

# HIVE-TC CMU: Comprehensive, Flexible and FAIR Tools for the HuBMAP HIVE



# Major focus in Year 1

- Our major focus in the first year has been on:
  - Pipelines for processing HuBMAP data
  - Working on the data portal and website
  - Connecting with other efforts (mainly HCA)
- Most important issues we have addressed so far:
  - Establishing and working with the DRT's to fully define the set of pipeline needed for the initial release
  - Initial implementation and testing of containerized pipelines for sequencing and imaging data
  - Imaging data workshop
  - Minimizing redundancy with work at HCA
  - Initial development and discussions about the HuBMAP portal

# Automatically collecting and processing scRNA-Seq data

NCBI Resources | How To | Sign In to NCBI

GEO DataSets | GEO DataSe | VA-seq[all] AND RNA[all] AND expression profiling by high throughput sequencing[GTPY] | Search

Entry type: DataSets (0), Series (501), Samples (0), Platforms (0)

Search results: Items: 1 to 20 of 501

Transcriptome analysis of high glucose loaded mouse embryo at 2-cell stage treated with ZGW

ArrayExpress | "RNA-seq of coding RNA from single cells" AND mc

Home | Browse | Submit | Help | About ArrayExpress | Contact Us | Login

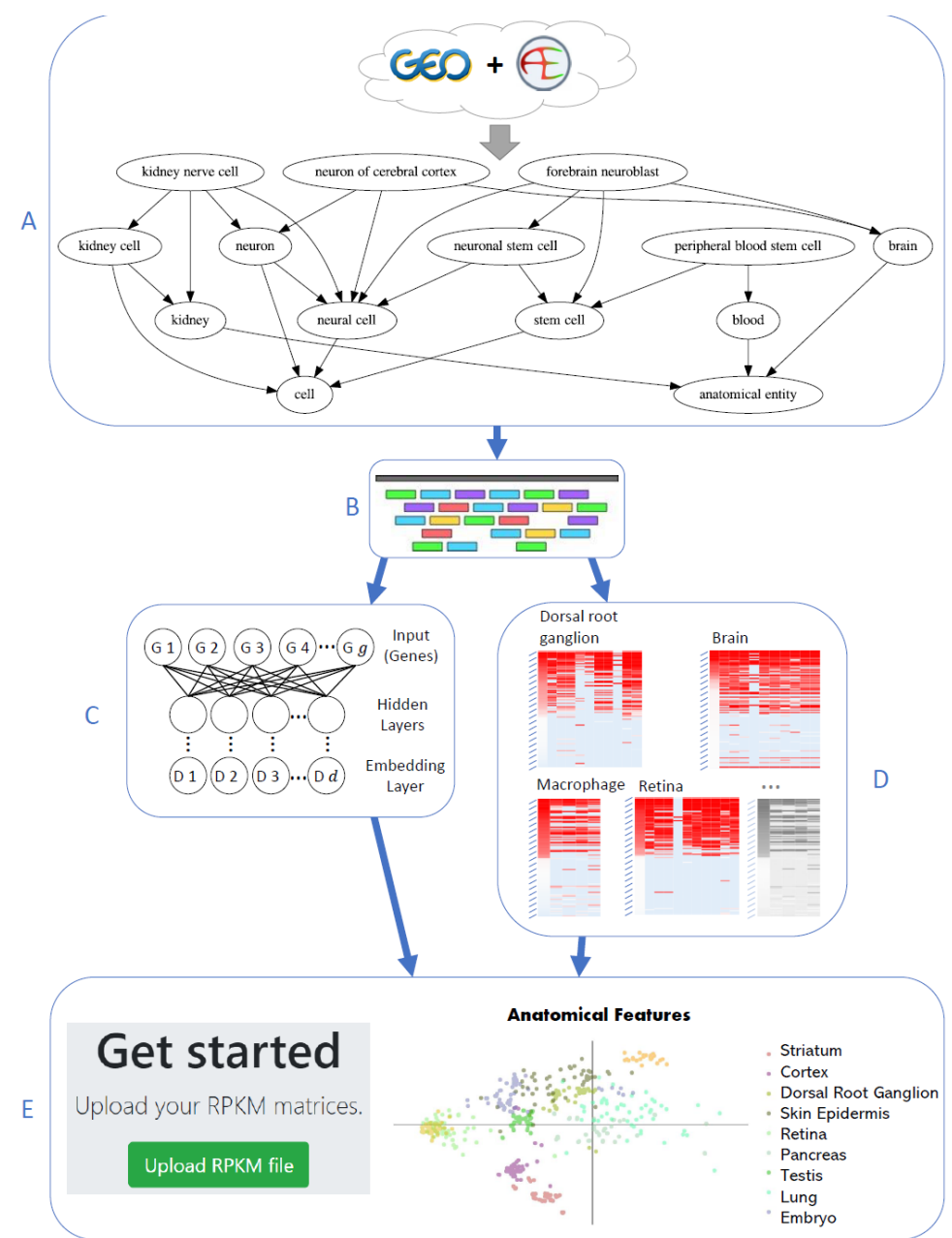
Filter search results | Show more data from EMBL-EBI

Search results for "RNA-seq of coding RNA from single cells" AND mouse

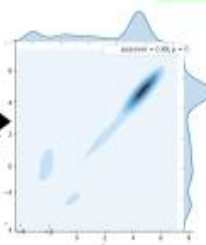
Filtered by AE only on

Page 1 of 2 | Showing 1-25 of 30 experiments | Page size 25 | 50 | 100 | 250 | 500

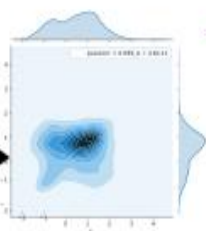
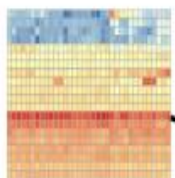
Accession	Title	Type	Organism	Assays	Released	Processed	Raw	Views	Atlas
E-MTAB-5466	Single-cell RNA-seq of motor	RNA-seq of	Mus musculus	202	01/02/2018	-	-	74	-



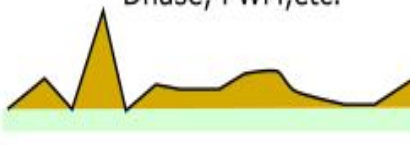
single cell RNA-seq



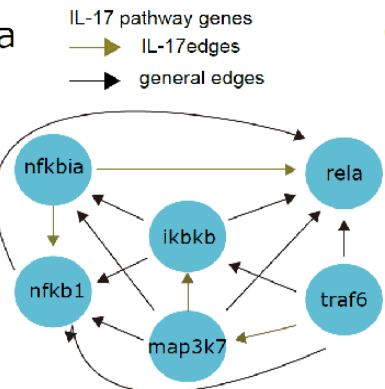
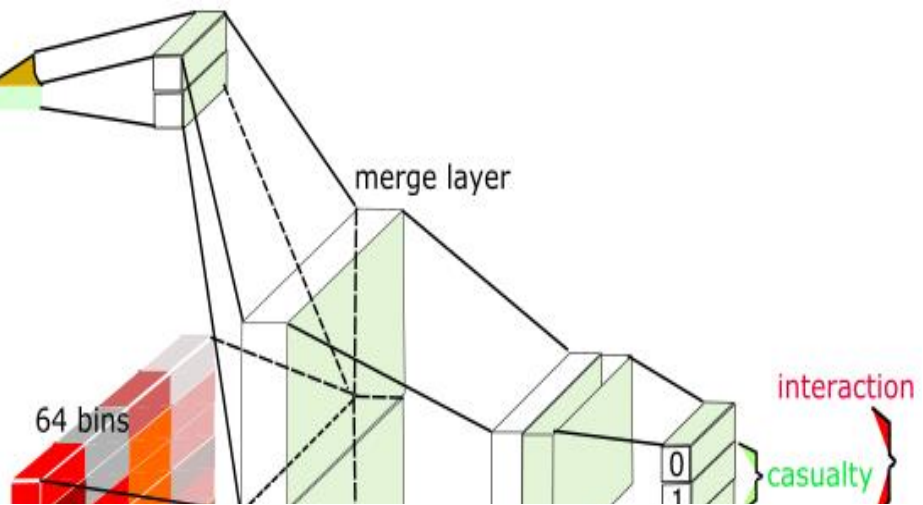
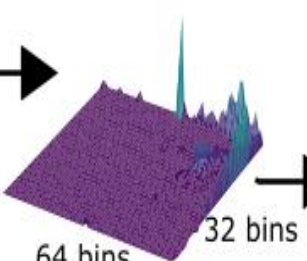
bulk RNA-seq



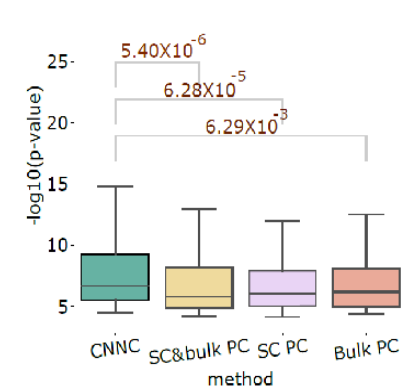
Dnase, PWM, etc.



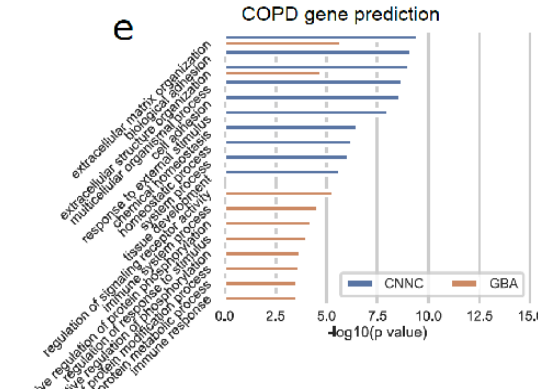
combined NEPDF



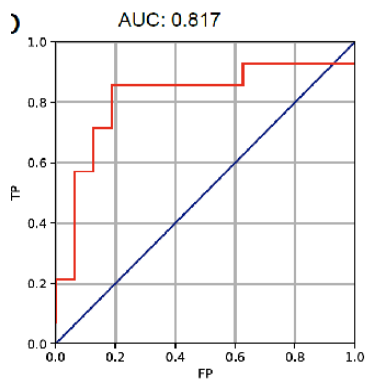
c



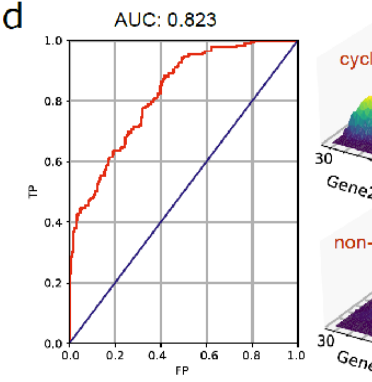
e



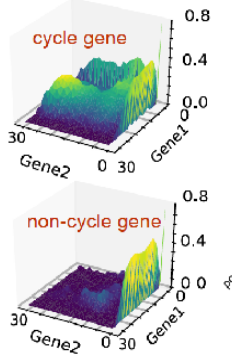
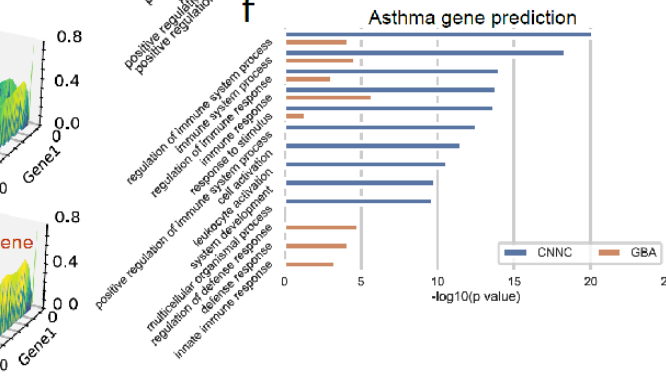
d



d

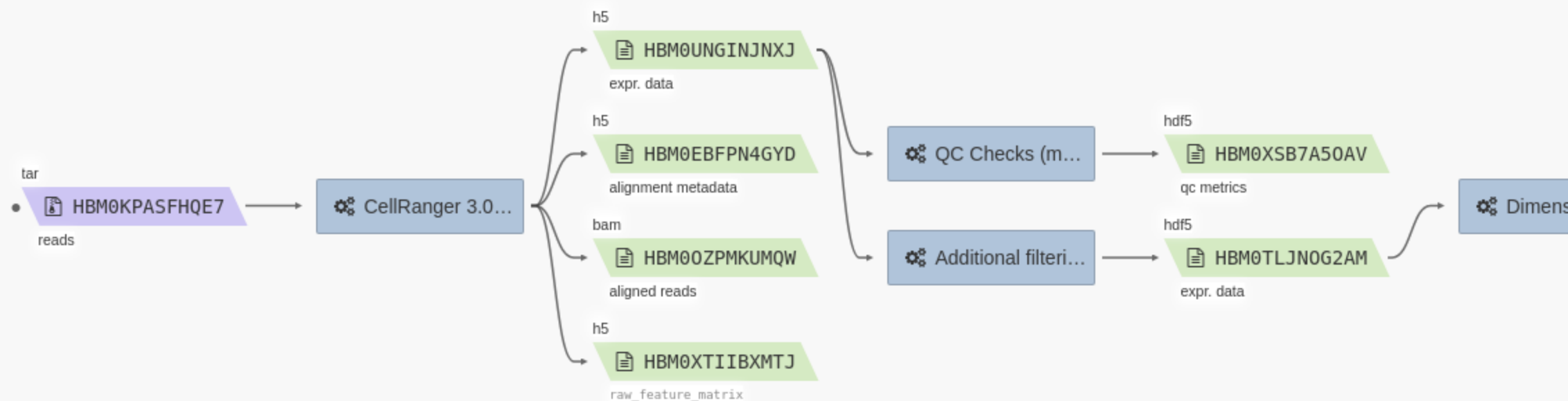


f





## Graph

[Show Reference Files](#) [Show Parameters](#)[Analysis Steps](#)[Center Nodes](#)

## Graph Legend

- Input File
- Output File
- Input Reference File
- Intermediate File

# Microscopy Pipelines

	A	B	C	D	E	F	G	H	I
1	Technology	Center	Tissue	Bench PoC	Data PoC	Project Manager	How do we get from initial raw data to usable analysis results? (For each pipeline step: name/description of method,		
2	CODEX	Florida	Thymus, Spleen, Lymph	Marda Jorgensen	Jesus Penalzoa	<a href="#">Marda Jorgensen</a>	Full tissue, stain, 10 micron thick section, raw data as ~10,000 raw images in TIFF format	For each stack, find best plane of focus, stitch together (performed by machine at time of processing, but probably upload results to HIVE also)	
3	CODEX	<a href="#">Stanford</a>	Colon	Sarah/Christian	Vishal (vgautham)	<a href="#">(? ) Aaron Horning</a>	(ahorning@stanford.edu) OR Cynthia Nogoy (cnogoy@stanford.edu)		
4	CODEX	Vanderbilt	Kidney	<a href="#">Jeff Spraggins?</a>	<a href="#">Heath Patterson</a>	<a href="#">Danielle Gutierrez</a>			
5	Lightsheet	Florida	Thymus, Spleen, Lymph	Seth	Seth	<a href="#">Marda Jorgensen</a>	Light sheet microscopy imaging (2x2x5mm*3 volume max; LEVEL 0 Raw data: *.czi (Zeiss) and *.sis (Arivis). TMC performing	Volumetric immunohistochemistry analysis. LEVEL 1: "Blob finder" segmentation -- identify objects based on thresholding and object size. Raw data: *.czi (Zeiss) and *.sis (Arivis). Are these proprietary formats? If so can we also have data in open ones? Software package(s) used? TMC performing.	Volumetric immunohistochemistry analysis. LEVEL 2: Segmentation analysis -- Object data (count, surface area, volume, co-localization). *.sis files. Software package(s) used? Method(s)? TMC performing
6	Autofluorescence	Cal Tech				<a href="#">Dana Jackson</a>			
7	Autofluorescence	Stanford				<a href="#">(? ) Aaron Horning</a>	(ahorning@stanford.edu) OR Cynthia Nogoy (cnogoy@stanford.edu)		
8	Autofluorescence	UCSD		Richard Q	Richard Q	Richard Q			
9	Autofluorescence	Vanderbilt		<a href="#">Jeff?</a>	Jeff	<a href="#">Danielle Gutierrez</a>	Raw 2D autofluorescence data -- LEVEL 0. ~1um pixel size, 4 channels (red/green/blue/violet) X100 sections. Zeiss CZI format, instrumental data as text file. TMC performing.	Data conversion to processed 2D data -- LEVEL 1. ~1um pixel size. 4 channel fluorescence microscopy (red/green/blue/violet) X 100 sections. OME-TIFF format. TMC performing	Register with mass spectrometry data. Image registration text file. Software, methods used? TMC performing
10	SeqFISH	Cal Tech	Blood vessel, Breast, He	Nina	Nico	<a href="#">Dana Jackson</a>	Raw data: LEVEL 0 -- raw TIFF images x 10-20 fields of view x 4.2 megapixels x 20-25 sections. Tag antibodies with heavy metals	Preprocessing -- processed data, LEVEL 1. Processed TIFF images. *.csv offsets for DAPI alignment. Software, methods, code used? TMC performing.	Decoding -- decoded data, LEVEL 2. *.csv gene locations/counts/cell. Software, methods, code used? TMC performing

# What is Dockstore?

Dockstore is a free and open source platform for sharing scientific tools and workflows. It is a registry of Docker-based resources described using popular workflow languages CWL, WDL, and Nextflow.

- **Portability**
  - Run workflows in any environment that supports Docker
- **Interoperability**
  - Standardize computational analysis through GA4GH APIs
- **Reproducibility**
  - Create, Share, Use
  - Containers + Popular descriptor languages

Search Docker Tools and Workflows for the Sciences:

Enter Keyword...

*Dockstore, developed by the Cancer Genome Collaboratory, is an open platform used by the GA4GH for sharing Docker-based tools described with the Common Workflow Language (CWL), the Workflow Description Language (WDL), or Nextflow (NFL)*

Sign up to Contribute >

Quick Start >

News and Events >

Discuss >

Browse Tools | Browse Workflows

A tool is a docker container with an associated descriptor describing how to run it.

Name	Author	Format	Project Links	Stars ↓
✓ pancancer/pcawg-dkzf-workflow	Brian O'Connor	CWL	GitHub Quay.io	4★
✓ pancancer/pcawg-sanger-cgp-workflow	Keiran Raine	CWL,WDL	GitHub Quay.io	3★
wtsicgp/dockstore-biobambam2/bamtobfastq	Keiran Raine	CWL	GitHub Quay.io	2★

Tweets by @DockstoreOrg

Dockstore Retweeted

Mike Lin @DNAdmin

miniwdl v0.3.0 is the first version with independent capability to run #OpenWDL workflows on the local host. Early testing, eager for community to try it out & report interop problems. @WDL\_dev @DockstoreOrg @GA4GH @czscience github.com/chanzuckerberg...

chanzuckerberg/ml... A static analysis tool... github.com

Jul 19, 2019

Now on version 1.6.0, first presented version 1.25 at BOSC2017





## Organizations



# The Human BioMolecular Atlas Program

The vision for the Human BioMolecular Atlas Program (HuBMAP) is to catalyze development of a framework for mapping of the human body at high resolution to transform our understanding of tissue organization and function.

HuBMAP@mail.nih.gov

<https://commonfund.nih.gov/hubmap>



Collections

Members

Events

### HuBMAP Analysis Pipelines

Analysis pipelines developed by the Human BioMolecular Atlas Program



### HCA Skylab Pipelines

Secondary analysis pipelines for the Human Cell Atlas.



### nf-core Pipelines

Bioinformatics analysis pipelines used for RNA sequencing data



## HuBMAP

The overall goals of the NIH Common Fund Human Biomolecular Atlas Program (HuBMAP) are to (1) accelerate development of the next generation of tools and techniques for constructing high-resolution spatial tissue maps that quantify multiple types of biomolecules; (2) generate foundational 3D tissue maps using validated high-content, high-throughput imaging and omics assays; (3) establish an open data platform that will develop novel approaches to integrating, visualizing, and modeling imaging and omics data to build multidimensional maps, and making data rapidly findable, accessible, interoperable, and reusable by the global research community; (4) coordinate and collaborate with other funding agencies, programs, and the biomedical research community to build the architecture and tools for mapping the human body



Search

Organizations

Docs

My Tools

My Workflows

mruffalo

Collections



# The Human BioMolecular Atlas Program / HuBMAP Analysis Pipelines

Analysis pipelines developed by the Human BioMolecular Atlas Program



[github.com/mruffalo/cellranger-pipeline](https://github.com/mruffalo/cellranger-pipeline)

Last updated Sep 17, 2019

Remove

Select edit to add a markdown description



[github.com/mruffalo/chromvar-cwl/chromvar](https://github.com/mruffalo/chromvar-cwl/chromvar)

Last updated Sep 3, 2019

Remove

# Deliverable for Year 2

## Goals for Year 2:

- We will work with the TMCs to make sure that pipelines needed to process HuBMAP data are available
- We will containerized pipelines for processing HuBMAP sequencing, imaging and mass spec data
- We will attempt to standardize pipelines currently used by the different groups
- We will work with others to make sure that the HuBMAP data is accessible, searchable and downloadable
- We will attempt to test joint analysis pipelines for HuBMAP data
- Current needs
  - Cooperation from TMCs and DRTs

# Collaborations

- We are starting to work on a project for the analysis of spatial transcriptomics and proteomics data
- Very early stages, no updates yet

# What should HuBMAP Do ?

- Upload sample data for all expected modalities in the initial release ASAP
- Make sure that all the data you expect to provide is covered by one of the DRTs
- Better discussion and collaboration with external efforts (NIH or others)
- Start uploading actual data collected