does not provide methods for automating data retrieval or submitting sequences for annotation.

*cazy\_webscraper* retrieves user-specified datasets from CAZY, producing a local SQL database enabling thorough interrogation of the data. *cazy\_webscraper* can also retrieve protein sequences from GenBank [2] and download structure files from RCSB PDB [3].



**Fig.1 CAZY database structure**  
 CAZY catalogues proteins into classes that are divided into families, some of which are divided into subfamilies.

### Method

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

**Scraping** is invoked using the command `python3 cazy_webscraper`. All optional flags can be found in the GitHub repository README.

**Expanding** the dataset beyond CAZY is achieved using the `expand` module.

### 1. GenBank

Each unique CAZyme is identified by its **primary** GenBank accession, consolidating duplicate CAZY entries in the local database.

Retrieve all CAZY family annotations for a given protein by querying the local CAZyme database by its GenBank accession.

*cazy\_webscraper* automates retrieving **protein sequences from GenBank**.

*cazy\_webscraper* can update sequences in the local CAZyme database if a newer sequence is available in NCB, **keeping the dataset up to date**.

### 2. CAZY Families

*cazy\_webscraper* automates and quickly scrapes CAZY. Scraping CAZY family GH1, containing **43,649 proteins**, takes **44 minutes**, instead of users manually reading **44 webpages**.

Eukaryota					
Protein Name	EC#	Organism	GenBank	UniProt	PDB/3D
unknown (W501228_P13) (fragment)		<i>Proteus trichosarca</i>	ARK95221.1	ASPT53	
lactase, partial (LcT)		<i>Pristionchus tritrichosarca</i>	AUD47938.1		
lactase, partial (LcT)		<i>Proteus trichosarca</i>	AUD47906.1		
prunasin hydrolase		<i>Prunus artemisiaca</i>	AHE74128.1		
prunasin hydrolase		<i>Prunus artemisiaca</i>	AHE74133.1		
prunasin hydrolase		<i>Prunus artemisiaca</i>	AHE74135.1		
prunasin hydrolase		<i>Prunus artemisiaca</i>	AHE74129.1		
β-glucosidase (Pa BG)	3.2.1.21	<i>Prunus avium</i>	AAA91166.1	Q83014	
Prudu_015238		<i>Prunus dulcis</i>	BBH04164.1		
Prudu_016574		<i>Prunus dulcis</i>	BBH05239.1		
Prudu_014650		<i>Prunus dulcis</i>	BBH03706.1		
Prudu_016581 (fragment)		<i>Prunus dulcis</i>	BBH05246.1		

**Fig.2 CAZY database structure**  
 An HTML table users had to previously parse manually to retrieve data from CAZY

Unlike previous scrapers [4], *cazy\_webscraper* can retrieve data for **specific CAZY classes and (sub)families**, reducing waiting times from **hours to minutes**.

### 3. EC Numbers

Use *cazy\_webscraper* to collate quickly CAZymes having similar activity by scraping by EC number or querying the local CAZyme database.

### 4. Taxonomy

Scrape specific taxa. Apply a combination of **kingdoms, genus, species**, and /or **strain** filters. Use the taxonomy data to track the evolution of functions through **phylogenetic analysis**.

### 5. CAZomes

Automate retrieving the CAZome (all CAZymes within a genome) of species of interest from CAZY.

Or quickly retrieve CAZomes by querying the local CAZyme database.

With one command, retrieve all protein sequences of a CAZome, ready for homolog searchers.

### 6. UniProt

Expand the dataset beyond CAZY by incorporating data *via* UniProt accessions. For example, retrieve CAZyme subcellular localisation data from UniProt, to **elucidate the functions** of uncharacterised CAZymes.

### 7. RCSB PDB

Automate rapid retrieval of **all PDB** structures for the dataset of interest in CAZY using *cazy\_webscraper*.

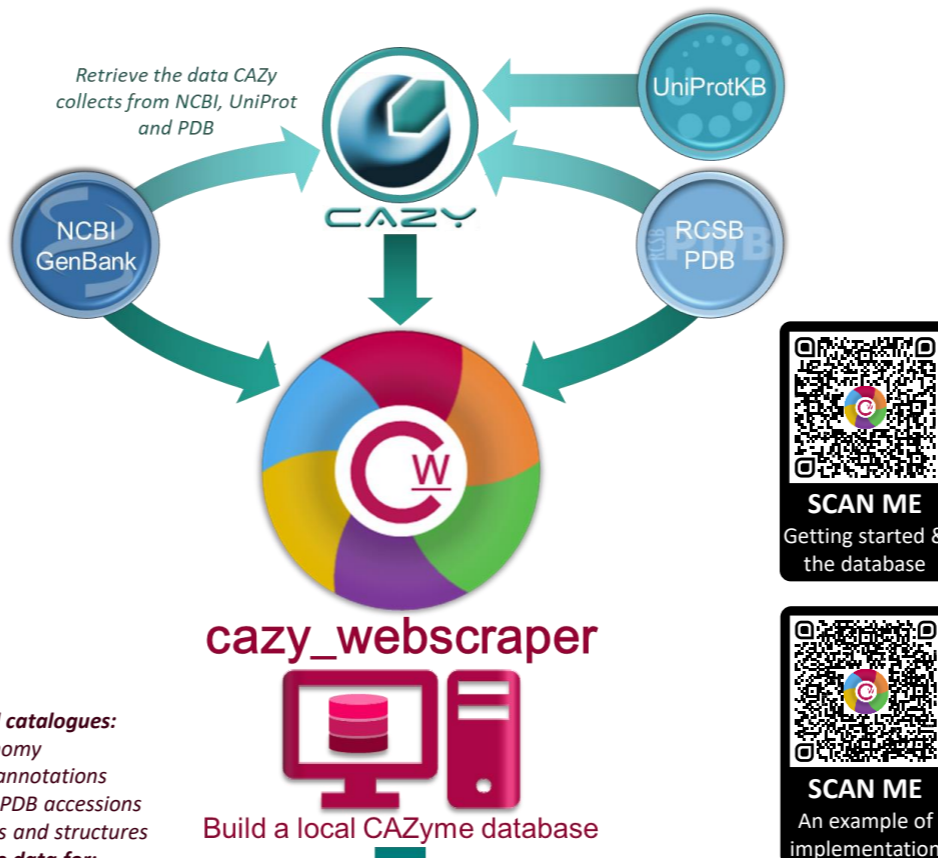
Query using a combination of taxonomy, CAZY (sub)family, CAZY class and EC number filters.

### 8. SQL Database

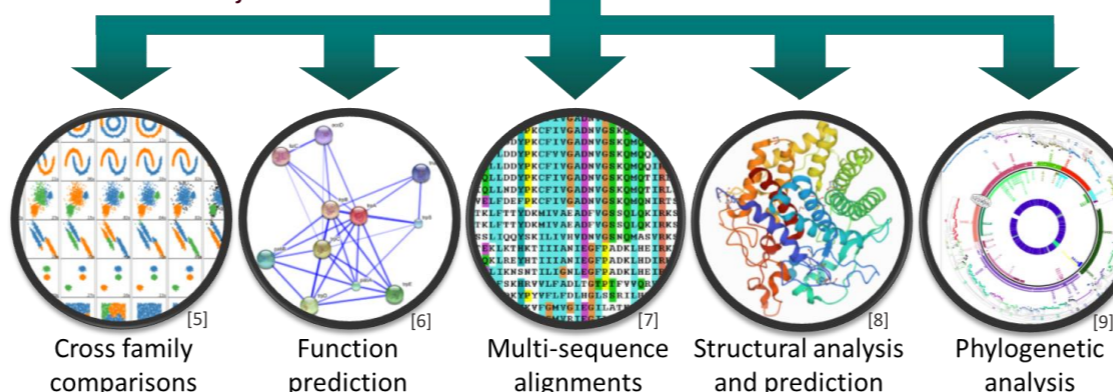
Building an SQL database instead of a plaintext [5], enables thorough interrogation of the data *via* complex queries using SQL.

Perform complex queries that cannot be performed on the CAZY website.

For example, retrieve all species with at least one CAZyme in GH1 and at least one CAZyme in PL9.



**Retrieves and catalogues:**  
 Taxonomy  
 CAZY family annotations  
 NCBI, UniProt & PDB accessions  
 Protein sequences and structures  
**Utilise these data for:**



**Fig. 2 Sources and application of data stored in the CAZyme database created by cazy\_webscraper**  
 Numbers in brackets indicate the source of the image.

### Reproducibility

Use *cazy\_webscraper* to generate **reproducible and shareable datasets**, facilitating reproduction of downstream analyses.

Optional configuration by a YAML file and generation of a log file, generates **shareable documentation** to bolster reproducibility.

### Conclusions

*cazy\_webscraper* provides new, **previously unachievable** access to the proteomic data within CAZY. This facilitates inclusion of CAZY data in functional, evolutionary, structural, genomic and metabolic studies. Thus, *cazy\_webscraper* opens up numerous new avenues of investigation.

- **Automate** retrieving CAZY annotations, protein sequences and structure files
- **Expand** the dataset beyond that stored in CAZY
- **Thoroughly** interrogate the dataset using complex queries in SQL

### References

- Lombard, V. et al. (2014) 'The carbohydrate-active enzymes database (CAZY) in 2013', *Nucleic Acids Research*, 42, pp.D490-D495
- Sayers, E. W. et al. (2020) 'GenBank', *Nucleic Acids Research*, 49(D1), pp.D92-96
- Berman, Helen M. et al. (2000) 'The Protein Data Bank', *Nucleic Acids Research*, 28(1), pp.235-242
- Honorato, R. V. (2016) 'CAZY-parser a way to extract information from the Carbohydrate-Active enZymes Database', *The Journal of Open Source Software*, 1(8), 53
- Chilamakuri, C. S. R. et al. (2011) 'Cross-genome comparisons of newly identified domains in *Mycoplasma gallisepticum* and domain architectures with other mycoplasma species', *International Journal of Genomics*, 2011, pp. 878973
- Wikipedia (2009) 'Protein function prediction', accessed 2021.03.27
- Andrade, M (2006) 'Multiple sequence alignment', *Wikipedia*, accessed 2021.03.27
- Parsiegla, G. (2002) 'Crystal structure of the cellulase Cel9M enlightens structure/function relationships of the variable catalytic modules in glycoside hydrolases', *Biochemistry*, 41(37), pp.11134-11142
- Barrett, K., Lange, L. (2019) 'Peptide-based functional annotation of carbohydrate active enzymes by conserved unique peptide patterns (CUPP)', *Biotechnology for biofuels*, 12, 102

### Acknowledgements

We thank the EASTBIO Doctoral Training Partnership for funding our work.