

MCDB 187 Lab Introduction

Colin Farrell / Feiyang Ma

MCDB 187

April 3, 2019

Contact Info

Lab Notes: <https://nuttylogic.github.io/teaching/>

- Colin Farrell
- 3rd Year PhD Student, Human Genetics
- colinpatfarrell@g.ucla.edu
- Feiyang Ma
- 3rd Year PhD Student, Molecular Biology
- mafeiyang@g.ucla.edu



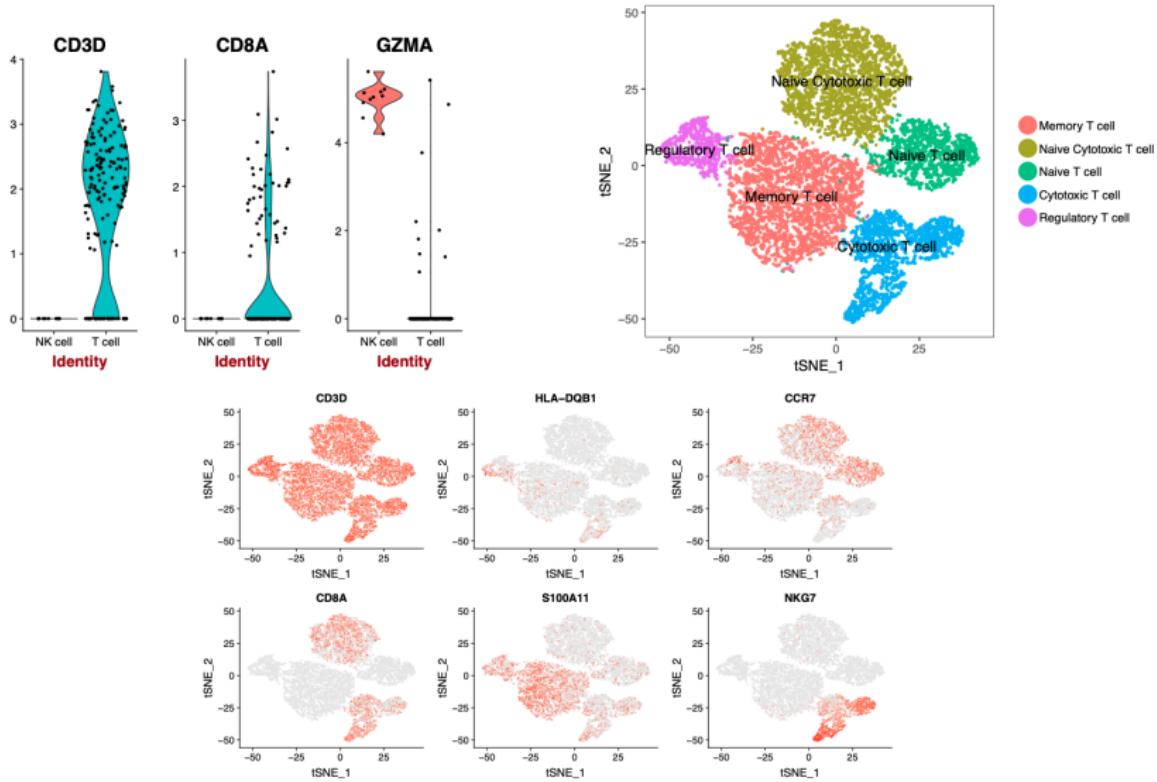
Feiyang Background



Chongqing, China
中国 重庆



Feiyang Research

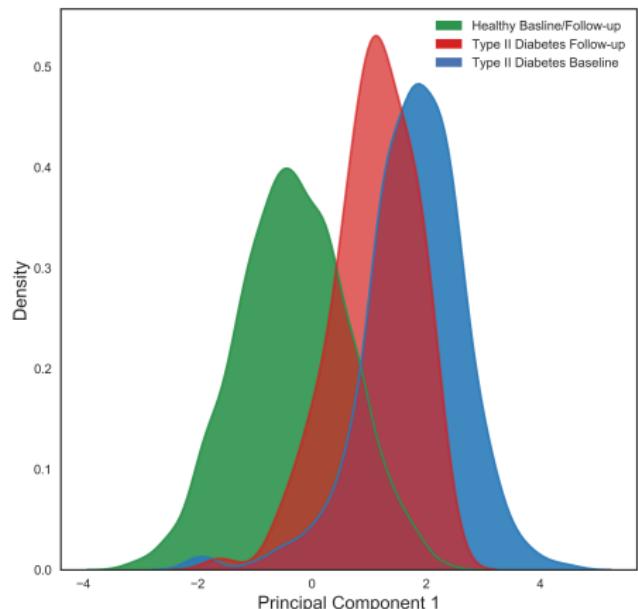
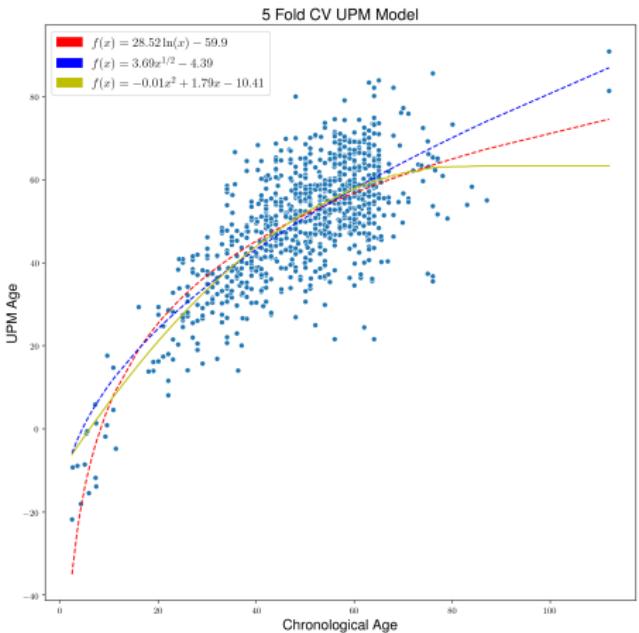


Feiyang Interests

Big fan of basketball (76ers)



Colin Research

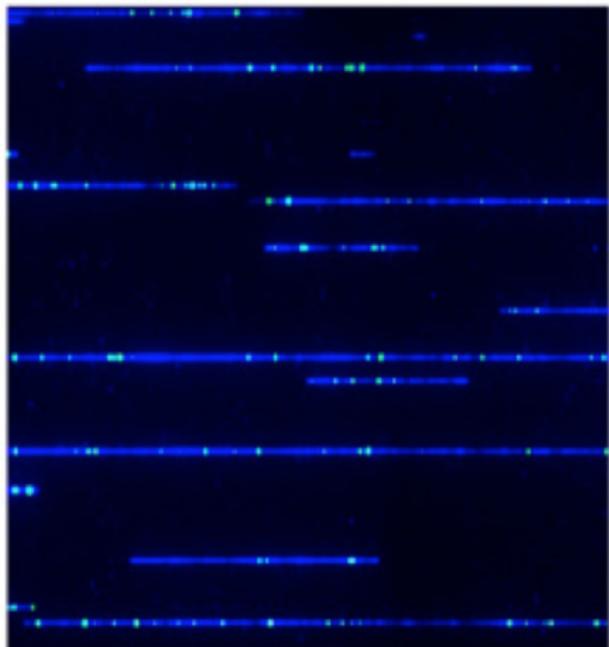


Colin Interests



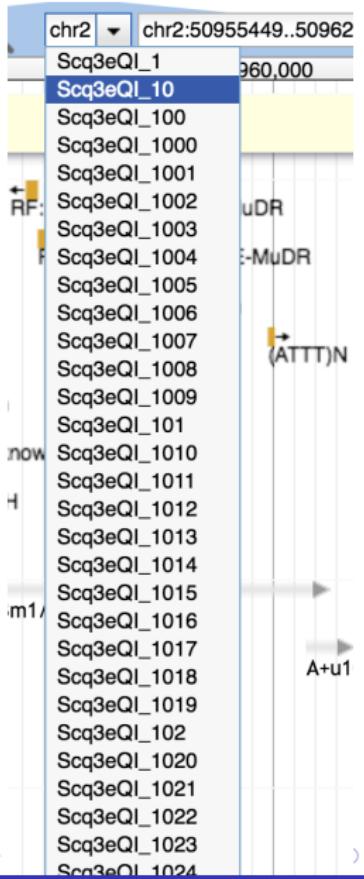
Genome Assembly

- Pull together overlapping sequence
- Consensus sequence between long stretches of highly repetitive DNA sequence is hard to link together
- This issue is complicated by the technology of choice, short read parallel sequencing
- Older methods get around this issue (BAC) as well as emerging technologies (optical mapping / long read sequencing)



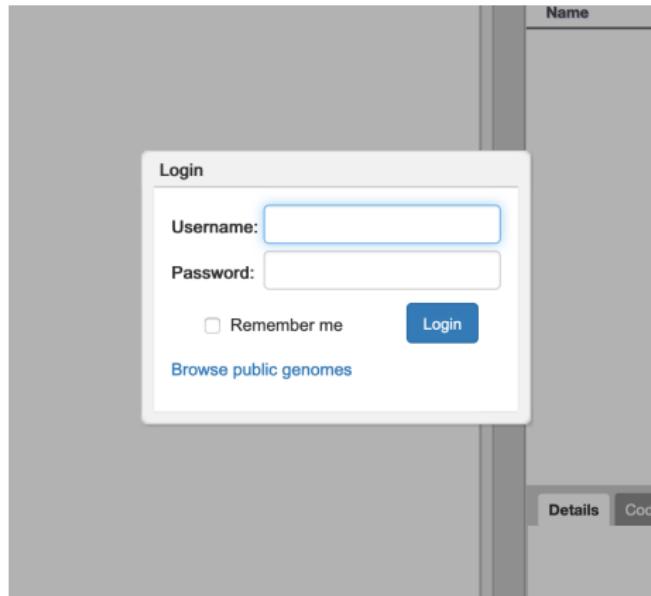
Fucus Genome

- A contig is a set of overlapping DNA sequence that represent a consensus region of DNA
- 110,288 Contigs
- 118579 - 203 base pairs
- Current genome assembly has many more contigs than chromosomes, why?
- 60 (2X) largest contigs randomly assigned to individuals in class roster
- Assigned contigs will vary in quality



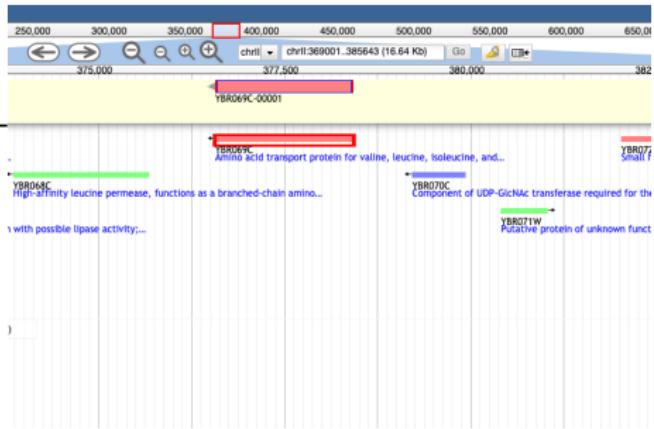
Apollo - Genome Annotation Web Server

- Apollo hosted at 159.89.132.226
- Login = My.UCLA Email
(lowercase)
- Password = UCLA ID Number
- Test Server Load
- 2 Test Data Sets
 - Yeast
 - Volvex (Globe Algae)



Apollo - Yeast Example

- Apollo user manual,
<http://genomearchitect.org/users-guide/>
- Drag feature to user space to begin annotation
- User features are automatically saved



Apollo - Validating Models

- BLAST
 - blast.ncbi.nlm.nih.gov
- Multiple Sequence Alignment

The screenshot shows the NCBI BLAST search interface. At the top, it says "BLAST® > Blastp suite". Below that, the search form is titled "Standard Protein BLAST".
The "Enter Query Sequence" field contains a placeholder "Enter accession number(s), gbk(s), or FASTA sequence(s)".
Below the query field, there are options for "Or, upload file" and "Job Title".
Under "Choose Search Set", the "Database" dropdown is set to "Non-redundant protein sequences (nr)".
The "Program Selection" section shows the "Algorithm" dropdown set to "blastp (protein-protein BLAST)".
At the bottom, there is a large blue "BLAST" button.

BLAST

What happens if the input sequence is a truncated version of the target sequence?

Input

TGATTCACTGTAGTGTGACATATGATAGGAAGGCCATTATC
ATAAAATGGGCAAAGGAGAATTCTCTCTGGCTCTATACA
TTACCCCCAGAACACTCCTGAGATGTGGGAGGATCTGATAC
AGAAAAGCAAAAGATGGAGGTTAGATGTGGTTGAGACCTAT
ATTTTTGGAATGTTCATGAGCCTCTCCTGGCAATTACAA

Target

ATGGAAATCAAACTCAGTTCCAAGTTGTTCTTTTATTTGGT
TTAGCGTGGTTCATTGGTTCTTCCAGCTGATTCACTGTAGT
GTGACATATGATAGGAAGGCCATTATCATAAAATGGGCAAAGG
AGAATTCTCTCTGGCTCTACATTACCCCCAGAACACT
CCTGAGATGTGGGAGGATCTGATACAGAAAAGCAAAAGATGGA
GGTTTAGATGTGGTTGAGACCTATATT

BLAST

Output

100% Match

```
TGATTCACTGTAGTGTGACATATGATAGGAAGGCCATTATC  
ATAAAATGGGCAAAGGAGAATTCTCTTCTGGCTCTATACA  
TTACCCCCAGAACACTCCTGAGATGTGGGAGGATCTGATAC  
AGAAAAGCAAAAGATGGAGGTTTAGATGTGGTTGAGACCTAT  
ATTTTTTGGAAATGTTCATGAGCCTCTCCTGGCAATTACAA
```

Problem?

Multiple Sequence Alignment

- Clustal Omega <https://www.ebi.ac.uk/Tools/msa/clustalo/>
- Protein or Nucleotide sequence

>sequence_label_1

```
TGATTCA GTGTAGTGTGACATATGATAGGAAGGCCATTATC  
ATAAATGGGCAAAGGAGAATTCTCTCTGGCTCTATACA  
TTACCCCCAGAACACTCCTGAGATGTGGGAGGATCTGATAC  
AGAAAGCAAAAGATGGAGGTTAGATGTGGTTGAGACCTAT  
ATTTTTTGGAAATGTTCATGAGCCTTCTCCTGGCAATTACAA
```

>sequence_label_2

```
TGATTCA GTGTAGTGTGACATATGATAGGAAGGCCATTATC  
ATAAATGGGCAAAGGAGAATTCTCTCTGGCTCTATACA  
TTACCCCCAGAACACTCCTGAGATGTGGGAGGATCTGATAC  
AGAAAGCAAAAGATGGAGGTTAGATGTGGTTGAGACCTAT  
ATTTTTTGGAAATGTTCATGAGCCTTCTCCTGGCAATTACAA
```

Clustal Omega Example

>Quercus_lobata

ATGGAAATCAACTCAGTTCCAAGTTGTTCTTTTATTGGTTAGCGT
TGATTCACTGTAGTGTGACAT...CAACCATTAAAGAACTGGCA

>Quercus_suber_beta-galactosidase_3

GACCTTGAGTGTGGAGTCAAAGCACAAATGGACCCTGTAGTAGTCTAC
TAGACCAAAT....CCGACCTTTAACCGGAAATGGGC

Clustal Omega Example

Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

Important note: This tool can align up to 4000 sequences or a maximum file size of 4 MB.

STEP 1 - Enter your input sequences

Enter or paste a set of

DNA

sequences in any supported format:

```
AGTGATGTTTGTAGATGGTAGTACCTCTGGCCAAAAGAAGTGTTCAGCCGTAAGGCCATTGGT  
TCTTAATTAAAGTTGGTAATGGACTAATGAAATATAGTGATTGTAGGTAGATAAGTTGTGAATTCCC  
GGGGTCCCTATAATTGGGGCCTAGGTGAATTTCGAGGCCAGTCCCCGACCTTTAACCGGAATGGC  
ACTGGGAAAAGAGGTGATTCTGTAGTGATTTATATATTATGTTGTATCATTGGGGAAAATC  
AAACCCAGTTGAGTTAGGCTATGTTGGGAGCGGATCTACTATGACCATTGTTATGTGCAGTTCT  
CAATTGAACTCCATTCAATTACATCAAAGTGGCATGTTTTGTC
```

Or, upload a file: [Choose File](#) No file chosen

STEP 2 - Set your parameters

OUTPUT FORMAT

ClustalW with character counts

The default settings will fulfill the needs of most users.

[More options...](#) (Click here, if you want to view or change the default settings.)

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

Clustal Omega Example

[Alignments](#)[Result Summary](#)[Phylogenetic Tree](#)[Submission Details](#)[Download Alignment File](#)[Send to Simple Phylogeny](#)[Send to MView](#)

CLUSTAL O(1.2.4) multiple sequence alignment

Quercus_lobata	-----	0
Quercus_suber_beta-galactosidase_3	GACCTTGAGTGTGGAGTCCAAAGCACAAATGGACCCTGTAGTAGTCTAGTTCTGTC	60
Quercus_lobata	-----	0
Quercus_suber_beta-galactosidase_3	TCTGCACTAGTAGACCAAAACCTCCATTAAAGACCTCTCTCTCTCTCTCTCTCT	120
Quercus_lobata	-----	4
Quercus_suber_beta-galactosidase_3	CTCTCTTTGCATTTCACAAAAGATTGTAACCTTAGTGTACACAGAGAGAACATGG****	180
Quercus_lobata	-----	64
Quercus_suber_beta-galactosidase_3	AAATCAACTCAGTTCCAAGTTCTTTTATTGGTTAGCGTGGTCATGGTTCTAAATGGGCAAATGGTTCT*****	240
Quercus_lobata	-----	124
Quercus_suber_beta-galactosidase_3	TCCAGCTGATTCACTGAGATGTGACATATGATAGGAAGGCCATTATCATAAATGGGCAAATGGTTCT*****	300
Quercus_lobata	-----	184
Quercus_suber_beta-galactosidase_3	GGAGAAATTCTCTCTGGCTCTACATTACCCAGAACGACTCCTGAGATGTGGGAGGGAGAATTCCTCTGGCTCTACATTACCCAGAACGACTCCTGAGATGTGGGAGG*****	360
Quercus_lobata	-----	244
Quercus_suber_beta-galactosidase_3	ATCTGATACAGAAAAGCAAAAGATGGAGGTTAGATGTGGTTGAGACCTATATTGGAGAATTCCTGGCTCTACATTACCCAGAACGACTCCTGAGATGTGGGAGG*****	420

Multiple Sequence Alignment

Solanum_tuberosum(potato)_beta-galactosidase_3	GCAGCTTCTGTCAAACAT-----GACTGGAAATCTGCTGC-TAGAGTAAATGTTCAAT	2610
Arabidopsis_thaliana_beta-galactosidase_3	AGCCCTTTGAAATGAGCACAGAG-----GCACTGGACCGCT-----GGTAGAGATCTCTCAT	1833
Quercus_lobata	-----ATGGAATCAACTCAGTTCCAAGTGTCTT	32
Quercus_suber_beta-galactosidase_3	A-ACCTTGTGTGACACAGAGAGCAATGGAATCACTCAGTTCCAAGTGTCTT	208
	*****	*
Solanum_tuberosum(potato)_beta-galactosidase_3	AACATGCACATAATTTCGCTCCCTGGTCATTAGCACTCTTCTGATGGAGAAATCTA	2670
Arabidopsis_thaliana_beta-galactosidase_3	TGCTTCTCTTATATCTGTGATCTGGATTGTTCTCTCAATCGG-----TAACCTTT	1888
Quercus_lobata	TTTATTGTTAGCGTGGTCTAGGTTCTCAGCTGATTCAG-----TGTAGTGTG	87
Quercus_suber_beta-galactosidase_3	TTTATTGTTAGCTAGTGTTGTCATGGTTCTCAGCTGATTCAG-----TGTAGTGTG	263
	*****	**
Solanum_tuberosum(potato)_beta-galactosidase_3	GTCITCAACACGGCCAAAGTACGTAATACACCTGTGTAGGAACCTCTGTACATTCAATAG	2730
Arabidopsis_thaliana_beta-galactosidase_3	TTATCTGTAATTTGGAAATGCTTACACAG-----GTTACACTGAGTCGGGCCAATG	1943
Quercus_lobata	ACATA----TGATAGGAAGGCCATTATCATAAATGGCAAAAGGAGAATTCTCTCTCTGG	143
Quercus_suber_beta-galactosidase_3	ACATA----TGATAGGAAGGCCATTATCATAAATGGCAAAAGGAGAATTCTCTCTGG	319
	*	*
Solanum_tuberosum(potato)_beta-galactosidase_3	TACTCTAGCTGCTTACGCTTCCTGGAGAAGATATGTTCTGTTATCTGCTACTCGCTAAC	2790
Arabidopsis_thaliana_beta-galactosidase_3	CACCATAGACAGCTTCAAGGCTTCGGCTGGTGTGACG-----TTCTACATACA	1993
Quercus_lobata	CTCTACATACATCCCCAGAACGACTCTCAGAGTGTGGAGGA-----TGTGATACAC	194
Quercus_suber_beta-galactosidase_3	CTCTACATACATCCCCAGAACGACTCTCAGAGTGTGGAGGA-----TGTGATACAC	370
	*	***
Solanum_tuberosum(potato)_beta-galactosidase_3	ACTCTGAA-AGCA--AGGATTTAAAGCTTTAAACCATATGTAGTTGGA-----	2840
Arabidopsis_thaliana_beta-galactosidase_3	AAAAG-----GAGGATCTTCTTAACTACTCATGGTGTCTACACTTACGATTGTTG	2047
Quercus_lobata	GAAAGCAAAGATGGAGTTTAGAGTGTG-----TTGAGACCTATATTTTTGGAGATGTTCA	251
Quercus_suber_beta-galactosidase_3	GAAAGCAAAGATGGAGTTTAGAGTGTG-----TTGAGACCTATATTTTTGGAGATGTTCA	427
	***	***
Solanum_tuberosum(potato)_beta-galactosidase_3	-----GTTCAA-----CATCAAAGATGGAAATGTCACCGACTAACTCAGAGA	2883
Arabidopsis_thaliana_beta-galactosidase_3	TTAACTCATCTAGCTCTTACACACTGTAACGCTTGGTTATCATCGACATCTTAAATA	2107
Quercus_lobata	TGAGCCTTCTCTGGCAATTACATTGAGGGAGATA-----	291
Quercus_suber_beta-galactosidase_3	TGAGCCTTCTCTGGCAATTACATTGAGGGAGATA-----	467
	**	***
Solanum_tuberosum(potato)_beta-galactosidase_3	TGCTATCTGGAGACTTACAGGATGAGATCTGCTAGATAGCTGAGCTCATCAATTA	2943
Arabidopsis_thaliana_beta-galactosidase_3	TAAGATCCATTCTCTGGAGACAGATCTGAGGAACTACT-----TTGAGAAC	2165
Quercus_lobata	---GATTGGTGGAGATTCTTACAGGACCATAGAAAAGCTGGGCTTATGCTCATATTGCC	348
Quercus_suber_beta-galactosidase_3	---GATTGGTGGAGATTCTTACAGGACCATAGAAAAGCTGGGCTTATGCTCATATTGCC	524
	***	***
Solanum_tuberosum(potato)_beta-galactosidase_3	GGTCATTTGGCTCTTGGCAAAATAAGTAACTGAGA-----GATAACAAGTGATTACTTG	3000
Arabidopsis_thaliana_beta-galactosidase_3	AG-----CCGGAGCTTCATTTGGTCAACACAGCTATGATGATGCTCCGA	2212
Quercus_lobata	AT-----TGGACCTTATGTTGTCAGAGTGGATTGGAGATTTCTCTG	395
Quercus_suber_beta-galactosidase_3	AT-----TGGACCTTATGTTGTCAGAGTGGATTGGAGATTTCTCTG	571
	*	***
Solanum_tuberosum(potato)_beta-galactosidase_3	GGTACATAACCAGGTGAGTTGGTTAGTATATAACCTGTTCTCGATTGCTCTC	3060
Arabidopsis_thaliana_beta-galactosidase_3	TCGCA-----GAAATGTGAGTGGCAACAAA-----ACTCTCTA	2254
Quercus_lobata	AT-----AGACGGCGAACAAA-----ACTCTCTA	449
Quercus_suber_beta-galactosidase_3	AT-----TTGGCTCAAGTATGCTCC--AGGCACTCAGT-----TTCAAGACAGCAATGAG-----CCTTTCAA	625
	*	***
Solanum_tuberosum(potato)_beta-galactosidase_3	AGCAGGAAATGGATTGTTGGCAGTCGAAATGTTTATGCTTA-----TATCTTTT-ACT	3116
Arabidopsis_thaliana_beta-galactosidase_3	CACATTCTGAAATAATGAGTGCAATGCTTATGCTAACTAAGTGAGTCGTGTTAAA	2314
Quercus_lobata	GAGGGCAATGCAAAAGGTTCTGACTGAA-----AGATTGTTGAGCTGA-----TGAAGAG	497
Quercus_suber_beta-galactosidase_3	GAGGGCAATGCAAAAGGTTCTGACTGAA-----AGATTGTTGAGCTGA-----TGAAGAG	673

Phylogenetic Tree

Results for job clustalo-l20180412-222345-0717-10122482-pg

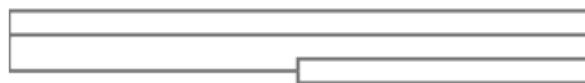
[Alignments](#) [Result Summary](#) [Phylogenetic Tree](#) [Submission Details](#)

Phylogenetic Tree

This is a Neighbour-joining tree without distance corrections.

[Download Phylogenetic Tree Data](#)

Branch length: Cladogram Real



Solanum_tuberosum_potato__beta-galactosidase_3 0.28303
Arabidopsis_thaliana_beta-galactosidase_3 0.24119
Quercus_ilobata 0.00441
Quercus_suber_beta-galactosidase_3 0.0085

Questions?