# Collaborative Genomic Data Analyses in the Cloud

**Steven B. Roberts**
Associate Professor
School of Aquatic and Fishery Sciences
University of Washington

robertslab.info

# Open Science

- You are free to Share!
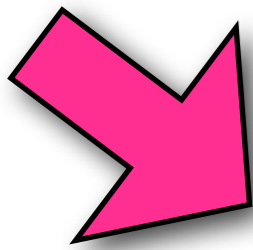
- Our lab practices open notebook science
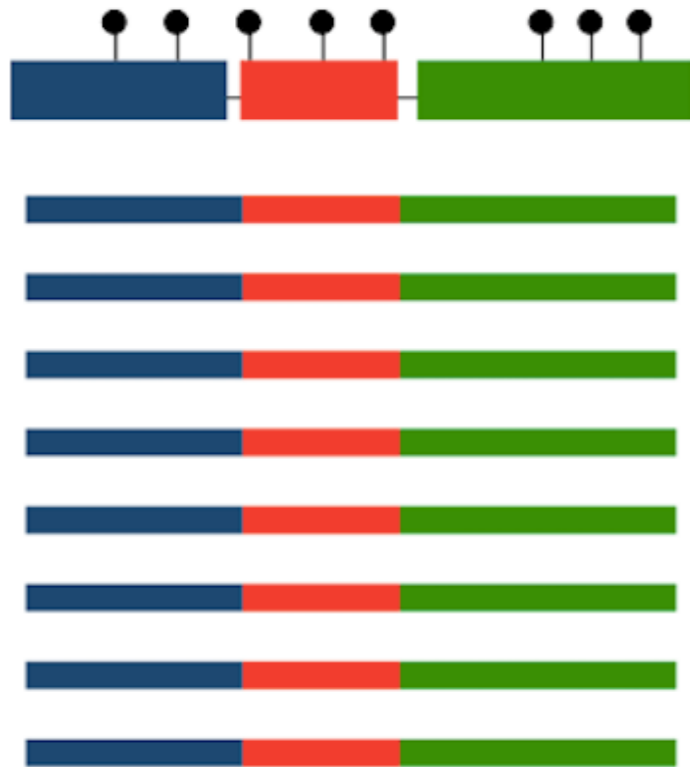
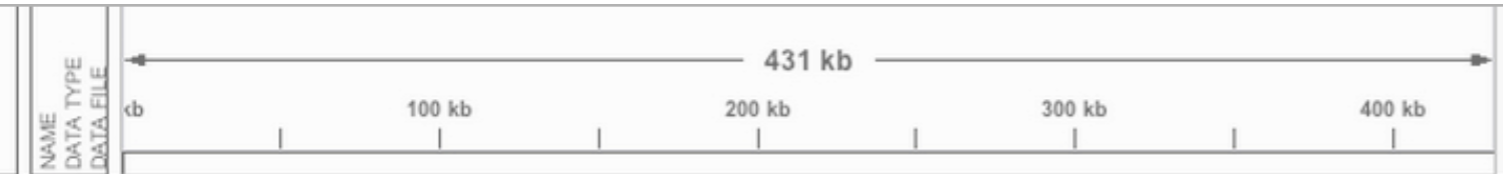robertslab.info          sr320@uw.edu

SQLSHARE

Galaxy

iPlant Collaborative™

Hyak

(GIGA)$^n$ DB    UniProt

NCBI    GO the Gene Ontology    EBI

# Stochastic Variation

431 kb

100 kb    200 kb    300 kb    400 kb

NAME
DATA TYPE
DATA FILE

[0 - 1.00]

[0 - 31]

CCCFFFFFHHHHGIEHEGJIJJHJGJJJJGIIJJJI
@HWI-ST700693:193:C05B7ACXX:3:1101:1342:2141
GTAAGAATTTAGGGTTGTTTTTAAAAATTGAAGTAA
+
BB@FFFFFHHHHHJHIIJJJJJJJJJJJJJJJJGII
@HWI-ST700693:193:C05B7ACXX:3:1101:1458:2156
TTGGTATGTATCATTTATCAATGACAGTATTGTTTT
+
@C@FFFFFHFHFFIIJIJHJJJIJGIGHIJJJHHJI
@HWI-ST700693:193:C05B7ACXX:3:1101:1496:2166
TATATTTAATATTTTGTGTTATTTAGTTTTGTAAAG
+
CCCFFFFFGHGHHJJJGHIJJJJJGIJJJJIJJJJ
@HWI-ST700693:193:C05B7ACXX:3:1
TTTAAAAGAATGTTTTTTTTTTTATAAATAAA
+

CGI_10015536        CGI_10015541

CGI_10015536        CGI_10015541

HWI-ST700693:193:C05B7ACXX:3:2308:17487:198950  272     C135321
011     255     36M     *       0       0       TAAAAAAAAAATTGAA
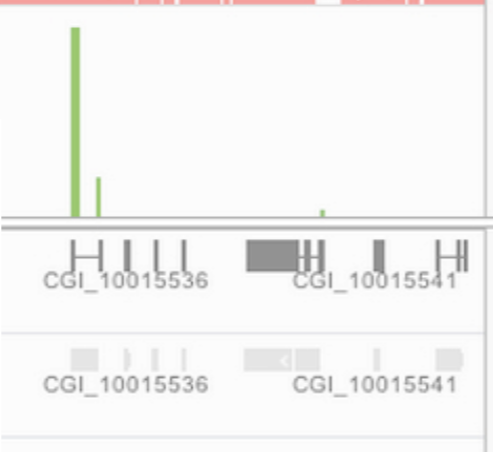AAAATACTACTTAAAATATAA    JJJJJJJJJJJJJJJJJJIHHHIFHHHHFFFFFCCC  N
M:i:1   ZS:Z:-+
HWI-ST700693:193:C05B7ACXX:3:2308:17728:198909  272     scaffol
d124    1441    255     36M     *       0       0       CAAATTT
TAACGAATTTTCATTAAATATATACCAAA   JJIJJJIJJJJJJJJJIJJJJIJJHHHHGFFF
FFCCC   NM:i:0   ZS:Z:-+
HWI-ST700693:193:C05B7ACXX:3:2308:17749:198922  272     C149943
572     255     36M     *       0       0       TACTATATACTCATT
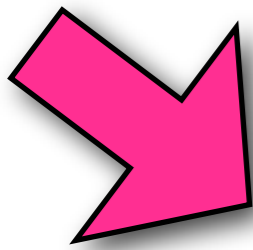ACGTCATTAATTATCAAAAAC   GEJJJJIHFJIJJJIHIIIJJHGHHAHGFFFFFBBB   N
M:i:0   ZS:Z:-+
HWI-ST700693:193:C05B7ACXX:3:2308:17857:198878  256     C149349
54      255     36M     *       0       0       TTAAGTTTGGTTGAA
ATTGGTTTAGTGATTTTGGAG   @@@DDDDFHHDHH@FGHIHHEFHHIHHIIEIGHIDG   N
M:i:0   ZS:Z:++
HWI-ST700693:193:C05B7ACXX:3:2308:18471:198 50  16.     C138363
705     255     36M     *       0       0       ACGCATATTATTTAATCTCCA   HGGGEHHJHHGIHFJIGII
M:i:0   ZS:Z:-+
HWI-ST700693:193:C05B7ACXX:3:2308:18302:198
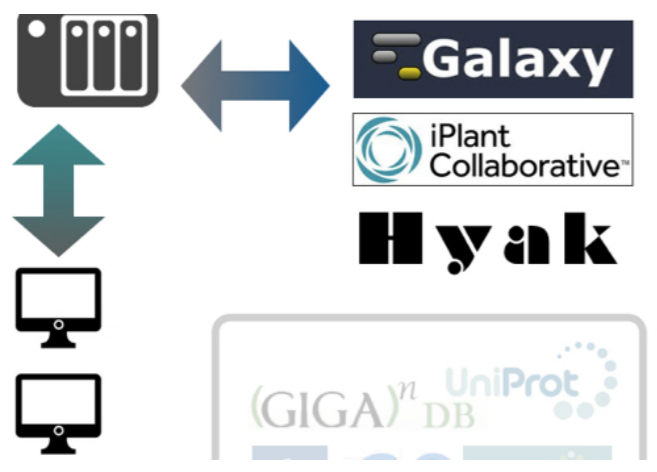
C10093  129     +       TTCAT   0.000   2.00    0       2       0       0       0.000   0.658
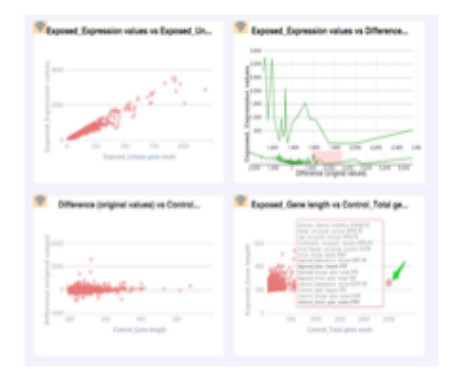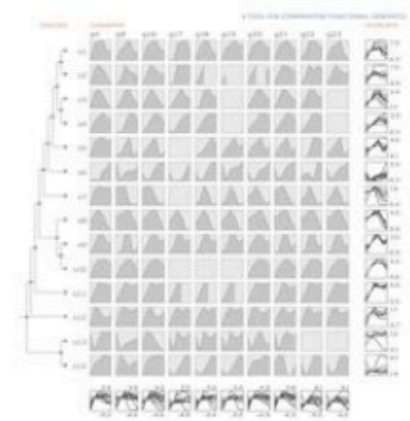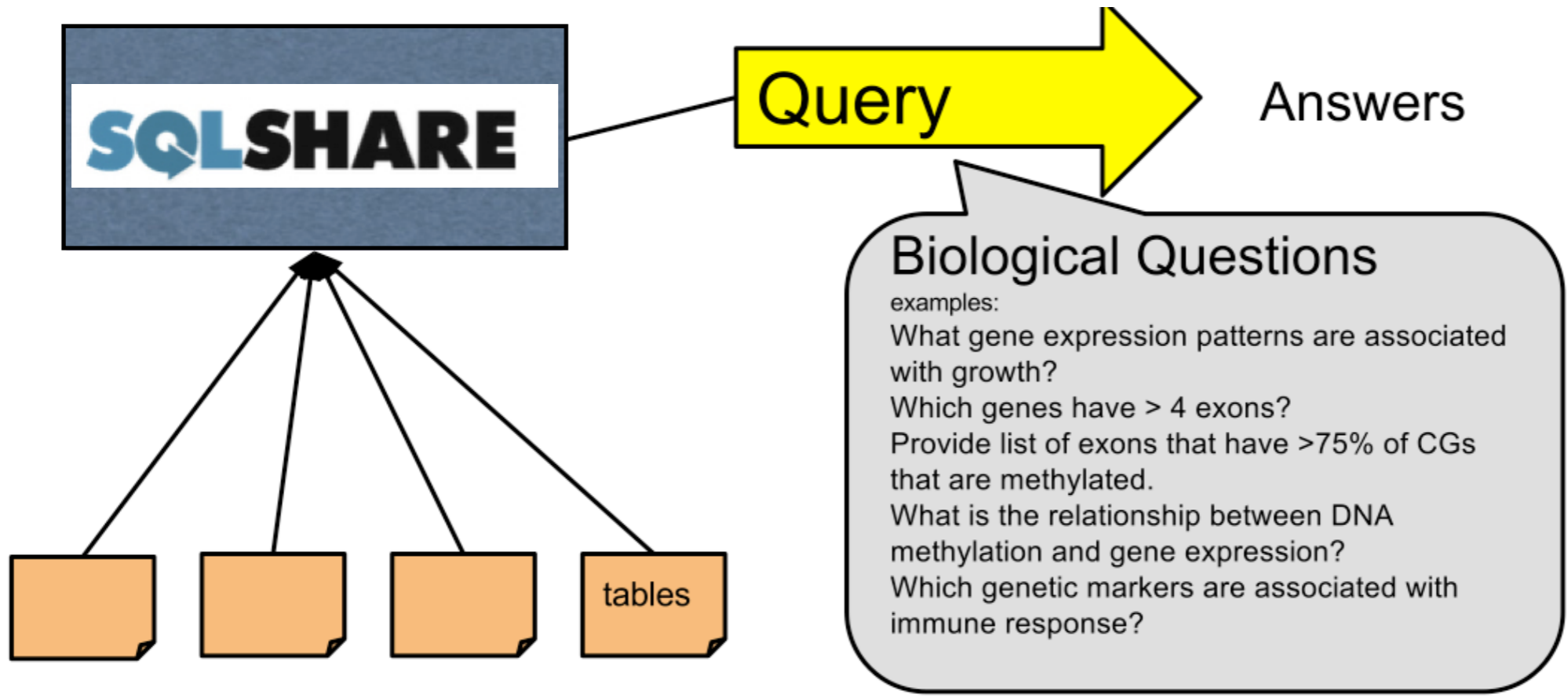C10093  133     +       TTCTC   0.000   2.00    0       2       0       0       0.000   0.658
C10093  135     +       CTCTA   0.000   2.00    0       2       0       0       0.000   0.658
C10093  139     +       AGCTA   0.000   2.00    0       2       0       0       0.000   0.658
C1011   55      -       CAGGC   0.000   1.00    0       1       0       0       0.000   0.793
C1011   73      +       TCCGG   1.000   1.00    1       1       1       1       0.207   1.000
C1011   75      -       CGGCA   0.000   1.00    0       1       0       0       0.000   0.793
C1011   78      -       CAGGG   0.000   1.00    0       1       0       0       0.000   0.793
C1011   79      -       AGGGA   0.000   1.00    0       1       0       0       0.000   0.793
C1011   80      -       GGGAT   0.000   1.00    0       1       0       0       0.000   0.793
C1011   88      -       TTGCT   0.000   1.00    0       1       0       0       0.000   0.793
C10153  106     -       CTGTA   0.000   1.00    0       1       0       0       0.000   0.793
C10153  115     +       TTCGT   0.000   1.00    0       1       1       1       0.000   0.793
C10153  121     -       CAGAT   0.000   1.00    0       1       0       0       0.000   0.793
C10153  124     -       ATGCT   0.000   1.00    0       1       0       0       0.000   0.793
C10153  127     -       CTGTA   0.000   1.00    0       1       0       0       0.000   0.793
C10153  132     -       AAGGA   0.000   1.00    0       1       0       0       0.000   0.793
C10153  133     -       AGGAT   0.000   1.00    0       1       0       0       0.000   0.793
C10195  119     +       AACGG   0.000   1.00    0       1       1       1       0.000   0.793
C10195  121     -       CGGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  123     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  125     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  127     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  129     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  131     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  133     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  135     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  137     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  139     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  141     -       GAGAG   0.000   1.00    0       1       0       0       0.000   0.793
C10195  143     -       GAGAT   0.000   1.00    0       1       0       0       0.000   0.793
C10195  146     -       ATGCA   0.000   1.00    0       1       0       0       0.000   0.793
C10197  33      +       AACAT   0.000   1.00    0       1       0       0       0.000   0.793
C10197  44      +       TTCCA   0.000   1.00    0       1       0       0       0.000   0.793
C10197  45      +       TCCAT   0.000   1.00    0       1       0       0       0.000   0.793
C10197  50      +       TTCAT   0.000   1.00    0       1       0       0       0.000   0.793

Why cloud?

**big, big, compute intensive, education**

raw - 70G

mapping - 60G

tables - 40G ........

**SQLSHARE**

Query → Answers

tables

**Biological Questions**

examples:
What gene expression patterns are associated with growth?
Which genes have > 4 exons?
Provide list of exons that have >75% of CGs that are methylated.
What is the relationship between DNA methylation and gene expression?
Which genetic markers are associated with immune response?

Galaxy

iPlant Collaborative™

Hyak

$(GIGA)^n{}_{DB}$   UniProt

# OYSTERGEN.ES

## eScience Institute

## SQLSHARE

**SQLShare is an easier way to store and share your data.** Get answers to your research questions right now.

**Log in using your account:**

UNIVERSITY of WASHINGTON

Google

**Don't have an account?**

Create a Google Account and start using SQLShare quickly.

...rone can view (and contribute) using the ...ides you with a screenshot of the data in ...en the table in a new webpage.

**Upload**

Upload any tabular data and start analyzing instantly. No need to install, configure, or design a database.

**Modify**

Exercise the full power of SQL even with zero programming experience: joins, subqueries, set operations.

**Share**

Analyze and compare your data collaboratively. Derive new datasets and share them with your colleagues.

https://sqlshare.es...
https://sqlshare.es...
...seq
https://sqlshare.es...
...seq

Logged in

Your datasets

All datasets
Shared datasets
Recent activity... 18
Recently viewed »

Upload dataset
New query

YOUR TOP VIEWED

qDOD Cgigas ...
BiGo_Larvae_j...
TJGR_CCD_d...
BiGill_RNAseq...

**Your Datasets**

Filter dataset by keyword: qdod

| Name | Sharing / Owner | Modi |
|---|---|---|
| qdod_proteome_blast_mouse | 🔒 sr320@washington.edu | Jan 2 |
| qDOD_v9_gene   GFF format file of oyster genes ~28k  gene | sr320@washington.edu | Nov |
| _qdod_goslim_graphtest | sr320@washington.edu | Oct |
| SNP_RNAseqLibary_SB_BiGill   SNP table from RNA-seq library - SB gill tissue pool (BiGill complement)  qdod2 | sr320@washington.edu | Oct |
| BiGill_meth_Zhang_exp   Gene-centric data including length, CG, percent methylation (gill) and tissue specific RPKM data from Zhang et a  qdod2 | sr320@washington.edu | Oct |
| qDOD_Cgigas_gene_fasta   Tabular format of Cgigas gene sequence fasta file Derived using Dataset: Genomic data from the Pacific oyst  qdod2 | sr320@washington.edu | Oct |
| qDOD Cgigas Gene Descriptions (Swiss-prot)   Description and evalues associated with the ~28k genes Derived us  blast | sr320@washington.edu | Oct |

*Use Cases*
• **Joining on Annotations**
• **File Conversion**
• **Querying Gene Tables**

Secondary stress: proteomics
Original input file had some peptides of charge state >2, so had to redo everything with fixed input file.
SR discovered that for some proteins, a peptide was sequenced multiple times and so had multiple expression values. From the unique protein associations file in SQLshare, I summed the expression values for all identical peptides.
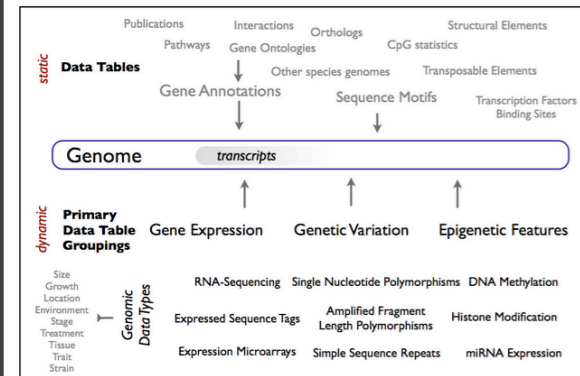
```
SELECT [peptide sequence], SUM([2_01 TotalArea]) AS CG2_01, SUM([2_02 TotalArea]) AS CG2_02, SUM([2_03
TotalArea]) AS CG2_03, SUM([5_01 TotalArea]) AS CG5_01, SUM([5_02 TotalArea]) AS CG5_02, SUM([5_03
TotalArea]) AS CG5_03, SUM([8_01 TotalArea]) AS CG8_01, SUM([8_02 TotalArea]) AS CG8_02, SUM([8_03
TotalArea]) AS CG8_03, SUM([11_01 TotalArea]) AS CG11_01, SUM([11_02 TotalArea]) AS CG11_02, SUM([11_03
TotalArea]) AS CG11_03, SUM([26_01 TotalArea]) AS CG26_01, SUM([26_02 TotalArea]) AS CG26_02, SUM([26_03
TotalArea]) AS CG26_03, SUM([29_01 TotalArea]) AS CG29_01, SUM([29_02 TotalArea]) AS CG29_02, SUM([29_03
TotalArea]) AS CG29_03, SUM([32_01 TotalArea]) AS CG32_01, SUM([32_02 TotalArea]) AS CG32_02, SUM([32_03
TotalArea]) AS CG32_03, SUM([35_01 TotalArea]) AS CG35_01, SUM([35_02 TotalArea]) AS CG35_02, SUM([35_03
TotalArea]) AS CG35_03, SUM([221_01 TotalArea]) AS CG221_01, SUM([221_02 TotalArea]) AS CG221_02,
SUM([221_03 TotalArea]) AS CG221_03, SUM([224_01 TotalArea]) AS CG224_01, SUM([224_02 TotalArea]) AS
CG224_02, SUM([224_03 TotalArea]) AS CG224_03, SUM([227_01 TotalArea]) AS CG227_01, SUM([227_02
TotalArea]) AS CG227_02, SUM([227_03 TotalArea]) AS CG227_03, SUM([230_01 TotalArea]) AS CG230_01,
SUM([230_02 TotalArea]) AS CG230_02, SUM([230_03 TotalArea]) AS CG230_03,
SUM([242_01 TotalArea]) AS CG242_01, SUM([242_02 TotalArea]) AS CG242_02, SUM([242_03 TotalArea]) AS
CG242_03, SUM([245_01 TotalArea]) AS CG245_01, SUM([245_02 TotalArea]) AS CG245_02, SUM([245_03
TotalArea]) AS CG245_03, SUM([248_01 TotalArea]) AS CG248_01, SUM([248_02 TotalArea]) AS CG248_02,
SUM([248_03 TotalArea]) AS CG248_03, SUM([251_01 TotalArea]) AS CG251_01, SUM([251_02 TotalArea]) AS
CG251_02, SUM([251_03 TotalArea]) AS CG251_03
```

20@washington.edu   Aug
20@washington.edu   Aug
20@washington.edu   Aug
20@washington.edu   Jul
20@washington.edu   Jul
20@washington.edu   Jul
20@washington.edu   Jul
20@washington.edu   Jun
20@washington.edu   Jun
20@washington.edu   May
20@washington.edu   May
20@washington.edu   May
20@washington.edu   May
20@washington.edu

# SQLSHARE

Logged in:

**Your datasets**
All datasets
Shared datasets
Recent activity...  18
Recently viewed »

Upload dataset
New query

YOUR TOP VIEWED

qDOD Cgigas ...
BiGo_Larvae_j...
TJGR_CCD_d...
BiGill_RNAseq...

## Your Datasets

| Name |
| --- |
| qdod_proteome_blast_mouse |
| qDOD_v9_gene  GFF format file of oyster genes ~28k |
| gene |
| _qdod_goslim_graphtest |
| SNP_RNAseqLibary_SB_BiGill  SNP table from RNA-seq library - SB gill tissue pool (BiGill complement) |
| qdod2 |
| BiGill_meth_Zhang_exp  Gene-centric data including length, CG, percent methylation (gill) and tissue specific RPKM data from Zhang et a |
| qdod2 |
| qDOD_Cgigas_gene_fasta  Tabular format of Cgigas gene sequence fasta file Derived using Dataset: Genomic data from the Pacific oyste |
| qdod2 |
| qDOD Cgigas Gene Descriptions (Swiss-prot)  Description and evalues associated with Cgigas 28k genes Derived using Dataset: Genomi |
| blast |

```sql
SELECT cgslim.CGI_ID,Description,evalue,SPID,GOID,term,GOSlim_bin,sequence
  FROM [sr320@washington.edu].[qDOD_Cgigas_GO_GOslim] cgslim

LEFT JOIN [sr320@washington.edu].[qDOD_Cgigas_gene_fasta] cgf
  on cgslim.CGI_ID = cgf.CGI_ID

Where term LIKE '%methyl%'
OR
term LIKE '%histone%'
```

sr320@washington.edu    Oct 2
sr320@washington.edu    Oct 2
sr320@washington.edu    Oct 2
sr320@washington.edu    Oct 2

## Secondary stress: proteomics

Original input file had some peptides of charge state >2, so had to redo everything with fixed input file.

SR discovered that for some proteins, a peptide was sequenced multiple times and so had multiple expression values. From the unique protein associations file in SQLshare, I summed the expression values for all identical peptides.

```sql
SELECT [peptide sequence], SUM([2_01 TotalArea]) AS CG2_01, SUM([2_02 TotalArea]) AS CG2_02, SUM([2_03
TotalArea]) AS CG2_03, SUM([5_01 TotalArea]) AS CG5_01, SUM([5_02 TotalArea]) AS CG5_02, SUM([5_03
TotalArea]) AS CG5_03, SUM([8_01 TotalArea]) AS CG8_01, SUM([8_02 TotalArea]) AS CG8_02, SUM([8_03
TotalArea]) AS CG8_03, SUM([11_01 TotalArea]) AS CG11_01, SUM([11_02 TotalArea]) AS CG11_02, SUM([11_03
TotalArea]) AS CG11_03, SUM([26_01 TotalArea]) AS CG26_01, SUM([26_02 TotalArea]) AS CG26_02, SUM([26_03
TotalArea]) AS CG26_03, SUM([29_01 TotalArea]) AS CG29_01, SUM([29_02 TotalArea]) AS CG29_02, SUM([29_03
TotalArea]) AS CG29_03, SUM([32_01 TotalArea]) AS CG32_01, SUM([32_02 TotalArea]) AS CG32_02, SUM([32_03
TotalArea]) AS CG32_03, SUM([35_01 TotalArea]) AS CG35_01, SUM([35_02 TotalArea]) AS CG35_02, SUM([35_03
TotalArea]) AS CG35_03, SUM([221_01 TotalArea]) AS CG221_01, SUM([221_02 TotalArea]) AS CG221_02,
SUM([221_03 TotalArea]) AS CG221_03, SUM([224_01 TotalArea]) AS CG224_01, SUM([224_02 TotalArea]) AS
CG224_02, SUM([224_03 TotalArea]) AS CG224_03, SUM([227_01 TotalArea]) AS CG227_01, SUM([227_02
TotalArea]) AS CG227_02, SUM([227_03 TotalArea]) AS CG227_03, SUM([230_01 TotalArea]) AS CG230_01,
SUM([230_02 TotalArea]) AS CG230_02, SUM([230_03 TotalArea]) AS CG230_03,
SUM([242_01 TotalArea]) AS CG242_01, SUM([242_02 TotalArea]) AS CG242_02, SUM([242_03 TotalArea]) AS
CG242_03, SUM([245_01 TotalArea]) AS CG245_01, SUM([245_02 TotalArea]) AS CG245_02, SUM([245_03
TotalArea]) AS CG245_03, SUM([248_01 TotalArea]) AS CG248_01, SUM([248_02 TotalArea]) AS CG248_02,
SUM([248_03 TotalArea]) AS CG248_03, SUM([251_01 TotalArea]) AS CG251_01, SUM([251_02 TotalArea]) AS
CG251_02, SUM([251_03 TotalArea]) AS CG251_03
```

l20@washington.edu    Aug
l20@washington.edu    Aug
l20@washington.edu    Aug
l20@washington.edu    Jul 2
l20@washington.edu    Jul 2
l20@washington.edu    Jul 1
l20@washington.edu    Jul 1
l20@washington.edu    Jun 2
l20@washington.edu    Jun 2
l20@washington.edu    May
l20@washington.edu    May
l20@washington.edu    May
l20@washington.edu    May
l20@washington.edu

## qdod: Querying Disparate Oyster Datasets

This respository provides access to genomic data and workflows (IPython notebooks) that are being integrated as part of effort to increase effeciency of biological discovery. The wiki associated with this repository will serve as the *primary means for documentation*. Currently the documentation is focused on 1) **describing current datasets** and 2) **providing workflow tutorials**.

In brief, data in the form of delimited text files is aggregated into SQLShare where they can be easily queried. Below is schematic representation of the different types of datasets.



During the initial phases the focus is on the Pacific oyster and primary data from the Roberts Lab.

### Select IPython Notebooks

- Static Data Tables - Universal
- Static Data Tables - Annotations

---

Currently the documentation is focused on 1) **describing current datasets** and 2) **providing workflow tutorials**.

### A. Raw Data

- Select NGS Data via Roberts Lab

### B. Datasets in SQLShare

- Universal
- Generic Oyster Datasets
- Tissue Specific Oyster Datasets

### C. Tutorials

- Simple Gene Search
- Standard SQLShare Queries
- Annotating Genes
- File Format Conversions

### D. Genome Browser Feature Tracks

- Canonical Tracks
- Bisulfite sequencing (gill tissue)
- Reference Genome Files

Please use GitHub's Issue feature to ask question, report problems, or suggest features.

Last edited by sr320, 9 days ago

github.com/sr320/qdod/wiki

# Sharing
# Collaboration*

Open Notebook Science

SQLSHARE
Galaxy
iPlant Collaborative™
Hyak

Reproducible

Collaboration          Open

# Open Notebook Science

... there is a URL to a laboratory notebook that is freely available and indexed on common search engines. It does not necessarily have to look like a paper notebook but it is essential that all of the information available to the researchers to make their conclusions is equally available to the rest of the world.

—Jean-Claude Bradley

# Open Notebook Science



genefish.wikispaces.com

# Open Notebook Science



**Automating a Workflow: Beyond Blast - to GO Slim**

The concept is that you can take a fasta file in a working directory and end up with GO slim information all within a single notebook that is automated. Currently this work by writing (and overwriting) as scratch file to SQLShare. Assumptions are that you are working in a directory with fasta file named `query.fa`. And blast algorithms are in PATH.

```
In [13]: #allows plots to be shown inline
         %pylab inline

         Populating the interactive namespace from numpy and matplotlib
```

```
In [4]: #Setting Working Directory
        wd="/Volumes/web/whale/fish546/qpx_go_val"
        #Setting directory of Blast Databases
        dbd="/Volumes/Bay3/Software/ncbi-blast-2.2.29\+/db/"
        #Database name
        dbn="uniprot_sprot_r2013_12"
        #Blast algorithim
        ba="blastx"
        #Location of SQLShare python tools: you can empty ("") if tools are in PATH
        spd="/Users/sr320/sqlshare-pythonclient/tools/"
```

```
In [5]: cd {wd}

        /Volumes/web/whale/fish546/qpx_go_val
```

```
In [5]: !{ba} -query query.fa -db {dbd}{dbn} -out {dbn}_{ba}_out.tab -evalue 1E-50 -num_threads 4 -max_hsps_per_subject 1 -

        BLAST Database error: No alias or index file found for protein database [/Volumes/Bay3/Software/ncbi-blast-2.2.29+
        /db/uniprot_sprot_r2013_12] in search path [/Volumes/web/whale/fish546/pipeline_test_dir4::]
```

```
In [6]: !head -1 {dbn}_{ba}_out.tab

        QPX_transcriptome_v1_Contig_2    sp|P52712|CBPX_ORYSJ    43.75    416    213    12    2095    869    6    40
        7        3e-98    326
```

```
In [17]: #Translate pipes to tab so SPID is in separate column for Joining
         !tr '|' "\t" <{dbn}_{ba}_out.tab> {dbn}_{ba}_out2.tab
```

```
In [18]: !head -1 {dbn}_{ba}_out2.tab
```

```
In [8]: #Uploads formatted blast table to SQLshare; currently has generic name and meant to be temporary: Warning will over
        !python {spd}singleupload.py -d scratchblast_out {dbn}_{ba}_out2.tab
        ...
```

```
In [9]: !python {spd}fetchdata.py -s "SELECT * FROM [sr320@washington.edu].[scratchblast_out]blast Left Join [sr320@washing
```

```
In [10]: !head -2 {dbn}_join2goslim.txt
         ...
```

```
In [11]: !python {spd}singleupload.py -d scratchjoin_slim {dbn}_join2goslim.txt

         processing chunk line 0 to 18037 (0.0718240737915 s elapsed)
         pushing uniprot_sprot_r2013_12_join2goslim.txt...
         parsing 9A18D989...
         finished scratchjoin_slim
```

```
In [12]: #Sets GO aspect
         !python {spd}fetchdata.py -s "SELECT Distinct Column1 as query, Column3 as SPID, GOSlim_bin FROM [sr320@washington.
```

```
In [13]: !head justslim.txt
         ...
```

```
In [15]: from pandas import *
```

Set some variables

blast

convert file format

upload to SQLShare
(python client)

join in SQLShare -
download

read in pandas

matplotlib generates
graph of GOsllim

Open Notebook Science

*Comparison*

Wiki - collaboration, versioning, search, publishing
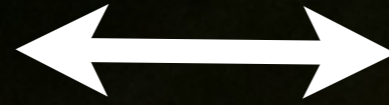
Evernote - simple, multi-platform

IPython - executable, versioning*

**no perfection solution**

Challenges: versioning, provenance, collaboration, simple sharing, discoverability

# Acknowledgements

Mackenzie Gavery
Claire Olson

Sam White
Brent Vadopalas
Jake Heare

Bill Howe
Dan Halperin

*DNA methylation*

Aquaculture Program

eScience Institute

robertslab.info        sr320@uw.edu