

# Advancing the Big Data Genomics Analysis using Cloud

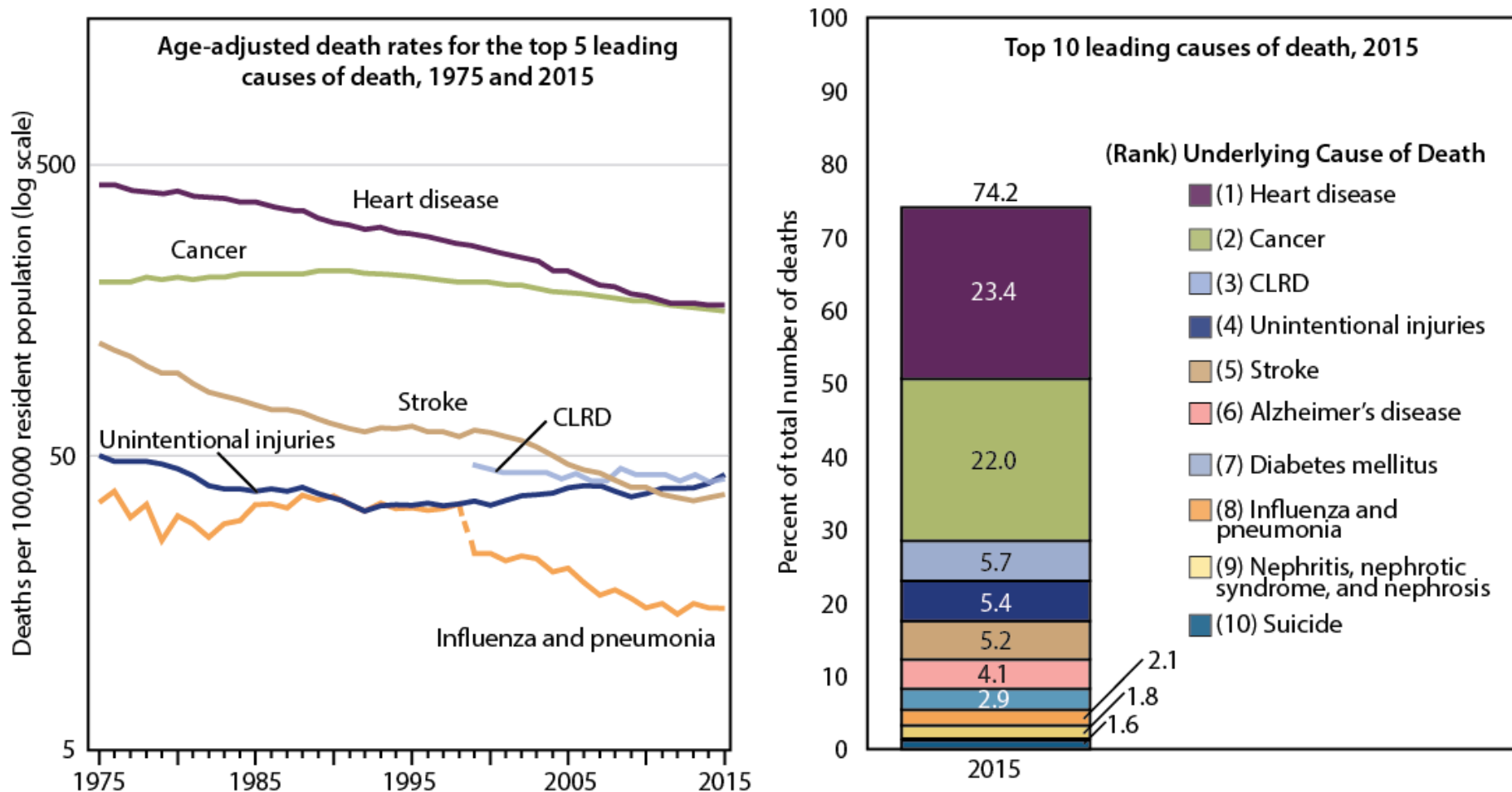


Segun C. Jung

Biomedical Informatics Consultant

University of Chicago & Argonne National Laboratory

# Leading causes of death is the US



NOTE: Due to coding changes for chronic lower respiratory diseases (CLRD) between ICD-9 and ICD-10, which prevent the direct comparison of trends prior to 1998 and after 1999, rates for CLRD are only shown for 1999 onwards.

SOURCE: NCHS, *Health, United States, 2016*, Figure 8. Data from the National Vital Statistics System (NVSS).

More than **90%** of cancer patients carry a mutation that may be responsive to a known drug

Mark Rubin, Weill Cornell Medical College and NewYork-Presbyterian Hospital in New York in *Nature*, April, 2015

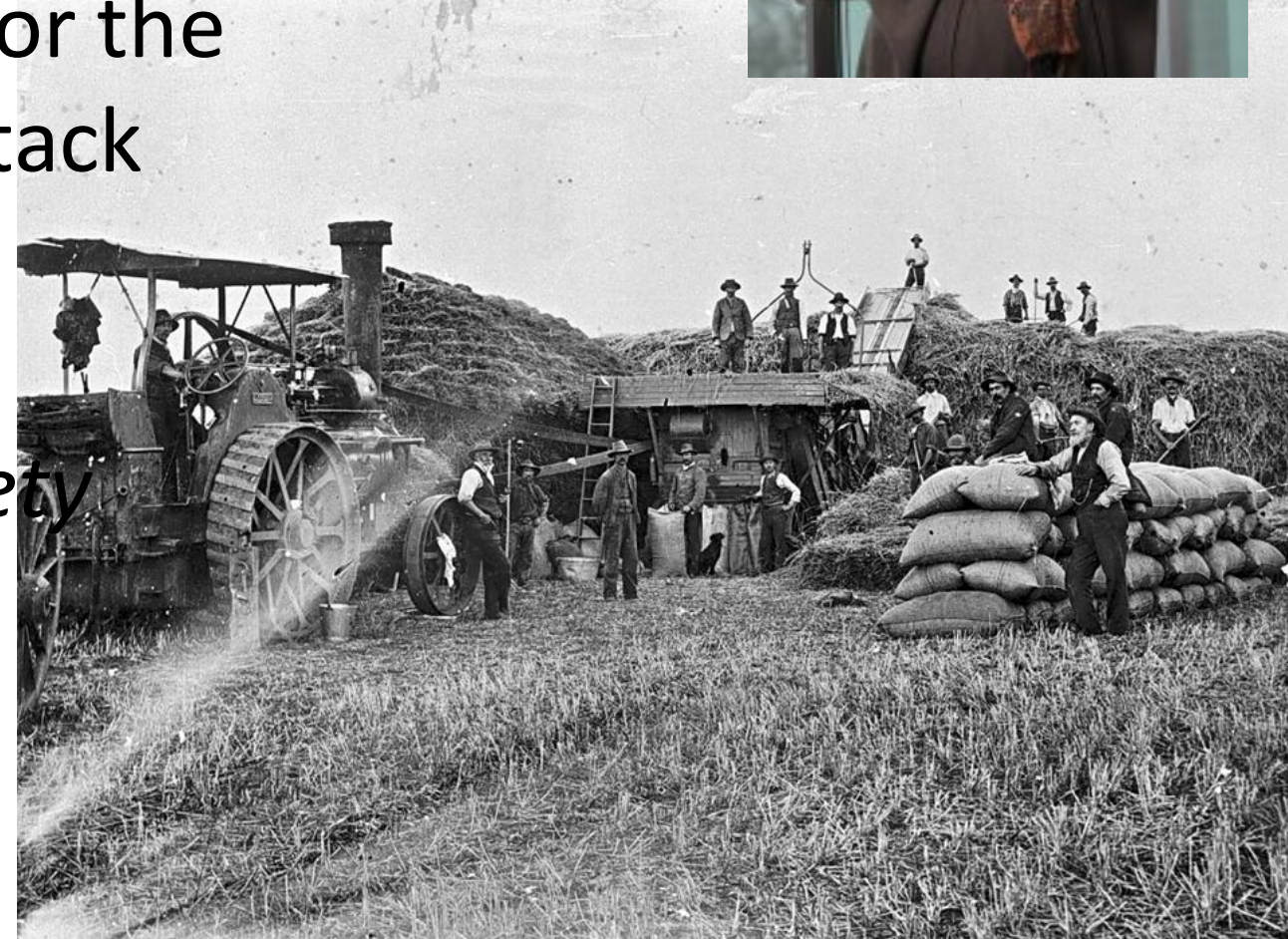




- Trying to find a single causative gene for diseases with a complex genetic background is like looking for the proverbial needle in a haystack



- Dr. Nancy Cox  
*President of American Society  
of Human Genetics*

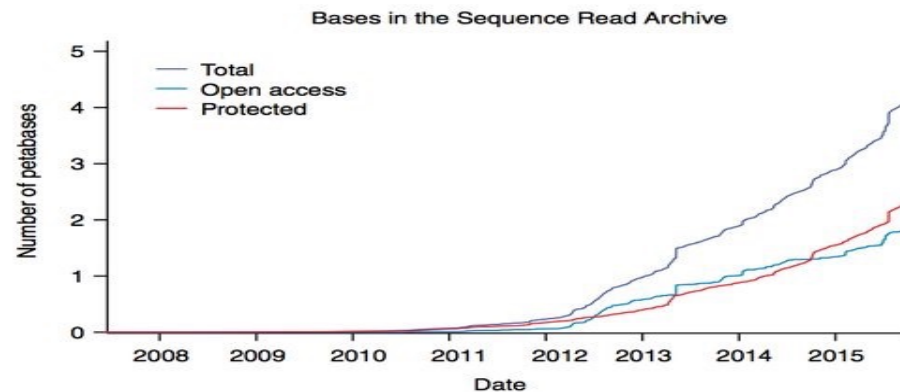




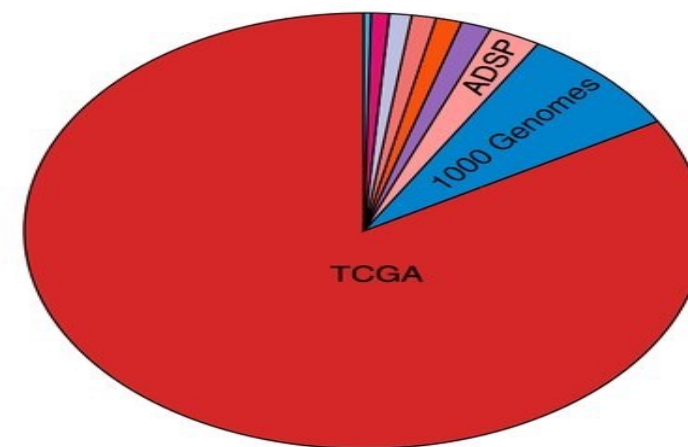


# Big Biomedical Data Sources

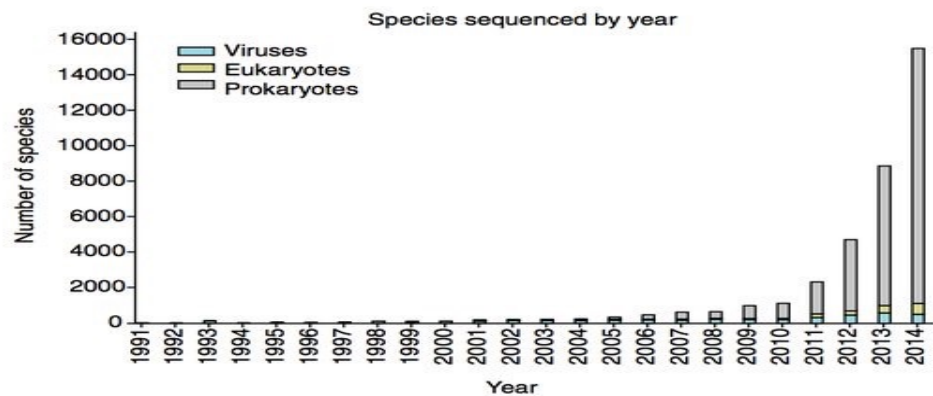
b



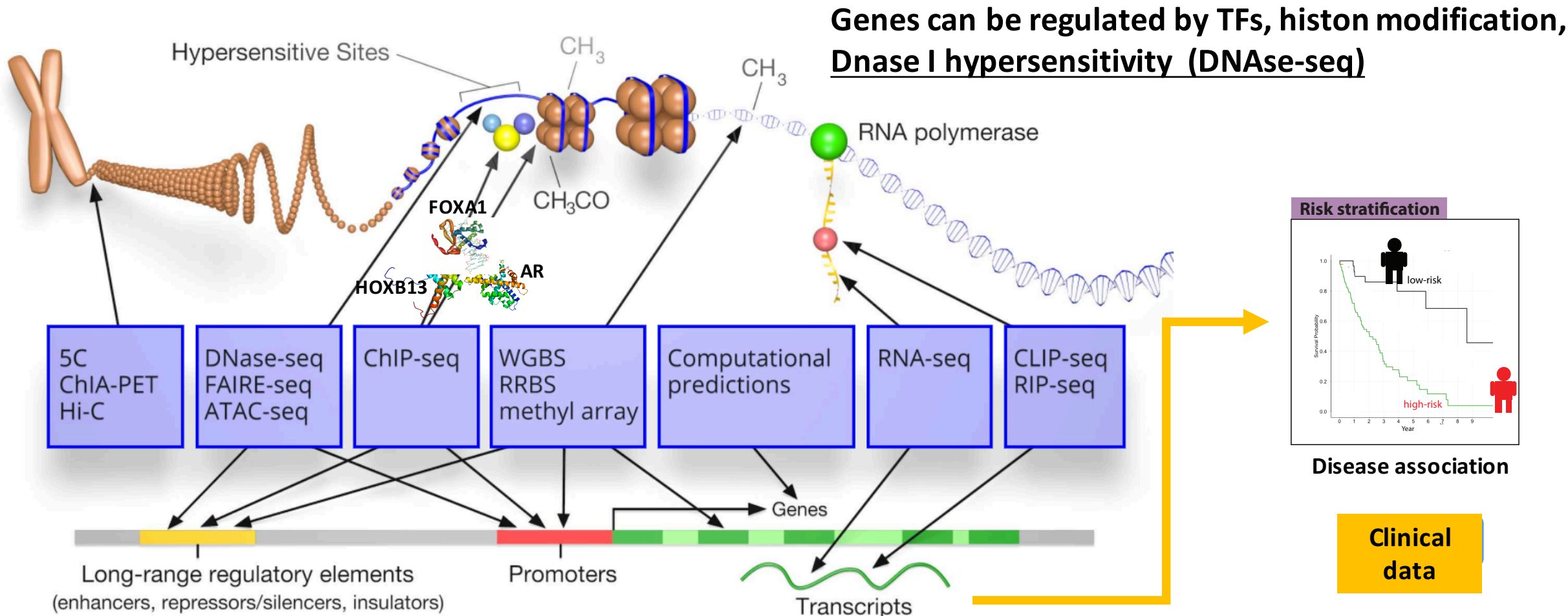
c



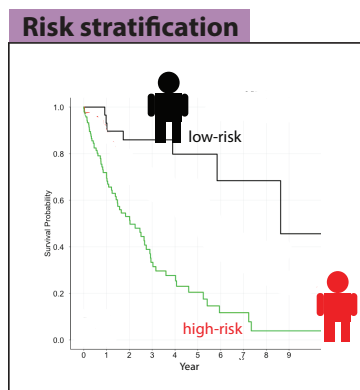
d



TCGA	- 2300 TB
1000 Genomes*	- 222 TB
ADSP	- 68 TB
NHGRI LSSP*	- 40 TB
GTEx	- 34 TB
NHLBI ESP	- 32 TB
HMP*	- 29 TB
ENCODE*	- 9 TB



Genes can be regulated by TFs, histon modification, Dnase I hypersensitivity (DNase-seq)



Disease association

Clinical data



Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

<https://www.encodeproject.org/>

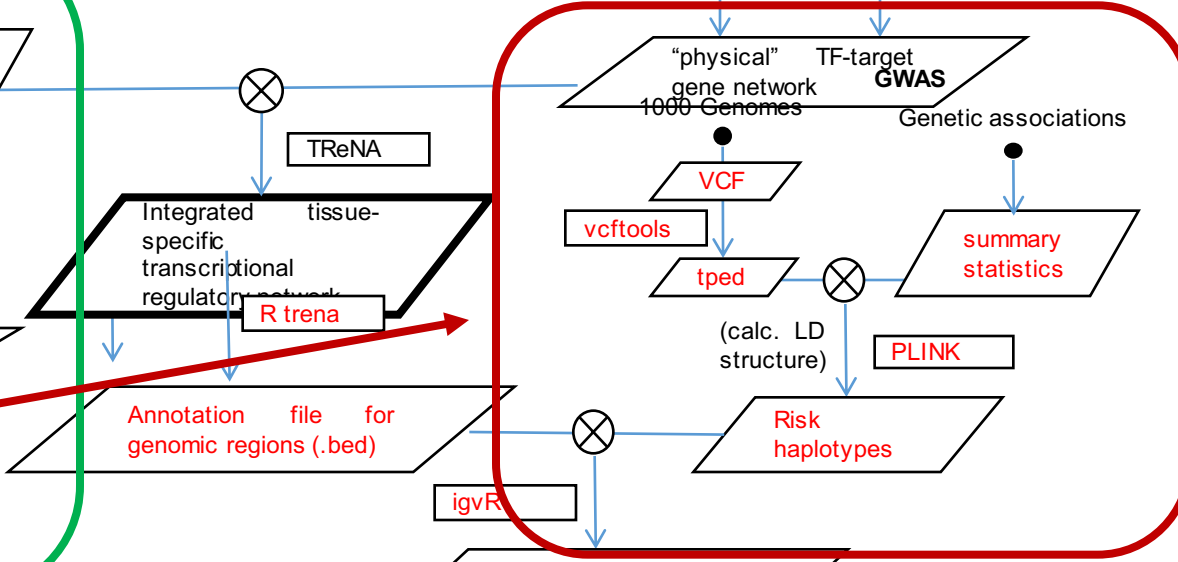
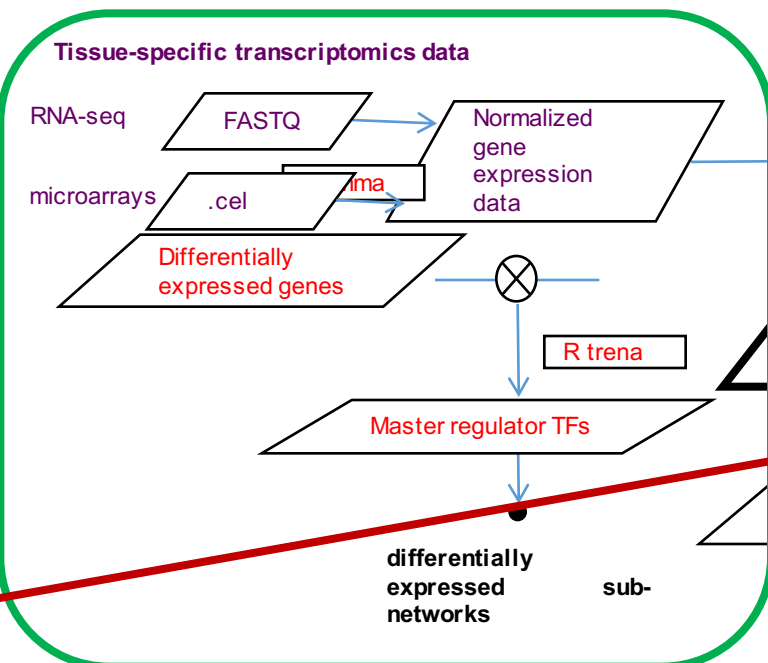
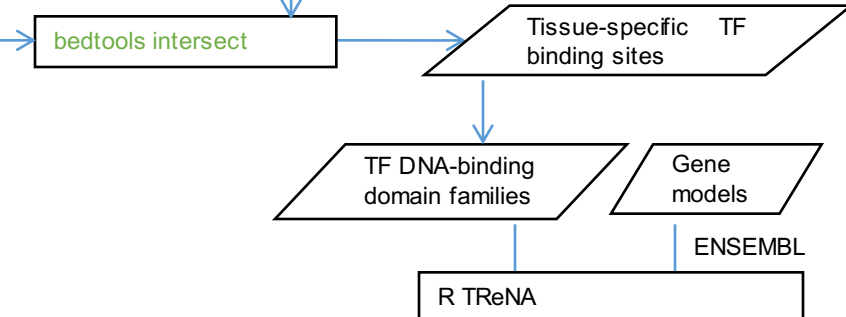
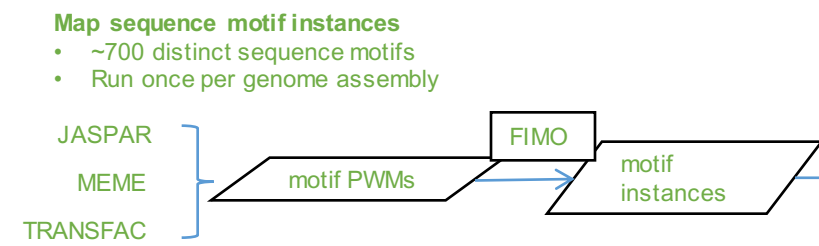
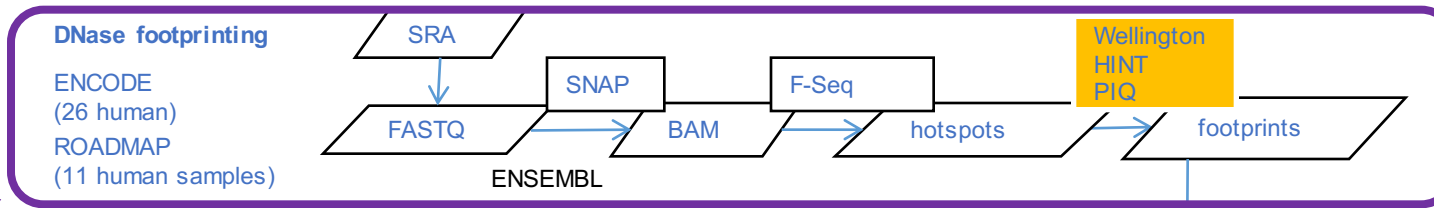
# Big Data Challenges in Transcriptional Regulatory Network Analysis (TReNA)

- Identifying and transferring raw data objects
  - BDBag and minid
- Building and optimizing workflows
  - Shell scripting to drag/drop tool box, File I/O
- Scalable data analysis on cloud
  - Batch submit, Monitor jobs



# What is TReNA?

Identifying sites for DNA-Protein interactions



Gene expression

Disease association

**regulatory genetic variation associated with disease risk**

# DNA-binding data in ENCODE

- DNase-seq, FAIRE-seq, and ATAC-seq for the hypersensitive site
- Total number of tissues: 27 (lymphoblast, brain, skin, etc)
- Total number of patient samples: 206
- Total number of fastq files: **1379**
  - Each patient sample has a few to many replicates
- Total size of the raw data: **2.5 TB**







# BDDS Solutions: Enabling TReNA – BDBag



Create a BDBag from an ENCODE search.

For example enter the following search:

[https://www.encodeproject.org/search/?type=Experiment&assay\\_title=RNA-seq&replicates.library.biosample.biosample\\_type=stem+cell](https://www.encodeproject.org/search/?type=Experiment&assay_title=RNA-seq&replicates.library.biosample.biosample_type=stem+cell)

Or paste in an Encode metadata file.

Encode Search Query

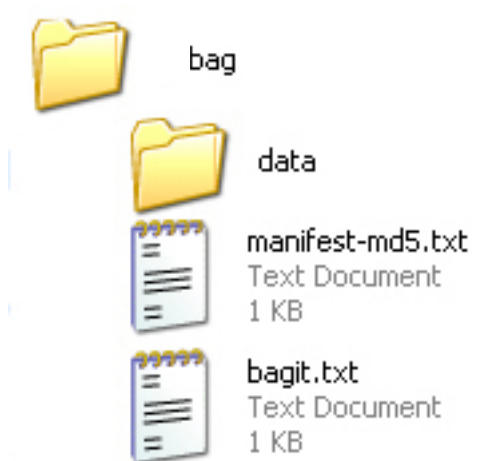
Encode Metadata File

File accession	File format	Output type	Experiment accession	Assay	Biosample term id	Biosample term name	Biosample type	Biosample life stage
Biosample sex	Biosample organism	Biosample treatments	Biosample subcellular fraction term name	Biosample phase	Biosample synchronization stage			
Experiment target	Antibody accession	Library made from	Library depleted in	Library extraction method	Library lysis method	Library crosslinking method		
Experiment date released	Project	RBNS protein concentration	Library fragmentation method	Library size range	Biosample Age	Biological replicate(s)		
Technical replicate	Read length	Run type	Paired end	Paired with	Derived from	Size Lab	md5sum	File download URL
					Assembly	Platform		

Create BDBag

BDBag created: <ark:/99999/fk40294z6m>

or you can access, transfer, and share the complete, materialized BDBag with [Globus](#)



**Lymphoblast  
metadata**

<https://github.com/ini-bdds/bdbag>

# Enabling TReNA – BDBag



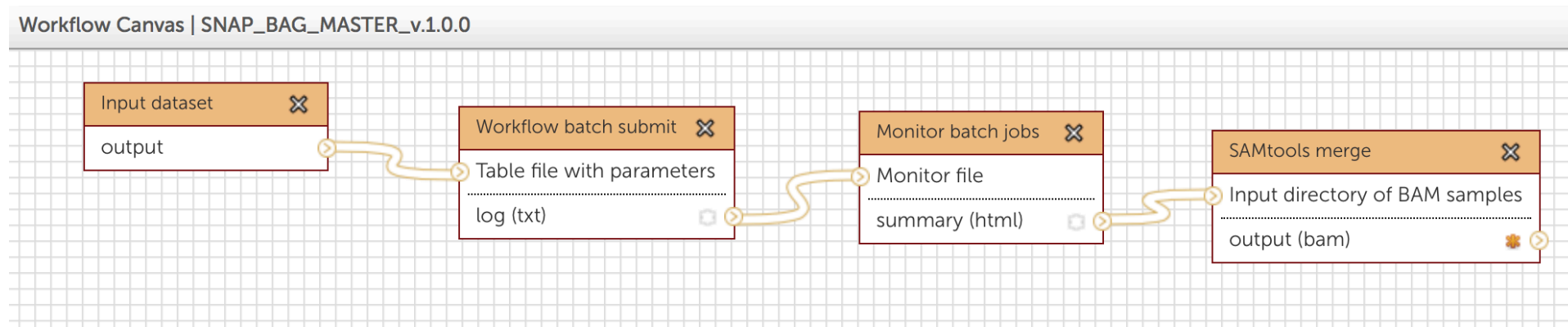
The screenshot shows the "globus Genomics" interface. The top navigation bar includes "Analyze Data", "Workflow", "Shared Data", "Visualization", "Admin", "Help", and "User". On the left, a "Tools" sidebar contains a search bar and a list of tools under the "DATA TRANSFER" category: "Globus Data Transfer", "Get Data", and "File Transfer Checksum". The main area displays the configuration for the tool "Get BDBag from MINID transfer data given a MINID to a bag dataset object (Galaxy Tool Version 1.0.0)". The "MINID" field is populated with the value "ark:/999999/fk4dj5pv64". An "Execute" button is visible below the input field.



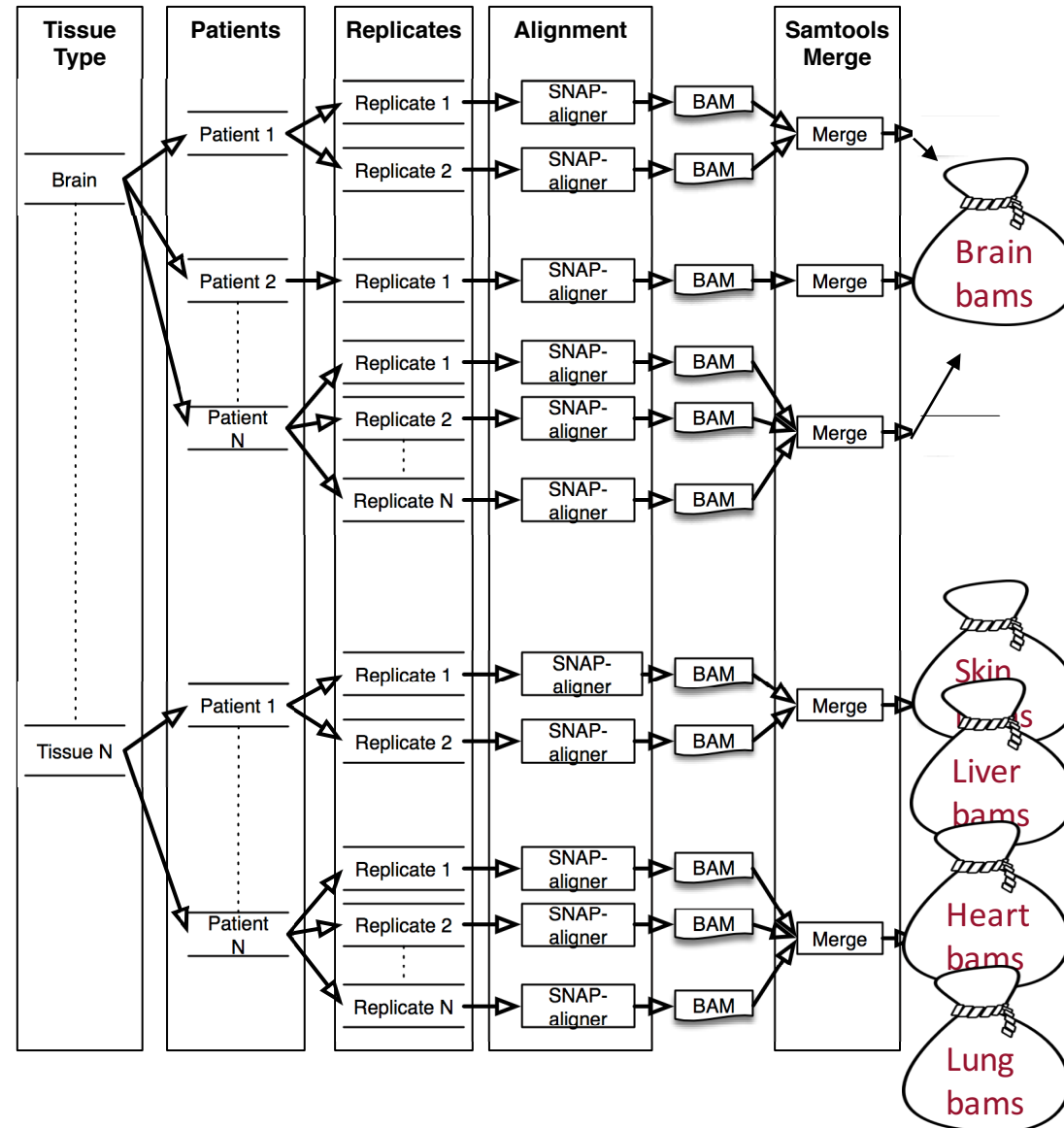
<https://github.com/ini-bdds/bdbag>



# BDDS Solutions: Enabling TReNA – Analysis pipelines

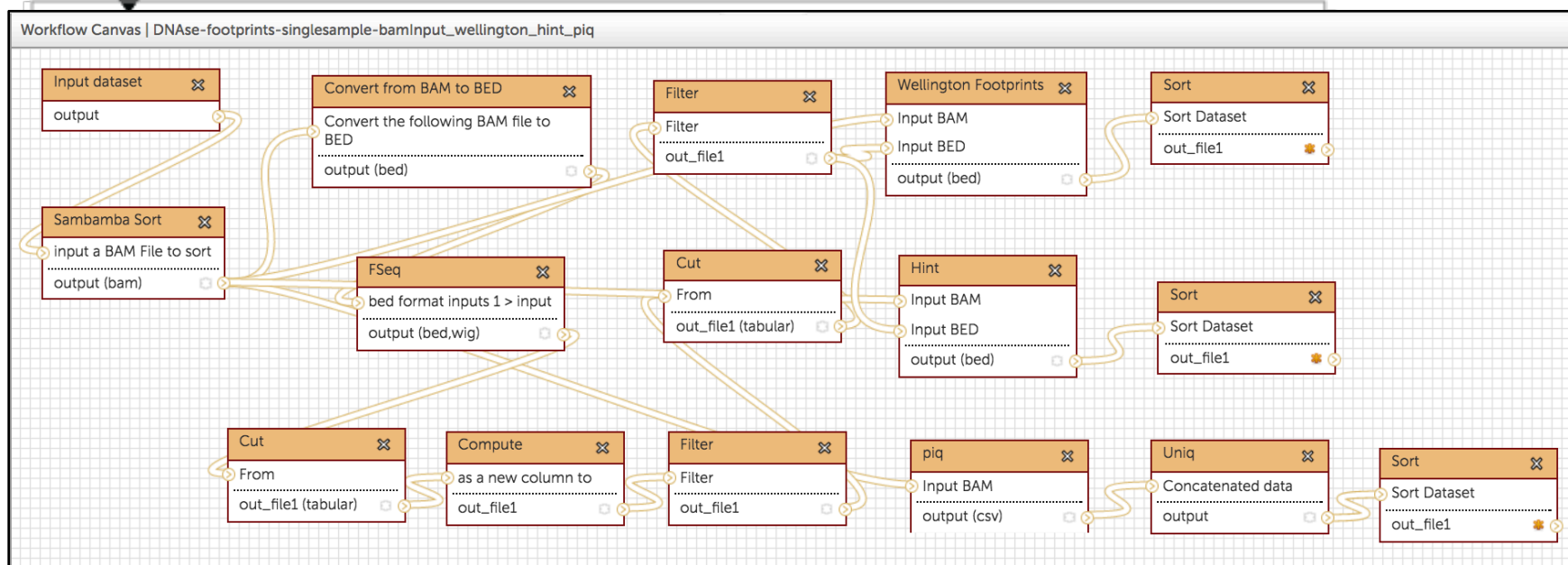


- Take a tissue sample bag (i.e. brain, skin, etc)
- Submit each bag to the alignment with the latest human reference genome GRCh38
- Merges samples that are from the same patient group





27 bags



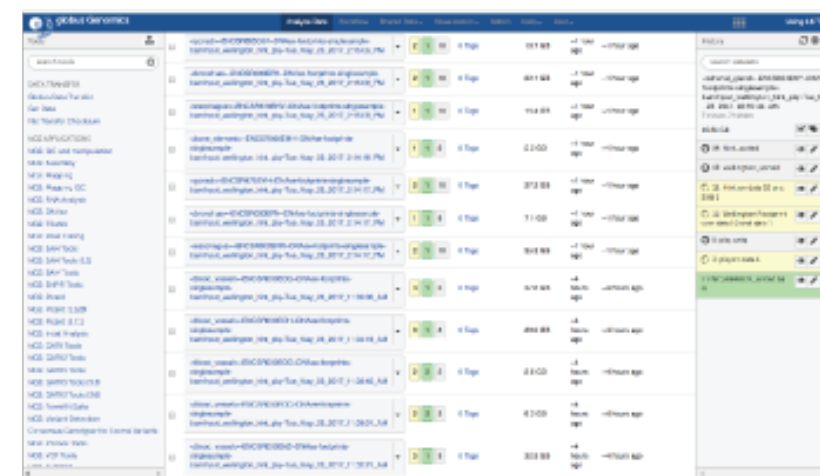
Translation of command lines to a galaxy workflow

```

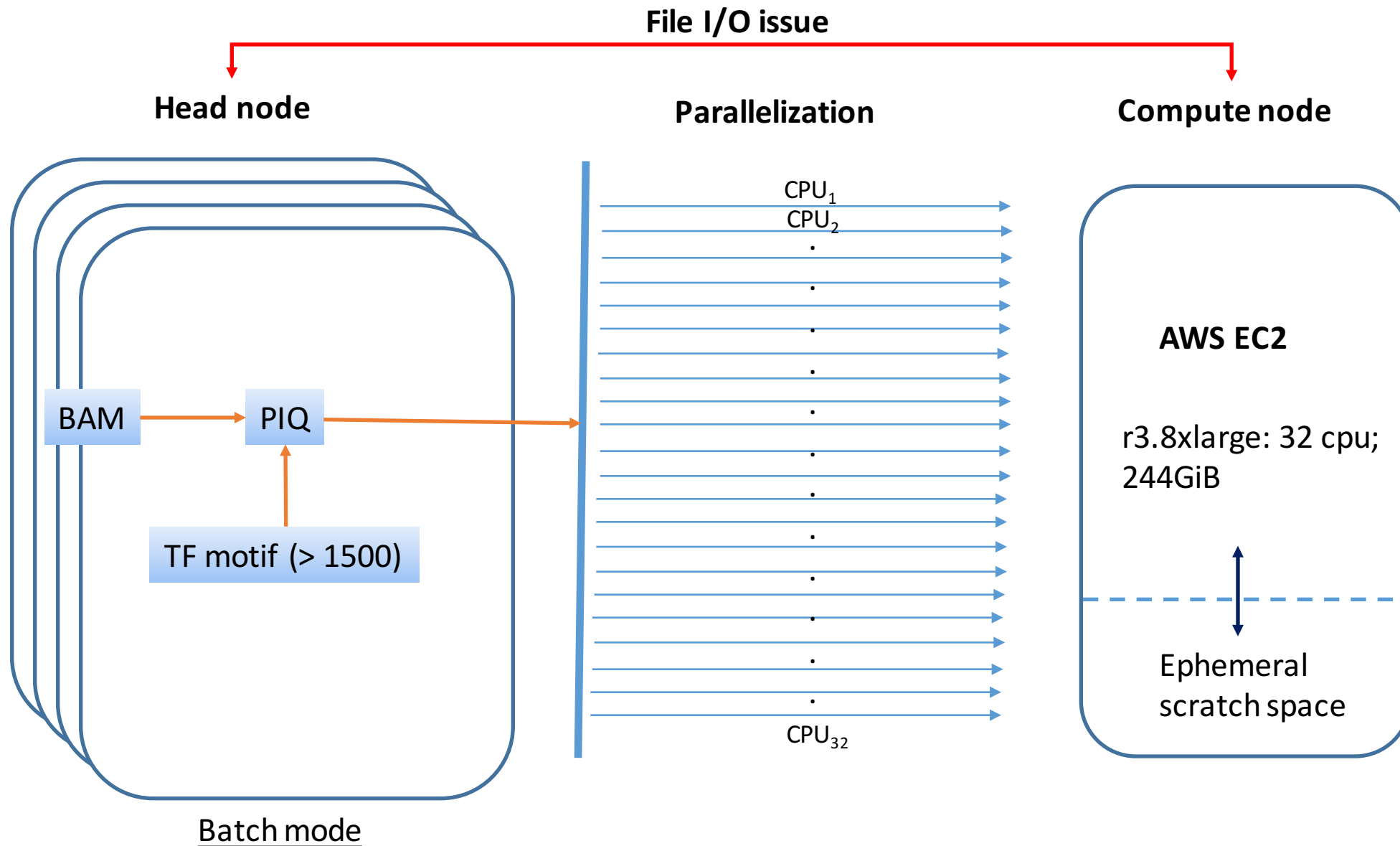
# This creates an LDHm file for the merged bam files, converts them to bed, runs Fseq and Wellington
#!/bin/bash
ls -l |> /usr/bin/awk -F= {print $1} |> /usr/bin/sort -n -k1,1
cd /usr/bin/sort -n -k1,1
mkdir -p /usr/bin/sort -n -k1,1
for i in $(ls -l |> /usr/bin/awk -F= {print $1} |> /usr/bin/sort -n -k1,1); do
  echo $i
done
# ... (rest of the script)

```

Running at scale with batch submission









BDDS / BDDS
Analyze Data Workflow Shared Data Admin Help User
Using 982.5 GB

Tools

ENCODE Tools

Lift-Over

Text Manipulation

Filter and Sort

Join, Subtract and Group

Convert Formats

Extract Features

Fetch Sequences

Fetch Alignments

Get Genomic Scores

Operate on Genomic Intervals

**STATISTICAL TOOLS**

Statistics

Wavelet Analysis

Graph/Display Data

Regional Variation

Multiple regression

Multivariate Analysis

**FASTA TOOLS**

Evolution

Motif Tools

Multiple Alignments

Metagenomic analyses

FASTA manipulation

NCBI BLAST+

Ontology services

**DATA MANAGEMENT**

History Management

Data Compression

Batch Management

Optimized Workflows

- DNA Exome Variant Analysis Optimized Workflow (BDDS)
- DNase Analysis Optimized Workflow (BDDS)

WORKFLOWS

- All workflows

### Saved Histories

Advanced Search

Name	Datasets	Tags	Sharing	Size on Disk	Created	Last Updated ↑	Status
<input type="checkbox"/> BDDS-GM19240~DNase_Optimized~Tue_Aug_09_2016_2:33:35_AM	5	0 Tags		42.8 MB	~ 12 hours ago	~ 10 hours ago	current history
<input type="checkbox"/> ~GM19239~DNase_Optimized~Tue_Aug_09_2016_2:33:19_AM	5	0 Tags		36.6 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM19238~DNase_Optimized~Tue_Aug_09_2016_2:33:03_AM	5	0 Tags		50.8 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM18507~DNase_Optimized~Tue_Aug_09_2016_2:32:45_AM	5	0 Tags		12.3 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM13976~DNase_Optimized~Tue_Aug_09_2016_2:32:27_AM	5	0 Tags		31.8 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM12892~DNase_Optimized~Tue_Aug_09_2016_2:32:11_AM	5	0 Tags		46.9 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM12891~DNase_Optimized~Tue_Aug_09_2016_2:31:55_AM	5	0 Tags		39.9 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM12878~DNase_Optimized~Tue_Aug_09_2016_2:31:38_AM	2	3	0 Tags	0 bytes	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM12865~DNase_Optimized~Tue_Aug_09_2016_2:31:23_AM	5	0 Tags		217.4 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM12864~DNase_Optimized~Tue_Aug_09_2016_2:31:09_AM	5	0 Tags		26.9 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM10248~DNase_Optimized~Tue_Aug_09_2016_2:30:55_AM	5	0 Tags		79.1 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> ~GM06990~DNase_Optimized~Tue_Aug_09_2016_2:30:38_AM	5	0 Tags		163.4 MB	~ 12 hours ago	~ 12 hours ago	
<input type="checkbox"/> batch	4	0 Tags	Shared	24.8 KB	May 19, 2016	~ 12 hours ago	

Your History

BDDS-GM19240~DNase\_Optimized~Tue\_Aug\_09\_2016\_2:33:35\_AM

42.8 MB

5: DNase Analysis Optimized Workflow (BDDS) on data 2 and data 1: minid  
6 lines  
format: txt, database: hg19

checksum for /scratch/galaxy/files/007/d...  
the TEST entity 03448f61c633791a93e7aa...  
HTTP connection (1): minid.bd2k.org  
identifier  
HTTP connection (1): minid.bd2k.org  
ed minid: ark:/99999/fk4wh2ts9q

4: DNase Analysis Optimized Workflow (BDDS) on data 2 and data 1: log

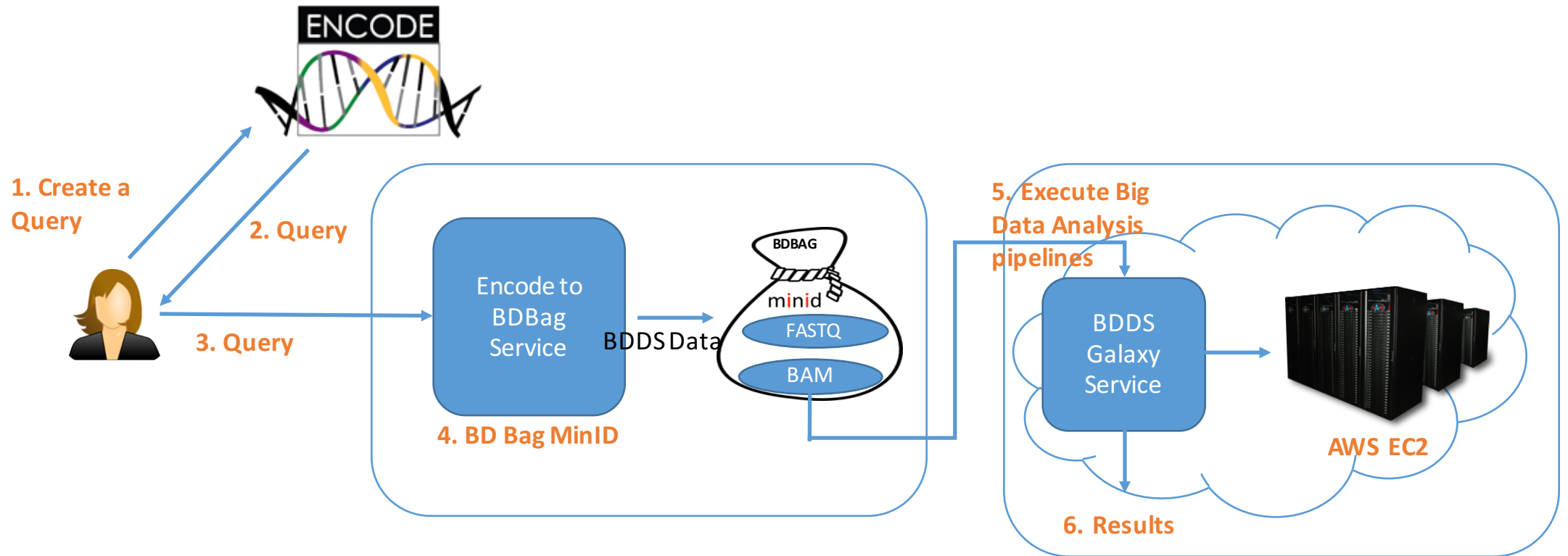
3: DNase Analysis Optimized Workflow (BDDS) on data 2 and data 1: bed

2: all\_motifs.meme

1: Homo\_sapiens.GRCh38.dna.primary\_assembly\_sorted.fa



# Dnase Hypersensitivity Analysis



## Results from TReNA analysis

- Total number of tissues: 27
- Total number of patient samples: 206
- Total number of fastq files: 1379
- Total size of the raw data: 2.5 TB
- Number of new tools added: 20
- Number of HPC workflows created: 11
- Number of compute hours for alignment: 24,000 CPU hours
- Number of compute hours for footprinting: 150,000 CPU hours
- Number of databases created: 106
  
- Notably, all the work was completed within two weeks which was originally expected for several months