



June 11, 2025



Building Data Workflows for Neuroscience and AI

Infrastructure × Operations × Collaboration

Dimitri Yatsenko, PhD
Thinh Nguyen, PhD





The Future of Research Operations

operational excellence - data integrity - governance - quality control - AI readiness

DataJoint Platform - computational workflows, future developments

DataJoint Elements - standardized workflows for neurophysiology

Data Publishing - DANDI & EMBER

MICrONS Project

Data access and analysis - DANDI, DataJoint, VORTEX

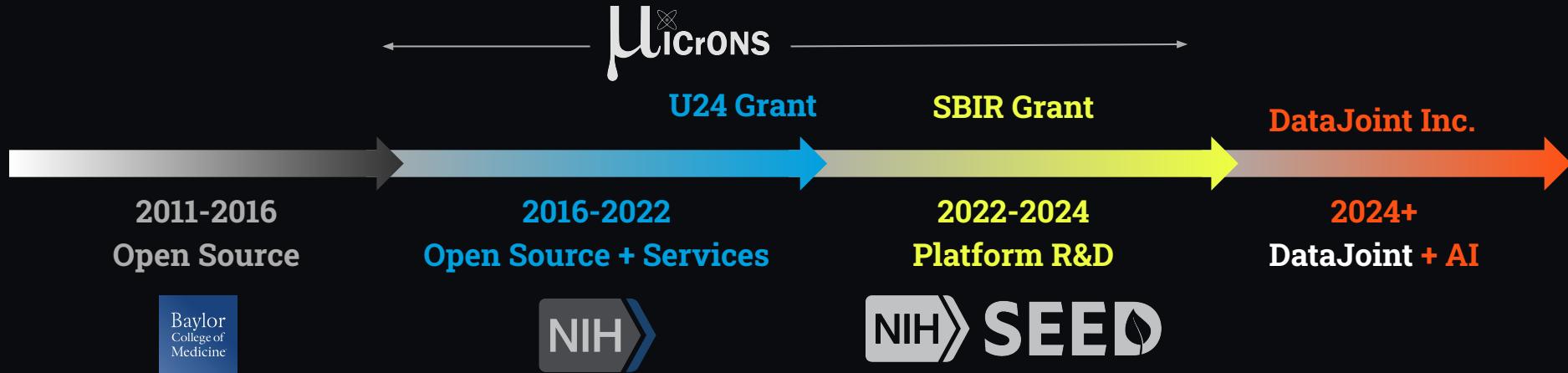
Analysis collaboration with Neuromatch, Stanford, led by Nima Dehghani MIT



AI in Research Operations

Co-Pilot → Co-Scientist → Co-Strategist

DataJoint: A Journey of Collaboration



*See [The U.S. Government Launches a \\$100-Million “Apollo Project of the Brain”](#) Scientific American (Mar. 8, 2016); [A milestone map of mouse-brain connectivity reveals challenging new terrain for scientists](#), Nature (Apr. 15, 2024).



OUR MISSION

To create a healthy operating environment for science that yields rapid, reproducible findings and aggregates data into a foundation for breakthrough discoveries.



Based on first principles of science & engineering.

SciOps Requirements

"SciOps" ≈ DevOps for Science

Reproducibility

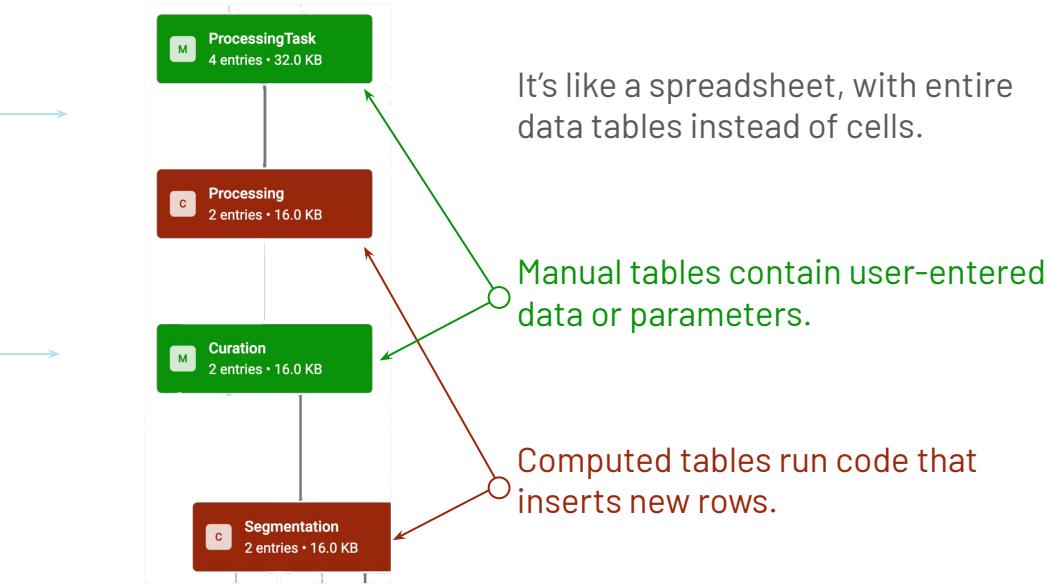
- ✓ Requires **relational integrity** among data, code, and process.

Flexibility + Resilience to Change

- ✓ Requires **unified management** of data, code, and process.

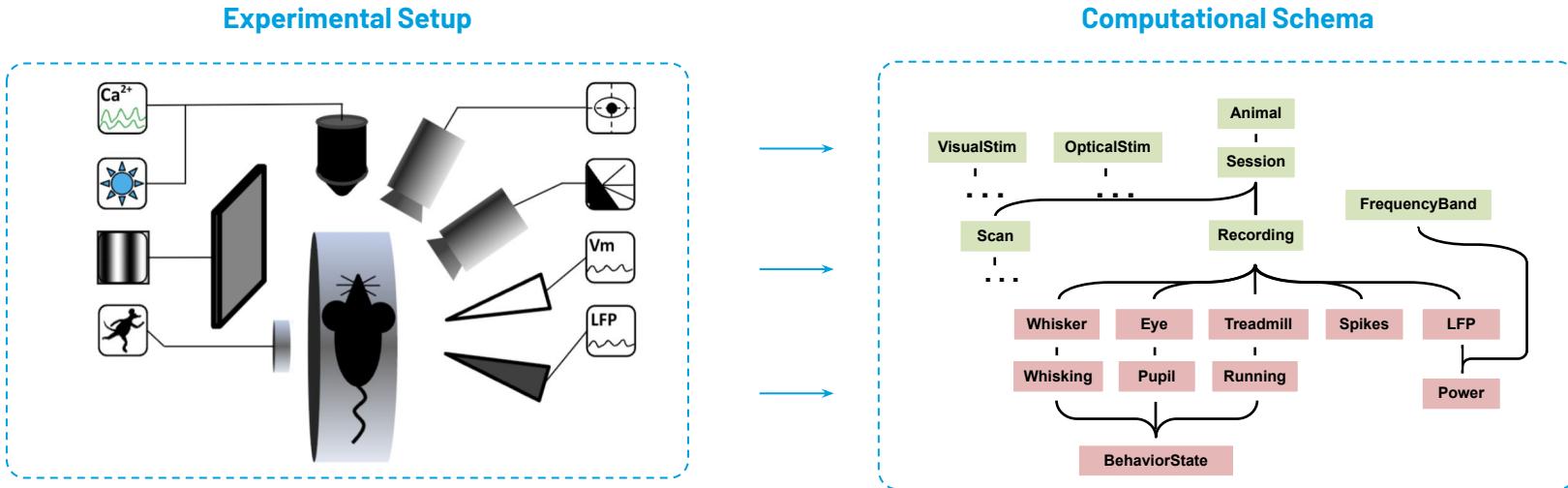
The Computational Database

A new database model





Express a study as a *pipeline* of data transformations.



Yatsenko et al (2015)
<https://arxiv.org/abs/1807.11104>



Demo



Open Source + Infrastructure + Operations

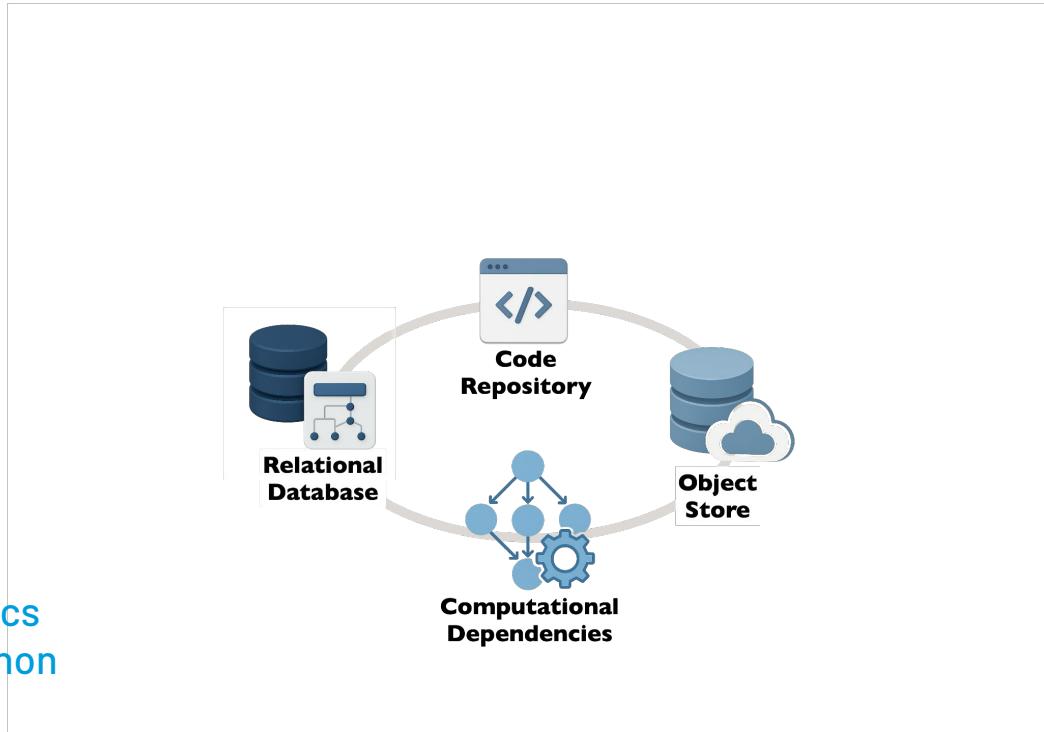
Open-source core

A flexible open standard for scientists to define all aspects of a study – so it can be *understood, validated, shared, and automated.*

Captures: code + data + dependencies

Operating in 100+ labs

<https://github.com/datajoint/datajoint-specs>
<https://github.com/datajoint/datajoint-python>





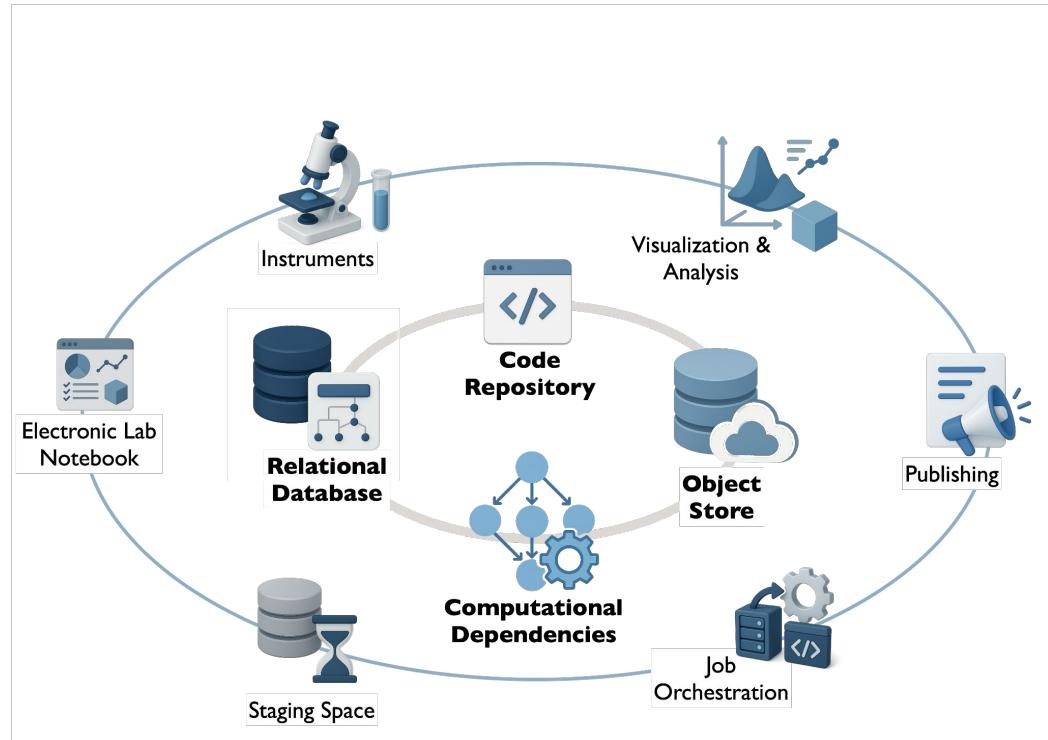
Open Source + Infrastructure + Operations

Open-source core

A flexible open standard for scientists to define all aspects of a study – so it can be **understood, validated, shared, and automated.**

Functional extensions for efficiency

Integration, automation, interfaces, security & compliance, sharing and publishing.





Open Source + Infrastructure + Operations

Open-source core

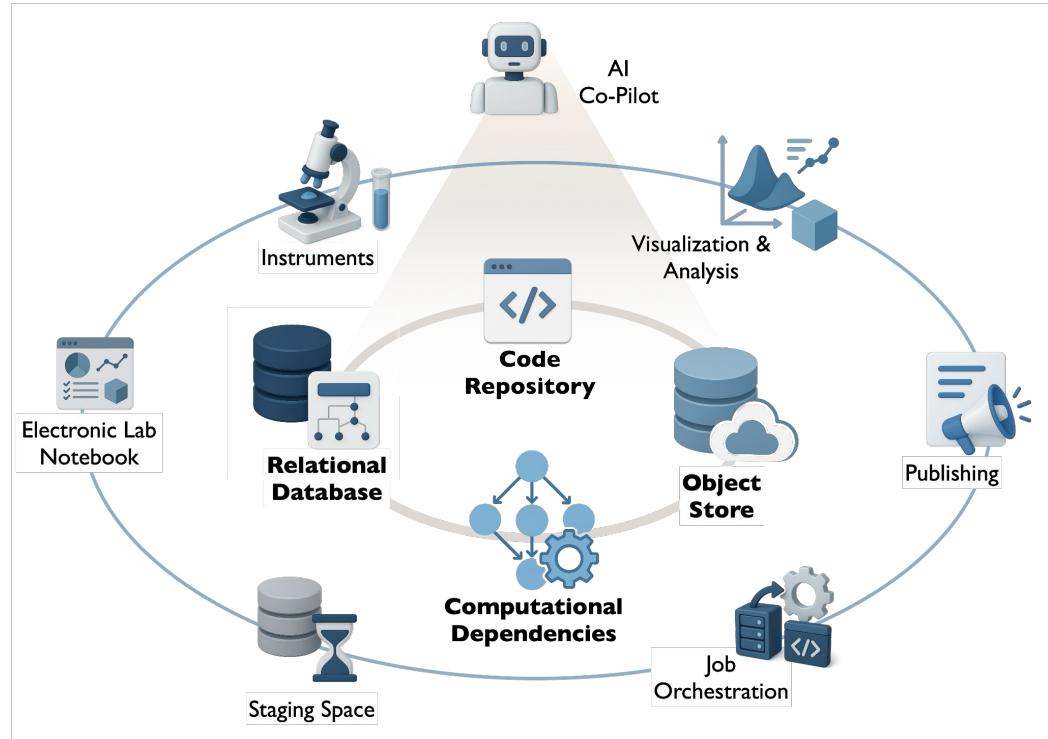
A flexible open standard for scientists to define all aspects of a study – so it can be **understood, validated, shared, and automated.**

Functional extensions for efficiency

Integration, automation, interfaces, security & compliance, sharing and publishing.

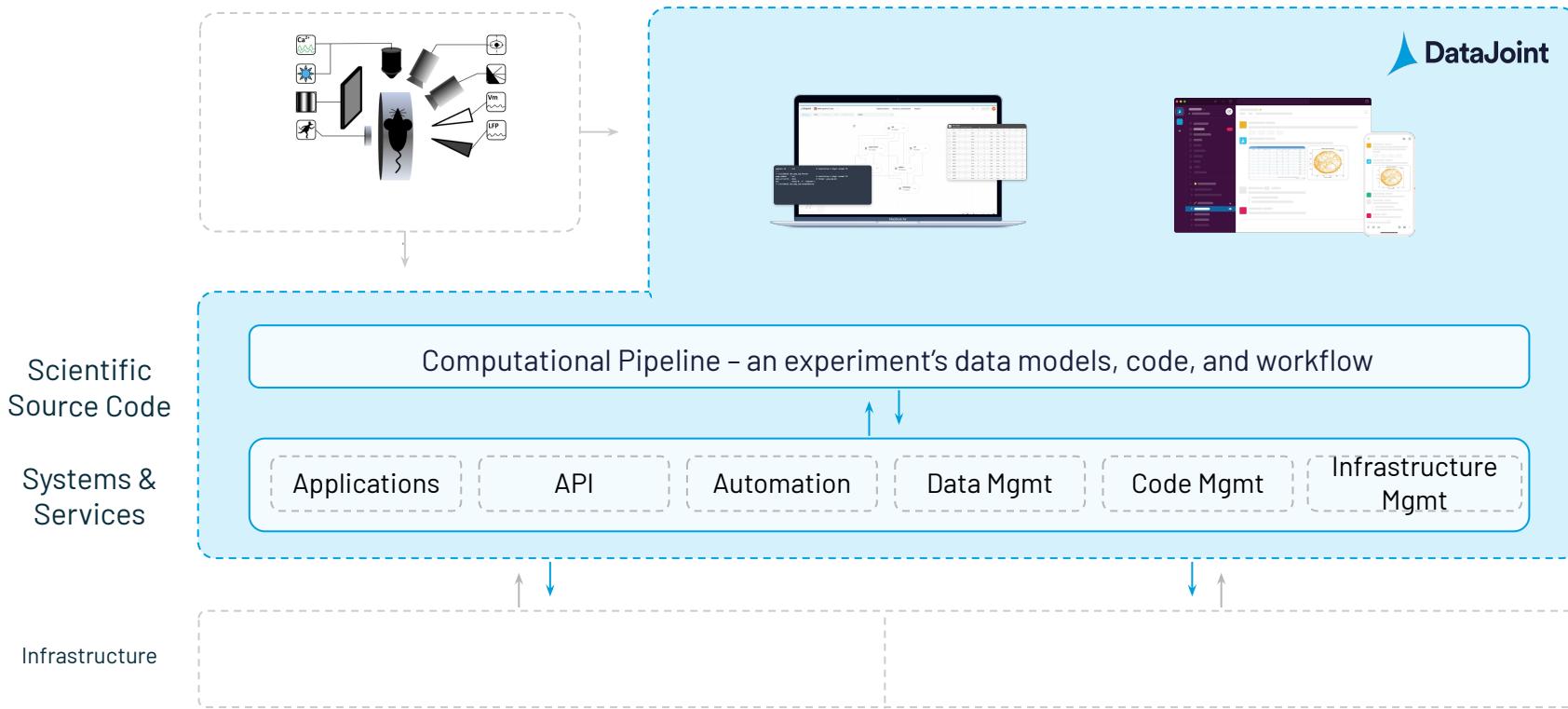
AI integration to expand capability

Data intelligence, coding assistance, knowledge tracking, writing assistance, strategy

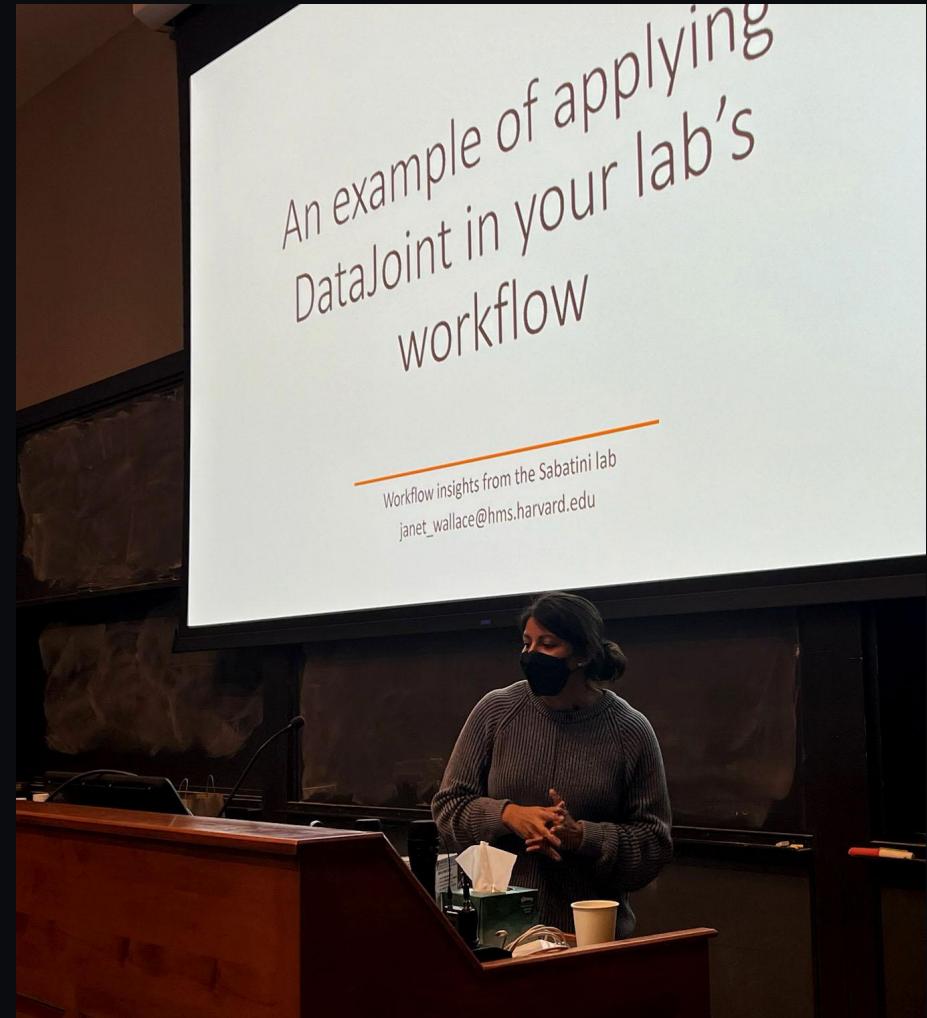




A full-stack operating platform for science.



Research teams
define their pipelines
and workflows





Modular pipelines

<https://docs.datajoint.com/elements/>

Common language

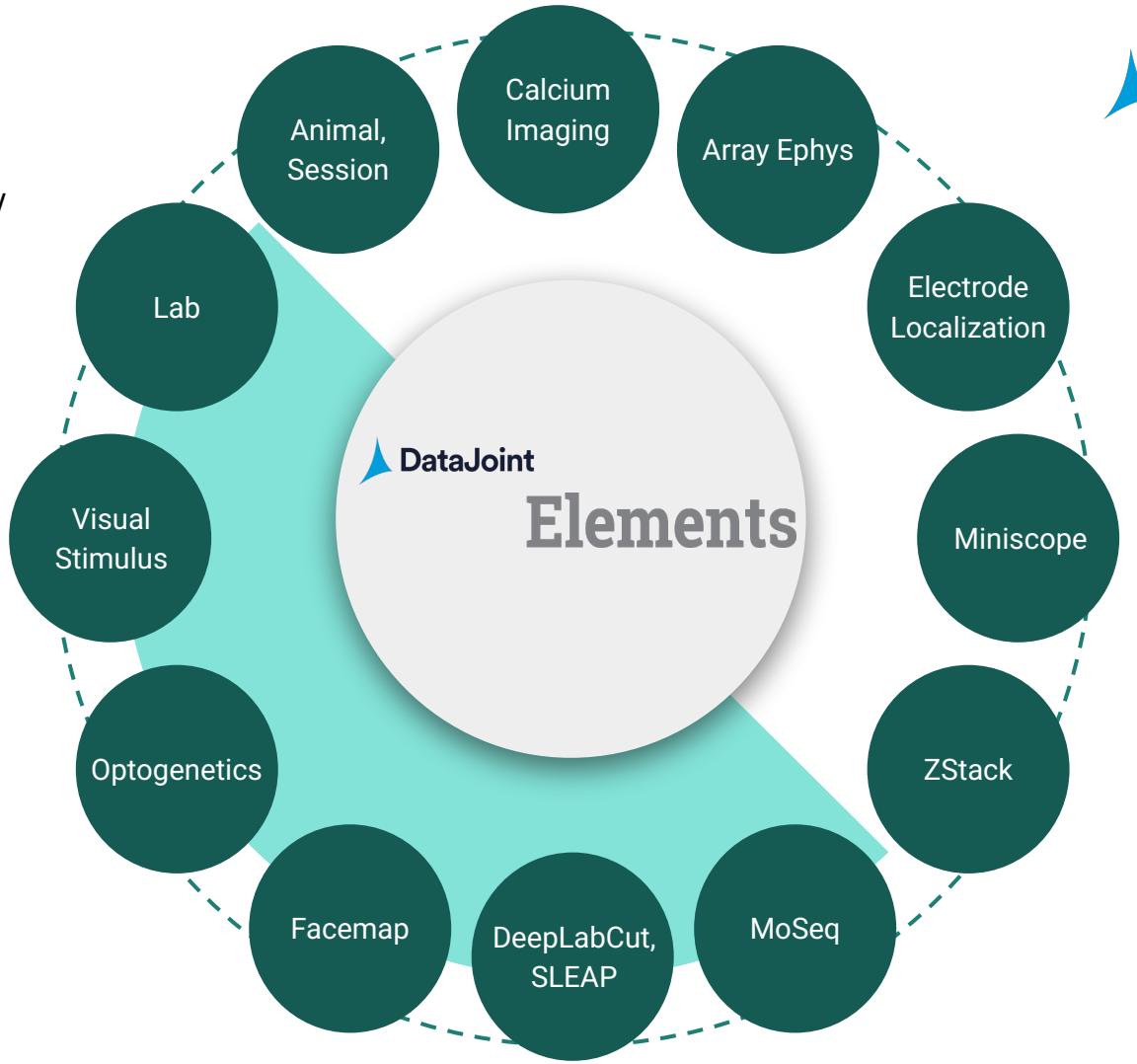
A flexible open standard for scientists to define all aspects of a study – so it can be **understood, validated, shared**, and **automated**.

Standardized modules

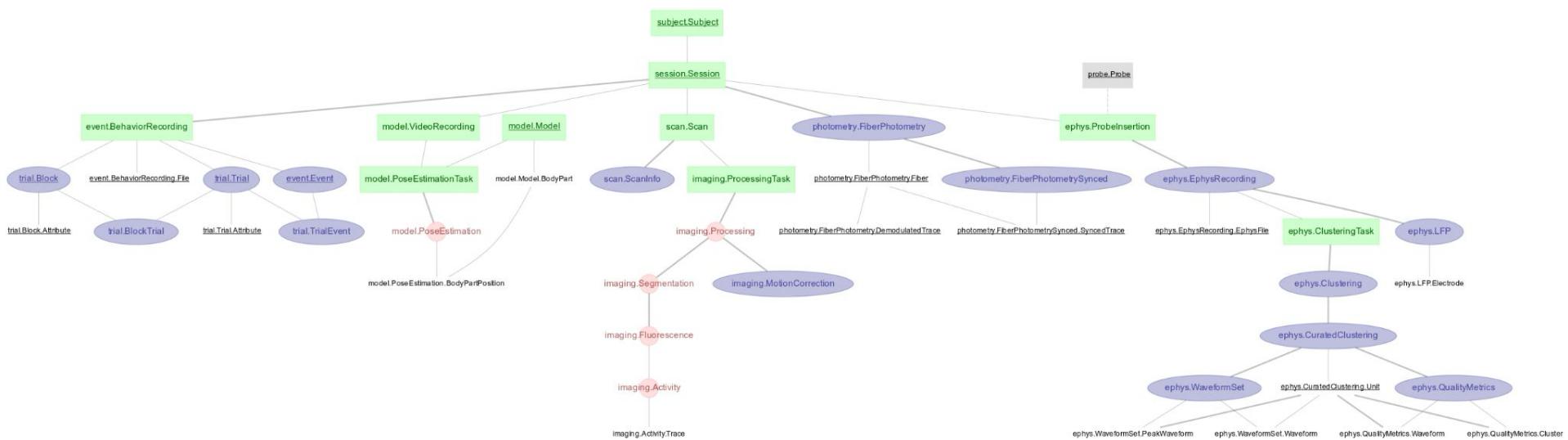
Validated open-source modules.
Integration, interfaces, customization

Uniform processes

Navigation, automation, queries,
visualization, sharing, publishing



Multimodal data pipeline



Operant Behavior
(trials, events)

Pose Estimation
(DLC)

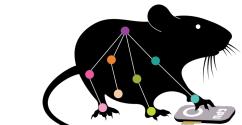
Calcium Imaging
(Suite2p)

Fiber
Photometry

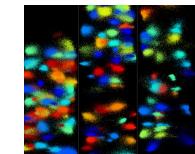
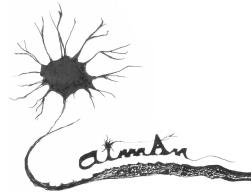
Ephys
(SpikeGLX,
Kilosort)

Open-source tools and informatics resources

Community, licensing, governance.

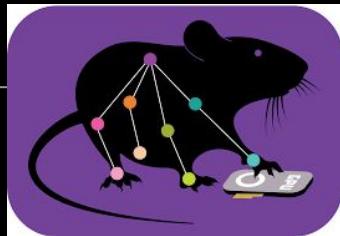


Allen Brain Observatory



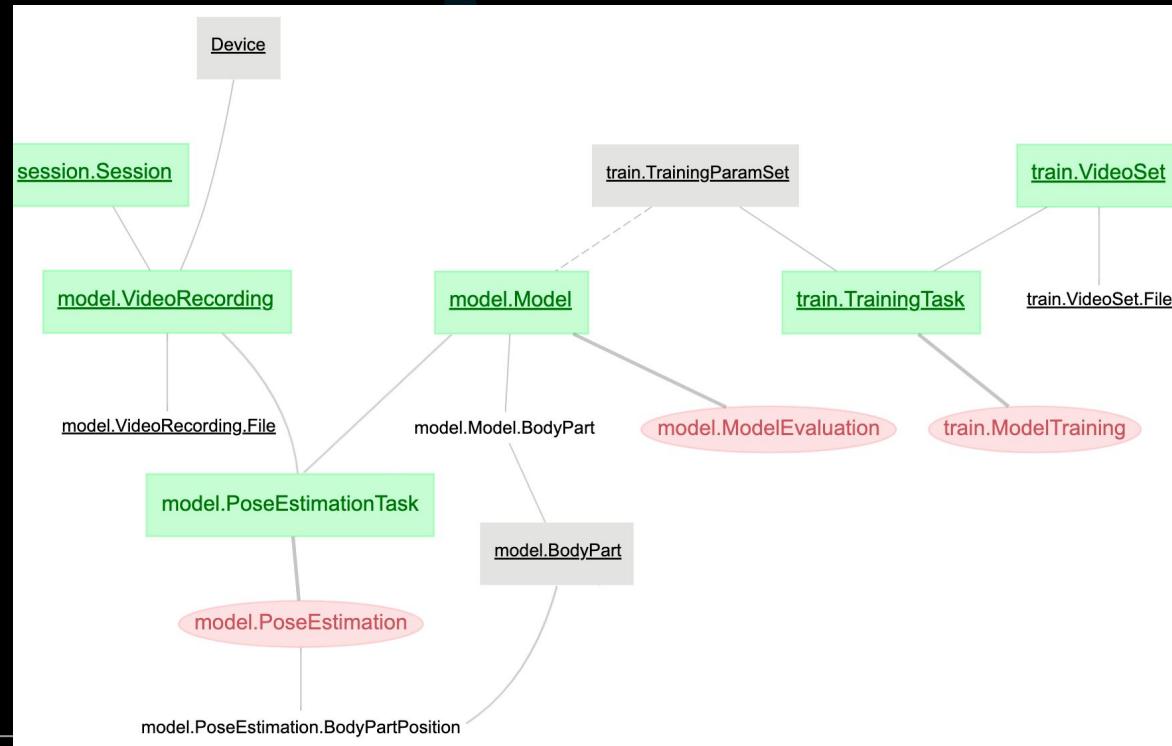
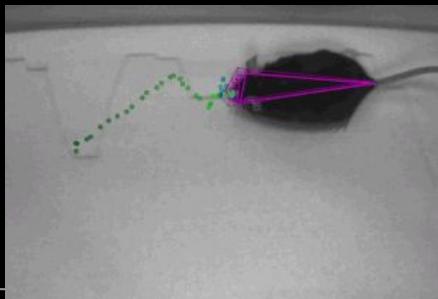
Element DeepLabCut

<https://github.com/DeepLabCut/DeepLabCut>



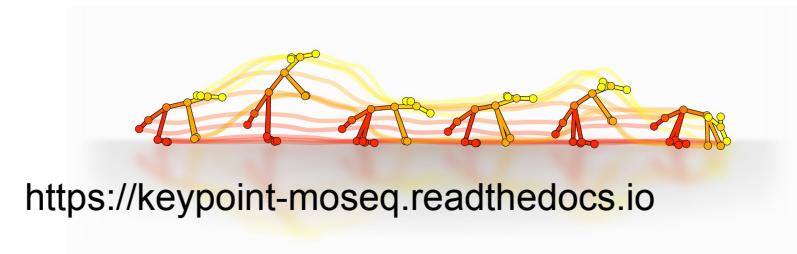
Analysis
Video Management
Model Training
Pose Estimation

Projects
Mesoscale Activity Project
Mathis Lab @ EPFL
Lu Lab @ Indiana U
Rose Lab @ Bonn U
Moser Group



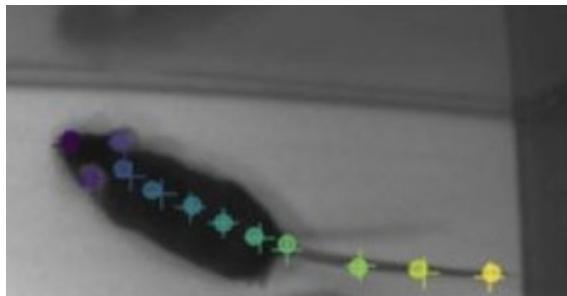


Keypoint-MoSeq

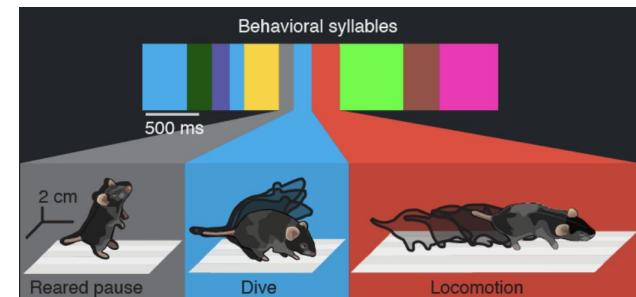


<https://keypoint-moseq.readthedocs.io>

Unsupervised machine learning algorithm segmenting continuous behavior into "syllables": e.g. rear, turn and pause



- 1. Model Training
- 2. Model Inference



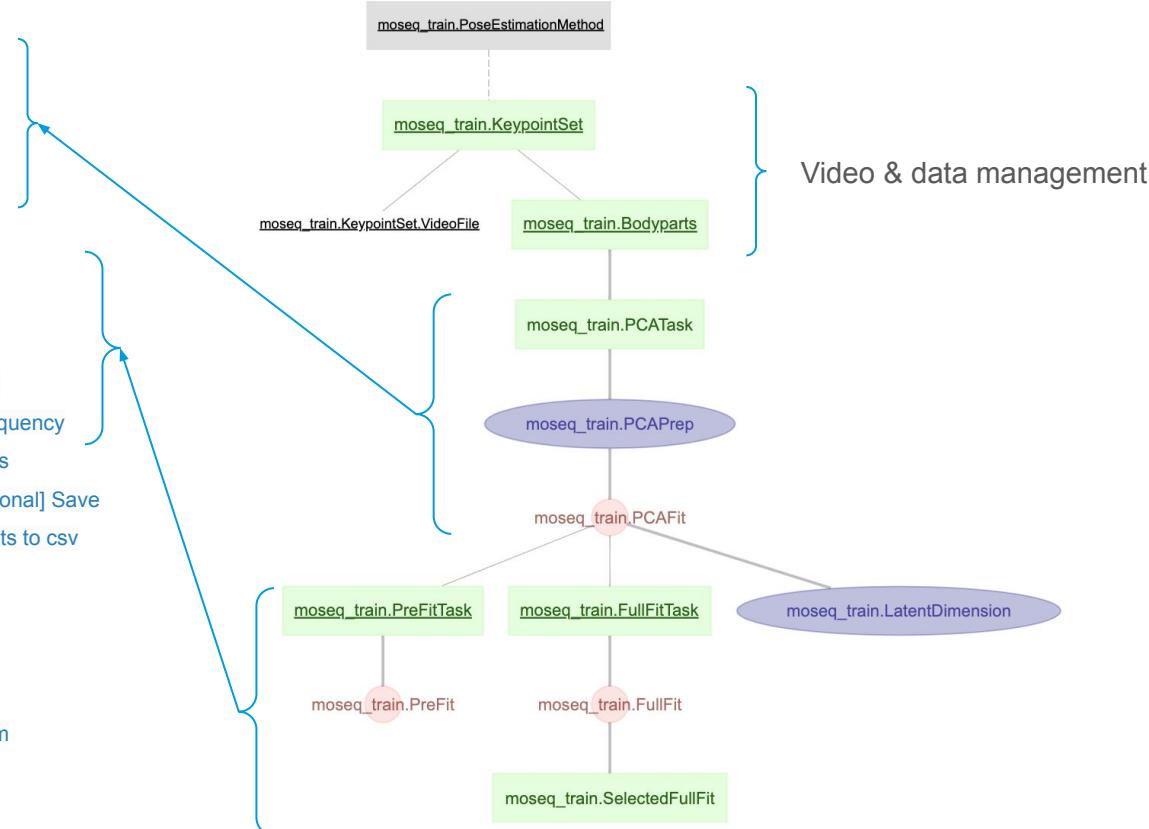
Article | [Open access](#) | Published: 12 July 2024

Keypoint-MoSeq: parsing behavior by linking point tracking to pose dynamics



MoSeq Model Training pipeline

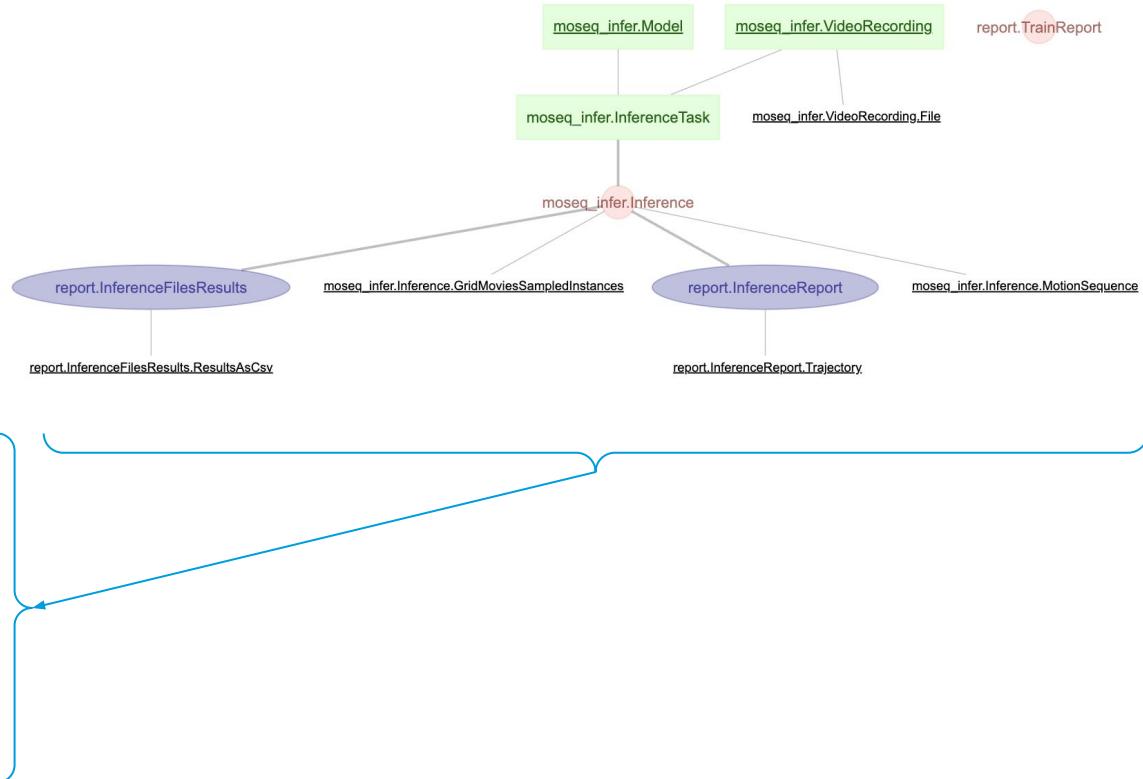
- Project setup
 - Edit the config file
 - Load data
 - Fit PCA
- Model fitting
 - Setting kappa
 - Initialization
 - Fitting an AR-HMM
 - Fitting the full model
 - Sort syllables by frequency
 - Extract model results
 - [Optional] Save results to csv
 - Apply to new data
- Visualization
 - Trajectory plots
 - Grid movies
 - Syllable Dendrogram





MoSeq Model Inference pipeline

- Project setup
 - Edit the config file
 - Load data
 - Fit PCA
- Model fitting
 - Setting kappa
 - Initialization
 - Fitting an AR-HMM
 - Fitting the full model
 - Sort syllables by frequency
 - Extract model results
 - [Optional] Save results to csv
 - Apply to new data
- Visualization
 - Trajectory plots
 - Grid movies
 - Syllable Dendrogram

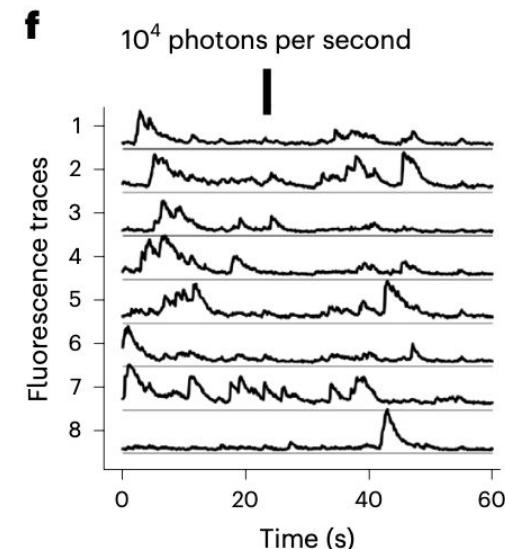
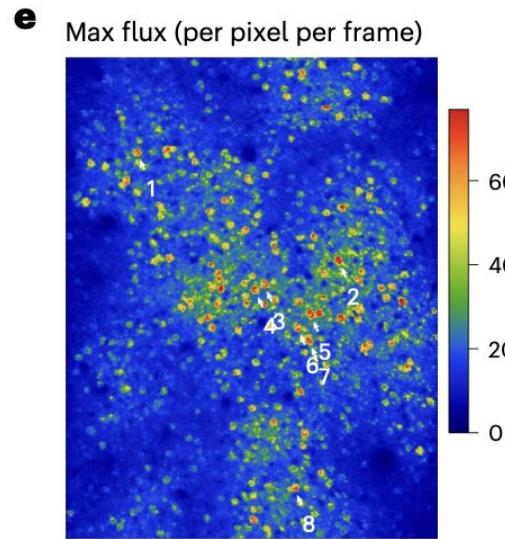
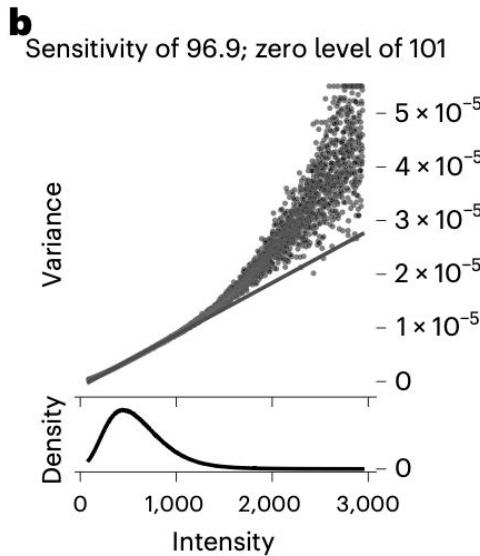




Protocol | Published: 17 March 2025

Uniform QC metrics for data quality across all data modalities

Standardized measurements for monitoring and comparing multiphoton microscope systems



Open Data In Neurophysiology: Advancements, Solutions & Challenges

Colleen J. Gillon^{†,1}, Cody Baker^{†,2}, Ryan Ly^{†,3}, Edoardo Balzani,⁴ Bingni W. Brunton,⁵ Manuel Schottdorf,⁶ Satrajit Ghosh,⁷ and Nima Dehghani^{7,8}

Open Data in Neurophysiology (2023) Ecosystem

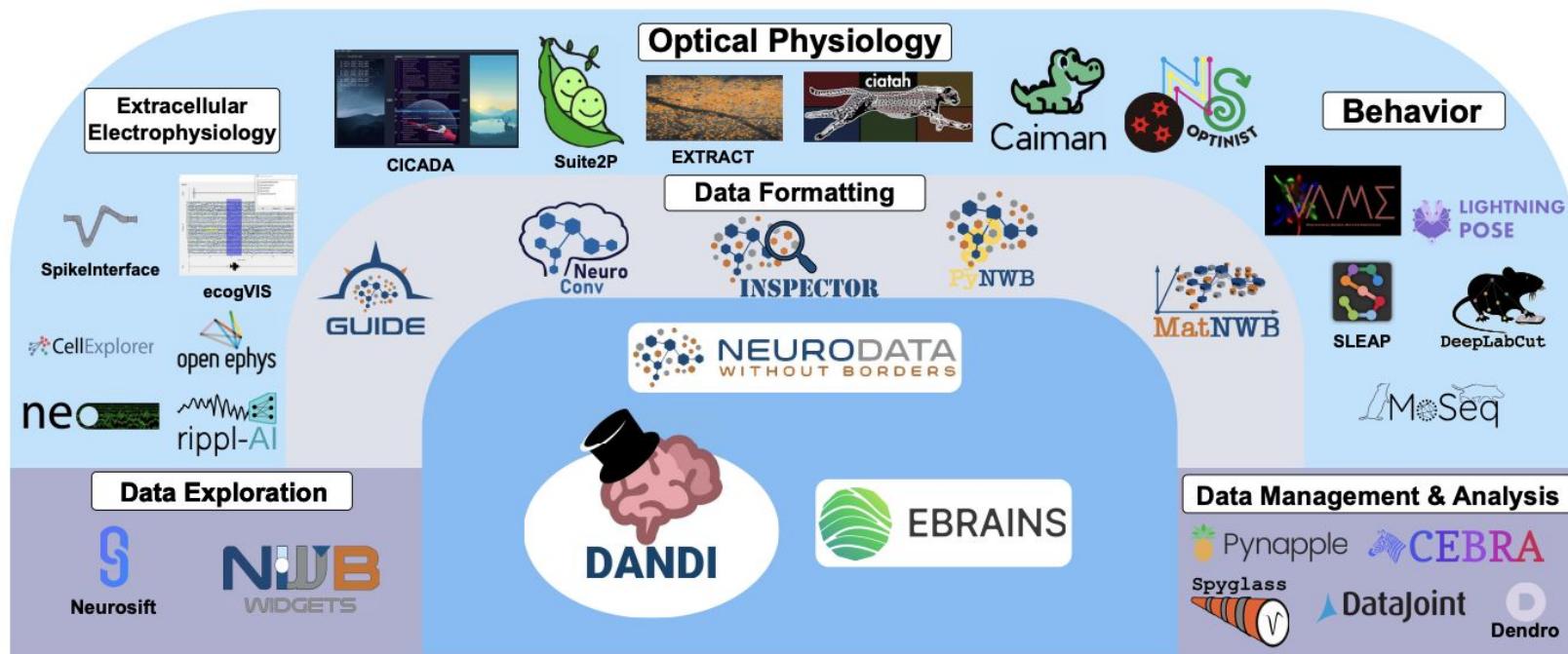


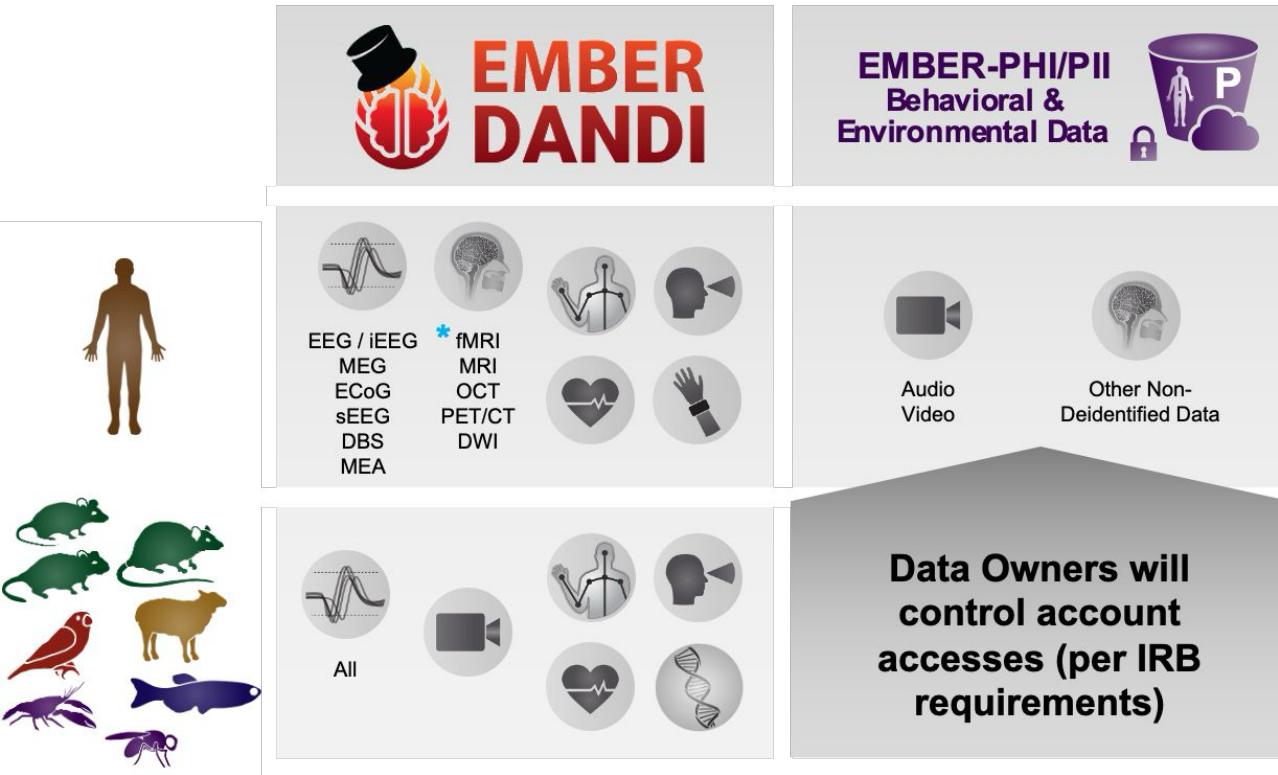
FIG. 1: The ecosystem of open source neurophysiology toolkits presented or discussed during ODIN 2023. See Table I for more information about each toolkit.

EMBER Multimodal Data Storage Strategy

Updated Plan

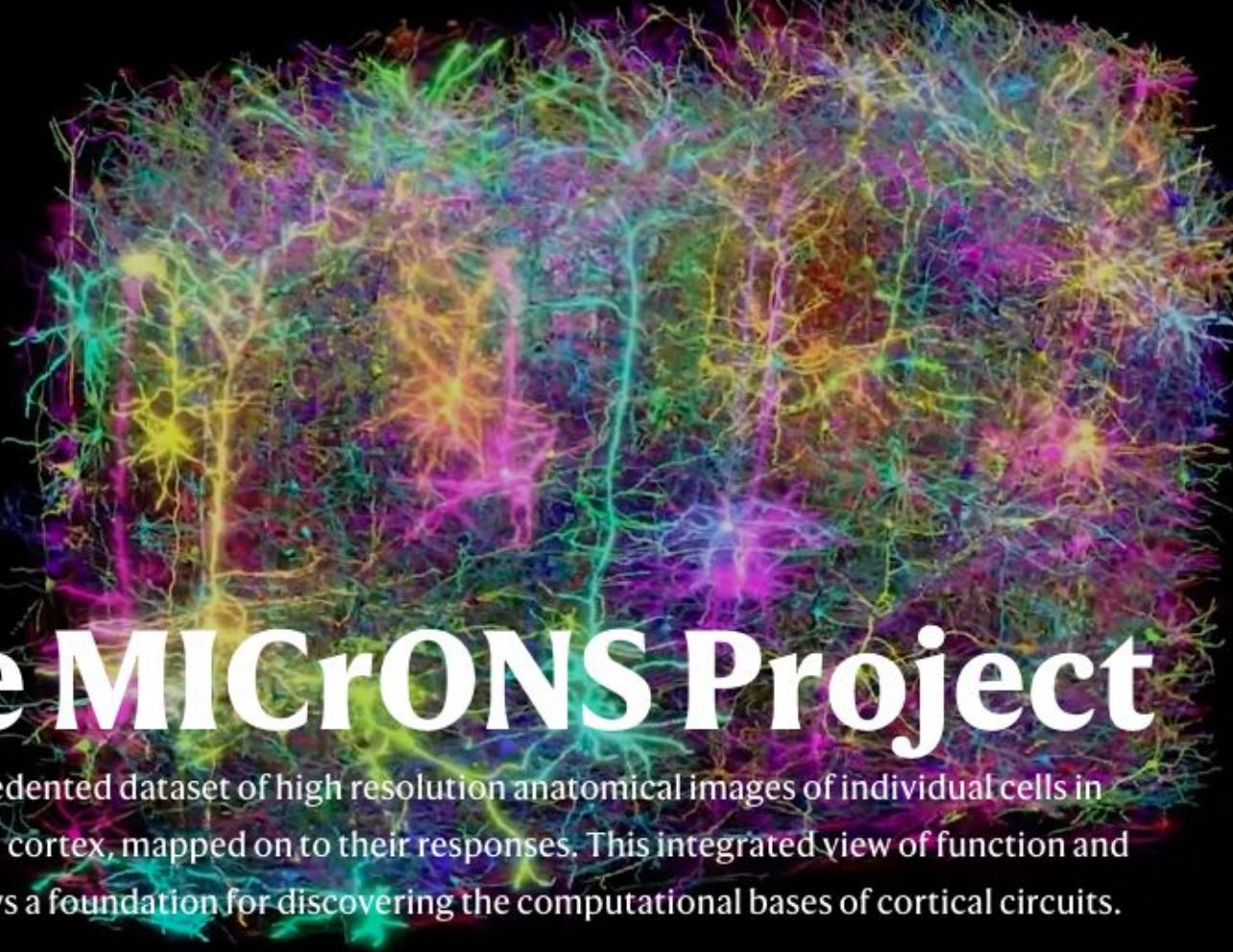


JOHNS HOPKINS
APPLIED PHYSICS LABORATORY



*Deidentified

We will be working directly with teams to assess and guide what data can be placed in DANDI



The MICRONS Project

An unprecedented dataset of high resolution anatomical images of individual cells in mouse visual cortex, mapped onto their responses. This integrated view of function and structure lays a foundation for discovering the computational bases of cortical circuits.



MICrONS: \$100M “Apollo Project of the Brain.”

2016



The “Machine Intelligence from Cortical Networks” program reverse-engineers the algorithms of the brain to develop less artificial AI.

✓ Collaboration

✓ Data Collection

✓ Analysis

2025



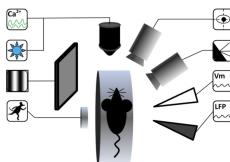
↗ Publication

“We now have the technology to look into the detailed organization of the brain.”

Andreas Tolias

DataJoint: the system of record

- In vivo SciOps
- Structure / function link
- Neural Data Access hosting



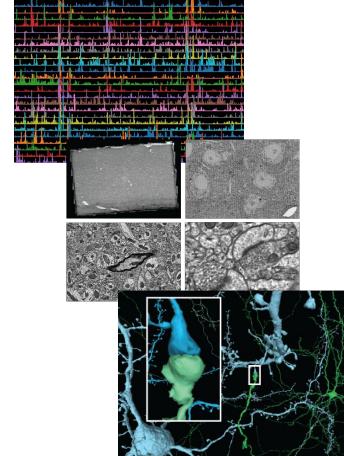
Experimental Methods



Function. Two-photon imaging records neuronal responses to visual stimuli in a cubic millimeter of mouse visual cortex.

Neuroanatomy. Serial-section electron microscopy images the same cubic millimeter of visual cortex.

Connectome. Convolutional nets align the image slices to reconstruct neurons and synapses in 3D and segment cells by class.



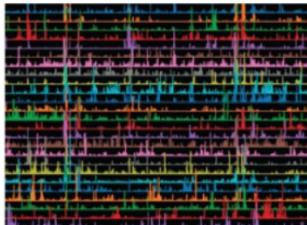
75,000 functional neurons

500,000 synapse junctions

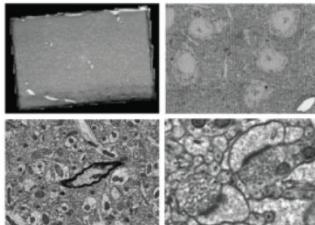
2 petabytes of data

a

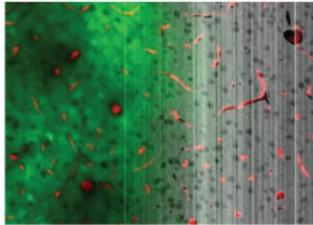
Data types available as data resource



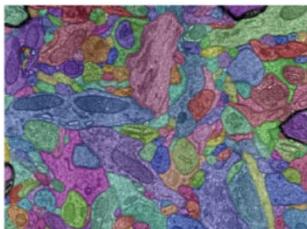
FUNCTIONAL DATA



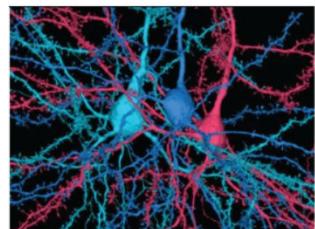
ELECTRON MICROSCOPY
IMAGERY



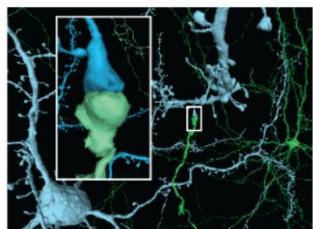
FUNCTIONAL -STRUCTURAL
CO-REGISTRATION



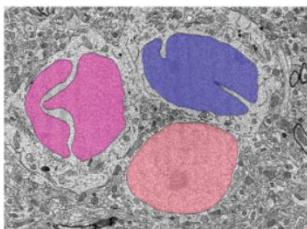
CELL SEGMENTATION



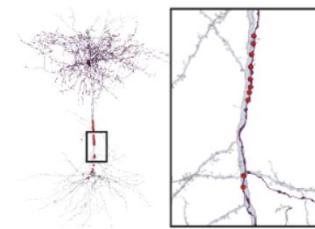
CELL MESHES



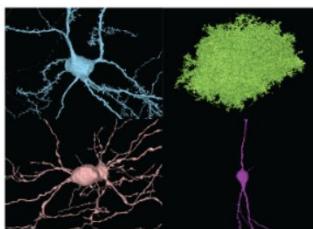
SYNAPSE CONNECTIVITY



NUCLEUS
SEGMENTATION



PROOFREADING
STATUS



CELL TYPES

MICrONS Data

<https://www.microns-explorer.org>

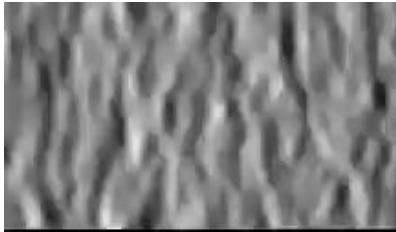
<https://datajoint.com/microns>

Visual Observatory of Cortex (VORTEX)
<https://www.microns-explorer.org/vortex>



MICrONS Visual Stimuli

"Monet"



Receptive fields
Orientation tuning
Direction tuning

"Trippy"



Responsiveness
Signal correlations

Videos



ML models
“Digital twin”

Modeling
Validation
Ethological relevance



Accessing Physiology Data

1. (Difficult) Download the SQL data

- Set up your own SQL server, load the data
- Access with the DataJoint API

2. (Easier) Connect to the database instance hosted by DataJoint.com

- Access with public credentials using the DataJoint API
- Complete tutorial in DevContainer - Codespaces:
https://github.com/datajoint/microns_phase3_ndu/

3. (Easiest) Download NWB files from DANDI Archive

<https://doi.org/10.48324/dandi.000402/0.230307.2132>



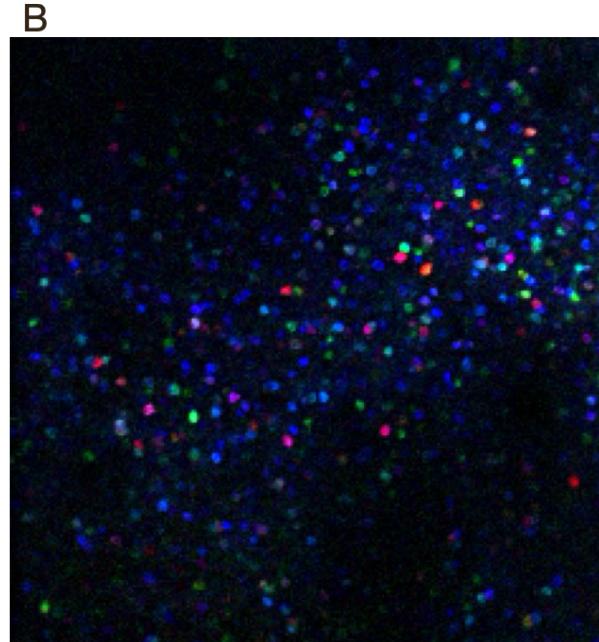
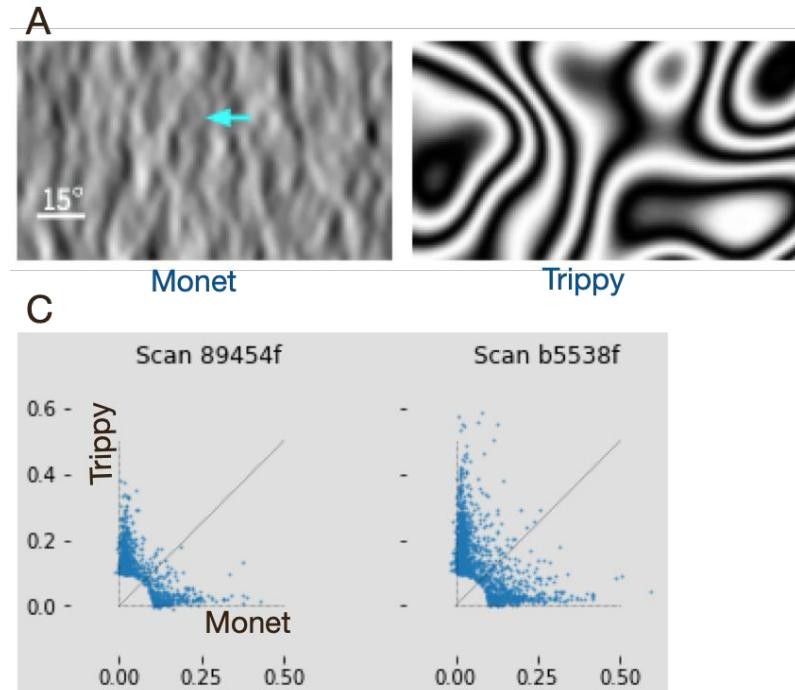
4. (Pro) Work on a complete data pipeline from raw data

Analysis collaboration: DataJoint + Neuromatch + Stanford

Analysis Collaboration

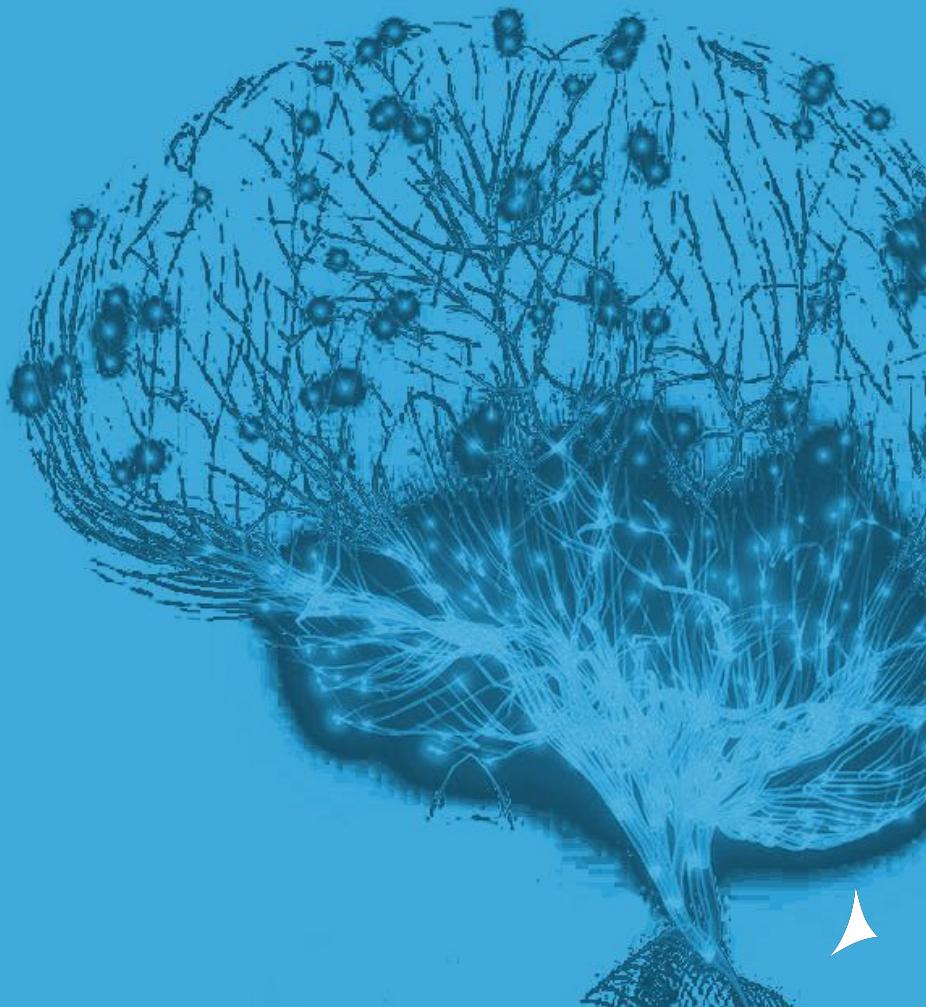


*D. YATSENKO¹, P. FAHEY², T. MUHAMMAD³, M. SHAKIBA⁴, R. ROKNI⁴, M. MOHAMMADI⁴, N. DEHGHANI⁵, A. S. TOLIAS²



Synthetic stimuli divide neuronal populations into two functional circuits

The NeuroAI Future



New Discipline: SciOps



Skills + Standards + Systems



DataJoint



JOHNS HOPKINS
APPLIED PHYSICS LABORATORY

SCIOPS: ACHIEVING PRODUCTIVITY AND RELIABILITY IN DATA-INTENSIVE RESEARCH

A PREPRINT

Erik C. Johnson^{1,a} Thinh T. Nguyen² Benjamin K. Dichter³ Frank Zappulla⁴
 Montgomery Kosma² Kabilar Gunalan^{2,5} Yaroslav O. Halchenko⁵ Shay Q. Neufeld⁷
 Kristen Ratan⁵ Nicholas J. Edwards⁹ Susanne Ressl¹⁰ Sarah R. Heilbronner¹¹
 Michael Schirmer^{12,16} Petra Ritter^{13,16} Brock Wester¹ Satrajit Ghosh^{5,16}
 Maryam E. Martone¹⁷ Dimitri Yatsenko^{2,b} Franco Pestilli¹⁸

¹ Research and Exploratory Development Department, Johns Hopkins University Applied Physics Laboratory, Laurel, MD, USA
² DataJoint Inc., Houston, TX, USA
³ CatalystNeuro, Benicia, CA, USA
⁴ Digital R&D Creation Center, Pfizer Inc., USA
⁵ McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA, USA
⁶ Center for Open Neuroscience, Department of Psychological and Brain Sciences, Dartmouth College, New Hampshire, USA
⁷ Inscopix, a Bruker company, Mountain View, CA, USA
⁸ Strategies for Open Science (Stratos), Santa Cruz, CA, USA
⁹ Potato Inc., San Diego, CA, USA
¹⁰ Department of Neuroscience, University of Texas at Austin, Austin, TX
¹¹ Neurosurgery, Baylor College of Medicine, Houston, TX
¹² Berlin Institute of Health (BIH) at Charité – Universitätsmedizin Berlin, Berlin, Germany
¹³ Department of Neurology with Experimental Neurology, Charité – Universitätsmedizin Berlin, Corporate Member of Freie Universität Berlin and Humboldt Universität zu Berlin, Berlin, Germany
¹⁴ Bernstein Focus State Dependencies of Learning and Bernstein Center for Computational Neuroscience, Berlin, Germany
¹⁵ Einstein Center for Neuroscience Berlin, Berlin, Germany
¹⁶ Einstein Center Digital Future, Berlin, Germany
¹⁷ Department of Otolaryngology, Harvard Medical School, Boston, MA, USA
¹⁸ Department of Neurosciences, University of California San Diego, La Jolla, CA, USA
^aerik.c.johnson@jpl.nasa.gov ^bdimitri@datajoint.com

September 22, 2024

ABSTRACT

Scientists are increasingly leveraging advances in instruments, automation, and collaborative tools to scale up their experiments and research goals, leading to new bursts of discovery. Various scientific disciplines, including neuroscience, have adopted key technologies to enhance collaboration, reproducibility, and automation. Drawing inspiration from advancements in the software industry, we present a roadmap to enhance the reliability and scalability of scientific operations for diverse research teams tackling large and complex projects. We introduce a five-level Capability Maturity Model describing the principles of rigorous scientific operations in projects ranging from small-scale exploratory studies to large-scale, multi-disciplinary research endeavors. Achieving higher levels of operational maturity necessitates the adoption of new, technology-enabled methodologies, which we















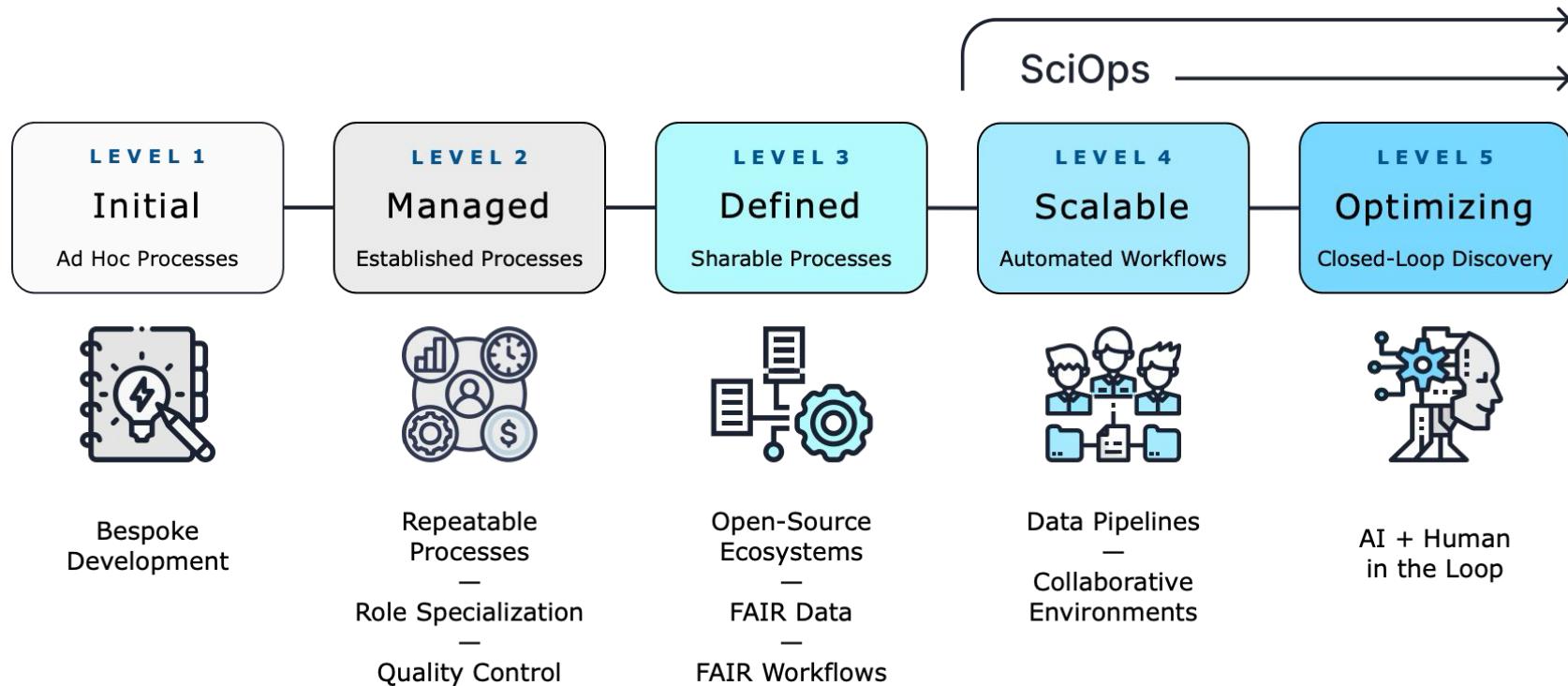
TEXAS

The University of Texas at Austin

Our [SciOps paper](#) is under review by *Nature Methods*.



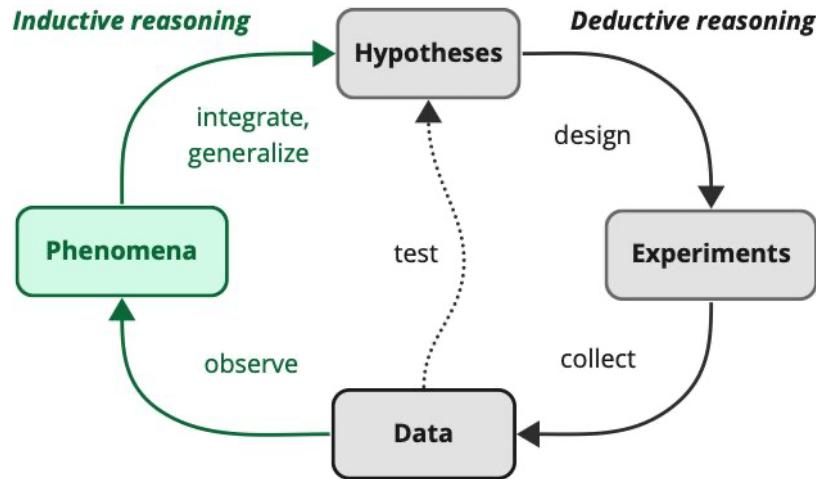
NeuroAI requires SciOps maturity.



<https://doi.org/10.48550/arXiv.2401.00077>



Closed-loop system from hypothesis to verification



AI requires mature scientific operations

SciOps: A Comprehensive Platform

Productivity through automation

- Ingest and store
- Prepare and standardize
- FAIR
- Flexible queries
- Collaborate, controlled sharing, and publish

AI Enablement

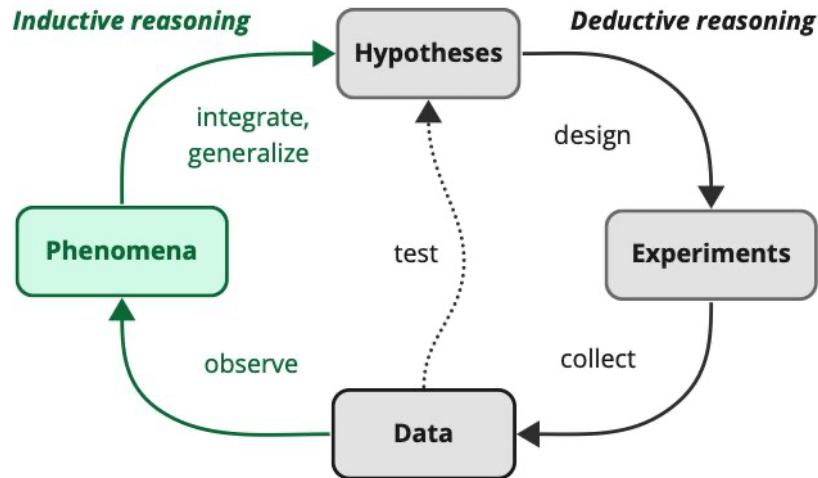
- Intuitive UI and common language
- Analyze and build models
- Integrate pipelines, tools (existing and 3rd party)
- Dynamic and flexible
- Compute orchestration

Strong Governance

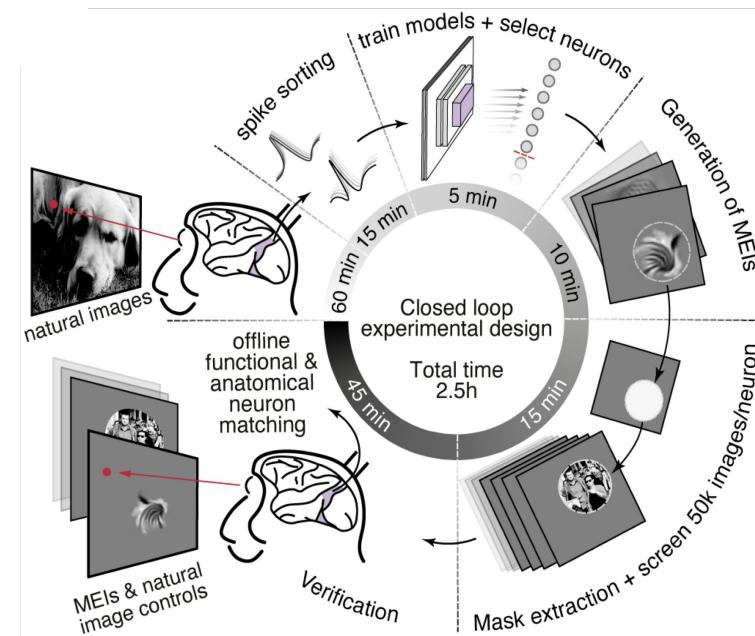
- System of record
- Code + data + dependencies + environment
- Data Integrity enforced and transparent
- Reproducibility
- Protect IP
- Well documented processes and transformations



Closed-loop neuroscience with live ML modeling.



AI embedding requires mature operations



courtesy Prof. Andreas Tolias

The Three AI Partners

A synergistic framework of AI collaborators. Click on any layer to learn more.

Level 2: AI Co-Visionary

Strategy & Vision

Guides long-term research trajectory, identifies funding and partnership opportunities, and helps map the strategic impact of scientific discoveries.

Level 1: AI Co-Scientist

Creativity & Exploration

Acts as an intellectual partner, generating novel hypotheses, synthesizing literature, and assisting in manuscript preparation through a closed-loop discovery process.

Level 0: AI Co-Pilot

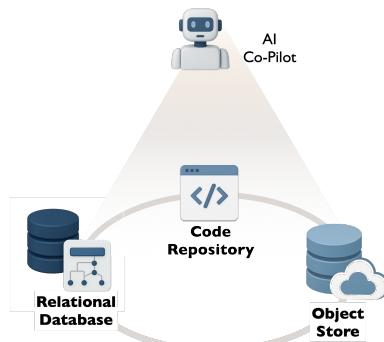
Rigor & Precision

Ensures absolute data integrity by translating natural language into precise, validated queries and orchestrating complex analysis workflows with perfect reproducibility.

 Deepinvent



- FutureHouse.org
- NovelSeek
- ReadySetPotato.com
- Google Co-Scientist
- Microsoft Discovery Platform
- HypSynth





Areas for collaboration

Restoring Gold Standard Science

Executive Orders | May 23, 2025

NIH AI: [NOT-OD-25-117](#)

- ✓ Operational excellence:
governance + infrastructure + operations + dev support + training
- ✓ AI-powered closed-loop science: Grant Opportunities
- ✓ Industry sponsored research
- ✓ Research communication and publishing: data + computation
- ✓ Technology dissemination & commercialization. Software licensing.

THANK YOU

