

# DolphinNext: A graphical user interface for reproducible pipelines

Onur Yukselen, Artur Manukyan,  
Alper Kucukural  
May 6, 2021

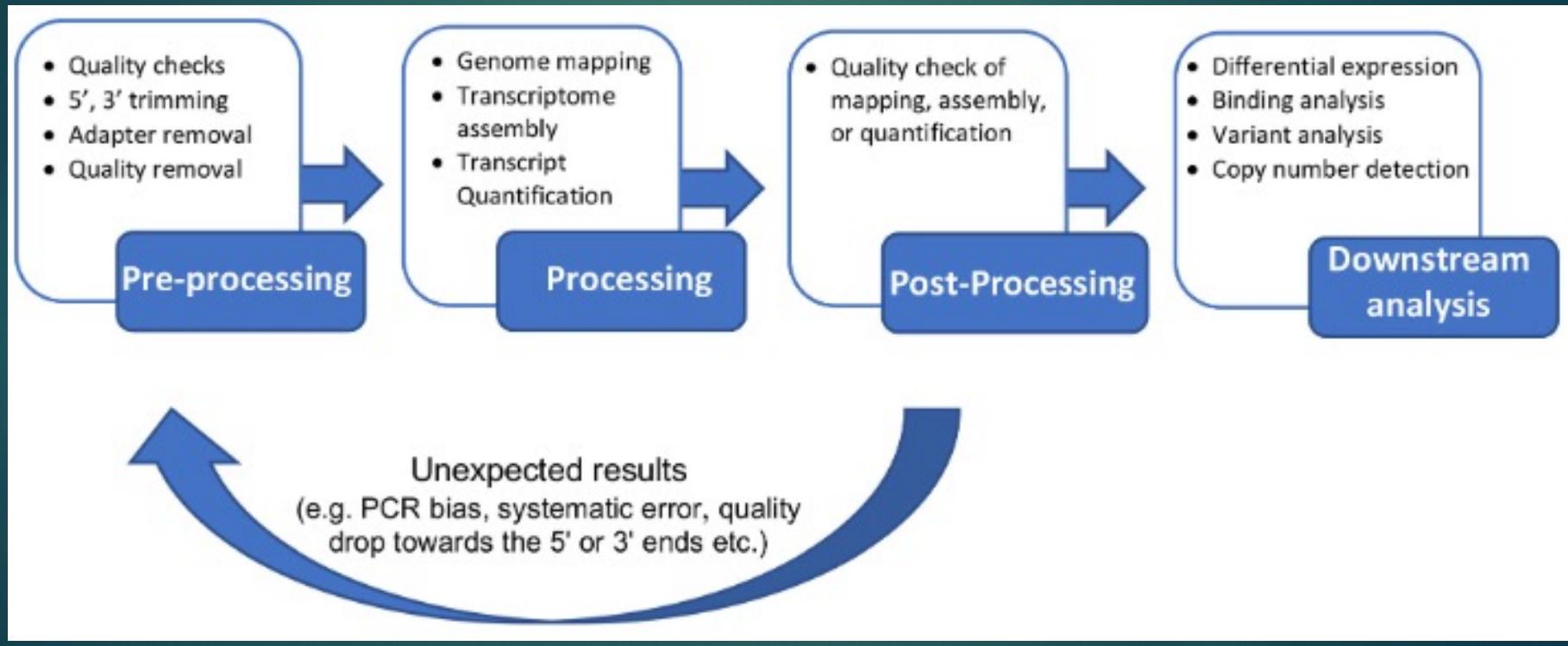


**nextflow**

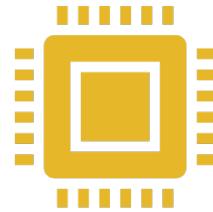


University of  
Massachusetts  
**Medical School**

# Sequence analysis: mostly reusable steps



# Large projects require

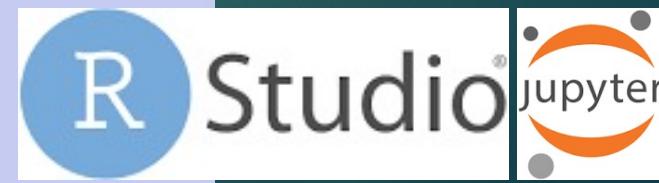
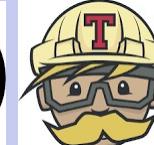
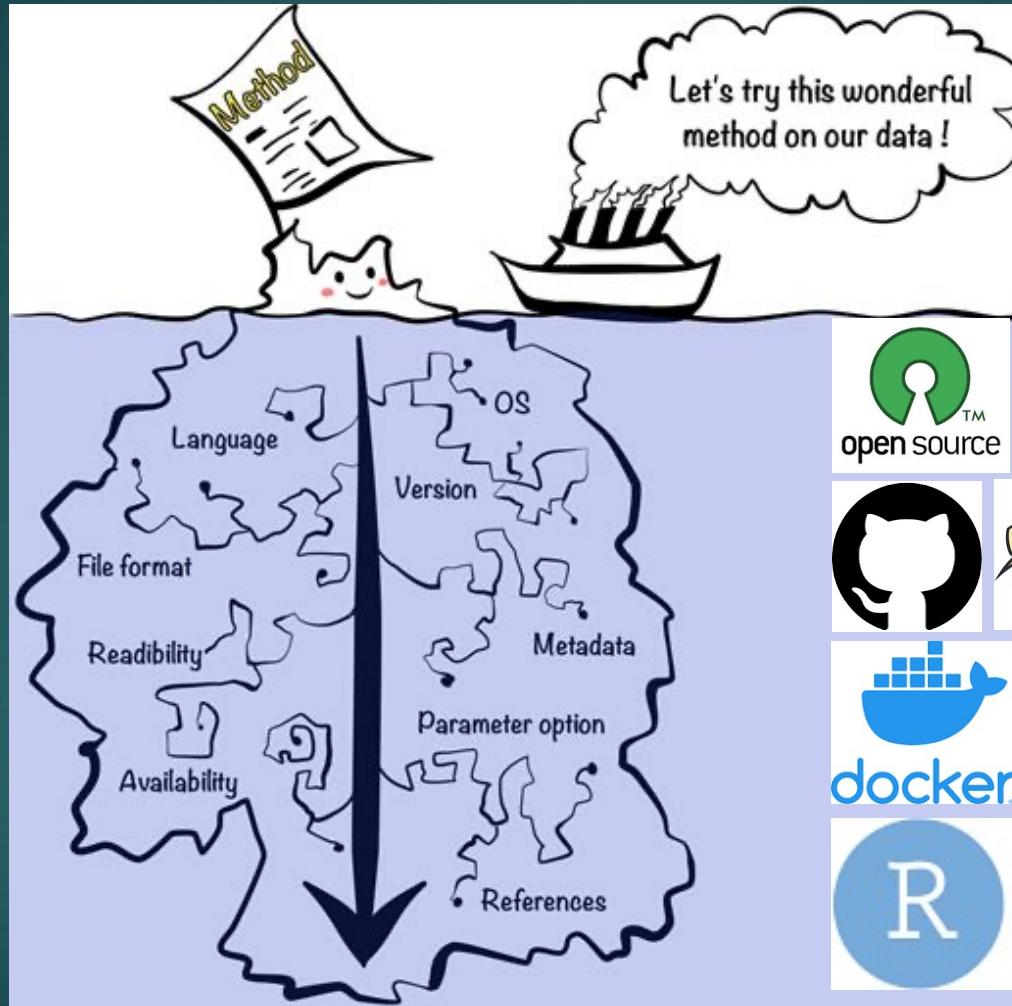


## Capture sample metadata

- Identify and correct batches
- Troubleshoot
- Drive processing

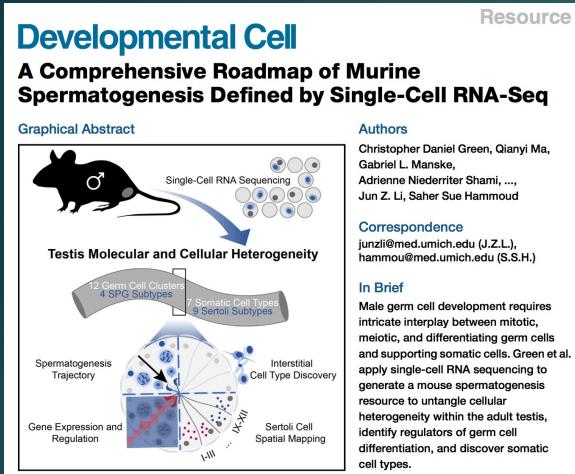
## Pipelines should be

- Highly parallel
- Self monitoring → restart on hardware failures
  - Easily restartable from mid points
- Built from reusable components
- Versioned
- Handle large number of samples



# Replicating a Method

## Dependency issues

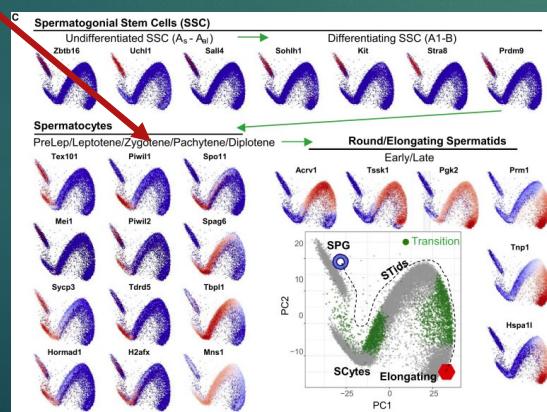


```

1  ## R script for pseudotemporal ordering of germ cells using Monocle in Nov 2017 by Qianyi
2  ### Related to Figure S2C right panel: germ cells (N=20,646) with >1k detected genes
3
4  ###
5  ### Read Data
6  home="/scratch/junzli_flux/qzm/Dropseq_analysis/"
7  file="figJul2017_10CellTypes_MouseAdultST24mergedclusters/monocle_"
8  load(file =paste0(home,"data_DGE/MouseAdultST24genesUMI20cell15.Robj"))
9  dge20=dge # Gene Filter 1 (N=24482): genes >20totalUMIsOverAllCells & >15Cells +31 genes
10 table(dge@data.info[,52:53])
11 # Cluster31Seriation0L0_GeneFilter1 is the reordered cluster ID 1-31
12 # Cluster 8-31 are germ cells

```

## Hardcoded Scripts



Software and Algorithms	
Drop-seq_tools (v1.12)	
Picard Tools (v2.6.0)	
Samtools (v1.2)	
STAR (v2.5.2b)	
R (v3.3.3)	
Seurat (v1.4.0.3)	
Seriation (v1.2-2)	
Monocle 2	



## R Markdown Code for Reproducing Clustering Analysis

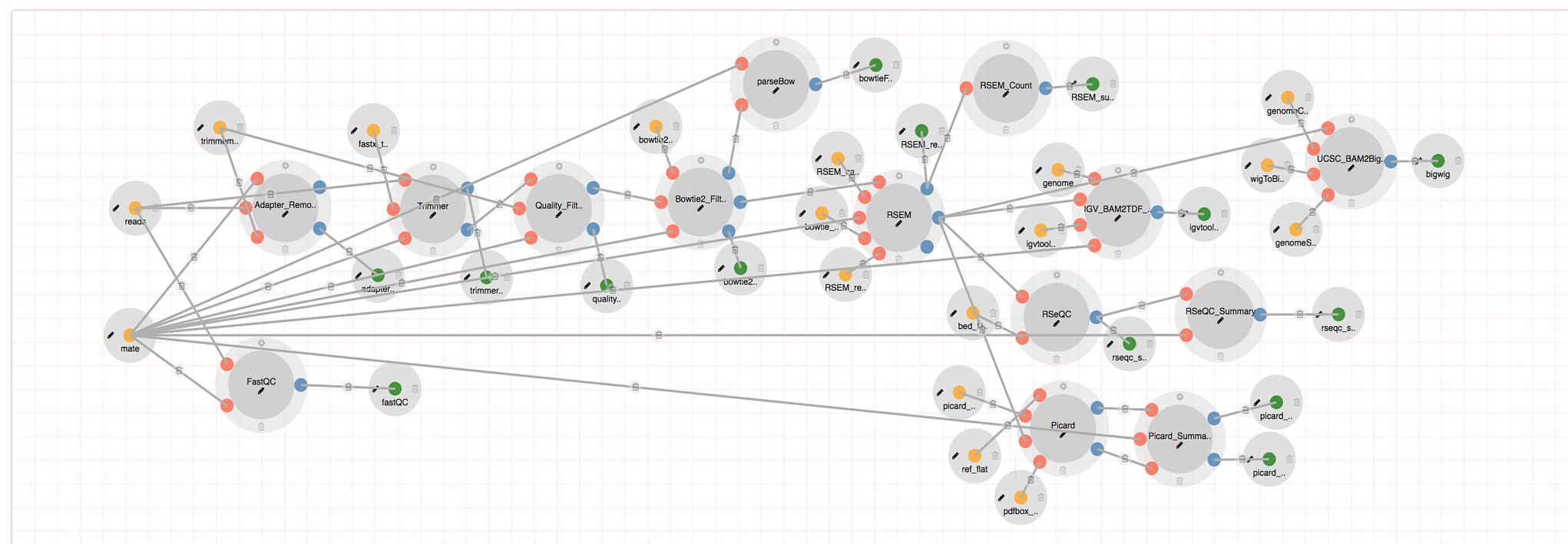
As an accompaniment to this paper, we provide an R markdown file that describes step-by-step procedures, including the loading and preprocessing of the Drop-seq digital expression matrix, PCA, Louvain-Jaccard clustering, data visualization, ordering by seriation, and differential expression tests. The R commands are provided at [https://github.com/qianqianshao/Drop-seq\\_ST](https://github.com/qianqianshao/Drop-seq_ST).

Created by onuryuksele on 2018-09-10 17:37:33 • Last edited on 2018-10-11 15:24:00

## Description

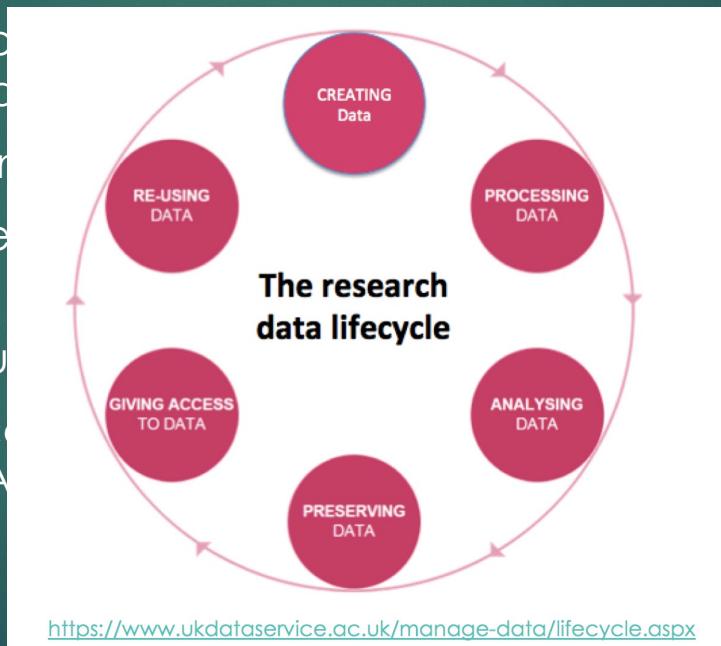
RSEM pipeline includes Quality Control, rRNA filtering, Genome Alignment using Bowtie, and estimating gene and isoform expression levels by RSEM.

## Steps:



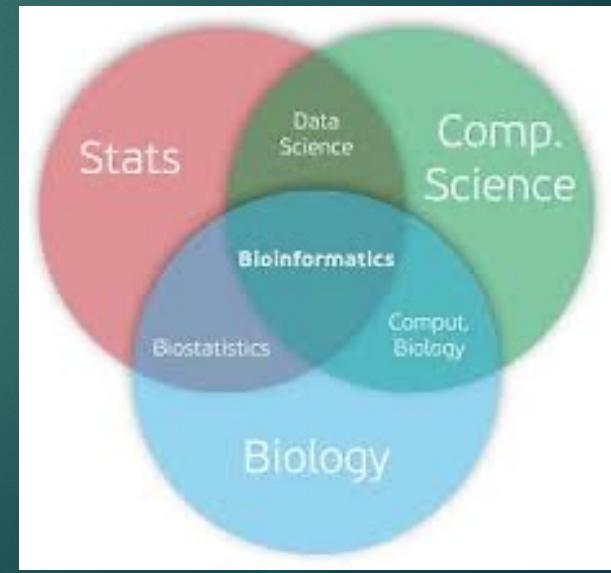
# Best practices for reproducible data processing

- ▶ For every result, keep how it was produced
- ▶ Avoid manual data reentry
- ▶ Archive the exact version of all programs used
- ▶ Version control all code
- ▶ Provide (public) access to code, runs, and results in FAIR ways

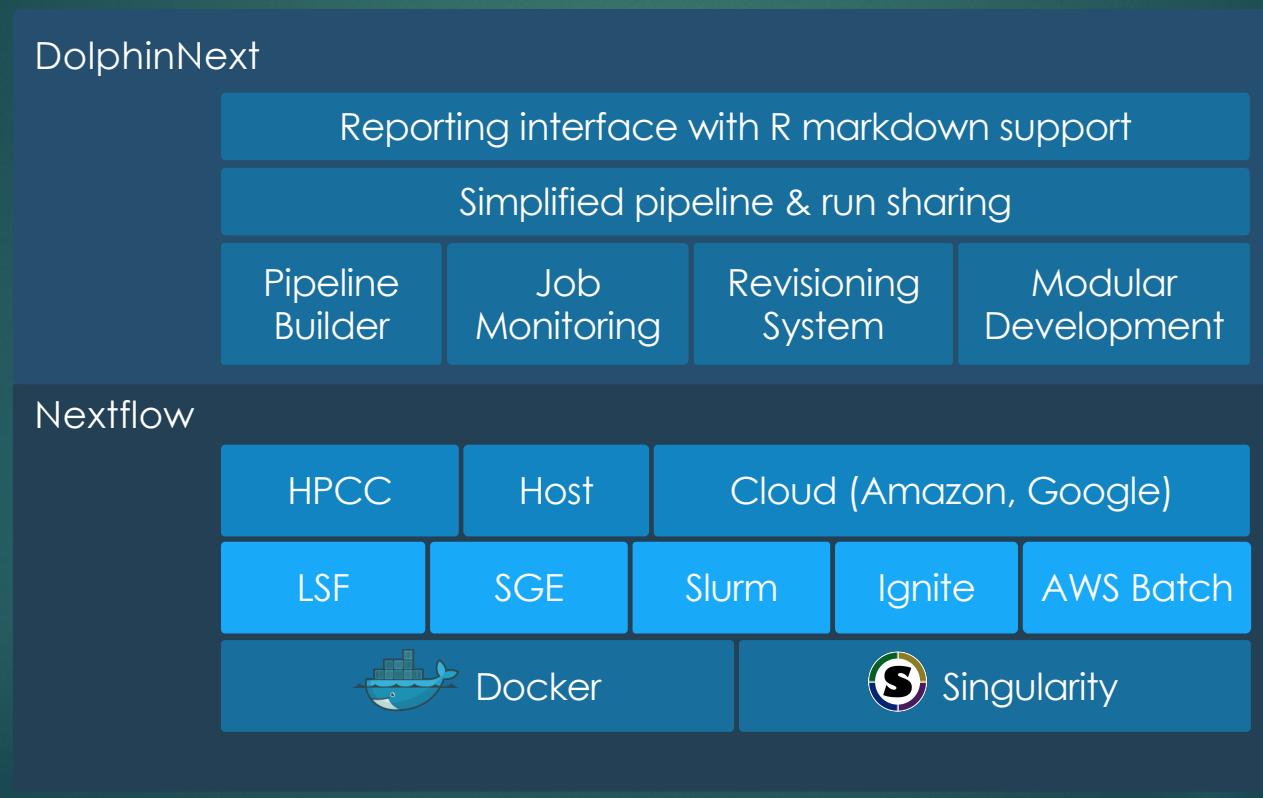


- ▶ DolphinNext is an easy-to-use web platform for creating, deploying, and executing complex pipelines for high throughput data processing and analysis.

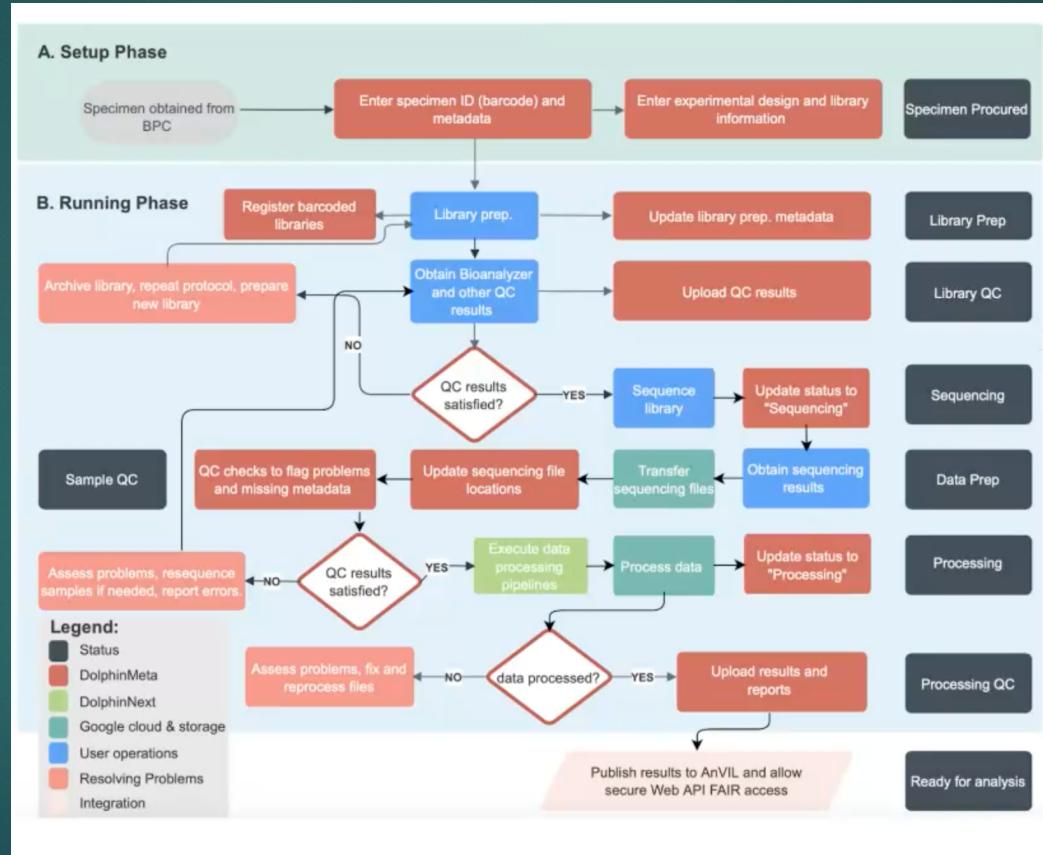
For all skill levels and areas to perform bioinformatics analysis



# Structure

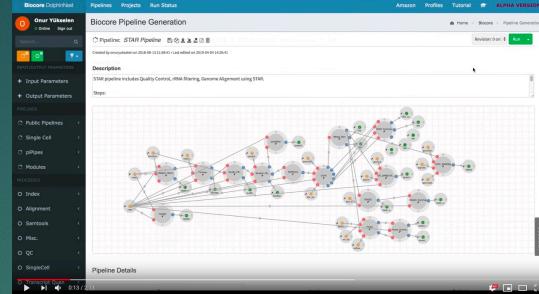


# DolphinLIMS





D-meta and D-portal



# Dnext



The screenshot shows the 'D-Meta' dashboard. At the top, there are tabs for 'Dashboard', 'Admin', and 'Import Data'. The main area has a search bar and a table titled 'Samples'. The table columns include: ID, DEB, Assumptions, Name, Contract, Disease Location, Number of Meters, Volume (ml) (mls), Cell Density (E03) (cells/mls), Total Cells, % Live Cells, and Status. There are 10 rows of sample data.

ID	DEB	Assumptions	Name	Contract	Disease Location	Number of Meters	Volume (ml) (mls)	Cell Density (E03) (cells/mls)	Total Cells	% Live Cells	Status
1	C000123456789		C000123456789			-	-	-	-	-	-
2	V9997123456789		V9997123456789			-	-	-	-	-	-
3	V9997123456789		V9997123456789			-	-	-	-	-	-
4	V9997123456789		V9997123456789			-	-	10000	-	-	-
5	V9997123456789		V9997123456789			-	-	10000	-	-	-
6	V9997123456789		V9997123456789			-	-	50000	-	-	-
7	C000123456789		C000123456789			-	-	50000	-	-	-
8	C000123456789		C000123456789			-	-	50000	-	-	-
9	V9997123456789		V9997123456789			-	-	20000	-	-	-
10	V9997123456789		V9997123456789			-	-	20000	-	-	-

# Dmeta



The screenshot shows the 'D-Portal' dashboard. At the top, there are tabs for 'Dashboard', 'Admin', and 'Import Data'. The main area includes three bar charts: 'Experiment Series', 'Experiments', and 'Samples'. Below the charts is a table titled 'Properties' with columns for Name, Status, Experiment, Patient, Aliquot, Clinical phenotype, and Skin. There is also a 'Patient Note' section.

Name	Status	Experiment	Patient	Aliquot	Clinical phenotype	Skin	Patient Note
C000123456789	Processed	V9997123456789	C0001	L2	Local Erythema	Loose	fat skin (2nd skin on SLE, rare)
V9997123456789	Processed	V9997123456789	C0001	L2	Local	Loose	-
V9997123456789	Processed	V9997123456789	C0002	N1	Non-Losomal	-	-
V9997123456789	Processed	V9997123456789	C0002	L2	Local	Loose	-
V9997123456789	Processed	V9997123456789	C0004	N1	Non-Losomal	-	-
V9997123456789	Processed	V9997123456789	C0005	N1	Non-Losomal	-	-
C000123456789	Processed	V9997123456789	C0007	L2	Local Erythema	Loose	fat skin (2nd skin on SLE, rare)

# Dportal



# Use Case

D-Portal Dashboard v0.0.22

1 Experiment Series | 1 Experiments | 102 Samples

Search: □

Properties

- + Status
- + Experiment series
- + Experiment
- + Patient
- + Disease
- + Clinical phenotype
- + Skin
- + Patient Note
- + Site
- + Project
- + Owner

Portal update

D-Meta Dashboard Admin Import Page v0.0.39

ID	Specimen ID	Name	Contract	Biopsy Location	Number of Blots	VOLUME (µL)	Cell Density (TC) (volumes)	Total Cells	% Live Cells	Notes
1	CB001_U1_V1_B01	CB001_U1_V1_B01	-	-	-	-	-	-	-	-
2	VB001_U1_V1_B01	VB001_U1_V1_B01	-	-	-	-	-	-	-	-
3	VB002_U1_V1_B01	VB002_U1_V1_B01	-	-	-	-	-	-	-	-
4	VB002_U1_V1_B02	VB002_U1_V1_B02	-	-	-	-	-	-	-	-
5	VB002_U1_V1_B03	VB002_U1_V1_B03	-	-	-	-	-	-	-	-
6	VB002_U1_V1_B04	VB002_U1_V1_B04	-	-	-	-	-	-	-	-
7	VB002_U1_V1_B05	VB002_U1_V1_B05	-	-	-	-	-	-	-	-
8	VB002_U1_V1_B06	VB002_U1_V1_B06	-	-	-	-	-	-	-	-
9	VB002_U1_V1_B07	VB002_U1_V1_B07	-	-	-	-	-	-	-	-
10	VB002_U1_V1_B08	VB002_U1_V1_B08	-	-	-	-	-	-	-	-

Metadata update

BioCore Dashboard Pipelines Projects Run Status Alpha Version

Pipeline: STAR Pipeline

Description: STAR pipeline includes Quality Control, RNA Filtering, Genome Alignment using STAR.

Steps:

Pipelines

- Public Pipelines
- My Pipelines
- Biocore
- Index
- Alignment
- Demultiplex
- QC
- Sample
- Import

Pipeline Details

Dportal



Dmeta-skin  
↑  
Users enter metadata

B	E	G	
Patient	Skin	S3Bucket	FastqName
CB17	Healthy	s3://biocorebackup/dolphin_import	CB17_H2.fastq.gz
CB19	Healthy	s3://biocorebackup/dolphin_import	CB19_HA.fastq.gz
CB19	Healthy	s3://biocorebackup/dolphin_import	CB19_H2.fastq.gz
CB20	Healthy	s3://biocorebackup/dolphin_import	CB20_H1.fastq.gz
CB20	Healthy	s3://biocorebackup/dolphin_import	CB20_H2.fastq.gz

Dnext

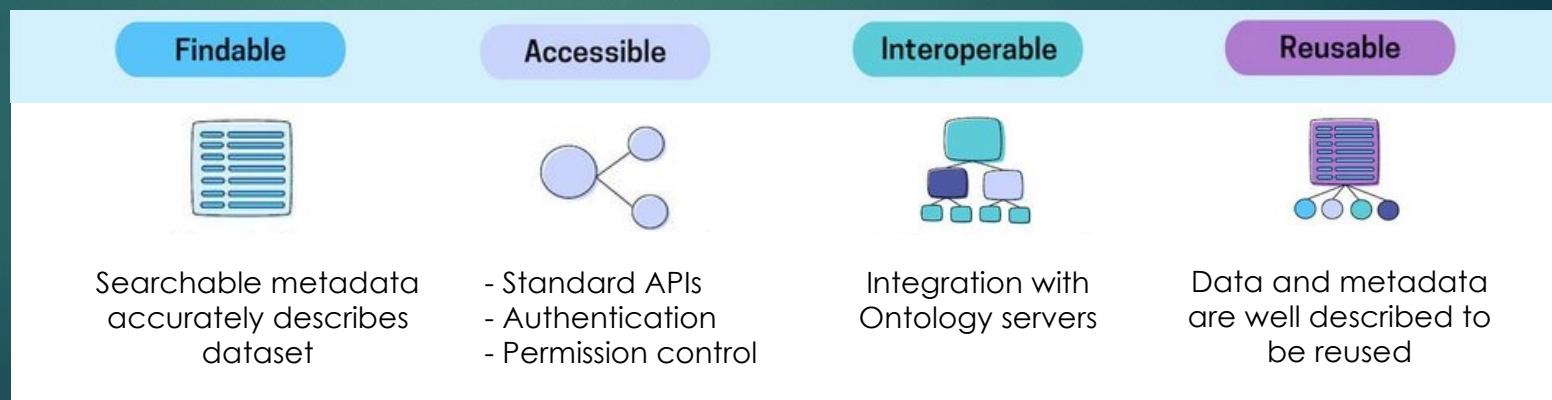
Automatic run submission

# Dmeta

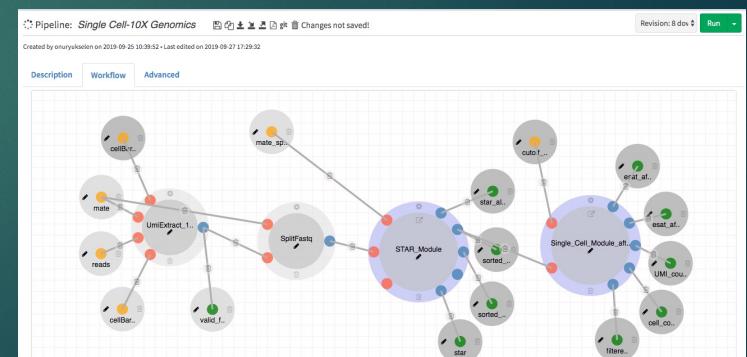
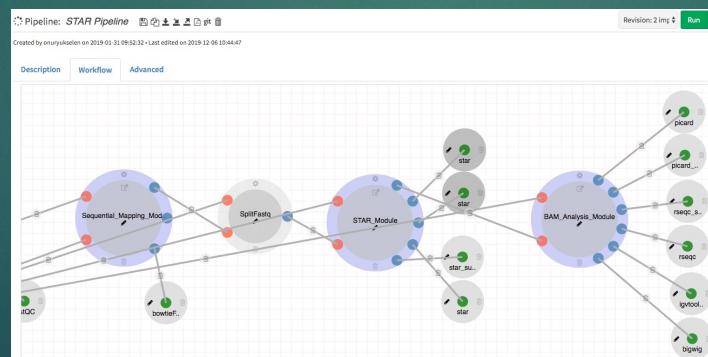
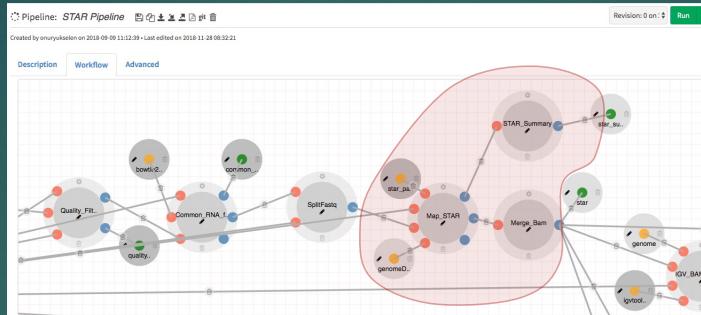
- ▶ Ontology based, metadata tracking system



- ▶ FAIR standards:



# Dnext



**Reproducibility**

**Versioning**

**Copy Run**

**Pipeline: Cell Ranger Pipeline**

Created by onuryskelen on 2019-08-03 11:46:26 • Last edited on 2019-08-28 15:18:34

Revision: 0 on 2019-04-30 21:34:34  
Revision: 1 mtkfastq added on 2019-08-03 13:01:51  
✓ Revision: 2 cellranger-v3.1.0 on 2019-08-28 15:18:34

Description Workflow Advanced

**Pipelines** **Projects** **Run Status** **Amazon** **Profiles** **Tutorial** **VERSION 0.3.62**

**Biocore Run**

**Run:** rsem

**Created by:** kucukura

**Run Settings**

**Run Description**

**Work Directory (Full path)**

/n/umw\_manuel\_garber/kucukura/data/vernia/process5

**Copy Run**

New revision of this pipeline is available. If you want to create a new run and keep your revision of pipeline, please click "Keep Existing Revision" button. If you wish to use same input parameters in new revision of pipeline then click "Use New Revision" button.

**Target Project:** rsem

**Buttons:** Cancel, Keep Existing Revision, Use New Revision



Github  
Permissions  
Shareability

Biocore Pipeline Generation

Description Workflow Advanced

Permissions to View

- Only me
- Only my group
- Everyone

Group Selection

- ✓ Choose group
- Admin
- test
- Zamore-Lab
- Demo

Publish

- No

Biocore Pipeline Generation

Description Workflow Advanced

GitHub onuruykselen / maseq /

Summary

build pending DOI: 10.1101/699539

RNA-seq pipeline includes Quality Control, rRNA filtering, Genome Alignment using HISAT2, STAR and Tophat2, and estimating gene and isoform expression levels by RSEM and featureCounts.

Steps:

- For Quality Control, we use FastQC to create qc outputs. There are optional read quality filtering (trimmmomatic), read quality trimming (trimmmomatic), adapter removal (cutadapt) processes available.
- Bowtie2/Bowtie/STAR is used to count or filter out common RNAs (eg. rRNA, miRNA, tRNA, piRNA etc.).
- RSEM is used to align RNA-Seq reads to a reference transcripts and estimates gene and isoform expression levels.
- HISAT2, STAR and Tophat2 are used to align RNA-Seq reads to a genome. Optionally, counting reads to genomic features such as genes, exons, promoters and genomic bins could be done by featureCounts.
- Genome-wide Bam analysis is done by Picard.
- Optionally you can create Integrative Genomics Viewer (IGV) and Genome Browser Files (TDF and Bigwig, respectively).

Program Versions:

- FastQC v0.11.8
- Star v2.6.1
- Hisat2 V2.1.0

Git Console

pushGit Account: onuruykselen

Repository: maseq

Branch: master

Push to GitHub

dolphinnext / maseq

RNA-Seq pipeline https://dolphinnext.umassmed.edu/index.html

80 commits 2 branches 0 releases All 2 contributors MIT

Branch: master New pull request Projects Wiki Security Insights

This branch is 4 commits ahead of onuruykselen:master.

nephante Merge pull request #18 from onuruykselen/master

Latest commit: 2 days ago

conf 10-09-2019 10:32:58 3 days ago

docs 10-09-2019 10:32:58 3 days ago

.travis.yml 10-09-2019 10:32:58 3 days ago

Dockerfile 10-09-2019 10:32:58 3 days ago

LICENSE 10-09-2019 08:48:07 3 days ago

README.md 10-09-2019 08:48:07 3 days ago

environment.yml 10-09-2019 10:37:26 3 days ago

main.rn 11-09-2019 17:32:02 2 days ago

main.nf 26-08-2019 15:57:19 18 days ago

nextflow.config 10-09-2019 08:48:07 3 days ago

README.md

build pending DOI: 10.1101/699539

RNA-seq pipeline includes Quality Control, rRNA filtering, Genome Alignment using HISAT2, STAR and Tophat2, and estimating gene and isoform expression levels by RSEM and featureCounts.

Steps:

- For Quality Control, we use FastQC to create qc outputs. There are optional read quality filtering (trimmmomatic), read quality trimming (trimmmomatic), adapter removal (cutadapt) processes available.
- Bowtie2/Bowtie/STAR is used to count or filter out common RNAs (eg. rRNA, miRNA, tRNA, piRNA etc.).
- RSEM is used to align RNA-Seq reads to a reference transcripts and estimates gene and isoform expression levels.
- HISAT2, STAR and Tophat2 are used to align RNA-Seq reads to a genome. Optionally, counting reads to genomic features such as genes, exons, promoters and genomic bins could be done by featureCounts.
- Genome-wide Bam analysis is done by Picard.
- Optionally you can create Integrative Genomics Viewer (IGV) and Genome Browser Files (TDF and Bigwig, respectively).

# Increase the team efficiency with

- ▶ Easy execution
- ▶ Easy monitoring
- ▶ Easy reporting

The screenshot shows a bioinformatics pipeline interface for running the RSEM pipeline. The top navigation bar includes 'Run: rsem demo run-vernia', 'Project: rsem', 'Pipeline: RNA-seq Pipeline', and status indicators like 'Completed' and 'Copy Run'. Below the navigation is a 'Run Description' field with placeholder text 'Enter run description here..'. The 'Work Directory (Full path)' is set to '/nl/umw\_manuel\_garber/kucukura/data/vernia/process4'. The 'Run Environment' section shows a dropdown menu with options: Cluster (Remote machine: ak97w@ghpcc06.umassrc.org), Local (Remote machine: kucukura@galaxy.umassmed.edu) (which is selected and highlighted in blue), Blocore (Remote machine: kucukura@blockcore.umassmed.edu), and Amazon (Amazon: Status:terminated Image id:ami-07659c1e5fe51b189 Instance type:5a.2xlarge). The main configuration area is divided into 'Inputs' and 'Rsem Settings'. The 'Inputs' section contains fields for 'Given Name' and 'Select Items' for various parameters: 'reads' (test dataset), 'mate' (pair), 'genome\_build' (mouse\_mm11), 'run\_RSEM' (yes), 'run\_HISAT2' (no), 'run\_STAR' (yes), and 'run\_Tophat' (no). The 'Rsem Settings' section includes fields for 'RSEM reference type' (bowtie), 'RSEM parameters' (-p 4), and 'Output Genome BAM' (true). A note states: 'If true is selected, RSEM will generate a BAM file, with alignments mapped to genomic coordinates and annotated with their posterior probabilities. (default=true)'. An 'Ok' button is at the bottom right of the settings panel.

# Increase the team efficiency with

- ▶ Easy execution
  - ▶ Easy monitoring
  - ▶ Easy reporting

Run Status								
Show <input type="text" value="10"/> entries		Search: <input type="text"/>						
ID	Run Name	Pipeline	Work Directory	Description	Status	Date Created	Owner	Options
1721	test.sample	RNA-seq Pipeline	/home/oy28w/nextflowruns/rsemNew2wq		Initializing	2019-09-13 10:46:11	onuryukseLEN	<button>Options</button>
1768	Y1lib4_2019_9_12	Cell Ranger Pipeline	/project/umw_nathan_lawson/DolphinOuts/Y1scRNaseqOut		Error	2019-09-12 16:01:25	lawsonN	<button>Options</button>
1754	IPSCs_R1_R2_R3-no_rmsk	RSEM Pipeline	/n1/umw_robert_brown/ozgun/RNAseq/analysis_IPSCs_R1_R2_R3_no_rmsk	3 replicates combined RNAseq samples of IPSCs_no_rmsk	Completed	2019-09-11 09:36:52	Ozgunuy	<button>Options</button>
1761	cell_ranger-HiSeq7-12	Cell Ranger Pipeline	/n1/umw_robert_brown/kit/cell_ranger_allseq		Error	2019-09-12 10:27:03	Katherine.Mocarski	<button>Options</button>
1767	RRP01 RNA-Seq run	RNA-seq Pipeline	/project/uma_ravil_ranjan/RRP01_pipeline		Completed	2019-09-12 00:03:10	Ravi.Ranjan	<button>Options</button>
1760	fh332_abcc9acta2_minISeq_1	Cell Ranger Pipeline	/project/umw_nathan_lawson/DolphinOuts/fh332scRNaseq		Completed	2019-09-11 08:50:45	lawsonN	<button>Options</button>
1759	dgcR8_smallRNA_set1-copy	RNA-seq Pipeline	/project/umw_oliver_rando/CC_smallRNA/bigwig_next		Completed	2019-09-10 17:07:17	conine	<button>Options</button>
1616	chip git version	ChIP-seq Pipeline	/project/umw_garberlab/yukseLEN/chipGit45		Completed	2019-09-10 11:29:58	onuryukseLEN	<button>Options</button>
1641	ataC git version	ATAC-seq Pipeline	/project/umw_garberlab/yukseLEN/nextflowruns/chiptestgit2		Completed	2019-09-10 10:53:06	onuryukseLEN	<button>Options</button>
1756	test	RNA-seq Pipeline	/home/oy28w/nextflowruns/rsemNew2edzz		Completed	2019-09-10 10:38:11	onuryukseLEN	<button>Options</button>

# Increase the team efficiency with

- ▶ Easy execution
- ▶ Easy monitoring
- ▶ Easy reporting

The diagram illustrates the workflow integration between Multiqc, RSEM, and DESeq2. A central 'Run Report 2' interface is connected by red arrows to three separate windows:

- Multiqc Read Distribution:** This window shows a stacked bar chart of mapped reads across various genomic features. The legend includes: ETS\_Exons (blue), 5'UTR\_Exons (green), Intronic (orange), TSS\_up\_5kb (red), TES\_down\_1kb (purple), TES\_down\_5kb (brown), and Other\_intergenic (yellow).
- RSEM Rmarkdown:** This window displays R Markdown code for DESeq analysis. It includes a scatter plot of log2FoldChange vs -log10(padj) with points colored by expression status (red for expressed, grey for non-expressed).
- DESeq Rmarkdown:** This window also displays R Markdown code for DESeq analysis, showing a similar scatter plot.

# DolphinNext Use Cases

# Walhout Lab use case

Pipelines Projects Run Status

## Biocore Pipeline Generation

Pipeline: CelSeq Pipeline-XL

Created by xl95w on 2019-12-09 12:10:48 • Last edited on 2020-03-04 22:29:50

Description Workflow Advanced

Revision: 5 debug on 2019-12-14 23:44:39  
 Revision: 6 debug on 2019-12-14 23:44:25  
 Revision: 7 new setting on 2019-12-14 23:44:12  
 Revision: 8 setting on 2019-12-14 23:43:58  
 Revision: 9 job config opt on 2019-12-14 23:43:47  
 Revision: 10 back to UmiExtractor on 2019-12-14 23:43:30  
 Revision: 11 add dedup after star on 2019-11-24 20:55:01  
 Revision: 12 add dedup after star on 2019-12-14 23:42:56  
 Revision: 13 add trim double count on 2019-12-20 12:46:15  
 ✓ Revision: 14 remove deduplication on 2020-03-04 22:29:50

ID	Sample	Pipeline	Project	Description	Status	Created	Last Run	Options
4144	metabolic_lipid_met4_lib1_and_lib6	CelSeq Pipeline-XL (Rev 14)	/project/umw_marian_walhout/data/xuhangLi/XL_HZ_CelSeq_02082021/dolphinNext/met4_lib1_and_lib6/	processing BGI-sequenced data of the metabolic RNAi library. Three bio replicates were pooled in the same sequencing li..	Completed	2021-02-08 11:48:36	xl95w	<button>Options</button>
4026	metabolic_lipid_met3_lib5_and_lib6	CelSeq Pipeline-XL (Rev 14)	/project/umw_marian_walhout/data/xuhangLi/XL_HZ_CelSeq_12112020/dolphinNext/met3_lib5_and_lib6/	processing BGI-sequenced data of the metabolic RNAi library. Three bio replicates were pooled in the same sequencing li..	Completed	2021-01-07 11:12:22	xl95w	<button>Options</button>
4025	metabolic_lipid_met3_lib3_and_lib4	CelSeq Pipeline-XL (Rev 14)	/project/umw_marian_walhout/data/xuhangLi/XL_HZ_CelSeq_12112020/dolphinNext/met3_lib3_and_lib4/	processing BGI-sequenced data of the metabolic RNAi library. Three bio replicates were pooled in the same sequencing li..	Completed	2021-01-06 21:15:47	xl95w	<button>Options</button>

Showing 1 to 10 of 140 entries (filtered from 2,864 total entries)

Previous 1 2 3 4 5 ... 14 Next



Xuhang “Hang” Li

# Reprocessed all data to use updated pipelines

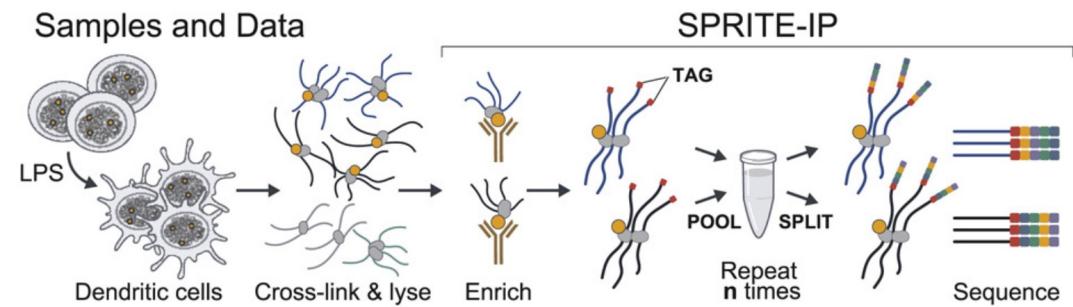
- ▶ Updated genome versions
  - ▶ mm9 to mm10
  - ▶ hg19 to hg38
- ▶ 300+ RNA-seq datasets to have same parameters and updated pipelines
- ▶ lncRNA quantification with Fitzgerald Lab
  - ▶ Easy updating with newer annotations

# Training and Sharing

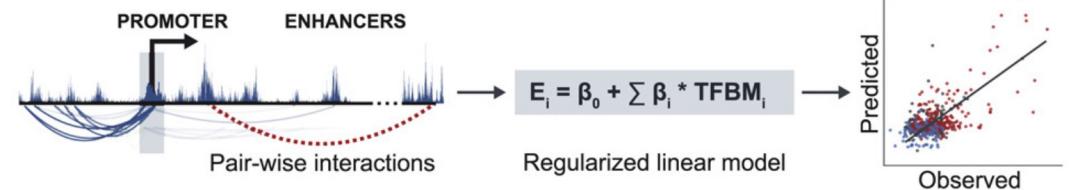
- ▶ Shared the runs and pipelines with
  - ▶ Summer interns (undergrads)
  - ▶ New grad students
  - ▶ Research associates
  - ▶ Wet-lab scientists

## Creating new pipelines using pre-built modules

- ▶ HiC
- ▶ SPRITE



### View point analysis



# DolphinNext Demo and Tutorials

<https://github.com/dolphinnext/dolphinnext-tutorial>



# Thanks to

- **Manuel Garber**
- **Onur Yukselen**
- **Artur Manukyan**
- **Pranitha Vangala**

- Alan Derr
- Elisa Donnard

All Garber Lab Members

- Craig Mello
- Ahmet Ozturk
- Osman Turkyilmaz

- Ruijia Wang
- Xuhang Li

- Nathan Lawson
- Julie Zhu

