**Brigham and Women's Hospital**
Founding Member, Mass General Brigham

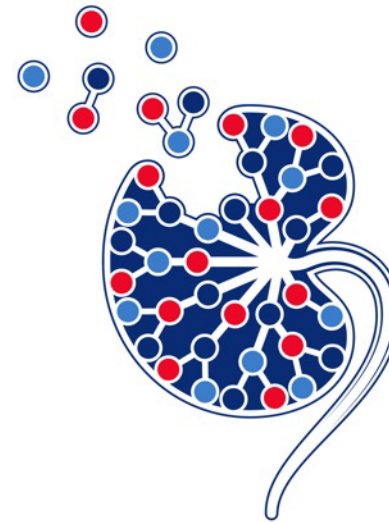# Metadata and Data Standards for NIDDK Research Data - The ATLAS-D2K Experience

M. Todd Valerius, Ph.D.

Brigham and Women's Hospital / Harvard Medical School

BWH - Renal Division

**HARVARD MEDICAL SCHOOL**
TEACHING HOSPITAL

# The ATLAS-D2K Center

A kidney and lower urinary tract-focused data discovery hub with access to visualizations and analysis tools.

*Bringing GUDMAP & RBK data under one ATLAS that embraces open science and provides links to related consortiums.*
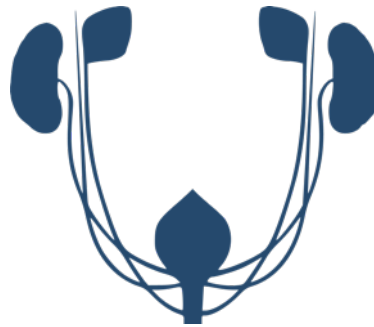
# Overall Aims:

Our long-term goal is to bring complex data into an accessible form for our research community.

Establish connections between molecular data of kidney and lower urinary tract present in GUDMAP, RBK, KPMP, HuBMAP, and the HCA.

Enable researchers of various levels of experience by providing tools to interact with the data.

# GenitoUrinary Development Molecular Anatomy Project (GUDMAP) & (Re)Building a Kidney (RBK): Overview

Overarching Program Goal:

- **GUDMAP:** high resolution molecular anatomy of the developing and mature genitourinary system (mouse, human, rat, dog)
- **RBK:** optimize differentiation of human kidney cell types in defined structures, and determine methods to promote kidney repair, to generate or repair nephrons that can function within the kidney (human, human iPSCs, zebrafish)

Number of investigators involved: GUDMAP: 9    RBK: 25

Technology Focus: array of gene expression techniques on tissues, differentiation of stem cells, a range of imaging techniques

Current Gaps? robust anatomical ontologies broadly implemented, metadata standardization, interactive tools for data analysis/data annotation (e.g., cluster data)

# Ontologies and controlled vocabularies

- Why is this important?
  - Gene expression and function occurs in tissues. Consistent use of names removes confusion amongst researchers and **enables** computation of complex queries.
  - Quickly apparent when trying to connect data
    - GUDMAP had generated thousands of wholemount & section *in situ* hybridizations, scored for expression, from two groups. -> **anatomical ontology** needed to connect.

Abler, L.L. *et al.* (2011) *Developmental dynamics : an official publication of the American Association of Anatomists*. https://doi.org/10.1002/dvdy.22730.
Georgas, K.M. *et al.* (2015) *Development (Cambridge, England)*. https://doi.org/10.1242/dev.117903.
Harding, S.D. *et al.* (2011) *Development (Cambridge, England)*. https://doi.org/10.1242/dev.063594.
Henry, G.H. *et al.* (2018) *Cell Rep*. https://doi.org/10.1016/j.celrep.2018.11.086.
Little, M.H. *et al.* (2007) *Gene expression patterns : GEP*. https://doi.org/10.1016/j.modgep.2007.03.002.

# Boolean Search on Scored Expression

# Anatomical annotation using established terms

Manual annotation of structures

Links to other data through anatomy

# Ontologies and controlled vocabularies

- Gene expression and function occurs in tissues. Consistent use of names removes confusion amongst researchers and *enables* computation of complex queries.

- Quickly apparent when trying to connect data
  - GUDMAP had generated thousands of wholemount & section *in situ* hybridizations, scored for expression, from two groups. -> *anatomical ontology* needed to connect.

- To accommodate cross-species data, we use Uberon and Cell Ontology
  - Uberon multi-species anatomy ontology
  - The Cell Ontology

- Healthy adult human tissue focus in **HuBMAP ASCT+B Tables**
  - Data-driven effort lead by Sanjay and enhanced by many.

Recommendation:

Select a source of anatomical and cell type terms that fit your research, use them as a standard, and capture the source of those terms from established ontologies.

# Data formats:
## Is anything as future proof as plain text?

**RAW sequence data** is, but there are privacy issues.

- Detailed *sequencing metadata* is produced and captured by computational tools.
- *Biosample metadata* needs to be captured well at the time of experiment.
- Protocols should be well referenced. (consortiums rely on self hosting OR commercial repositories like Protocols.io

**Image data** is a mature, poor example

- Center on "open" established standards like OME-TIFF and related formats (e.g., Zarr and OME-NGFF) that capture microscopy/imaging metadata.
  - The benefit: these formats work well with open-source software like ImageJ and QuPath.
  - Images adjusted for publication and presentation <u>are not useful for downstream reuse and quantitative analysis</u>.
- The "biosample metadata" associated with an image needs to be captured EARLY in the process.

# Data Curation/Interaction UI: Gene



Direct to expression data associated with this gene

Search for presence of different types of expression data, scored exp. region, etc.

Direct to all data associated with this gene

Direct to different specimens associated with this gene

# Normal anatomical and structural changes

## 3-D Mapping of Tissues

# Data formats:
Is anything as future proof as plain text?

## RAW sequence data is, but there are privacy issues.

- Detailed *sequencing metadata* is produced and captured by computational tools.
- *Biosample metadata* needs to be captured well at the time of experiment.
- Protocols should be well referenced. (consortiums rely on self hosting OR commercial repositories like Protocols.io

## Image data is a mature, poor example

- Center on "open" established standards like OME-TIFF and related formats (e.g., Zarr and OME-NGFF) that capture microscopy/imaging metadata.
  - The benefit: these formats work well with open-source software like ImageJ and QuPath.
  - Images adjusted for publication and presentation <u>are not useful for downstream reuse and quantitative analysis</u>.
- The "biosample metadata" associated with an image needs to be captured EARLY in the process.

Recommendations:
1. Develop a plan to capture biosample metadata and protocol with the sequencing data.
2. Capture and associate biosamples metadata and consolidate on a lossless image file standard.

# Data storage and levels of sharing, accessibility to open-source tools.

- **Sequence data – a layered approach**
  - "Processed" data is useful to a wider range of researchers.
    - Count files are ready for analysis without computationally intensive genome aligners in HPCs. More researchers can use such data immediately.
  - Privacy of participants - what is reasonable to share even with full consent?
    - The ability to de-identify participants from limited sequencing data expands before thoughtful policy will catch up. Think beyond the contractual protection to anticipate while maintaining data availability.
  - What intermediate products are available?
    - R objects like Seurat capture analysis decisions for a data generators fine analysis and can be used by less experienced researchers subsequently.
  - Processed data more freely shareable, but what happens when references change?
    - As data ages, re-alignment may be necessary as reference genomes change and improve.

# Direct linking to data for efficiency with large datasets



- New approach to scientific rigor and reproducibility
  - Data followed from slide to database image
  - "Largest possible" supplementary data

- Collections designed around specific structures

## Multiple layers with scRNA-seq

1. RAW sequencing files (fastq)
2. Processed gene expression matrix files (txt)
3. R objects of analysis (Rds, e.g., Seurat)
4. Static visualization tools
5. Interactive visualization tools

# Direct linking to data for efficiency with large datasets



- New approach to scientific rigor and reproducibility
  - Data followed from slide to database image
  - "Largest possible" supplementary data

- Collections designed around specific structures

## Multiple layers with scRNA-seq
1. RAW sequencing files (fastq)

# scRNA-seq Visualizations - Accessibility



- New approach to scientific rigor and reproducibility
  - Data followed from slide to database image
  - "Largest possible" supplementary data
- Collections designed around specific structures

Multiple layers with scRNA-seq
1. RAW sequencing files (fastq)
2. Processed gene expression matrix files (txt)
3. R objects of analysis (Rds, e.g. Seurat)
4. Static visualization tools
5. Interactive visualization tools

# Data storage and levels of sharing, accessibility to open-source tools.

- **Sequence data – a layered approach**
  - "Processed" data is useful to a wider range of researchers.
    - Count files are ready for analysis without computationally intensive genome aligners in HPCs. More researchers can use such data immediately.
  - Privacy of participants - what is reasonable to share even with full consent?
    - The ability to de-identify participants from limited sequencing data expands before thoughtful policy will catch up. Think beyond the contractual protection to anticipate while maintaining data availability.
  - What intermediate products are available?
    - R objects like Seurat capture analysis decisions for a data generators fine analysis and can be used by less experienced researchers subsequently.
  - Processed data more freely shareable, but what happens when references change?
    - As data ages, re-alignment may be necessary as reference genomes change and improve.

# mRNA-Seq Reanalysis Progress

| # Studies | # Experiments | # Replicates | # Files |
|---|---|---|---|
| 49 | 184 | 602 | 851 |

⬇ **First Round of execution**

| Execution Status | Count (# Replicates) | Description |
|---|---|---|
| Success | 110 (18.3%) | |
| Error | 492 (81.7%) | |
|   - Metadata | 311 (51.7%) | - Validate: Species, Paired-End, Strandedness, Spikes-in<br>- No metadata or mismatched<br>- Incorrect sequencing type (e.g. ChiP-Seq instead of mRNA-Seq) |
|   - File | 181 (30.0%) | - Mismatched #Reads of R1 and R2, multiple runs, missing files<br>- Not fastq structure (e.g. fastq+bam) |

⬇ **After a few rounds of resolution and execution**

| Success | - 549 mRNA-Seq replicates (45 studies)<br>- Re-labeled 3 mRNA-Seq replicates to ChiP-Seq (1 study) |
|---|---|
| Outstanding issues | - Missing files: 3 replicates<br>- Conflicting files: 50 replicates (3 studies) |

As of 08/10/2021

Malladi & Henry, UT Southwestern

# mRNA-Seq QC, Processed Files, and Visualization



QC data (including execution status)

User submitted seq files and hub-processed analysis files are all accessible

https://www.gudmap.org/id/Q-Y4GY

Interact visualization of TPM expression (group by Experiment, Anatomical Source, Stage, etc)

https://dev.gudmap.org/id/Q-Y4GY (dev server only)

Individual execution run contains workflow definition, version, source code URL, input and output for reproducibility

https://www.gudmap.org/id/17-BMCM

# Data storage and levels of sharing, accessibility to open-source tools.

- **Sequence data – a layered approach**
  - "Processed" data is useful to a wider range of researchers.
    - Count files are ready for analysis without computationally intensive genome aligners in HPCs. More researchers can use such data immediately.
  - Privacy of participants - what is reasonable to share even with full consent?
    - The ability to de-identify participants from limited sequencing data expands before thoughtful policy will catch up. Think beyond the contractual protection to anticipate while maintaining data availability.
  - What intermediate products are available?
    - R objects like Seurat capture analysis decisions for a data generators fine analysis and can be used by less experienced researchers subsequently.
  - Processed data more freely shareable, but what happens when references change?
    - As data ages, re-alignment may be necessary as reference genomes change and improve.

Recommendations:
1. Consider asking your core to use programmatically published base processing pipelines.
2. Publish these AND analysis code to a coding repository. (Contemporaneous documentation saves time later.)
3. Consider what intermediate layers should be protected (RAW versus counts).
4. Publish analysis intermediates.

# ATLAS-D2K Team



**Carl Kesselman**
Data Management, public cloud infrastructure, tools, coordination center

Data & Tool Integration

**Phil Blood** (*& Jonathan Silverstein, Co-I*)
Private cloud infrastructure, tools

**M. Todd Valerius, Sanjay Jain**
*Scientific Co-Directors*
Mission, Operation, Biocuration, Ontology

Tool Integration

**Katy Börner**
3D CCF tissue registration,
outreach/hackathons

*Matthias Kretzler (Co-I)*
Quality control processes, bioinformatic
analysis, UX expertise

DK135157, DK133090