

In silico identification of GH family 10 homologs in metagenomes

Jenny, Matt, Shannon, William

BIOS 700

April 30th, 2013

Goals!

- 1.) Practice the bioinformatics knowledge and computing skills by applying them to a research project, which might lead to novel research findings
- 2.) Learn how to design a bioinformatics workflow to answer biology/evolution questions
- 3.) Learn how to identify useful tools, datasets and existing knowledge from the research papers to help design the workflow

Let's find out!

DID WE ACCOMPLISH THIS??

What are we looking at?

RESEARCH BACKGROUND

Glycoside hydrolases

- These enzymes hydrolyze the glycosidic bond between two (or more) carbohydrates, or between a carbohydrate and a non-carbohydrate moiety
- These enzymes were then broken down into different families based on sequence similarities

Mostly

But some are...

**This family also includes
some cellobiohydrolases
([EC:3.2.1.91](#))

Glycoside Hydrolase Family 10

Known Activities	endo-1,4- β -xylanase (EC 3.2.1.8); endo-1,3- β -xylanase (EC 3.2.1.32)
Mechanism	Retaining
Clan	GH-A
3D Structure Status	(b/a)8
Catalytic Nucleophile/Base	Glu (experimental)
Catalytic Proton Donor	Glu (experimental)
Note	formerly known as cellulase family F.
External resources	CAZypedia ; HOMSTRAD ; PRINTS ; PROSITE ;
Commercial Enzyme Provider(s)	MEGAZYME ; NZYTech ; PROZOMIX ;
Statistics	GenBank accession (1849); Uniprot accession (932); PDB accession (120); 3D entries (25); cryst (1)
Summary	All (1516) Archaea (8) Bacteria (906) Eukaryota (288) unclassified (314) Structure (25 - 1 cryst) Characterized (260)

Last update: 2013-04-16 © Copyright 1998-2013
AFMB - CNRS - Université d'Aix-Marseille

Predominant GH10 enzymes are denoted as EC 3.2.1.8, therefore classified as xylanases.

Table 1 . Acronyms for genes and encoded enzymes

Enzyme	Gene	Protein	EC designations
Cellulase	<i>cel</i>	Cel	EC 3.2.1.4; EC 3.2.1.91
Xylanase	<i>xyn</i>	Xyn	EC 3.2.1.8
Mannanase	<i>man</i>	Man	EC 3.2.1.78
Lichenase	<i>lic</i>	Lic	EC 3.2.1.73; EC 3.2.1.58
Laminarinase	<i>lam</i>	Lam	EC 3.2.1.39

Domain organisations: Glycoside hydrolase, family 10 (IPR001000)

This family has the following domain organisations:

Showing 1 - 20 of 177

1 2 3 4 5 6 7 8 9 > >>

Domain organisation	Number of proteins
■	1777
■ ■	84
■ ■ ■	77
■ ■ ■	69
■ ■	66
■ ■ ■	49
■ ■ ■	46
■ ■ ■	42
■ ■ ■	41
■ ■ ■ ■	19
■ ■ ■ ■ ■	18
■ ■ ■ ■ ■	16
■ ■ ■ ■ ■	15
■ ■ ■ ■ ■ ■	15
■ ■ ■ ■ ■ ■ ■	11
■ ■ ■ ■ ■ ■ ■ ■	11
■ ■ ■ ■ ■ ■ ■ ■ ■	10
■ ■ ■ ■ ■ ■ ■ ■ ■	7
■ ■ ■ ■ ■ ■ ■ ■ ■	7

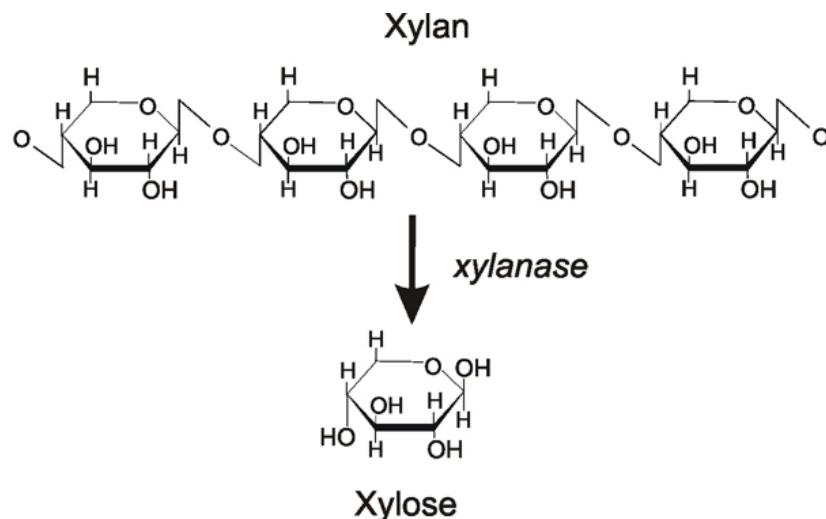
Colour key

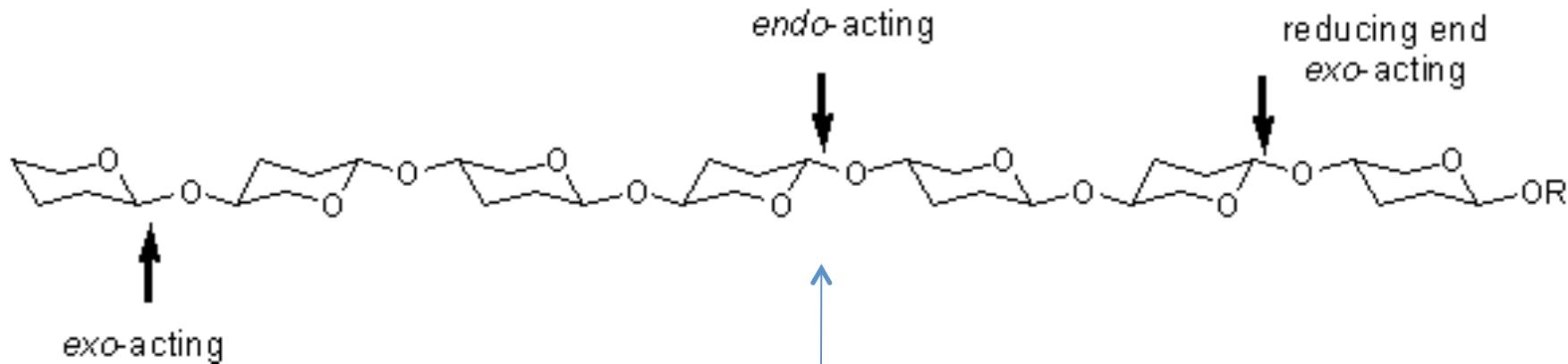
- Carbohydrate-binding, CenC-like (IPR003305)
- S-layer homology domain (IPR001119)
- Cellulosome enzyme, dockerin type I (IPR016134)
- Ricin B lectin domain (IPR000772)
- Glycoside hydrolase, catalytic domain (IPR013781)

- Carbohydrate-binding domain, family 9 (IPR010502)
- Immunoglobulin-like fold (IPR013783)
- Galactose-binding domain-like (IPR008979)
- Cellulose-binding domain, fungal (IPR000254)
- Cellulose-binding domain, family II, bacterial-type (IPR001919)
- Glycosyl hydrolase, five-bladed beta-propellor domain (IPR023296)

Endo-1,4- β -Xylanase

- Class of enzyme that degrades β -1,4-xylan into xylose
 - This breaks down one of the major component of plant cell walls, hemicellulose



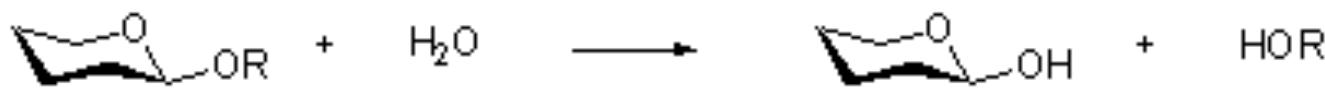


This enzyme cleaves β -1,4-xylan in an *endo* manner
(as seen above)

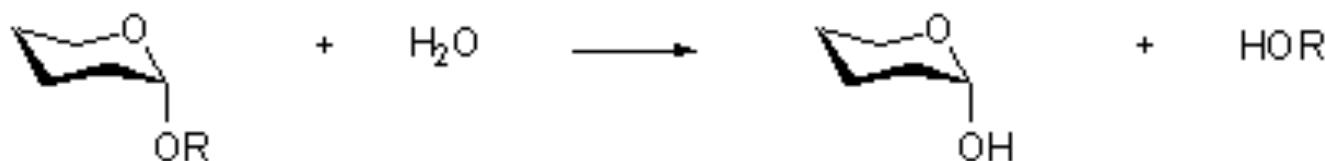
Retaining Mechanism

Retaining glycoside hydrolases:

1.) α



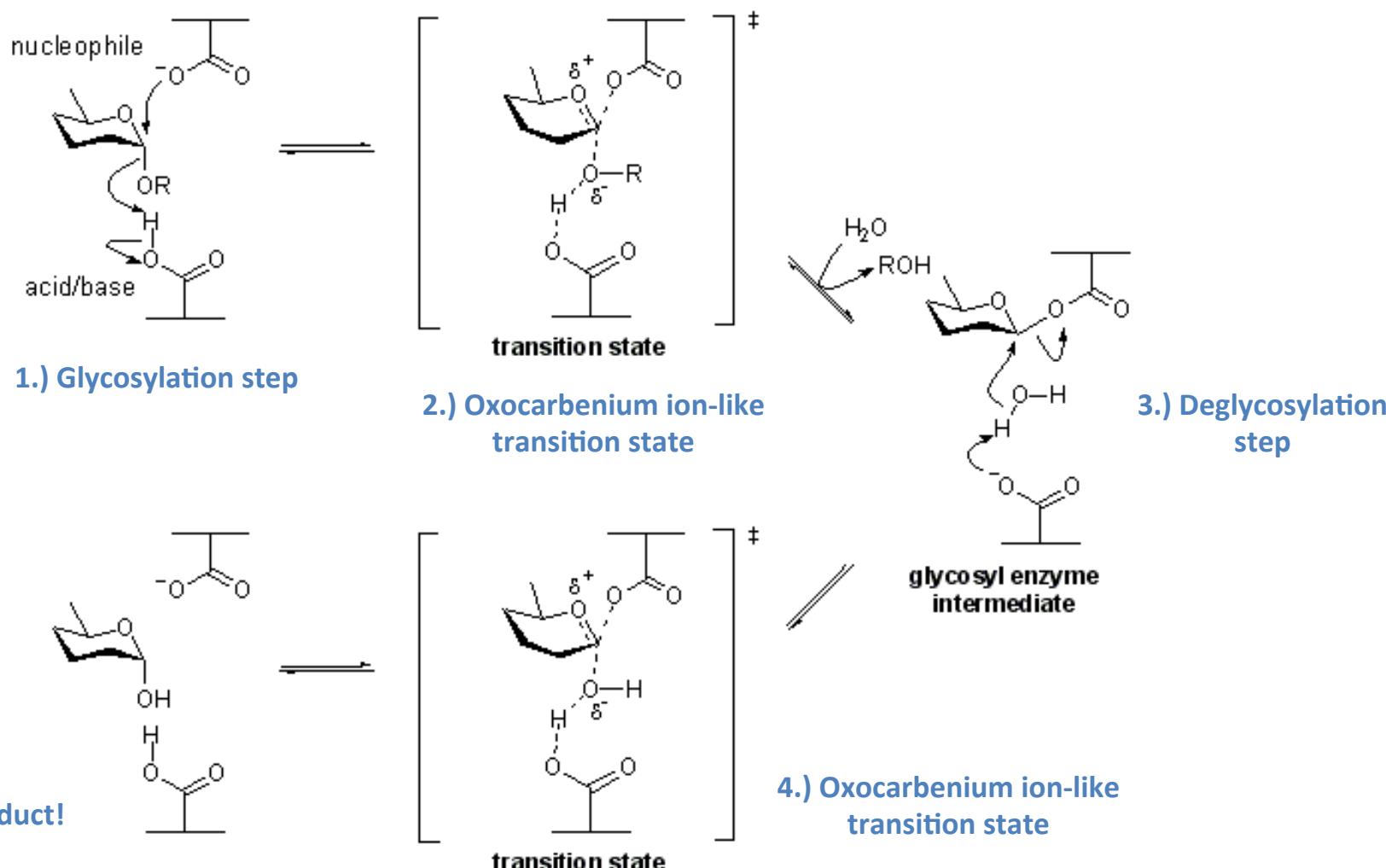
2.) β



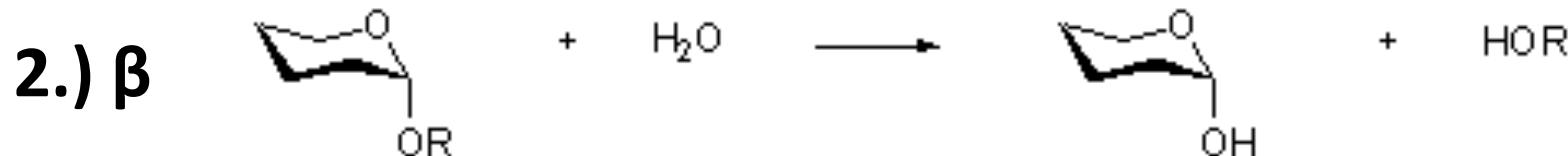
Retaining glycoside hydrolases:



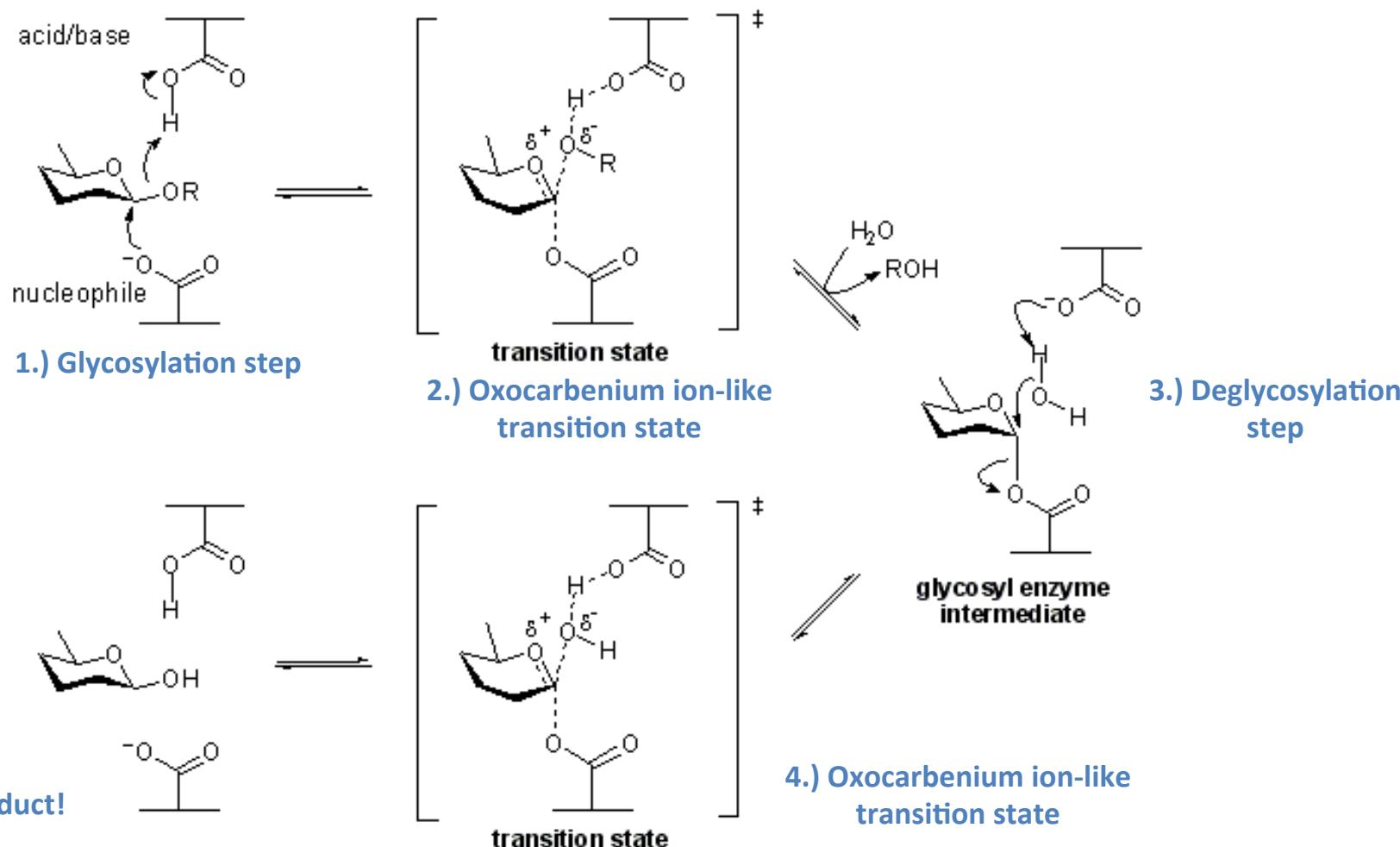
Retaining mechanism for an α -glycosidase:



Retaining glycoside hydrolases:

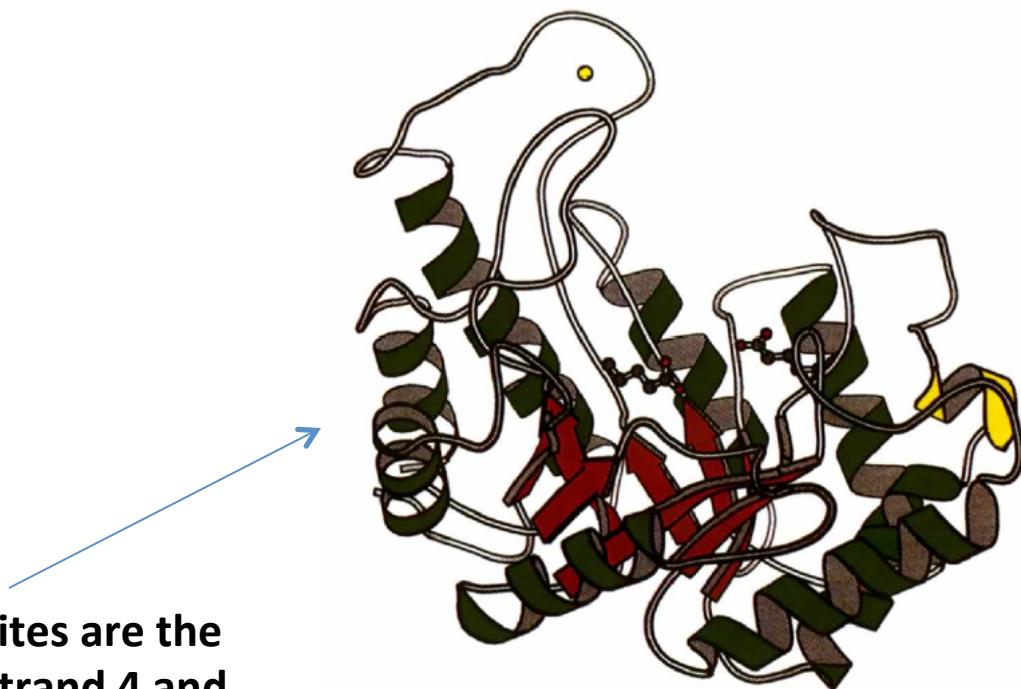


Retaining mechanism for a β -glycosidase:



3D Structure of GH10 / Active Sites

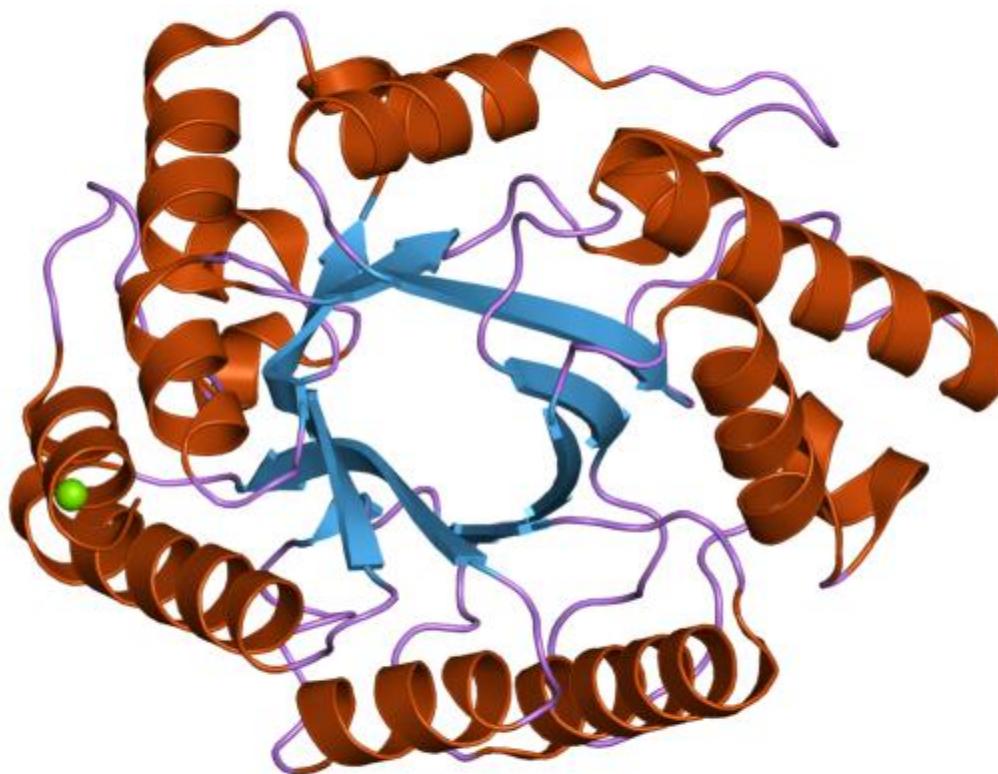
(the refined X-ray crystal structure at 1.8 Å, resolution of the catalytic core of the family F xylanase A from *Pseudomonas fluorescens* subspecies *cellulosa*)



*These active sites are conserved across all homologous GH10 proteins

Fig. 1. The eightfold β/α -barrel of xylanase A viewed perpendicular to the barrel axis showing the acid base Glu127 on strand 4 and the nucleophile Glu246 on strand 7. The Ca atom (yellow dot) stabilizes the long loop after strand 7. This figure was prepared using MOLSCRIPT (Kraulis, 1991). The β -strands are shown as arrows in red and the eight α -helices are shown in green; the additional α -helix ($\alpha 4\alpha$) is shown in yellow. The glutamates are 5.5 Å apart consistent with a retention of the configuration during catalysis.

Crystal structure of the catalytic domain of xylanase A from *Streptomyces halstedii* jm8



Again, note the
eightfold β/α -barrel, a
notable characteristic of
GH10 family members

Key Species

Key species	Number of proteins
 <i>Oryza sativa subsp. japonica</i> (Rice)	33
 <i>Arabidopsis thaliana</i> (Mouse-ear cress)	27

Taxa

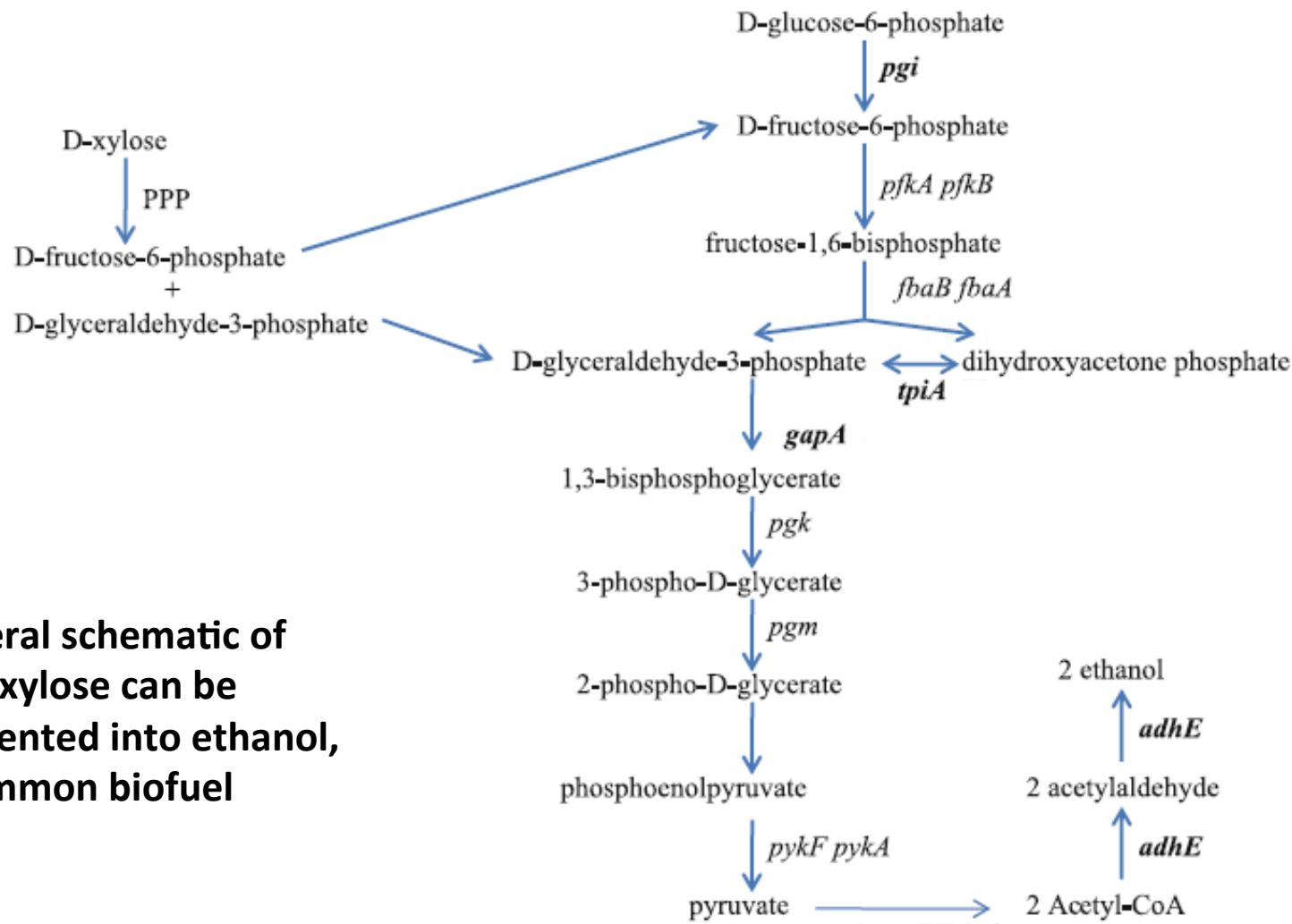
- **cellular organisms** [2110 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - Archaea [18 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Euryarchaeota [18 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - Bacteria (eubacteria) [1350 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Actinobacteria [261 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Bacteroidetes/Chlorobi group [245 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Chlamydiae/Verrucomicrobia group [42 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Chloroflexi [3 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Cyanobacteria (blue-green algae) [33 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Deinococcus-Thermus [3 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Dictyoglomi [4 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Fibrobacteres/Acidobacteria group [21 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Firmicutes (Gram-positive bacteria) [339 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Planctomycetes [23 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Proteobacteria [248 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Spirochaetes [16 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Thermotogae [28 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ environmental samples [82 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ unclassified Bacteria [2 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - Eukaryota (eucaryotes) [742 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Alveolata (alveolates) [8 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Opisthokonta [402 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Stramenopiles (heterokonts) [21 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ Viridiplantae [275 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ environmental samples [36 proteins](#) | [FASTA](#) | [Protein IDs](#)
- **unclassified sequences** [561 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ environmental samples [556 proteins](#) | [FASTA](#) | [Protein IDs](#)
 - ⊕ metagenomes [5 proteins](#) | [FASTA](#) | [Protein IDs](#)

Roles of Xylanase

- Important role for microorganisms that live off of plant sources
 - An example of which is seen with fungi that degrade plant sources for nutrients
- Human uses include...
 - Chlorine-free bleaching of wood pulp prior to papermaking
 - Digesting silage (for cattle/biofuel)

Why is this important?

- Xylan is a major component of the hemicellulose portion of plant cell walls and constitutes up to 35% of the total dry weight of higher plants (Badal, 1999).
- Therefore studying xylanases helps contribute to bioenergy research and the creation of biofuels! ☺



General schematic of how xylose can be fermented into ethanol, a common biofuel

What do we hope to accomplish?

GOALS

Compare the HMM domains of known GH10 proteins to determine if any JGI metagenome contain GH10 domains

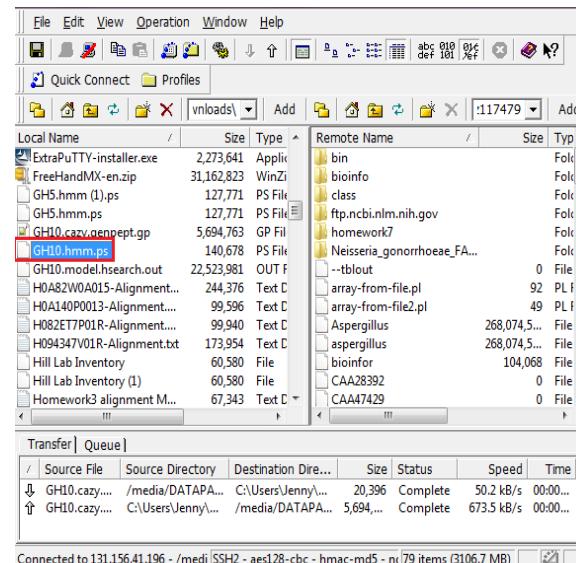
How on earth did we do this?!

METHODS

Step 2 – Download Datasets

2.2.2: HMM model

A screenshot of the dbCAN web server interface. At the top, there's a banner for 'Carbohydrate-active enzyme ANnotation' (dbCAN) with EU flags. Below the banner, the navigation menu includes Home, Browse, Genome, BLAST, Download, Help, and About us. A search bar is present. The main content area shows 'Basic information' and tabs for Sub-families, Phylogenetic tree, Signature sequences, and Component proteins. Under 'Activities', it says 'See annotated activities from CAZyDB'. The 'Statistics of sequences' section provides details about the source (BGI-gut, CAMERA, CAZY, Cow-Rumen, JGI, NCBI-Bacteria, NCBI-env-nr, NCBI-Fungi, NCBI-nr, Phytosome-Plant, Swiss-Prot, TrEMBL, WUSTL-bacteria) and taxonomy (Archaea, Bacteria, Eukaryota, Unclassified, Unclassified sequences). The 'Domain logo' section has a link to a logo image. The 'Download' section features a button for 'HMM model' (Multiple sequence alignment), which is highlighted with a red box. A red arrow points from this button to the terminal command in the screenshot below.



2.2.1: Known GH10 proteins

A screenshot of the CAZy database interface. The top navigation bar includes links for HOME, ENZYME CLASSES, ASSOCIATED MODULES, GENOMES, Glycoside Hydrolases, GlycosylTransferases, Polysaccharide Lyases, Carbohydrate Esterases, Auxiliary Activities, and a search bar. Below the navigation, a search result for 'GH10' is shown. The results table lists various entries with columns for Known Activities, Mechanism, Clan, 3D Structure Status, Catalytic Mechanophile/Base, Catalytic Proton Donor, Note, External resources, Commercial Enzyme Provider(s), Statistics, and Summary. A red arrow points from the 'GH10' search results table to the terminal command in the screenshot below.

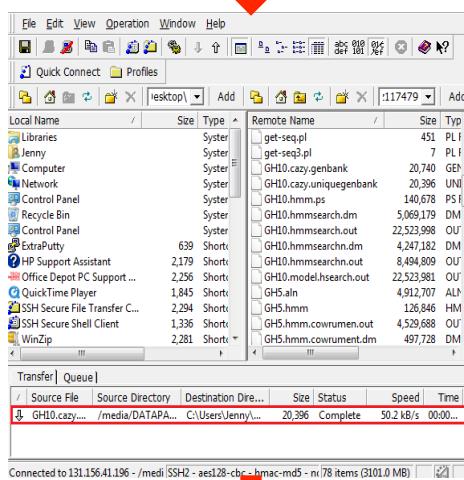
A screenshot of a terminal session on a Linux system. The command 'z117479@glu:~\$ vi GH10.cazy.genbank' is entered. A red arrow points from this command to a progress bar at the bottom of the terminal window, which shows the download of 'GH10.cazy.genbank' from 'http://cazy.scs.yale.edu/gh10.cazy.genbank' to '/media/DATAP...'. The progress bar shows the file size as 3106.7 MB and the speed as 673.3 kB/s.

```
ACP87428.1
ACW02056.1
ACP87436.1
ACW02064.1
ACP87439.1
ACW02067.1
ACP87441.1
ACW02069.1
ACP87442.1
ACW02070.1
ACP87443.1
ACW02071.1
ACP87447.1
ACW02075.1
ACP87450.1
ACW02078.1
ACP87451.1
ACW02079.1
ACP87452.1
ACW02080.1
ACP87455.1
ACW02083.1
ACP87458.1
-- INSERT --
```

Step 2.2.1 – Known GH10 Proteins

Parse for unique GenBank ID's

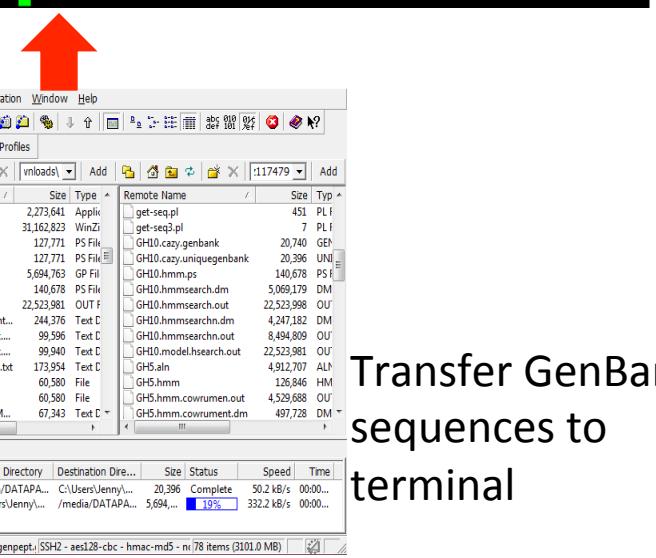
```
z117479@glu:~$ less GH10.cazy.genbank | sed '/^$/d' | sort | uniq | sort > GH10.cazy.uniquegenbank  
z117479@glu:~$ head GH10.cazy.uniquegenbank  
A43802  
AAA16427.1  
AAA17888.1  
AAA21812.1  
AAA23059.1  
AAA23062.1  
AAA23227.1  
AAA23286.1  
AAA56791.1  
AAA56792.1
```



Transfer unique ID file to local computer

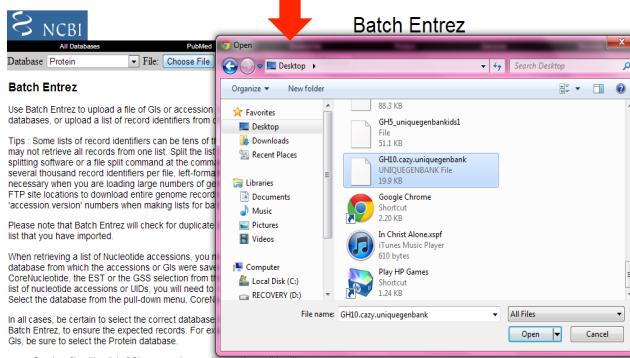
Reformat GenBank sequences to FASTA sequences

```
z117479@glu:~$ seqret -sequence GH10.cazy.genpept.gp -outseq GH10.cazy.fasta -sf  
ormat genbank -osformat fasta
```

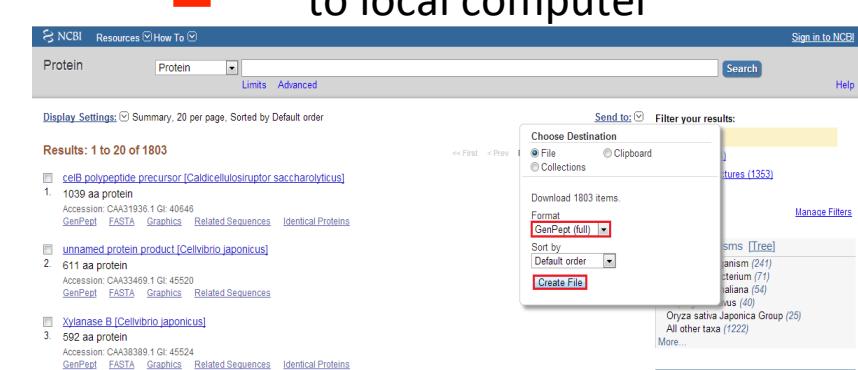


Transfer GenBank sequences to terminal

Download GenBank sequences to local computer



Search database for sequences



Step 3-4.1.1: Homology Search and Parsing Output Files

Hmmrsearch and parsing of JGI metagenome proteins against GH10 HMM

```
z117479@glu:~/project1$ hmmsearch --domtblout GH10.hmm.jgi.dm /media/DATAPART3/z  
117479/GH10.hmm.ps /home/yyin/db/jgi-v3.5.pr.fa > GH10.hmm.jgi.dm &  
z117479@glu:~/project1$ less GH10.hmm.jgi2.dm | grep -v '^#' | awk '{print $1,$3  
,$6,$7,$12,$13,$16,$17,$18,$19}' | awk '$6<1e-5&&($8-$7)/$3>.8' | sed 's/ /\t/g'  
> GH10.hmm.jgi.dm.ps
```

Hmmrsearch and parsing of GH10 known NCBI-nr proteins against GH10 HMM model

```
z117479@glu:~$ hmmsearch --noali --domtblout GH10.hmm.NCBI-nr.dim GH10.hmm.ps GH10.cazy.fasta > GH10.hmm.NCBI-nr.out &
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.dim | grep -v '^#' | awk '{print $1,$3,$6,$7,$12,$13,$16,$17,$18,$19}' | awk '$6<1e-5&&($8-$7)/$3>.8' | sed 's/ /\t/g' > GH10.hmm.NCBI-nr.dim.ps
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.dim.ps
```

CCA71122	1168	303	0	8.5e-102	8.5e-102	5
301	126	429				
CCA71122	1168	303	0	1.1e-100	1.1e-100	5
301	474	777				
CCA71122	1168	303	0	3.5e-101	3.5e-101	5
301	822	1125				
ADG76430	820	303	1.8e-214	1.2e-107	1.2e-107	
4	302	44	338			
ADG76430	820	303	1.8e-214	1e-107	1e-107	5
402	696					302
ADG75375	815	303	1.3e-205	4.1e-102	4.1e-102	
4	302	42	337			
ADG75375	815	303	1.3e-205	5.4e-104	5.4e-104	
3	303	395	691			
ABX44209	2457	303	1.9e-154	8.1e-85	8.2e-85	2
391	722					302
ABX44209	2457	303	1.9e-154	2.1e-68	2.1e-68	2
744	1083					302
ADB36415	366	303	3.7e-122	4.2e-122	4.2e-122	
2	302	37	365			
CCH02308	352	303	1.9e-120	2.2e-120	2.2e-120	
2	302	22	351			
AEE95681	380	303	6.9e-120	7.9e-120	8e-120	1
GH10.hmm.NCBI-nr.cm.psl						

Step 4.2-4.3: Parse Files to Extract Protein Sequences

Extract hit ID's

```
z117479@glu:~/project1$ less GH10.hmm.jgi.dm.ps | cut -f1 > GH10.hmm.jgi.dm.ps.id  
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.dm | grep -v '^#' | awk '{print $1,$3,$6,$7,$12,$13,$16,$17,$18,$19}' | awk '$6<1e-5&&($8-$7)/$3>.8' | sed 's/ /\t/g' | cut -f1 > GH10.hmm.NCBI-nr.dm.id
```

Parse for unique ID's

```
z117479@glu:~/project1$ less GH10.hmm.jgi.dm.ps.id | sort | uniq | sort > GH10.hmm.jgi.dm.ps.uniqid  
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.dm.id | sort | uniq | sort > GH10.hmm.NCBI-nr.dm.uniqid
```



Create perl script using bioperl modules to extract sequences from hit ID's

```
z117479@glu:~/project1$ vi step4.3.pl  
  
#!/usr/bin/perl -w  
  
open (ID,$ARGV[0]);  
while (<ID>) {  
    chomp $_;  
    $id_hash{$_}=1;  
}  
  
use Bio::SeqIO;  
  
$new=Bio::SeqIO->new(-file=>$ARGV[1], -format=>"fasta");  
  
while ($seq=$new->next_seq){  
    if (defined $id_hash{$seq->id}){  
        print ">",$seq->id,"\\n",$seq->seq,"\\n";  
    }  
}
```

CCA71122
CCA71122
CCA71122
ADG76430
ADG76430
ADG75375
ADG75375
ABX44209
ABX44209
ADB36415
CCH02308
AEE95681
ACP87343
ACW01971
ACX75889
ACX75889
ADL24782
ADL24782
AEV97130
ACP87323
ACW01951
AEE49315
BAM46412
GH10.hmm.NCBI-nr.dm.id

Step 4.3 – Extract Protein Sequences

Run perl script on parsed hit ID files and original databases, JGI metagenome and known GH10 proteins (NCBI)

```
z117479@glu:~/project1$ nohup perl step4.3.pl GH10.hmm.jgi.dm.ps.uniqid /home/yy  
in/db/jgi-v3.5.pr.fa > GH10.hmm.jgi.id.totalseq &  
[1] 50501  
z117479@glu:~/project1$ nohup perl step4.3.pl GH10.hmm.NCBI-nr.dm.uniqid GH10.ca  
zy.fasta > GH10.hmm.NCBI-nr.id.totalseq &  
[2] 50854  
>2000622610  
MHKRLLPLVLALVLVVTPFLLANPATPAAAAEIPWGDKSTPGMVSLGNRVRSDDEMPTMIRLMREAGVQWNREEIWWDQ  
VQFEPDGPFRWDGDSSRFYNDRAIQLQAEAGISILGLLDYNPAWFKGRNPHEEWLSDWGVDYVATVARYGRDRGQIKY  
WEVVNEPNLVPSGYESGLYNVDDFVRLQTASAAIRAADPEAKIVLGGVADIWSEIPEFAYDTPDYLQRLYALDAWPMFD  
ILGLHPYRPDAPEVPVLRDRSQTREQAAEIDALLAQFGNKPIWYTEVGWSTESDGIVSEDEQAALQRFYLLAMTHPG  
VEKIFWYDFRNDTGDNSNYTRPINDPNENQFHGFLRGNYPLNFDDQRMRKPSYSAYYHLHNHLGGVSWRAKYALP  
>2001264806  
MAFSKDKASFTRRSAIAAGLAAGVSACVTPQPGVPAMPSDLQLAKRSGRFGSAVGWGRPAGDRGSFANPAYAAILENEC  
SLLVPENELKWQWTRPGPGQFDFRQFDAIADYASRKGLLRGHTLFWLPEKWPWKWLVAHDFGAQPAAKAEAMLREHVQT  
VCRRYGT RIYSYDVVNEAVQPENGLIRDVTVKALGHEATLDVMFQTARAEEAPHAQLVYNDYMSWERNSEDETHMRGVLK  
LLEGFRKRGT PVDALGIQSHIRLLKPLTVAEIVRESGGPWRRFLDEVVAMGYKLEITEFDVNDRLLAPTNPLERDRMVADY  
AKAYLDVMLSYPQLRDIAMA GMVDRY SWLTGF DPKD KTLKRG TPYDKHFRPK LLRE AIAA FMGAA APA  
>2003420141  
MTPTRRAILAAPLALAACDRFASAEPSPPNVPLKSIAPAPFGTAIKASQIDDPDWVALARANVSQLTPEWEMKMEYILA  
DGLDRPNFDRSDRIAFAFARAEGMAMHHTLIWYAQGKEAFAGLSGAAFDRFDGYIATVAGRYRGKVRSDVVNEPILDD  
GSGMRDCHWSARYGHDTYILRAFEKARIADPDAVLFLNEYQNQESVPAKGAQFLKLVERLLKAGCPLQGVGLQSHLWIDIP  
EGVIAAYMREISQFGLPIHVSEL DCTLR TENRLDLRSQADRIAAQGARVTEL ASAF AALPKAQQFAFTVWGLRDTD  
>2003425276  
LGIGAQALDGMTTVLHIKVQDFPTAANLPVAMI PNCAATRHAFTLDGGCP RRHGGADAANAAGPAAPGCASPRSPLTM  
WFTRVSLANPVFATMLMLALMVLGVFSYQRLKVDQFPNVDFPVVVVTVDYPGASPEIVESEVTKKIEEGVN SIAGINALT  
SR SYESTSVIIIEFQLHVDGRRAEDVREKVASVRPTL RDEVKEPRV LRFDPASSPIWSAVLPQADAGAKAPDAVALTS  
WADQVLKKRLENVRGVGVAVNLVGATPREINIVLQPEALEAYAITPDVPPPLKSIAPAPFGTAIKAMQIDDPDWVALARANV  
SQLTPEWEMKMEYILANGLDRPNFDRTDRIAFAQAQGMAMHHTLIWYSQGQEAFAGMDDAAFDRFDGYIAAVAGRYR  
GH10.hmm.jgi.id.totalseq
```

Step 5: Retrieve HMM Domain Regions

Perl script using bioperl modules to extract HMM domain regions using first HMM parsed output files

```
z117479@glu:~/project1$ vi idsubseq-bioperl.pl
#!/usr/bin/perl

use Bio::SeqIO;

open (IN,$ARGV[0]);

while (<IN>){
    chomp$_;
    @col=split(/\t/,$_);
    push(@{$id_hash{$col[0]}},$_);
}

$new=Bio::SeqIO->new(-file=>$ARGV[1], -format=>"fasta");

while($seq=$new->next_seq()){
    if(defined $id_hash{$seq->id}){
        @id_array=@{$id_hash{$seq->id}};

        foreach(@id_array){
            @id_col=split(/\t/,$_);
            print ">",$id_col[0],"|",$id_col[-2],"-", $id_col[-1],"\n",
            $seq->subseq($id_col[-2],$id_col[-1]),"\n";
        }
    }
}
```

```
>2000622610|77-378
WWDQVQFEPDGPFRWDGDSSRFNYDRAIQLQAEAGISILGLLDYNPAWFGRNPHPEEFLSDWGDYVYATVARYGRDRG
QIKYWEEVNEPNLVPSEYESGLVNDDFVRVLQTASAIARADPEAKIVLGGVADIWEIPEFAYDTPDYLQRLYALDAW
PMFDILGLHPYRPDAPEPVLRDRSQTREQAAEIDALLAQFGNKPWYTEVGWSTESDGIVSEDEQAALQRFYLLAM
THPGVEKIFWYDFRNDTGDNNSYTRPINDPNENQFHGLLRGNYPLNFDDQRMRKPSYSAYY
>2001264806|42-381
QLAKRSGRRFGSAVGWGRPGADRGSFANPAYAAILENECSLLVPENELKWQWTRPGPGQFDFRQFDAIADYASRKGLLR
GHTLFWLPEWKWPWLVAHDFGAQPAAKAEAMLREHVTVTQCRRYGTRIYSYDVNEAVQPEGLIRDTVVTKALGHATEL
DVMFQTARAEEAPHAQQLVYNDYMSWERNSEDETHMRGVKLLEGFRKRGTQPVDALGIQSHIRLLKPLTVAEIVRESGGPWR
RFLDEVVAMYKLEITEFDVNDRLAPTNPLERDRMVADYAKAYLDVMLSYPQLRDIAMAWGMVDRYSLWTGFDPDKTLK
RGT PYDKHFRPKLLEAIAA
>2003420141|40-315
APFGTAIKASQIDDPDWALARANVSQLTPEWEMKMEYILADGLDRPNFDRSDRIAFAAREGMAMHHTLIWIYAQGKEA
FAGLSGAAFDRADFQYIATVAGRYRGKVRSDVVNEPILDGSGMRDCHWSARYGHGQYIILRAFEKARIADPDAVLFLNE
YNQESVPAKGAFQFLKLVERLLKAGCPLQGVGLQSHLWIDIPEGVIAAYMREISQFGLPIHVSLEDCTLRTENRLDMRNKA
DRIAAGARVTELASAFALA PKAQQFAFTVWGLRDT
>2003425276|296-603
APFGTAIKAMQIDDPDWALARANVSQLTPEWEMKMEYILANGLDRPNFDRSDRIAFAQAGMAMHHTLIWIYSQGQEA
FAGMDDAAFDRADFQYIATVAGRYRGKVRSDVVNEPILDGSGMRDCHWSARYGHGQYIILRAFEKARIADPDAVLFLNE
YNQESVPAKGVQFLKLVERLLKAGCPLQGLGLQSHLWIDIPEGVIAAYLREVAQFGLPIHVSLEDCTLRTENRLDMRNKA
DRIAAGARVTELASAFALA PKAQQFAFTAWGLRDTDSWYRQGEKDDGKDPLPFDSFGRPNPMAAAL
>2003454879|40-355
LAPFFVGATGMTGQLDDPTWVAADRHLQITPEWEILPERILGPGFAYDFSADRMVDWATARGRVGHALVWYAQGL
GH10.hmm.jgi.id.domseq
```

Run perl script using the sequences extracted from the hit ID's (step 4.3) and the first parsed hmmersearch parsed output files (step 4.1.1)

```
qz117479@glu:~/project1$ nohup perl idsubseq-bioperl.pl GH10.hmm.jgi.id.dm.ps GH10.
mm.jgi.id.totalseq > GH10.hmm.jgi.id.domseq &
[4] 53635
z117479@glu:~/project1$ nohup perl idsubseq-bioperl.pl GH10.hmm.NCBI-nr.dm.ps GH
10.hmm.NCBI-nr.id.totalseq > GH10.hmm.NCBI-nr.id.domseq &
```

Step 6: Combine JGI and NCBI files

Use “cat” to combine the JGI metagenome and known GH10 proteins (NCBI) files (step 5) containing the HMM domain sequence regions into one new file

```
z117479@glu:~/project1$ cat GH10.hmm.jgi.id.domseq GH10.hmm.NCBI-nr.id.domseq >
GH10.hmm.jgi.NCBI-nr.domseq
z117479@glu:~/project1$ less GH10.hmm.jgi.NCBI-nr.domseq
>2031944243|75-350
RGEYEPEHEGETQAHIRIRATAEWFAARAVRLKGHPLVWHTVKAPWMDSLSPDAAWHNTLDRIRREVGDFSGLIDTWDVVNE
AVIMPVFVNEPDGTPNVISRIAQQKGRIPLIAREAFDEARATNPGATLLNDFDLSSAYEILIEGLLEAGVRIDALGLQTH
MHQGYWGEETMLAMVDRFARYGLPLHLTENTLLSGDLMMPA HIVDLNDYQPASWPSTPEGEARQADEIARHYRSLVAHPAV
ASITYWGDISDAGAWLGAPAGLLRADGSRKPAYDALR
>CAA31936|45-372
lcesykddfmigvaiparclsndtdkrmvlkhfnsitaenemkpesllaggstglsyrfstadafvdfastnkigirgh
tlvwhnqtpdwffkdsgqrlskdallarlkqiydvvgyrkvyawdvvneaidenqpdssrrstwyeicgpeyieka
fiwaheadpnaklfyndynteiskkrdfiynmvknlkskgipihgigmqchinvnwpvsseiensiklfssipgieihit
eldmslynygssenystppqdllqkqsqkykeiftmlkkyknvvksvtfwglkddyswlrsfygkndwpplffedysakp
ayawviea
>CAA33469|270-607
adfpiwgavaasggnadiftssarqnivraefnqitaenimkmsymysgsnfstnsdrlvswaaqngqtvhghalvwhp
syqlpnwasdsnanfrqdfarhidtvaahfagqvkswdvvnealfdsaddpdgrgsangyrqsvfyrqfpgpeyideafr
rapradptaelyyndfnteengakkttalvnlvqrllnnvgvpidgvfqmhvmndypsianirqamqkivalsptlkikit
eldvrlnnpydgnssndytnrndcavscagldrqkarykeivqaylevvppgrggitvwgiadpdswlythqnlpdwpl
1fndnlqpkpayqgvvea
>CAA38389|321-562
akkfignittsgavrsdftrywnqitpeneskwgsvegtrnvynwapldriyayarqnnipvkahtfvwgaqspswlnnl
sgpevaveieqwirdycarypdtamidvvneavpghqpagyaqrafgnnwqrvfqlarqycpnsililndynnirwqhn
efialakaqgnyidavglqahelkgmtaaqvktaidniwnqvgkpiyiseydigtdndqvqlqnfqahfpvfynhphvhg
it
```

Step 8: Build Multiple Sequence Alignment

Use MAFFT program to build multiple sequence alignment using the newly combined file (step 6)

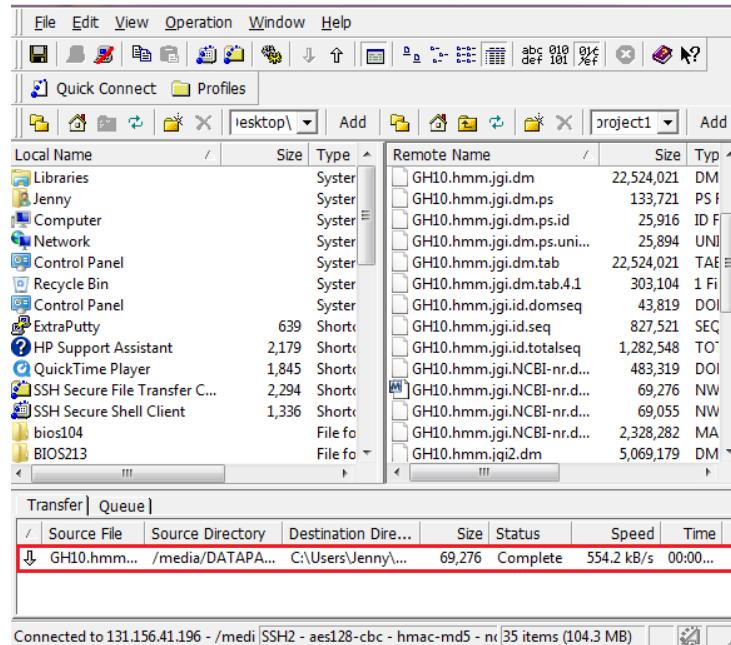
```
z117479@glu:~/project1$ mafft --anysymbol GH10.hmm.jgi.NCBI-nr.domseq > GH10.hmm.jgi.NCBI-nr.domseq.mafft &
```

```
>2000622610|77-378
-----
-----WWD-QVQF-EPDGPFRW-----
-----DGD-----SSRFY
NY-----DRA-----IQLQ---A-EAGI---S----I
LG--L---LDY----NPAWFKGRNPH-----P
EE-W-L
-----SDW-GDY
-----
-----
-----
-----
-----V-YATVA---R---Y-----GRDRG---Q-----
-----IKYW-----EVW-NEPNLVP-----
-----
-----
-----SGYESGLYNVDDFVRVLQTASAAIRAADPEAKIV-
LGG-V-----ADI-----WSE-----I PEFAY-DTP-----
-----DYLQRLYAL-----D-----AWPM
FD--ILGLHPYRPDAPEVPVLRRDRSQTYREQAAE-
-----IDALLAQF---G-----N-
--KPIWYTEVG-----WST---ESDGI-----VSEDE-
GH10.hmm.jgi.NCBI-nr.domseq.mafft
```

Step 8: Build Phylogenetic Tree

Use FastTree to build phylogenetic tree from the MAFFT alignment

```
z117479@glu:~/project1$ nohup /home/mrupani/Downloads/FastTree GH10.hmm.jgi.NCBI  
-nr.domseq.mafft > GH10.hmm.jgi.NCBI-nr.domseq.fasttree.nwk &
```



Step 9: Parse HMM Domain Sequence Files for hit ID's

To color code the JGI metagenome proteins and the known GH10 proteins, the exact ID's of the sequences needs to be obtained

```
z117479@glu:~/project1$ less GH10.hmm.jgi.id.domseq | grep ">" | sed "s/>//g" >  
GH10.hmm.jgi.id.domseq.id  
z117479@glu:~/project1$ less GH10.hmm.jgi.id.domseq.id  
  
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.id.domseq | grep ">" | sed "s/>//g"  
" > GH10.hmm.NCBI-nr.id.domseq.id  
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.id.domseq.id
```

```
2004144103|11-279  
2004144113|491-761  
2004144883|87-357  
2004146193|125-450  
2004146653|87-358  
2004175319|38-369  
2004380153|1-271  
2005561842|40-355  
GH10.hmm.jgi.id.domseq.id
```



Perl scripts were written to assign the hex blue color code to JGI metagenome proteins and the hex red color code to known GH10 proteins

```
z117479@glu:~/project1$ vi jgicolor.pl  
#!/usr/bin/perl  
  
while (<>){  
    chomp $_;  
    print $_, "\t#0000FF\n";  
}
```

```
z117479@glu:~/project1$ vi NCBI-nrcolor.pl  
#!/usr/bin/perl  
  
while(<>){  
    chomp$ _;  
    print $ _, "\t#F0000\n";  
}
```

Perl script was run to assign the color codes to the previously parsed files

```
z117479@glu:~/project1$ perl jgicolor.pl GH10.hmm.jgi.id.domseq.id > GH10.hmm.jg  
i.id.domseq.color  
z117479@glu:~/project1$ less GH10.hmm.jgi.id.domseq.color  
  
z117479@glu:~/project1$ perl NCBI-nrcolor.pl GH10.hmm.NCBI-nr.id.domseq.id > GH1  
0.hmm.NCBI-nr.id.domseq.color  
z117479@glu:~/project1$ less GH10.hmm.NCBI-nr.id.domseq.color
```

```
2004142063|245-619      #0000FF  
2004144103|11-279       #0000FF  
2004144113|491-761       #0000FF  
2004144883|87-357       #0000FF  
2004146193|125-450       #0000FF  
2004146653|87-358       #0000FF  
2004175319|38-369       #0000FF  
2004380153|1-271        #0000FF  
2005561842|40-355        #0000FF  
GH10.hmm.jgi.id.domseq.color
```



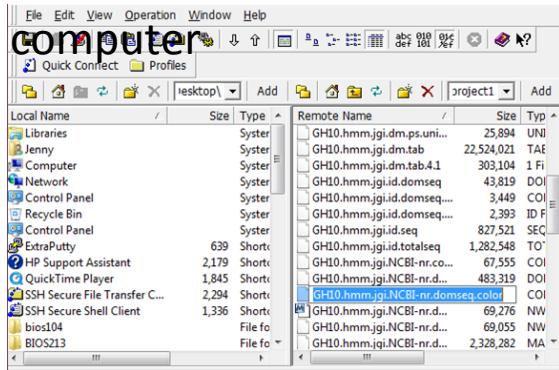
Step 9: Download Datasets to computer

Combine color files

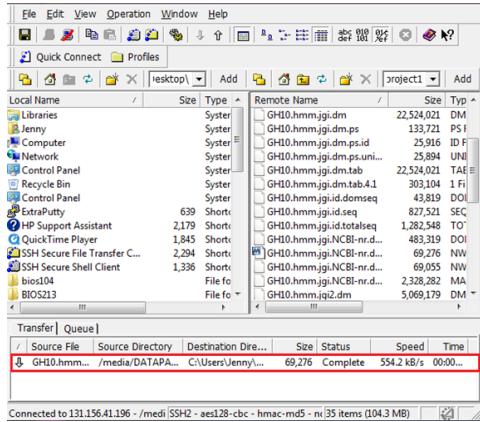
```
z117479@glu:~/project1$ cat GH10.hmm.jgi.id.domseq.color GH10.hmm.NCBI-nr.id.domseq.color > GH10.hmm.jgi.NCBI-nr.domseq.color  
z117479@glu:~/project1$ less GH10.hmm.jgi.NCBI-nr.domseq.color
```

2029678161 71-373	#0000FF
2029683389 3-296	#0000FF
2029724659 1-347	#0000FF
2030993216 7-327	#0000FF
2031944243 75-350	#0000FF
CAA31936 45-372	#FF0000
CAA33469 270-607	#FF0000
CAA38389 321-562	#FF0000
CAA43712 3-337	#FF0000
AAA23059 17-340	#FF0000
AAA23062 18-300	#FF0000

Transfer combined color files to local computer



Transfer FastTree Newick file to local computer



Connected to 121.156.41.196 - /media/SSH2 - aes128-cbc - hmac-md5 - nc(35 items (104.3 MB))

Step 9: Upload Datasets to iTOL

Upload JGI metagenome and known GH10 protein Newick file to iTOL

You can upload a list of NCBI taxonomy IDs and the tree of life will be pruned to include only the species from your file. download.

File with tax IDs: No file chosen

Upload your own tree

NEW! If you are using iTOL to upload your own trees, try [creating a personal account](#). More info about the iTOL pe-

Use this form to upload your own phylogenetic tree. It should be in plain text, in one of supported formats (Newick, Nexus available).

You can either paste your tree into the box, or upload a file using the file selector below. Don't forget to select the correct

Paste or type the tree:

Upload a file which contains your tree: GH10.hmm.jgi...sttree.nwk

Tree format:

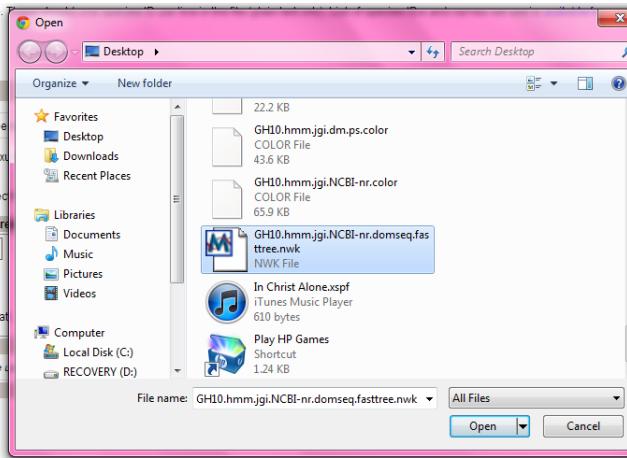
Make sure the file is plain text, and contains only trees in the selected format

Tree name: if you don't specify a name, a numeric ID will be used

Optional information

Advanced options (show)

(If you're uploading extra data with your tree, fill the dataset section below before clicking 'Upload')



Upload JGI metagenome and GH10 protein combined color file to iTOL

Upload datasets for your tree

TOL can annotate phylogenetic trees with several types of data. Please check our help pages for the detailed explanations.

Dataset 1 **Dataset 2** **Dataset 3** **Dataset 4** **Dataset 5** **Dataset 6** **Dataset 7**

Dataset 1 file: GH10.hmm.jgi...l-nr.color

please use plain text files only

Display label:

label will be used in the legends

Field delimiter: Space Tab Comma

make sure the correct delimiter is selected

Data type:

check the [help pages](#) for detailed descriptions

Bar size: Maximum: pixels

maximum value in the dataset will have this pixel size

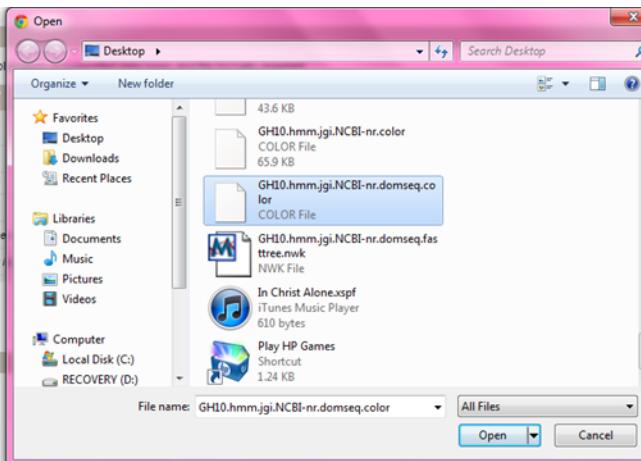
Dataset 1 color:

used in the legends and for datasets where the color is

Retrieve a previously uploaded tree

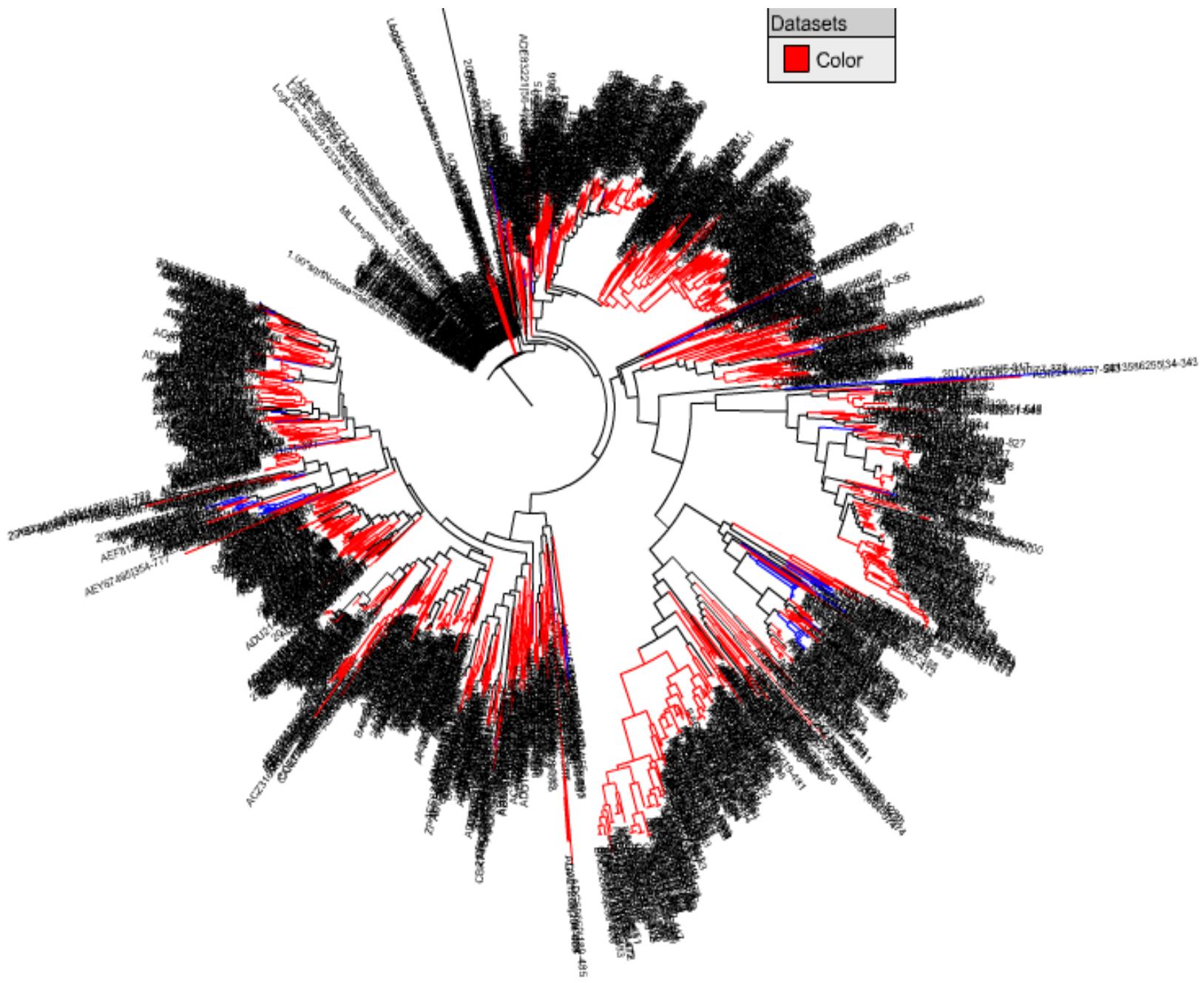
Tree ID: Retrieve

Your recently uploaded trees (1. is the most recent):

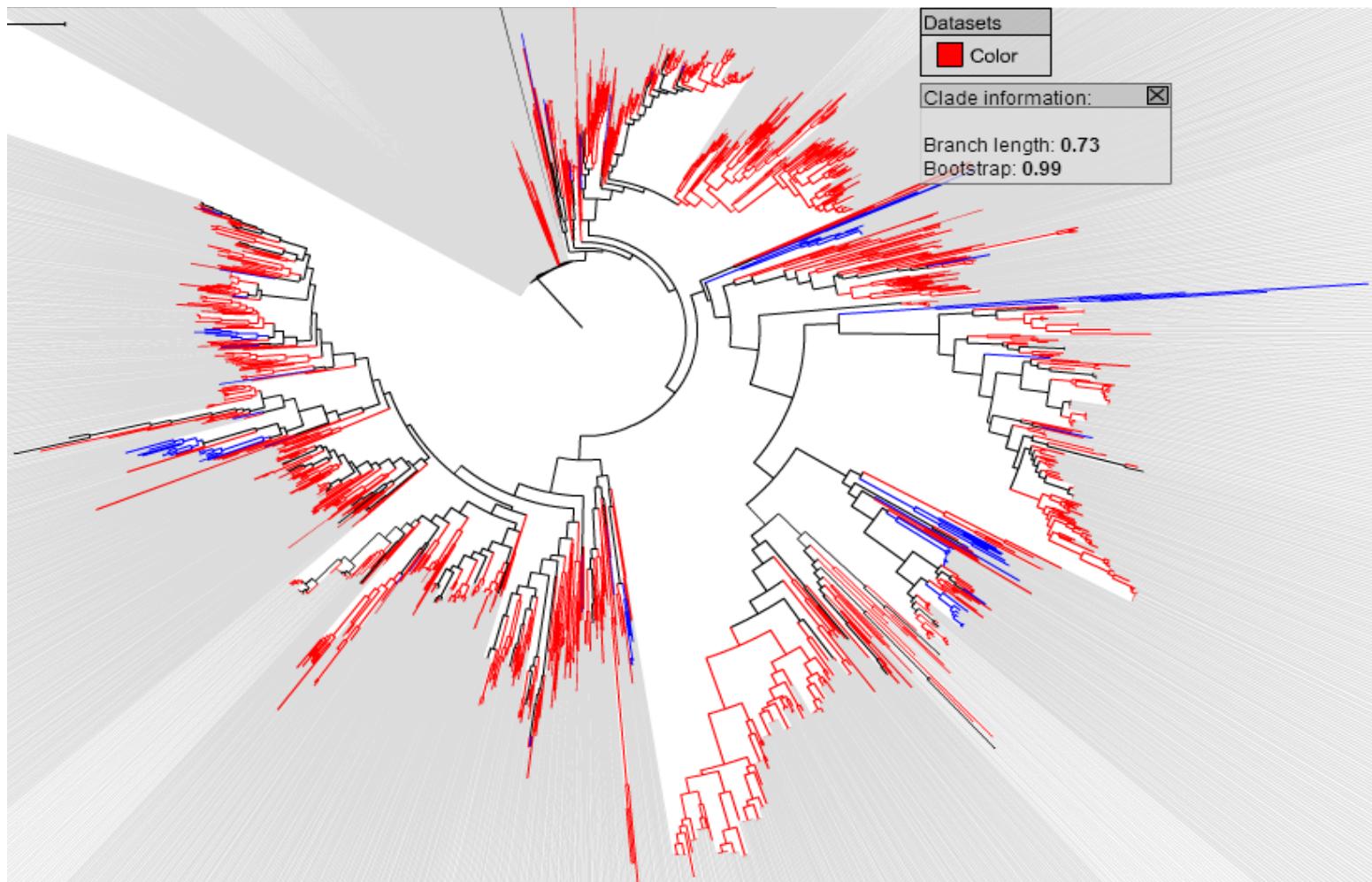


The fun stuff... ☺

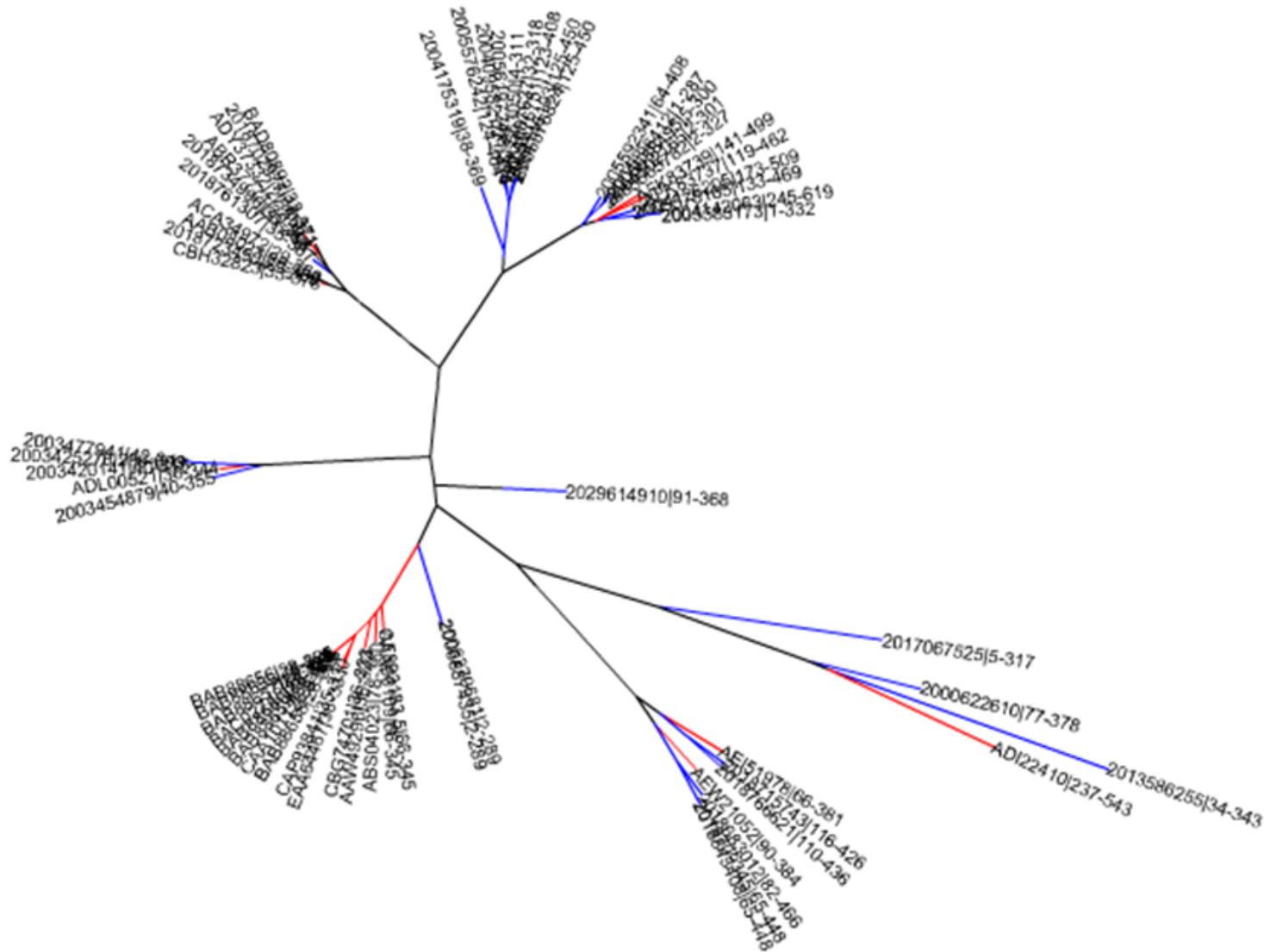
MAJOR FINDINGS



Pruning Phylogenetic Tree



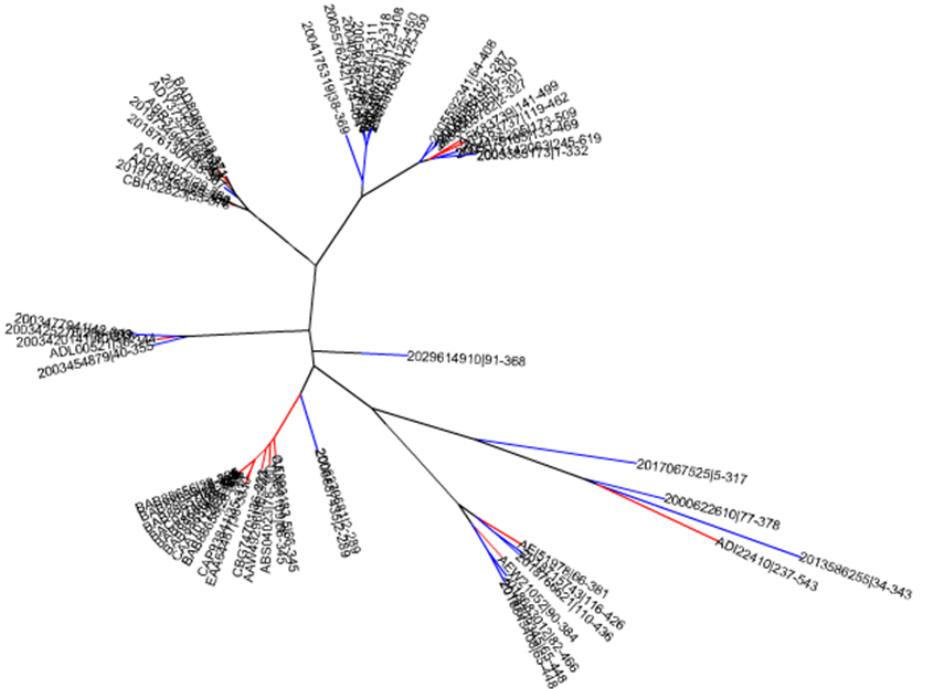
Pruned tree with JGI metagenome proteins



What do we take away from this?

RESULTS!

- A few metagenome proteins share stringent sequence homology with the HMM domains of known GH10 proteins



- What are the JGI metagenome proteins that share stringent sequence homology with known GH10 proteins?

Export pruned tree form iTOL

Export current tree (13115621110683413667469420) to other formats

Current display mode:	Other tree features:
Unrooted	Leaves visible: 64 Rotation: 0° Arc: 360° Branch lengths ignored: No Branch lengths displayed: No Bootstraps displayed: No



Please select the output format below and customize other export options as required. The default font size is likely wrong, click the 'suggest optimum' link to calculate the correct value for the current tree view.

Format	Encapsulated Postscript (eps) select the desired output format
Font size	<small>Font size for leaf labels</small>
Display leaf labels	<small>You can omit leaf labels from the exported trees using this option</small>
Line width	<small>Line width for tree branches, in pixels</small>
Branch colors	<small>branch colors are ignored, they will be black by default</small>
Leaf labels	<small>leaf labels were edited (or automatically assigned), use this option to show original tree IDs instead</small>
Text formats	Newick (txt)
Select the dataset(s)	Nexus (txt) PhyloXML (txt)
<small>Note: There are no datasets selected. You can include multiple datasets in the exported trees. iTOL will not check for possible overlaps, but datasets will be grouped into separate layers (in vector formats), making post processing easier.</small>	



Exporting the tree into newick format

Your file is ready for download. It will be kept on our server for 1 week.

Citations: Letunic and Bork (2006) *Bioinformatics* 23(1):127-8 and Letunic and Bork (2011) *Nucleic Acids Res* doi: 10.1093/nar/gkr201

Copy and paste newick file to terminal using vi

```
z117479@glu:~/project1$ vi GH10.hmm.jgi.NCBI-nr.itolprune.nwk
```

050|4-311:0.06525, ((2004087637|32-318:0.00910,2005576751|123-408:0.00321):0.0389
7[0.99], (2004146193|125-450:0.00962,2005576824|125-450:0.01532):0.02946[0.95]):0.
.06334[1.00]):0.06137[0.96], (2005576242|124-404:0.14891, (2004095821|21-315:0.000
14,2005619724|6-294:0.00014):0.08218[1.00]):0.04073[0.78]):0.31556[1.00]):0.0493
2[0.26]):0.09091[0.95], (((((2004126205|173-509:0.0,2005578165|133-469:0.0):0.20
382, (2004142063|245-619:0.00014,2005583173|1-332:0.00014):0.37234[1.00]):0.06424
[0.79], (AFK83739|141-499:0.30058,AFK83737|119-462:0.33129):0.02341[0.06]):0.0959
3[0.99], (2005592341|64-408:0.17191, ((2004096412|1-287:0.00853,2005588495|5-300:0.
.01095):0.09764[1.00], (2004102165|2-301:0.02596,2005565782|2-327:0.01088):0.0440
2[0.95]):0.11588[1.00]):0.05234[0.79]):0.03978[0.83]):0.40091[1.00]):0.16616[0.9
1]):0.03618[0.39]):0.13718[0.94]):0.18219[0.98]):0.08241[0.52]):0.18631[0.79]):0.

Extract tree ID's using bioperl modules

```
z117479@glu:~/project1$ vi get.treeid.pl
```

```
#!/usr/bin/perl -w

use Bio::TreeIO;

$treeio=Bio::TreeIO->new(-format=>"newick", -file=>$ARGV[0]);

while($tree=$treeio->next_tree){
    for $node ($tree->get_nodes){
        print $node->id, "\n";
    }
}
```

```
2029614910|91-368
AEW21052|90-384
2018683012|82-466
2018709345|65-448
2018843408|65-448
AEI51978|66-381
2018715743|116-426
2018766621|110-436
2017067525|5-317
2000622610|77-378
2013586255|34-343
ADI22410|237-543
2006270681|2-289
2006887435|2-289
AAW49296|36-333
CBG74701|36-333
BAB88655|28-325
BAB88656|28-325
BAB88657|28-325
BAB88660|28-325
BAC02742|28-325
CAA10112|28-325
BAB88659|28-325
GH10.hmm.jgi.NCBI-nr.itolprune.id
```

Run perl script on newly created newick file from pruned iTOL tree

```
z117479@glu:~/project1$ perl get.treeid.pl GH10.hmm.jgi.NCBI-nr.itolprune.nwk | sed '/^$/d' > GH10.hmm.jgi.NCBI-nr.itolprune.id
z117479@glu:~/project1$ less GH10.hmm.jgi.NCBI-nr.itolprune.id
```



Parse output file to remove sequence coordinates

```
z117479@glu:~/project1$ less GH10.hmm.jgi.NCBI-nr.itolprune.id | sed 's/|/\t/' | cut -f1 > GH10.hmm.jgi.NCBI-nr.itolprune.id.ps
z117479@glu:~/project1$ less GH10.hmm.jgi.NCBI-nr.itolprune.id.ps
```

```
2013586255
ADI22410
2006270681
2006887435
AAW49296
CBG74701
BAB88655
BAB88656
BAB88657
BAB88660
BAC02742
CAA10112
BAB88659
GH10.hmm.jgi.NCBI-nr.itolprune.id.ps
```

Works Cited

- Béguin, Pierre. "MOLECULAR BIOLOGY OF CELLULOSE DEGRADATION." *Annu. Rev. Microbiol* 4 (1990): 219-48. Print.
- HARRIS, GILLIAN W. "Refined Crystal Structure of the Catalytic Domain of Xylanase A from *Pseudomonas Fluorescens* at 1.8 .~ Resolution." *Acta Crystallographica Section D/ International Union of Crystallography* D52 (1996): 393-401. Web. <<http://journals.iucr.org/d/issues/1996/02/00/he0136/he0136.pdf>>.
- "InterPro." *Glycoside Hydrolase, Family 10*. EMBL-EBI, n.d. Web. <<http://www.ebi.ac.uk/interpro/entry/IPR001000>>.
- Manow, Ryan. "Partial Deletion of Rng (RNase G)-enhanced Homoethanol Fermentation of Xylose by the Non-transgenic Escherichia Coli RM10." *J Ind Microbiol Biotechnol* 39 (2012): 977-85. Print.
- Withers, Steven. "Glycoside Hydrolase Family 10." *CAZypedia*. N.p., 10 Sept. 2012. Web. <http://www.cazypedia.org/index.php/Glycoside_Hydrolase_Family_10#bibkey_HJG>.
- Withers, Steve, and Spencer Williams. "Glycoside Hydrolases." *CAZypedia*. N.p., 15 Jan. 2013. Web. <<http://www.cazypedia.org/index.php/Retaining>>.
- "Xylanase." *Wikipedia*. Wikimedia Foundation, 19 Apr. 2013. Web. <<http://en.wikipedia.org/wiki/Xylanase>>.