# Reproducible big data science: A case study in continuous FAIRness

## Kyle Chard

University of Chicago and Argonne National Laboratory

Ravi Madduri, Michael D'Arcy, Segun Jung, Alexis Rodriguez, Dinanath Sulakhe, Eric Deutsch, Cory Funk, Ben Heavner, Matthew Richards, Paul Shannon, Gustavo Glusman, Nathan Price, Carl Kesselman and Ian Foster

# Reproducibility requires continuous FAIRness

- Make all data findable, accessible, interoperable, reusable **at every stage**, via pervasive use of simple identifier and exchange format conventions

- Build on proven **security, data, identifier, and computation building blocks** that have large user communities inside and outside biomedicine

- Leverage industry best practices to meet **scalability, interoperability, sustainability, and reliability** needs
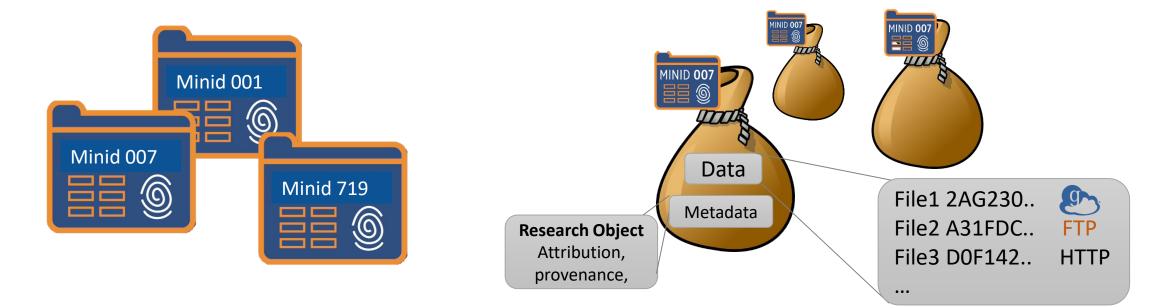
# Interoperability: naming and exchange

## Minid

- Lightweight identifiers for any product at any stage
- Easily created, dereferenced, validated
- Global integrity – validate content across the commons

## BDBag

- Self-describing and flexible format for exchange
- Extended BagIt Specification
- Standard manifest representation that supports different protocols
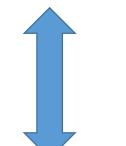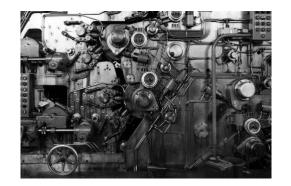- Research Object metadata

## Atlas of Transcription Factor Binding Sites from ENCODE DNase Hypersensitivity Data Across 27 Tissue Types

Cory C Funk, Segun Jung, Matthew A Richards, Alex Rodriguez, Paul Shannon, Rory Donovan, Ben Heavner, Kyle Chard, Yukai Xiao, Gustavo Glusman, Nilufer Erteskin-Taner, Todd Golde, Arthur Toga, Leroy Hood, John D Van Horn, Carl Kesselman, Ian Foster, Seth Ament, Ravi Madduri, Nathan D Price

## Reproducible big data science: A case study in continuous FAIRness

ⓘD Ravi K Madduri, Kyle Chard, Mike D'Arcy, Segun C Jung, Alexis Rodriguez, Dinanath Sulakhe, Eric W Deutsch, Cory Funk, Ben Heavner, Matthew Richards, Paul Shannon, Gustavo Glusman, Nathan Price, Carl Kesselman, ⓘD Ian Foster

# Generation of TFBS Atlas



Vernot et al.  doi: 10.1101/gr.134890.111

- Uniform processing of next generation sequencing data
  - Align to reference genome
  - Identify DNase hypersensitve regions
  - Apply multiple footprinting algorithms to locate putative transcription factor binding sites (TFBSs)
- Evaluate confidence in putative TFBSs
- Use TFBSs as features for machine learning approaches applied to disease-specific research

# Footprints Master Workflow

**A**

Dnase-seq

🔧 Get BDBag from MINID ✖
output1 (txt)

🔧 Create batch for patient SNAP alignment workflow
Bag location metadata
output (txt)

🔧 SNAP Workflow Batch Submit ✖
Table file with parameters
log (txt)

🔧 Generate SNAP BAG ✖
Monitor file
summary (html)

Align-ments

🔧 Create batch for patient footprints workflow ✖
BAM Data Object
output (txt)

🔧 Footprints Workflow Batch Submit ✖
Table file with parameters
log (txt)

Foot-prints

🔧 Generate Footprints ✖ BAG
Monitor file
summary (html)

**B**  Sub-Workflow for sequence alignment and BAM merge
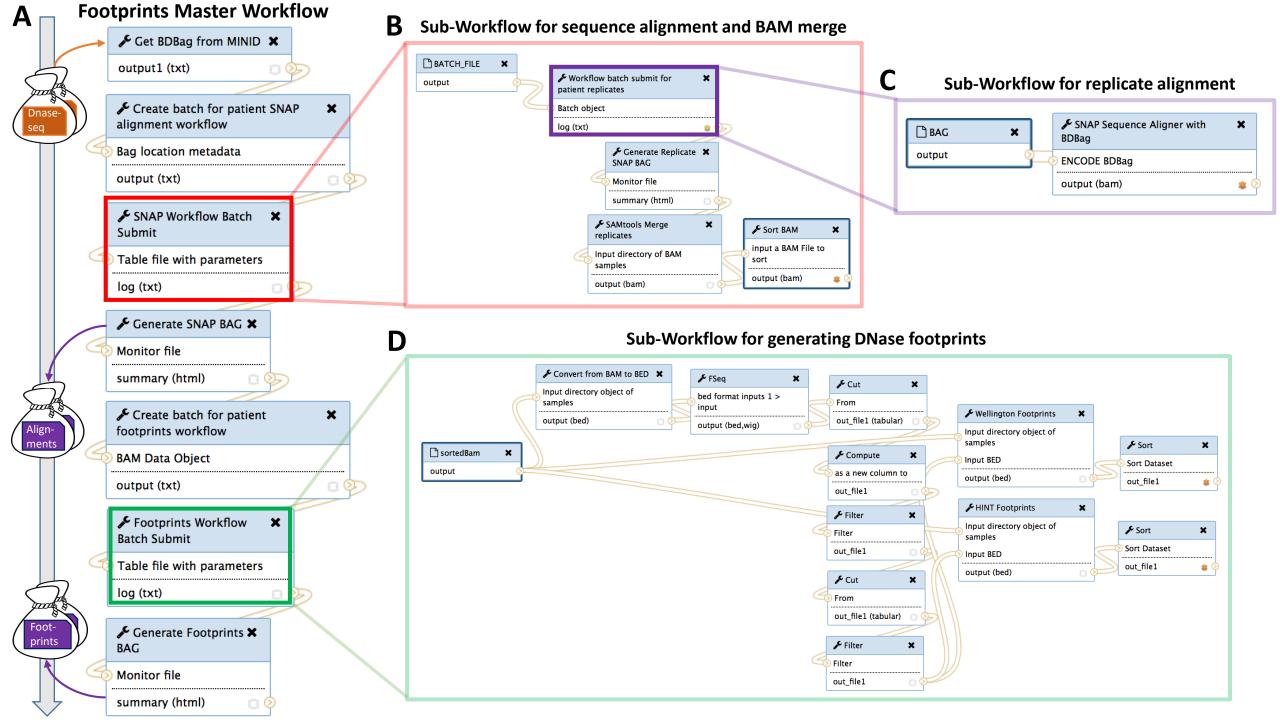
📄 BATCH_FILE ✖
output

🔧 Workflow batch submit for patient replicates ✖
Batch object
log (txt)

🔧 Generate Replicate SNAP BAG ✖
Monitor file
summary (html)

🔧 SAMtools Merge replicates ✖
Input directory of BAM samples
output (bam)

🔧 Sort BAM ✖
input a BAM File to sort
output (bam)

**C**  Sub-Workflow for replicate alignment

📄 BAG ✖
output

🔧 SNAP Sequence Aligner with BDBag ✖
ENCODE BDBag
output (bam)

**D**  Sub-Workflow for generating DNase footprints

🔧 Convert from BAM to BED ✖
Input directory object of samples
output (bed)

🔧 FSeq ✖
bed format inputs 1 > input
output (bed,wig)

🔧 Cut ✖
From
out_file1 (tabular)

📄 sortedBam ✖
output

🔧 Compute ✖
as a new column to
out_file1

🔧 Filter ✖
Filter
out_file1

🔧 Cut ✖
From
out_file1 (tabular)

🔧 Filter ✖
Filter
out_file1

🔧 Wellington Footprints ✖
Input directory object of samples
Input BED
output (bed)

🔧 Sort ✖
Sort Dataset
out_file1

🔧 HINT Footprints ✖
Input directory object of samples
Input BED
output (bed)

🔧 Sort ✖
Sort Dataset
out_file1

# Reproducibility

**1) Datasets**

| # | Name | Identifier | Role | Description | Size |
|---|------|-----------|------|-------------|------|
| D1 | DNase-seq | minid:b9dt2t | In | BDBag of 27 BDBags extracted from ENCODE by ❶, one per tissue: 1,591 FASTQ files in all. | 2.40 TB |
| D2 | Alignment | minid:b9vx04 | Out | BDBag of 54 BDBags produced by ❷, 1 per {tissue, seed}: 386 BAM files in all. | 5.30 TB |
| D3 | Footprints | minid:b9496p | Out | BDBag of 54 BDBags containing footprints computed by ❸, one per {tissue, seed}. Each BDBag contains two BED files per biosample, one per footprinting method. | 0.04 TB |
| D4 | Motifs | minid:b97957 | In | Database dump file containing the non-redundant motifs provided by Funk et al. [13]. | 31.5 GB |
| D5 | Hits | minid:b9p09p | Out | Database dump file containing the hits produced by ❹. | 0.04 TB |
| D6 | TFBSs | minid:b9v398 | Out | BDBag of 54 BDBags containing candidate TFBSs produced by ❺, one per {tissue, seed}. Each BDBag contains two database dump files, one per footprinting method. | 0.35 TB |

**1) Tools**

| # | Name | Identifiers for software |
|---|------|--------------------------|
| ❶ | Extract DNase-Seq | encode2bag service: https://github.com/ini-bdds/encode2bag-service  encode2bag client: https://github.com/ini-bdds/encode2bag |
| ❷, ❸ | Alignment, Footprints | Galaxy pipeline: minid:b93m4q  Dockerfile: minid:b9jd6f  Docker image: minid:b97x0j |
| ❹ | Hits | R script: minid:b9zh5t |
| ❺ | TFBSs | R scripts: minid:b9fx1s |

# Summary

- Complexity and level of effort increases with size of the data, computation, analysis team
  - TFBS example: big data (~10 TB), parallel computing (70K core hours), distributed data and processing, handful of researchers ...
- Continuous FAIRness requires various tools & services
  - Identifiers, BDBags, ROs, Containers,
  - Globus, Globus Genomics, GitHub
- User Study
  - 10/11 students and researchers were able to reproduce this study without assistance