



*Dimitri Yatsenko, Jacob Reimer, Edgar Walker  
(Andreas Tolias' Lab)*

---

# DataJoint

Intro

---

2016-12-09

# The DataJoint framework

---

Aim: Efficiency, integrity, and collaboration

- ❖ Large volumes of multimodal, dynamic data with evolving structure
- ❖ Explicit, systematic, self-documenting data organization
- ❖ Solid theoretical foundations: a strict and streamlined variant of the relational data model to make easy to understand and use
- ❖ MATLAB and Python interoperability
- ❖ Data integrity
- ❖ Precise, fast, interactive queries
- ❖ Data sharing with fine access control
- ❖ Distributed computation

# DataJoint history

---

- ❖ In continuous use and development since **2009** in Andreas Tolias lab but never as its own project.
- ❖ Used by several other labs since **2011**.
- ❖ **2016:** a major data organization tool for the MICrONS project.
- ❖ **Oct 2016:** Dr. Paul Schrater at U. of Minnesota was awarded \$50k supp funding to adopt DataJoint and to write tutorials
- ❖ **Nov 2016:** Dimitri Yatsenko et al founded Vathes LLC to apply for Phase I small-business funding for enhanced hosting and networking services based on DataJoint.

# DataJoint

data pipelines for the science lab  
with MATLAB and Python



Dimitri Yatsenko

In-depth guide to be released on  
March 1, 2017

<https://leanpub.com/datajoint>

Introductory tutorials online

# Big Data? NoSQL?

---

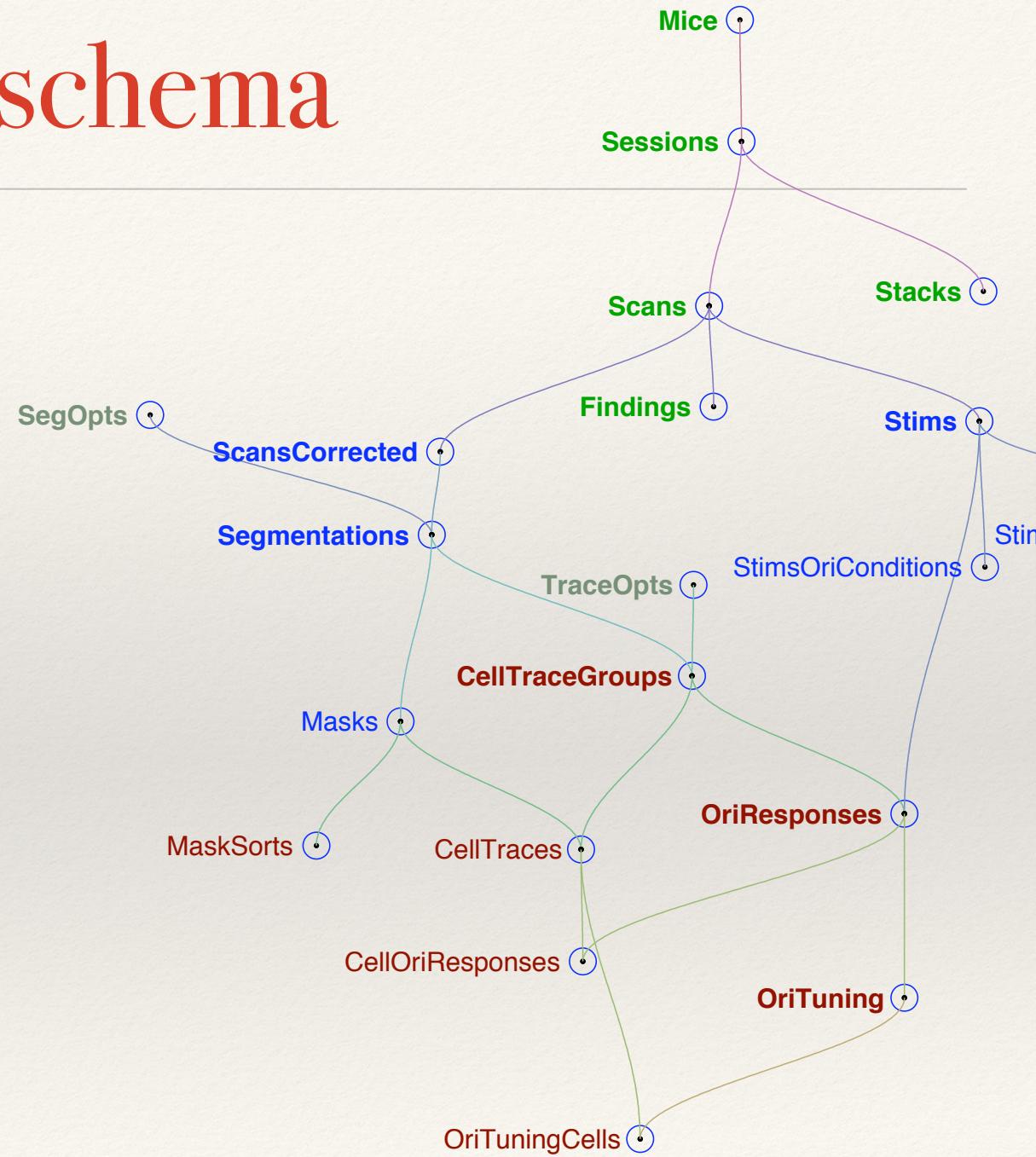
**Big Data** often describes data without predefined structure

- ❖ structure is learned from data itself
- ❖ in this sense, DataJoint is **not** Big Data as it enforces explicit structure

The trendy **NoSQL** databases abandon structure and consistency for the sake of flexible distributed hosting (scalability)

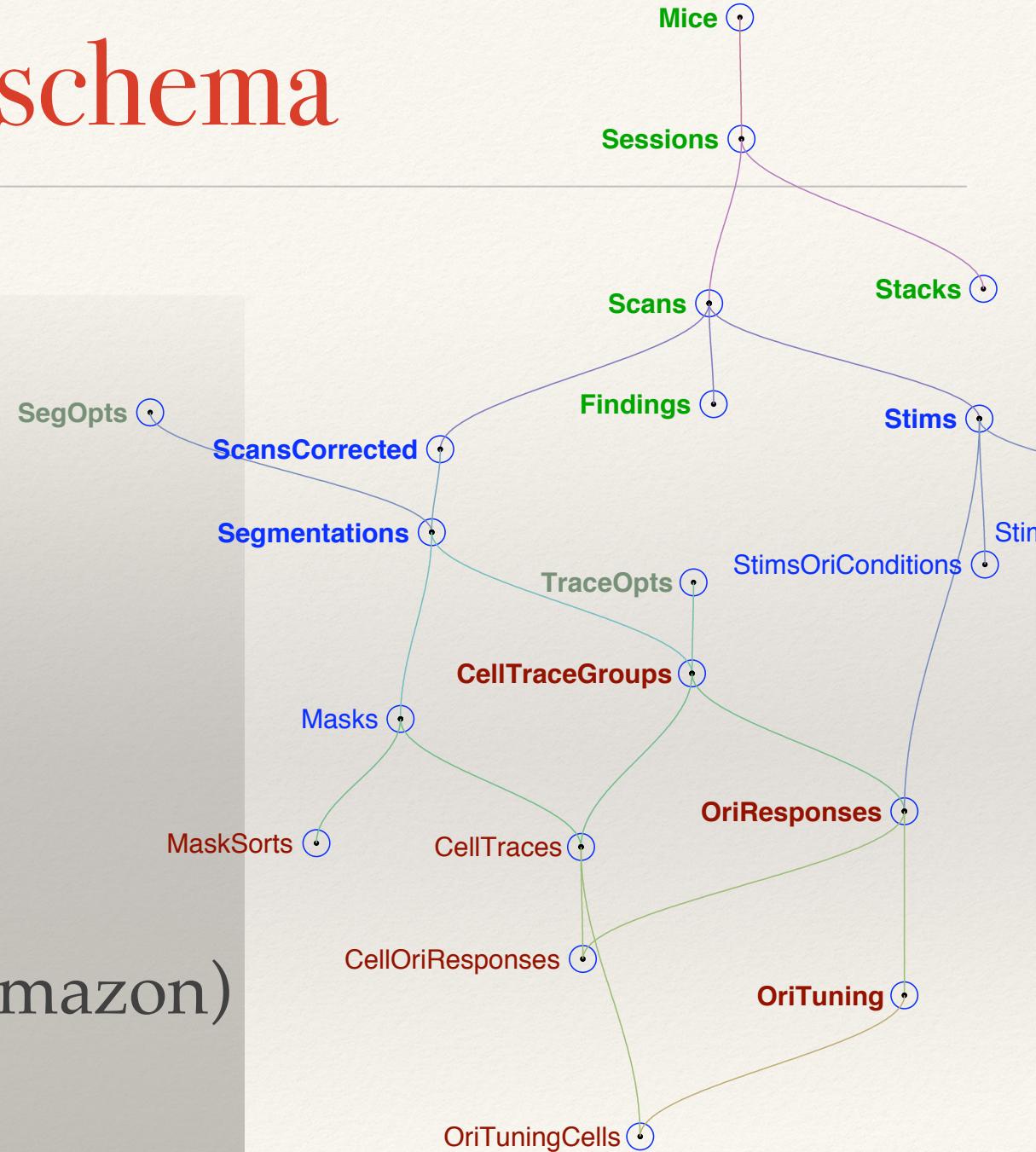
- ❖ in this sense, DataJoint is **not** NoSQL
- ❖ DataJoint strives to achieve scalability while maintaining consistency

# A DataJoint schema



# A DataJoint schema

- ❖ Backed by MySQL
- ❖ Hosting
  - personal computer
  - lab server
  - cloud server (e.g. Amazon)



# What is a DataJoint database?

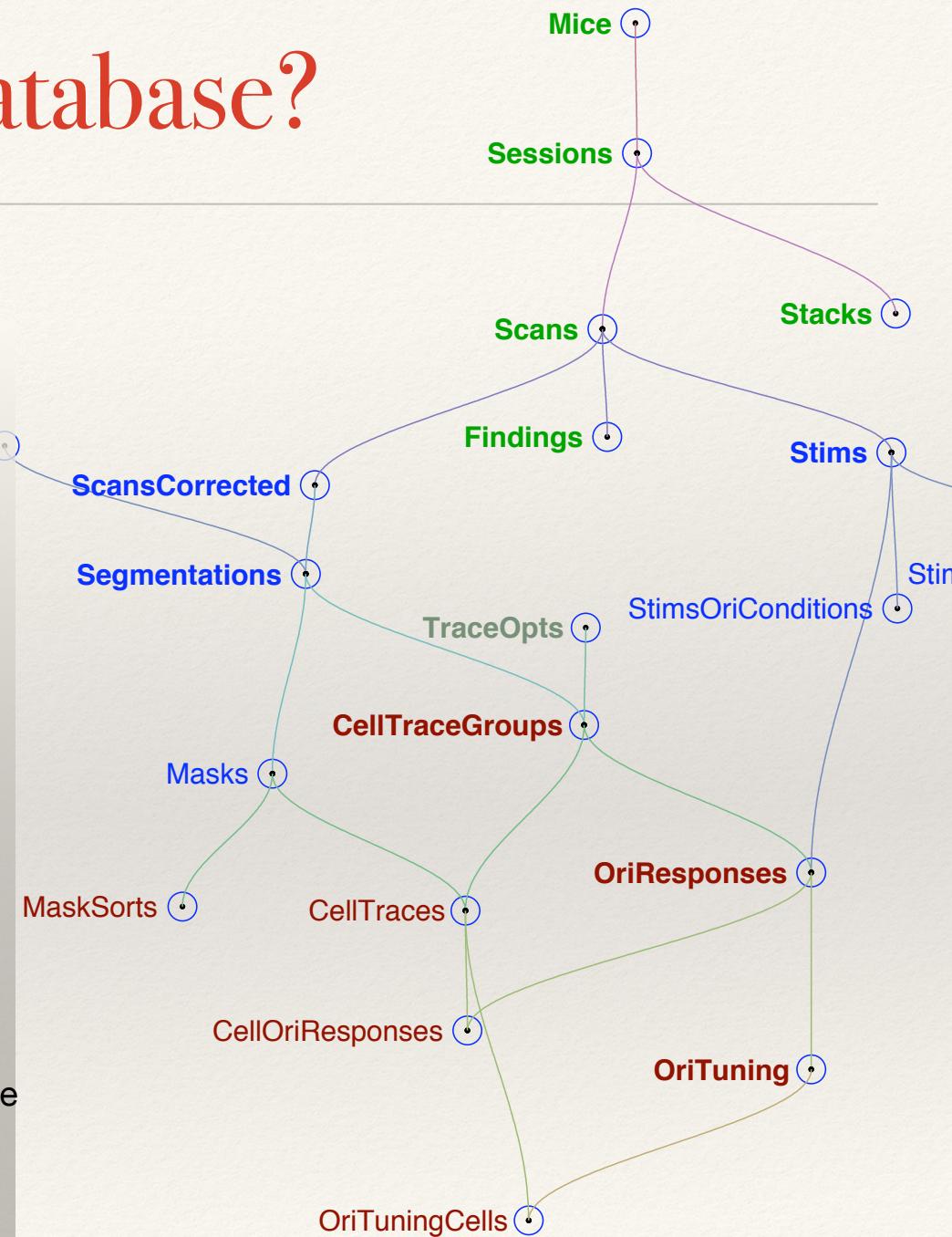
- ❖ Nodes are tables of data

	primary key			
	id	scan	scan_date	operator
2	1	2014-01-01	Edgar	
2	2	2014-06-01	Edgar	
2	3	2015-01-01	George	
3	1	2015-05-01	Cathryn	
3	2	2015-05-08	Shan	
3	3	2015-05-15	Shan	

attribute name

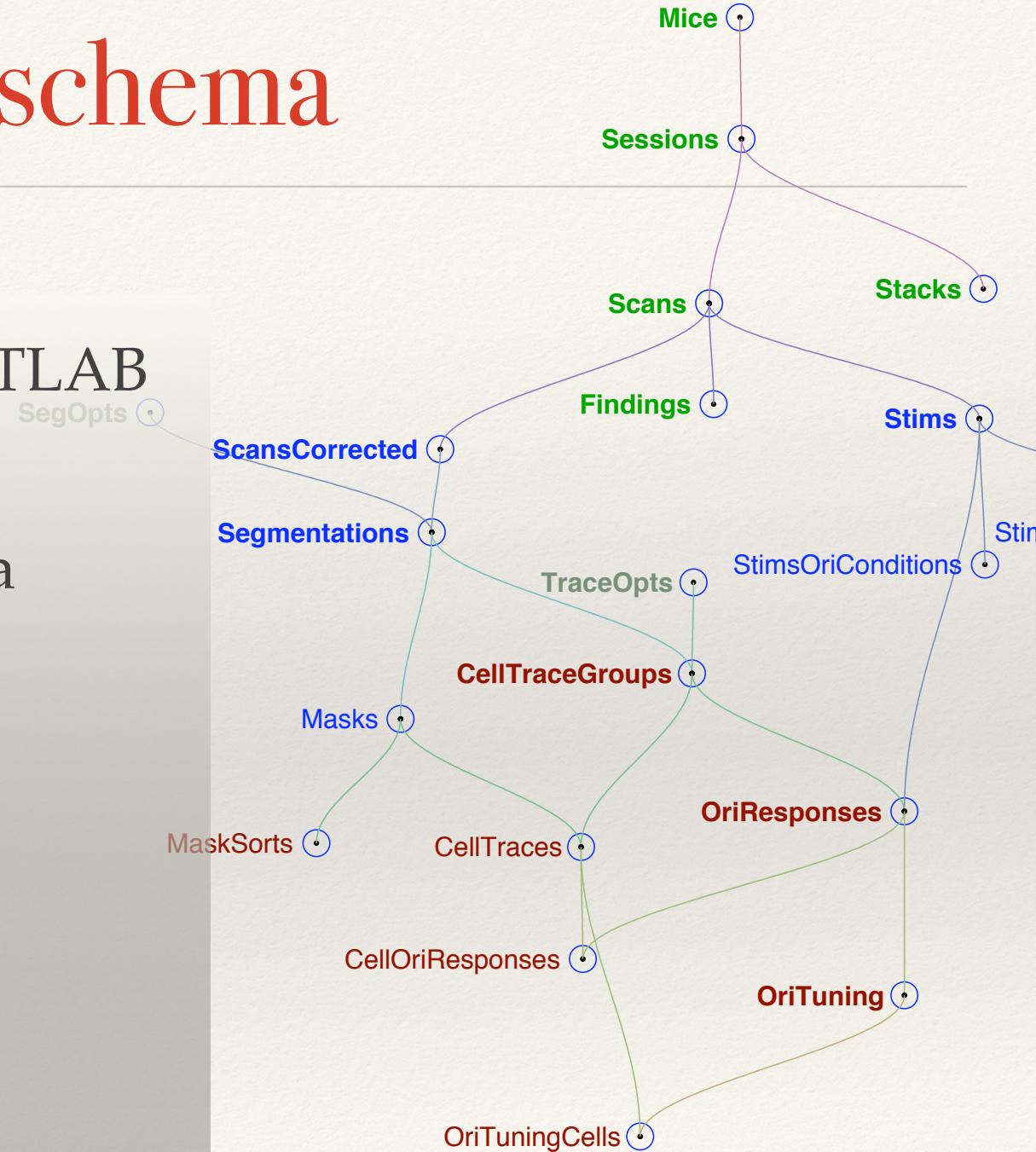
tuple

attribute value



# A DataJoint schema

- ❖ Nodes are classes in MATLAB or Python
- ❖ Nodes are also tables in a relational database
- ❖ Data echelons:
  - manual
  - imported
  - computed



# A DataJoint schema



# DataJoint infrastructure

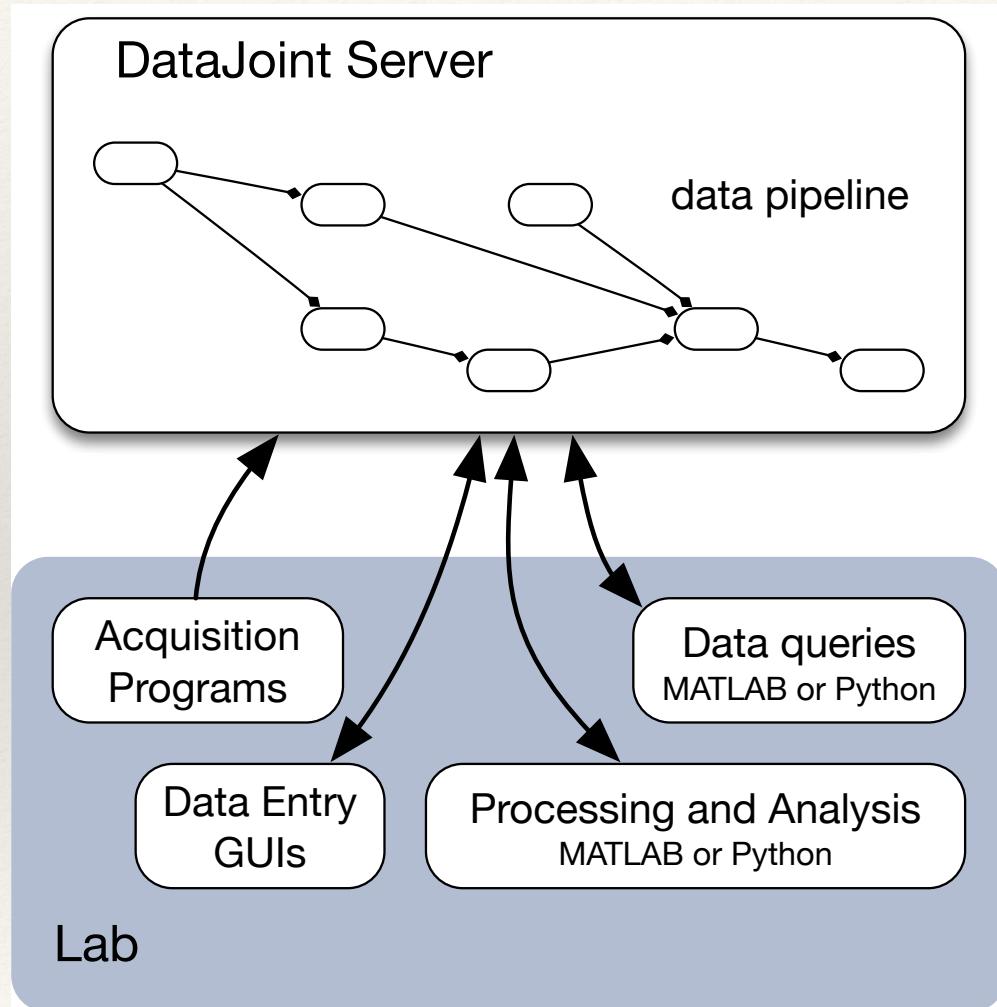
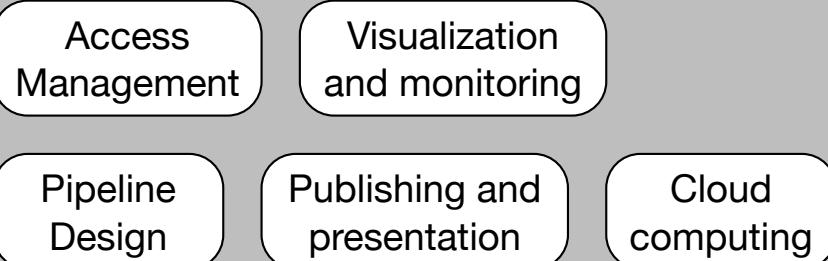
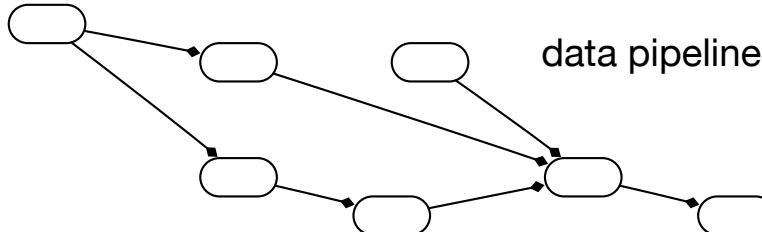


Figure 1: Basic DataJoint architecture.

## DataJoint Hub



## DataJoint Server



Acquisition  
Programs

Data Entry  
GUIs

Data queries  
MATLAB or Python

Processing and Analysis  
MATLAB or Python

Lab

# DataJoint infrastructure

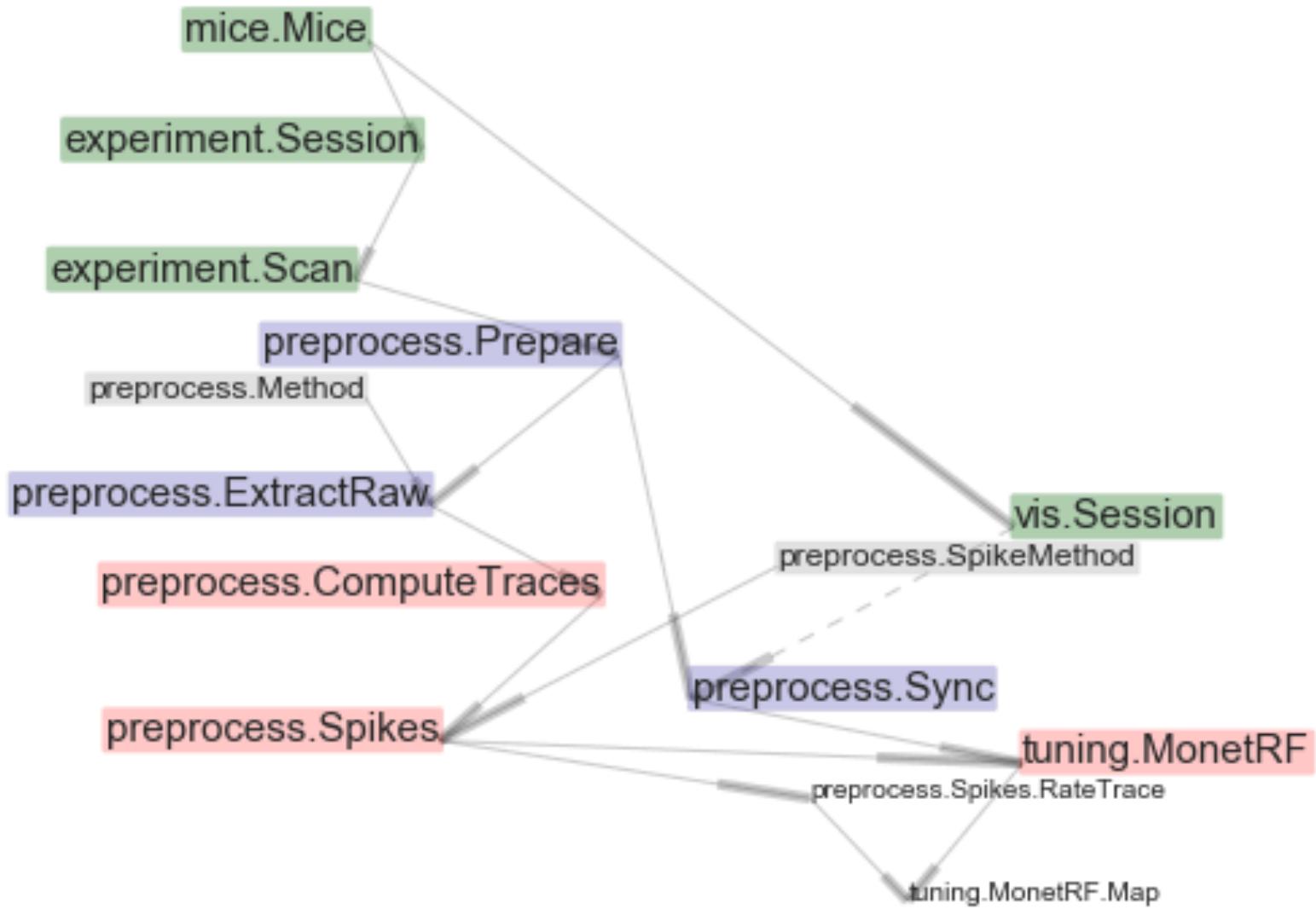
(planned in 2017)

## DataJoint Hub ([datajoint.io](http://datajoint.io))

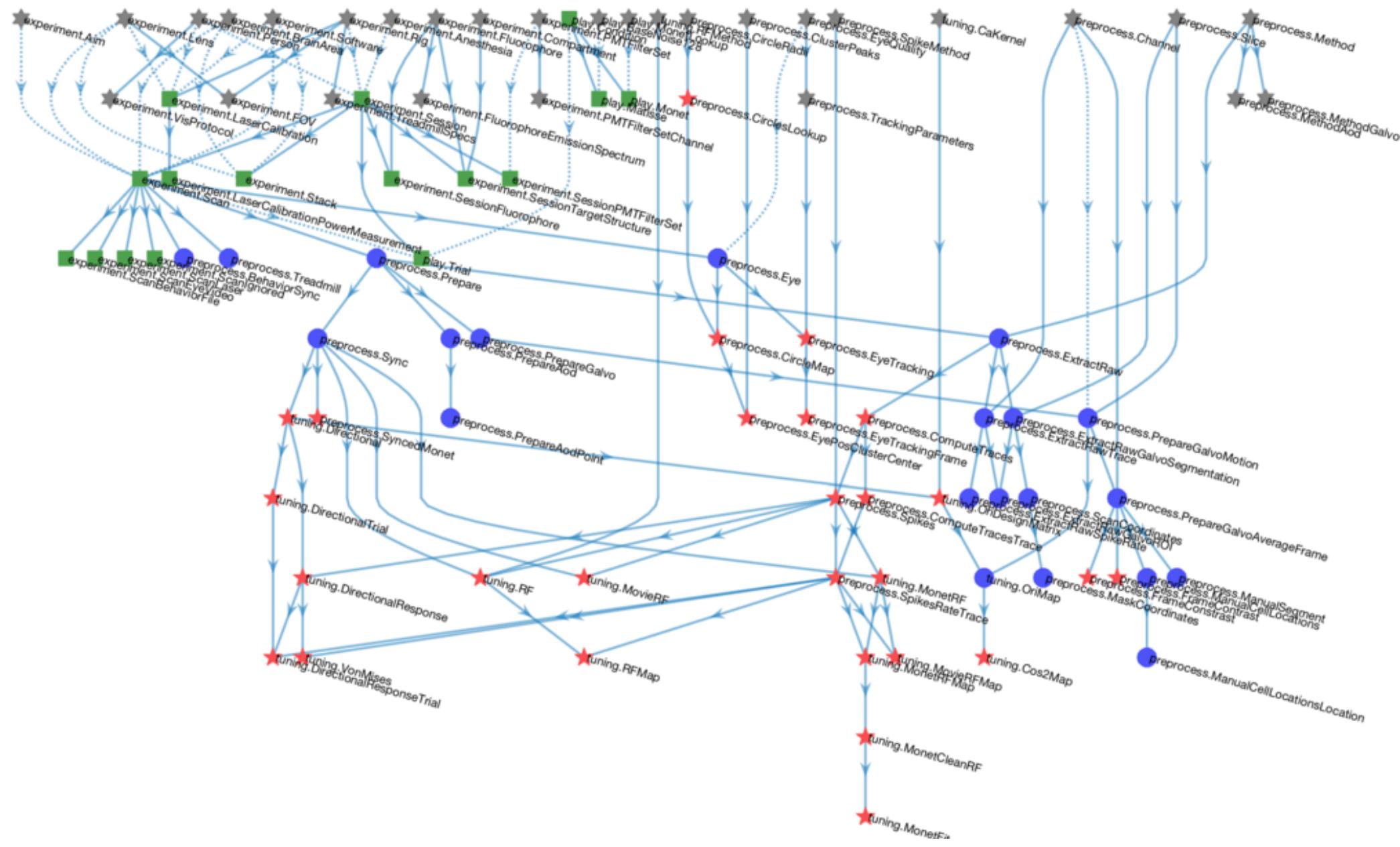
- pipeline hosting (free or subscription)
- distributed cloud computing

# MICrONS pipeline (fragment)

<https://github.com/cajal/pipeline>

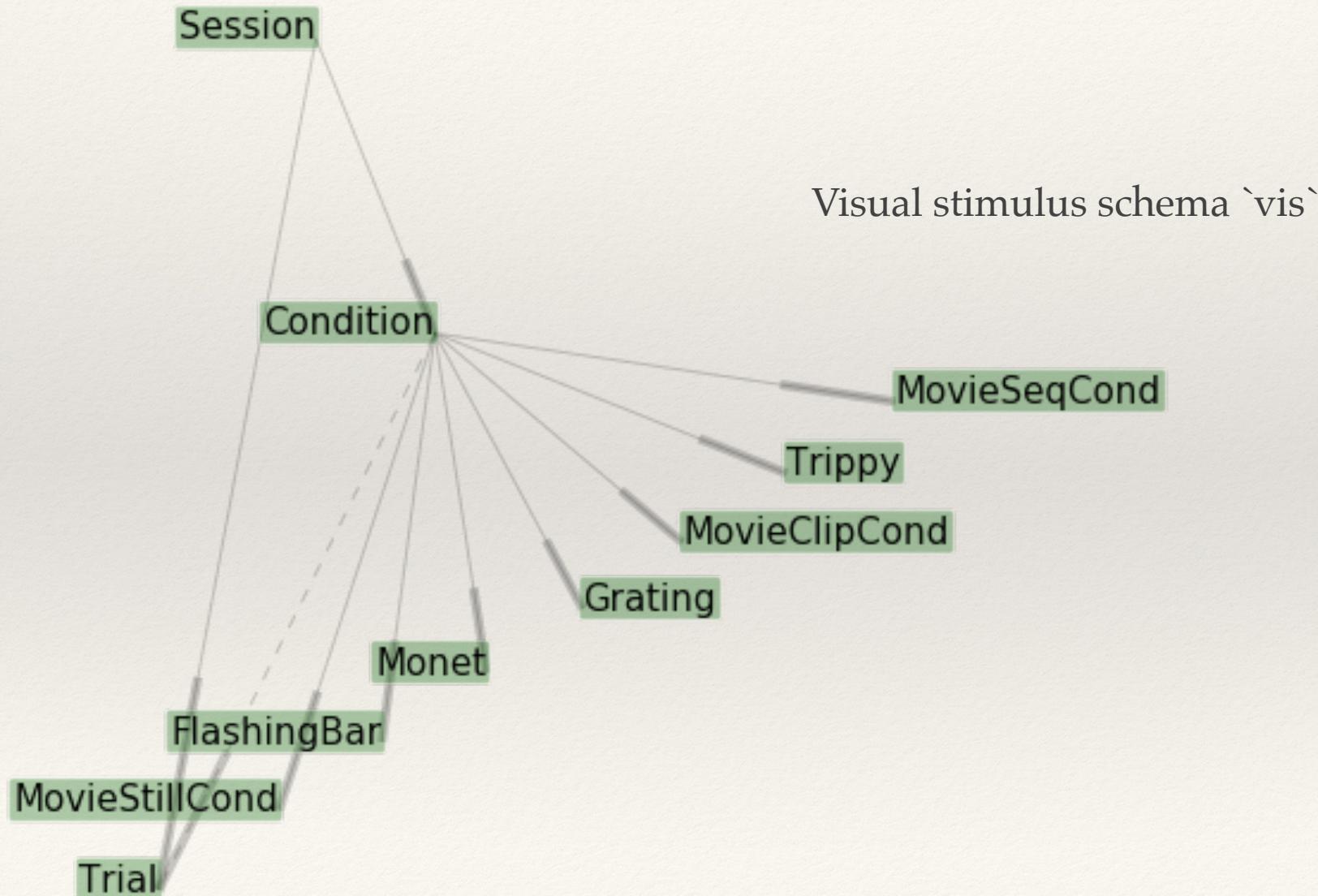


# MICrONS pipeline



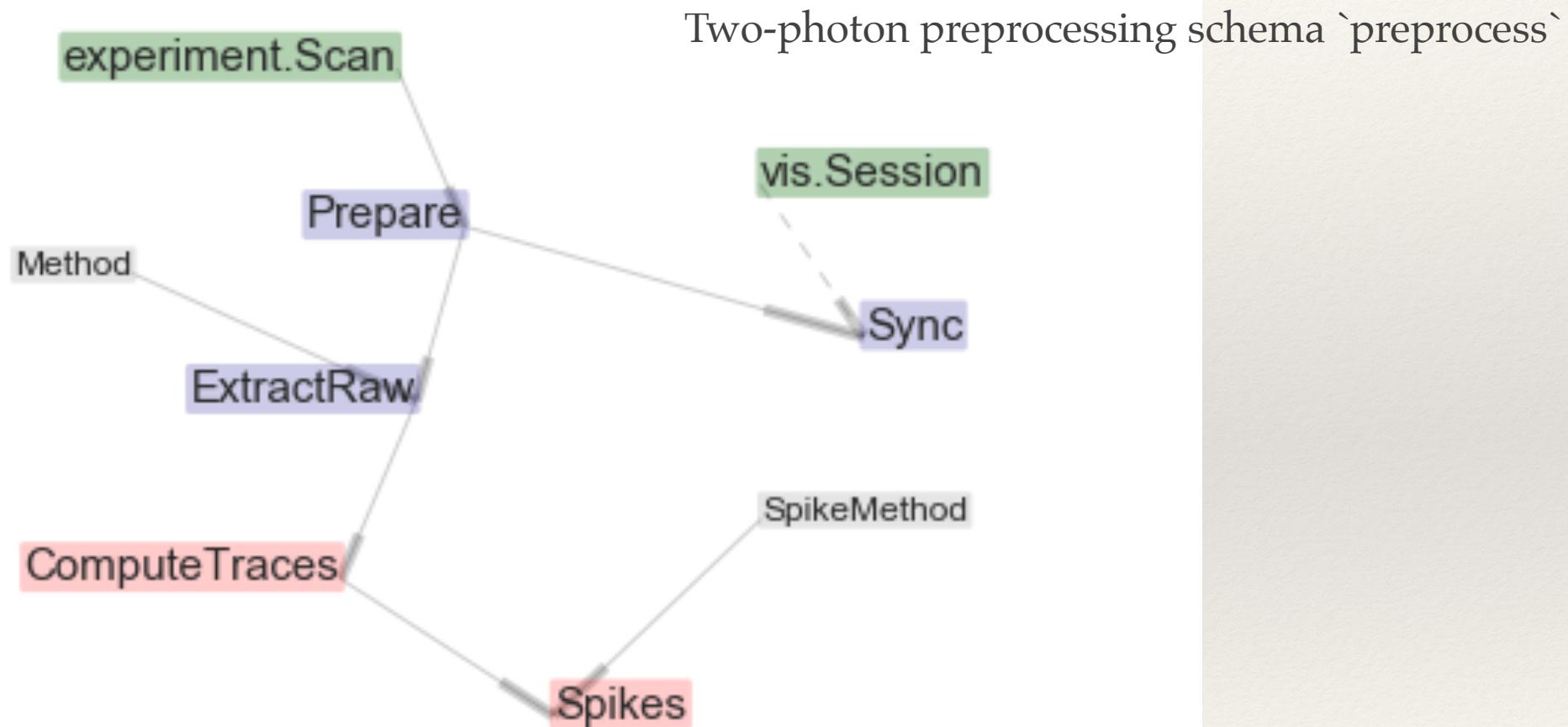
# MICrONS pipeline (fragment)

<https://github.com/cajal/pipeline>



# MICrONS pipeline (fragment)

<https://github.com/cajal/pipeline>



# A DataJoint node

```
@schema
class Scan(dj.Manual):
    definition = """ # scanimage scan info
-> Session
    scan_idx : smallint # scanimage-generated index
-----
    surfz          :float   # (um) z-coord at pial surface
    depth=0        :int     # manual depth measurement
    laser_wavelength :float   # (nm)
    laser_power    :float   # (mW) to brain
    cortical_area="V1" :enum('other','unknown','V1','LM','AL','PM')
    scan_notes = ""  :varchar(4095) # free-notes
    scan_ts = CURRENT_TIMESTAMP : timestamp   # don't edit
"""

```

Python

# A DataJoint node

```
%{  
common.Scan (manual) # scanimage scan info  
->common.Session  
scan_idx : smallint # scanimage-generated index  
----  
surfz :float # (um) z-coord at pial surface  
depth=0 :int # manual depth measurement  
laser_wavelength :float # (nm)  
laser_power :float # (mW) to brain  
cortical_area="V1" :enum('other', 'unknown', 'V1', 'LM', 'AL', 'PM')  
scan_notes = "" :varchar(4095) # free-notes  
scan_ts = CURRENT_TIMESTAMP : timestamp # don't edit  
%}  
  
classdef Scan < dj.Relvar  
end
```

MATLAB

# A DataJoint node

```
%{  
common.Scan (manual) # scanimage scan info  
->common.Session  
scan_idx : smallint # scanimage-generated index  
----  
surfz :float # (um) z-coord at pial surface  
depth=0 :int # manual depth measurement  
laser_wavelength :float # (nm)  


| animal_id | tp_session | scan_idx | surfz | depth | laser_wavelength | laser_power | cortical_ar |
|-----------|------------|----------|-------|-------|------------------|-------------|-------------|
| 3625      | 2          | 5        | 0     | 57    | 920              | 40          | V1          |
| 3626      | 1          | 1        | 0     | 80    | 920              | 34          | V1          |
| 3626      | 1          | 3        | 0     | 80    | 920              | 34          | AL          |
| 3626      | 1          | 5        | 0     | 120   | 920              | 34          | AL          |
| 3626      | 1          | 7        | 0     | 90    | 920              | 34          | V1          |

  
classdef Scan < dj.Relvar  
end
```

# A DataJoint node (computed)

```
@schema
class Power(dj.Computed):
    definition = """ # LFP power in each frequency band
-> LFP
-> FrequencyBand
-----
power :float # mV^2 Hz -- average power in the band
"""

def _make_tuples(self, key):
    # fetch required data
    lo, hi = (FrequencyBand() & key).fetch1['freq_lo', 'freq_hi']
    signal, dt = (LFP() & key).fetch1['voltage', 'dt']
    # compute
    signal = band_pass_filter(signal, lo, hi, dt)
    power = compute_average_power(signal, dt)
    # submit
    self.insert1(dict(key, power=power))
```

Python

# A DataJoint node (computed)

```
%{  
example.Power (computed) # LFP power in each frequency band  
-> example.LFP  
-> example.FrequencyBand  
----  
power :float # mV^2 Hz -- average power in the band  
%}  
-----
```

```
classdef Power < dj.Relvar & dj.AutoPopulate  
methods(Access=protected)  
    function makeTuples(self, key)  
        % fetch required data  
        [lo, hi] = fetch1(example.FrequencyBand & key, 'freq_lo', 'freq_hi');  
        [signal, dt] = fetch1(example.LFP & key, 'voltage', 'dt');  
        % compute  
        signal = band_pass_filter(signal, lo, hi, dt);  
        key.power = compute_average_power(signal, dt);  
        % submit  
        self.insert(key)  
    end  
end  
end
```

MATLAB

---

# Automated computation

---

matlab

```
>> populate(example.Power)
```

python

```
>>> example.Power.populate()
```

---

# Distributed computation

---

matlab

```
>> parpopulate(example.Power)
```

python

```
>>> example.Power().populate(reserve_jobs=True)
```

---

# Query language

---

- ❖ Formulate and refine your query before retrieving data
- ❖ Three operators:
  - restrict      A & cond
  - difference    A - cond
  - join            A \* B
  - union          A + B
  - project        A.proj(*attr-list*)
  - aggregate      A.aggr(B, *computations*)

# Query language

## Restrict A&B

A & B = all rows in A that have matching rows in B

<b>id</b>	<b>species</b>	<b>sex</b>	<b>date_of_birth</b>
1	mouse	F	2015-04-01
2	monkey	M	2011-12-01
3	mouse	M	2015-05-08
4	mouse	F	2015-05-08

&

<b>id</b>	<b>scan</b>	<b>image</b>
2	1	(BLOB)
2	2	(BLOB)
2	3	(BLOB)
3	1	(BLOB)
3	2	(BLOB)
3	3	(BLOB)

=

<b>id</b>	<b>species</b>	<b>sex</b>	<b>date_of_birth</b>
2	monkey	M	2011-12-01
3	mouse	M	2015-05-08

# Query language

## Join

$A * B = \text{all combinations of matching rows from } A \text{ and } B$

<b>id</b>	<b>species</b>	<b>sex</b>	<b>date_of_birth</b>
1	mouse	F	2015-04-01
2	monkey	M	2011-12-01
3	mouse	M	2015-05-08
4	mouse	F	2015-05-08

\*

<b>id</b>	<b>scan</b>	<b>image</b>
2	1	(BLOB)
2	2	(BLOB)
2	3	(BLOB)
3	1	(BLOB)
3	2	(BLOB)
3	3	(BLOB)

=

<b>id</b>	<b>scan</b>	<b>species</b>	<b>sex</b>	<b>date_of_birth</b>	<b>image</b>
2	1	monkey	F	2011-12-01	(BLOB)
2	2	monkey	F	2011-12-01	(BLOB)
2	3	monkey	F	2011-12-01	(BLOB)
3	1	mouse	M	2015-05-08	(BLOB)
3	2	mouse	M	2015-05-08	(BLOB)
3	3	mouse	M	2015-05-08	(BLOB)

# Query language

## Project:

Select, rename, or compute columns

<b>id</b>	<b>scan</b>	<b>date_of_birth</b>	<b>image</b>
2	1	2011-12-01	(BLOB)
2	2	2011-12-01	(BLOB)
2	3	2011-12-01	(BLOB)
3	1	2015-05-08	(BLOB)
3	2	2015-05-08	(BLOB)
3	3	2015-05-08	(BLOB)

**.proj('id -> animal\_id', 'image') =**

<b>animal_id</b>	<b>scan</b>	<b>image</b>
2	1	(BLOB)
2	2	(BLOB)
2	3	(BLOB)
3	1	(BLOB)
3	2	(BLOB)
3	3	(BLOB)

---

# Expressions

---

“all the two-photon sessions using the **25x lens** that have scans that have **not yet been synchronized** with the **visual stimulus**”

```
Session & 'lens="25x"' & (Scan - Sync)
```

# Relational algebra (combinations of operators)

“all the two-photon sessions using the **25x lens** that have scans that have **not yet been synchronized** with the **visual stimulus**”

```
Session & 'lens="25x"' & (Scan - Sync)
```

**translated into SQL**

```
SELECT *
FROM `tp`.`session`
WHERE lens="25x" AND (`animal_id`, `session_id`) IN (
    SELECT `animal_id`, `session_id`
    FROM `tp`.`scan`
    WHERE ((`animal_id`, `scan_id`, `session_id`) NOT IN (
        SELECT `animal_id`, `scan_id`, `session_id`
        FROM `tp`.`_sync`)))
```

# DataJoint in publications

<http://datajoint.github.com/publications/>

1. Baden T, Berens P, Franke K, Rezac M, Bethge M, and Euler T (2015). The functional diversity of retinal ganglion cells in the mouse. *Nature*, 529(7586), pp.345-350. [link](#)
2. Cadwell CR, Palasantza A, Jiang X, Berens P, Deng Q, Reimer J, Tolias K, Bethge M, Tolias AS (2015). Morphological, electrophysiological and transcriptomic profiling of single neurons using Patch-seq. *Nature Biotechnology* [link](#)
3. Jiang X, Shen S, Cadwell CR, Berens P, Sinz F, Ecker AS, and Tolias AS (2015). Principles of connectivity among morphologically defined cell types in adult neocortex. *Science*, 350(6264), aac9462. [link](#)
4. Yatsenko D, Josic K, Ecker AS, Froudarakis E, Cotton RJ, and Tolias AS (2015). Improved estimation and interpretation of correlations in neural circuits. *PLoS Comput Biol* 11(3): e1004083. doi:10.1371/journal.pcbi.1004083 [link](#)
5. Eriskin S, Vaicieliunaite A, Jurjut O, Fiorini M, Katzner S, and Busse, L (2014). Effects of Locomotion Extend throughout the Mouse Early Visual System. *Current Biology*, 24(24), 2899-2907.
6. Reimer J, Froudarakis E, Cadwell CR, Yatsenko D, Denfield GH, Tolias AS (2014). Pupil fluctuations track fast switching of cortical states during quiet wakefulness. *Neuron*, 84(2), 355-

aac9462. [link](#)

4. Yatsenko D, Josic K, Ecker AS, Froudarakis E, Cotton RJ, and Tolias AS (2015). Improved estimation and interpretation of correlations in neural circuits. *PLoS Comput Biol* 11(3): e1004083. doi: 10.1371/journal.pcbi.1004083. [link](#)
5. Erisken S, Vaicieliunaite A, Jurjut O, Fiorini M, Katzner S, and Busse, L (2014). Effects of Locomotion Extend throughout the Mouse Early Visual System. *Current Biology*, 24(24), 2899-2907. <http://datajoint.github.com/publications/>
6. Reimer J, Froudarakis E, Cadwell CR, Yatsenko D, Denfield GH, Tolias AS (2014). Pupil fluctuations track fast switching of cortical states during quiet wakefulness. *Neuron*, 84(2), 355-362.
7. Erisken S, Vaicieliunaite A, Jurjut O, Fiorini M, Katzner S, and Busse L (2014). Effects of Locomotion Extend throughout the Mouse Early Visual System. *Current Biology*, 24(24), 2899-2907.
8. Froudarakis E, Berens P, Ecker AS, Cotton RJ, Sinz FH, Yatsenko D, Saggau P, Bethge M, and Tolias AS. Population code in mouse V1 facilitates readout of natural scenes through increased sparseness. *Nature neuroscience* (2014). [link](#)
9. Ecker AS, Berens P, Cotton RJ, Subramanyan M, Denfield GH, Cadwell CR, Smirnakis SM, Bethge M, and Tolias AS (2014): State dependence of noise correlations in macaque primary visual cortex. *Neuron* 82(1). [link](#) [code](#) [data](#) [pdf](#)
10. Cotton RJ, Froudarakis E, Storer P, Saggau P, and Tolias AS (2013). Three-dimensional mapping of microcircuit correlation structure. *Frontiers in neural circuits*, 7. [pubmed 24133414](#)
11. Vaicieliunaite A, Erisken S, Franzen F, Katzner S, and Busse L (2013). Spatial integration in mouse primary visual cortex. *Journal of Neurophysiology*, 110(4), 964-972. [pubmed 23719206](#)

# DataJoint in publications

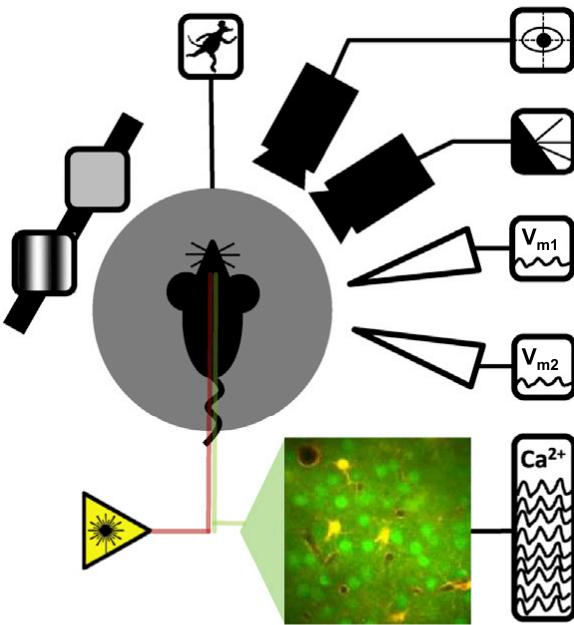


# Pupil Fluctuations Track Fast Switching of Cortical States during Quiet Wakefulness

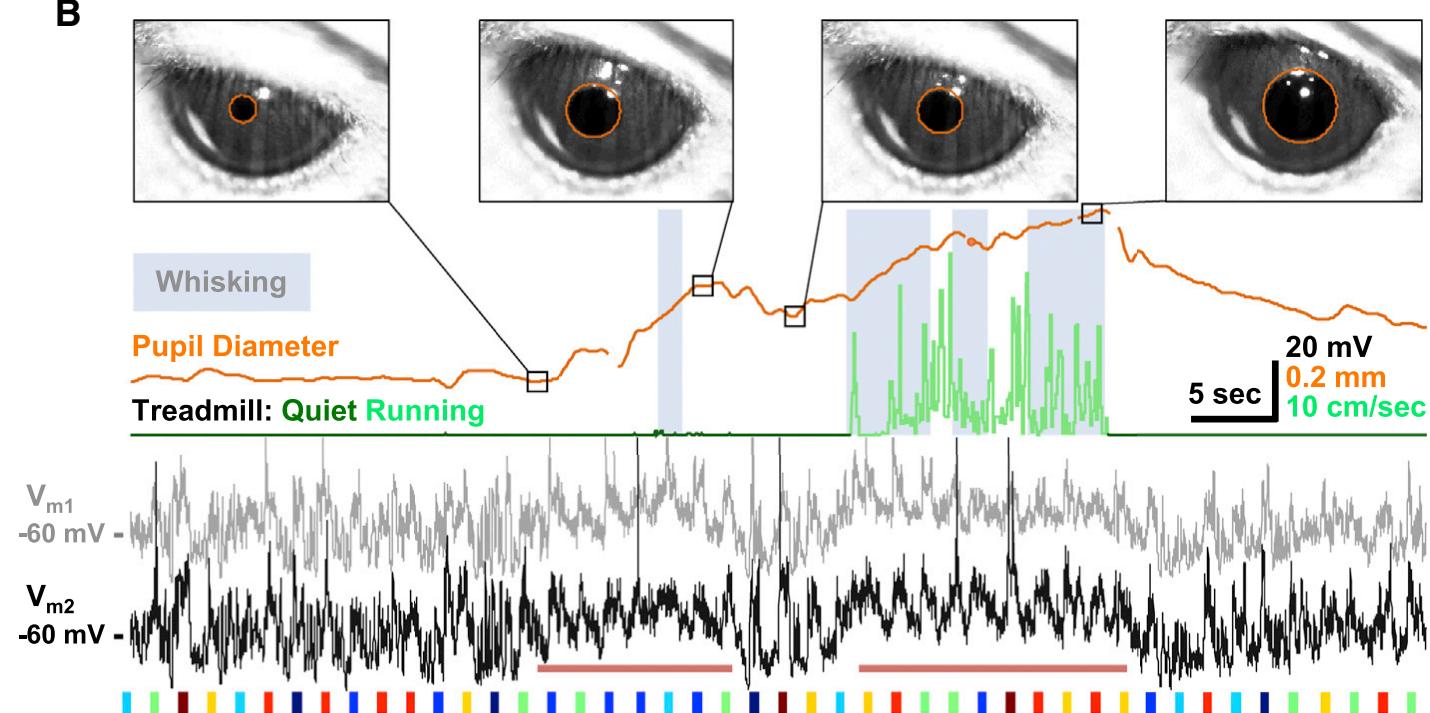
Jacob Reimer,<sup>1,\*</sup> Emmanouil Froudarakis,<sup>1</sup> Cathryn R. Cadwell,<sup>1</sup> Dimitri Yatsenko,<sup>1</sup> George H. Denfield,<sup>1</sup> and Andreas S. Tolias<sup>1,2,\*</sup>

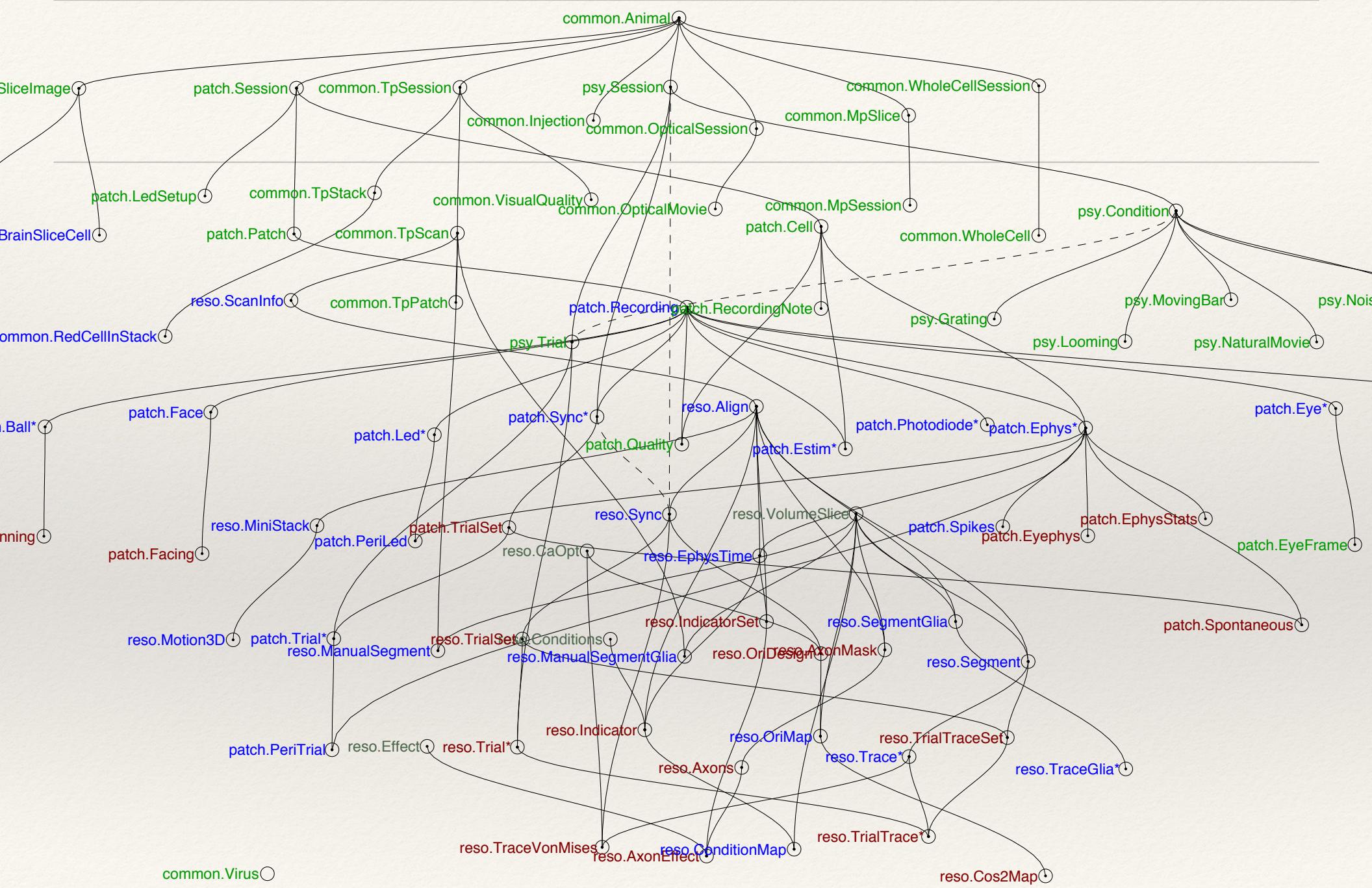
Neuron 84, 355–362, October 22, 2014

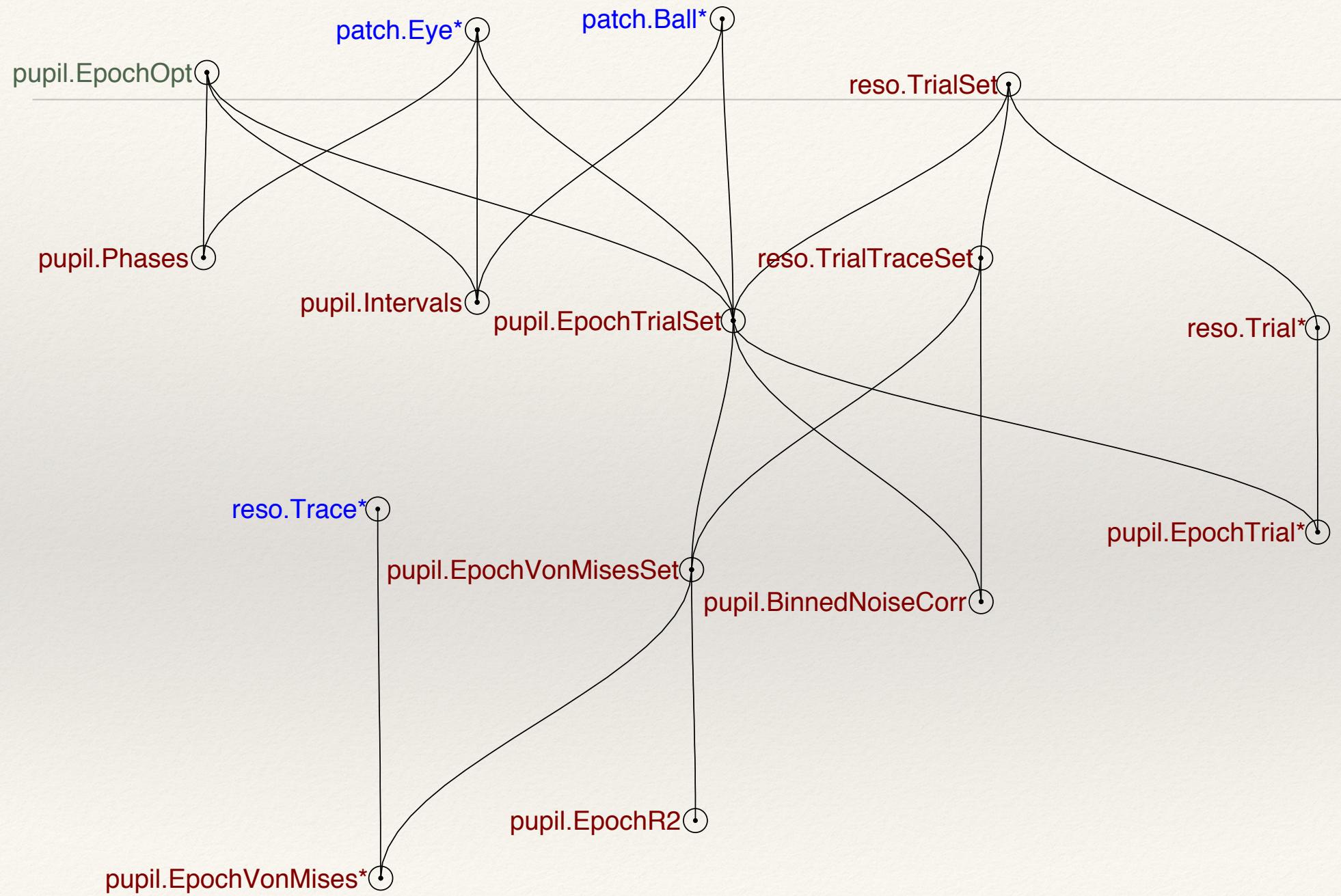
A



B





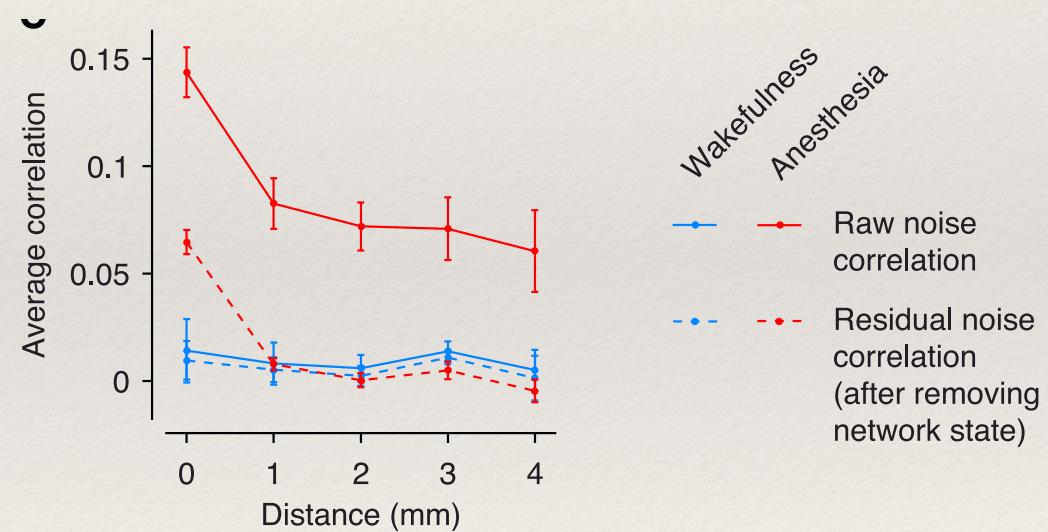
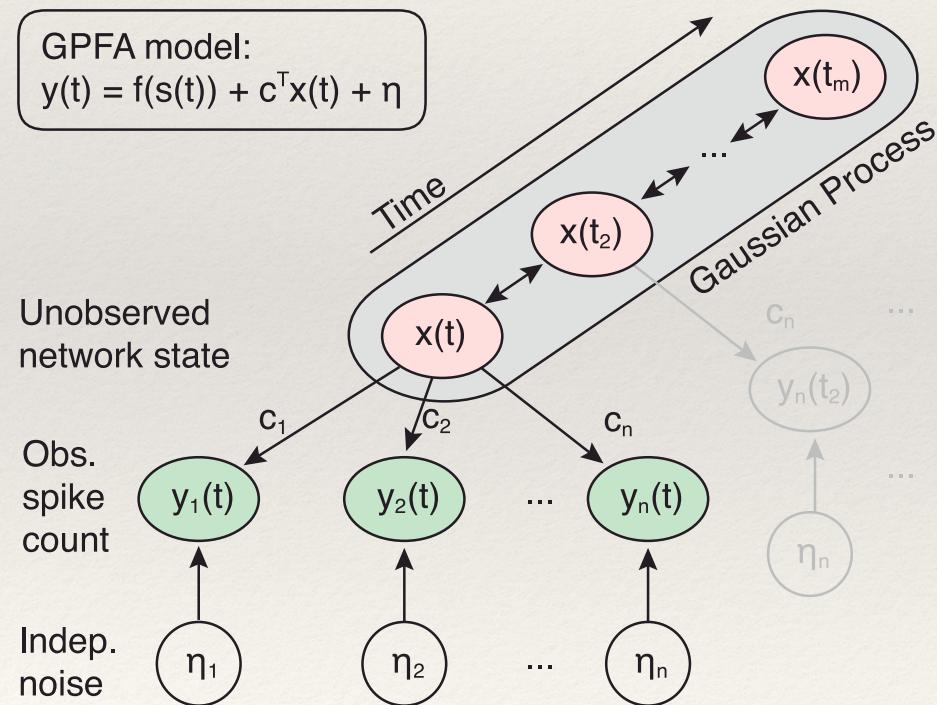




# State Dependence of Noise Correlations in Macaque Primary Visual Cortex

Alexander S. Ecker,<sup>1,2,3,5,\*</sup> Philipp Berens,<sup>1,2,3</sup> R. James Cotton,<sup>1</sup> Manivannan Subramaniyan,<sup>1</sup> George H. Denfield,<sup>1</sup> Cathryn R. Cadwell,<sup>1</sup> Stelios M. Smirnakis,<sup>1,4</sup> Matthias Bethge,<sup>2,3,5</sup> and Andreas S. Tolias<sup>1,3,6,\*</sup>

Neuron 82, 235–248, April 2, 2014 ©2014 Elsevier Inc. 235

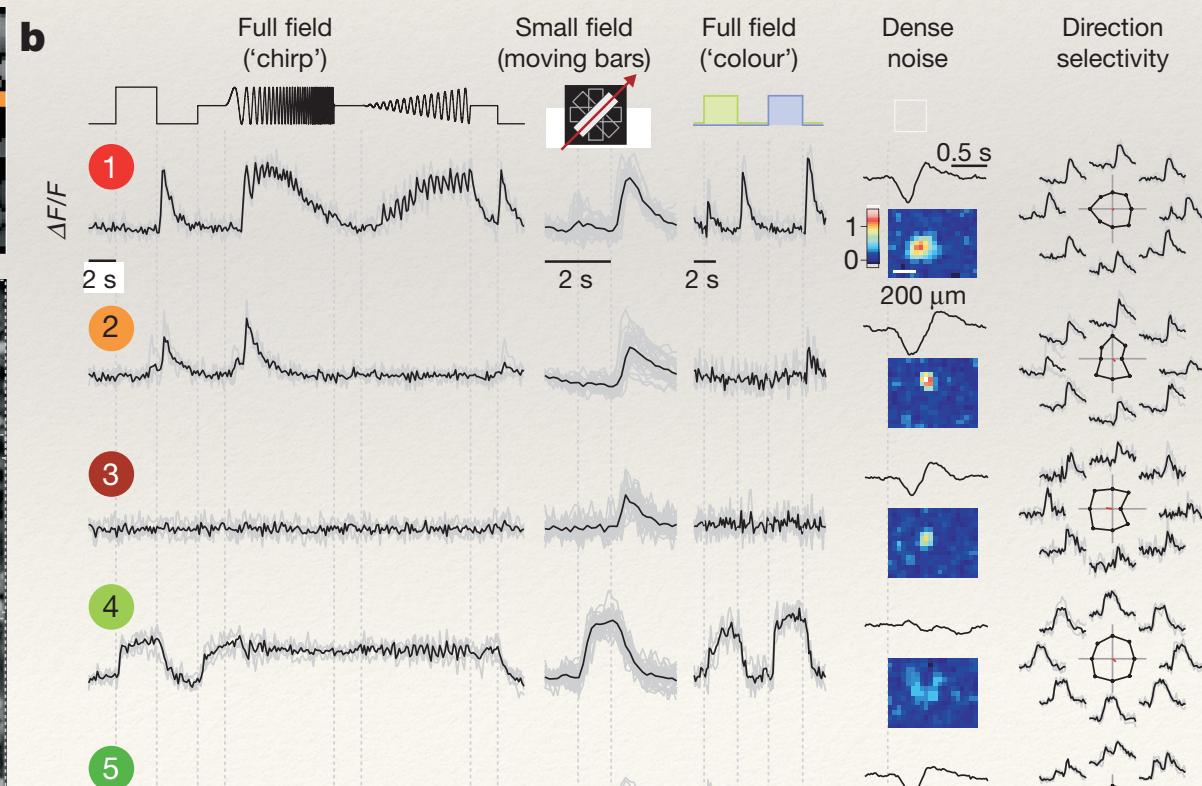
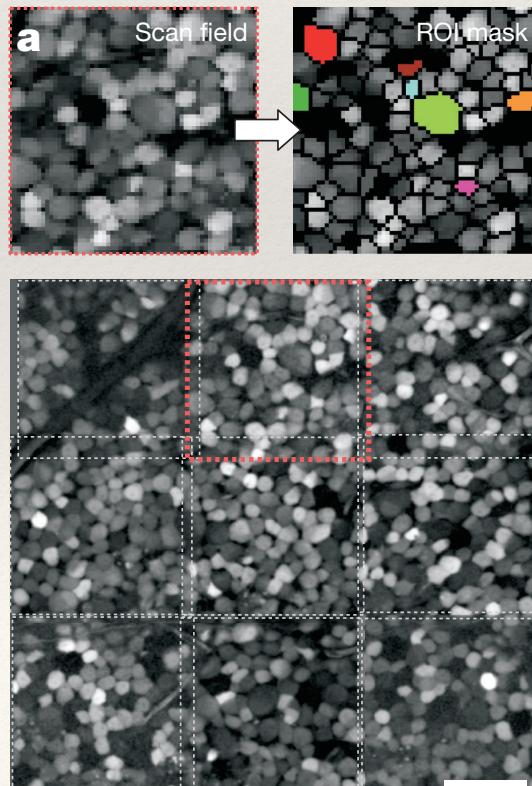


# ARTICLE

doi:10.1038/nature16468

# The functional diversity of retinal ganglion cells in the mouse

Tom Baden<sup>1,2,3\*</sup>, Philipp Berens<sup>1,2,3,4,5\*</sup>, Katrin Franke<sup>1,2,3,6\*</sup>, Miroslav Román Rosón<sup>1,2,3,6</sup>, Matthias Bethge<sup>1,2,5,7</sup> & Thomas Euler<sup>1,2,3</sup>

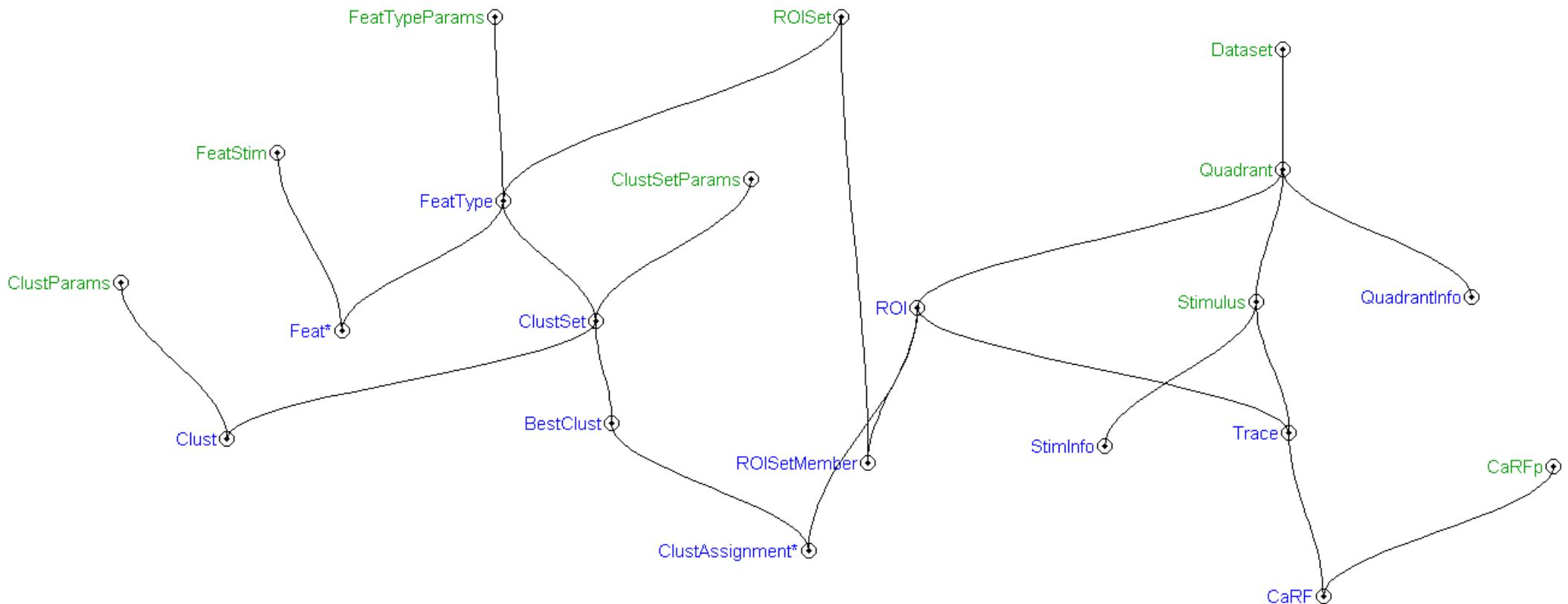


# ARTICLE

doi:10.1038/nature16468

# The functional diversity of retinal ganglion cells in the mouse

Tom Baden<sup>1,2,3\*</sup>, Philipp Berens<sup>1,2,3,4,5\*</sup>, Katrin Franke<sup>1,2,3,6\*</sup>, Miroslav Román Rosón<sup>1,2,3,6</sup>, Matthias Bethge<sup>1,2,5,7</sup> & Thomas Euler<sup>1,2,3</sup>



# Resources

---

- ❖ **Reference:**

Yatsenko, D., Reimer, J., Ecker, A.S., Walker, E.Y., Sinz, F., Berens, P., Hoenselaar, A., Cotton, R.J., Siapas, A.S. and Tolias, A.S., 2015.  
**DataJoint: managing big scientific data using MATLAB or Python.**  
bioRxiv, p.031658.

- ❖ **General information:** [datajoint.github.com](https://datajoint.github.com)

- ❖ **MATLAB:** [github.com/datajoint/datajoint-matlab](https://github.com/datajoint/datajoint-matlab)

- ❖ **Python:** [github.com/datajoint/datajoint-python](https://github.com/datajoint/datajoint-python)

---

# Acknowledgements

---

Alex Ecker

Edgar Walker

Fabian Sinz

Andreas Hoenselaar

Philipp Berens

Jacob Reimer

Andreas Tolias

---

# Thank you

---