

# Scalable, Distributed Factor Analysis in Spark

Daoyuan Jia, Tiancheng Liu, Danielle Rager

12/6/2016

# Overview

- Scientific Motivation
- Dimensionality Reduction Models:
  - PCA / sPCA / FA
- Allen Institute Dataset and Data Pre-processing
  - Distributed Implementation
- Distributed Factor Analysis (FA) Algorithm
- Conclusion

# Scientific Motivation

One research question in neuroscience:

*Neurons exhibit highly variable electrical responses, even for the same stimuli. Neuroscientists want to understand the structure of trial-to-trial variability in neural responses. Are there global effects in variability across neurons in the network?*



## Brain-data gold mine could reveal how neurons compute

Allen Brain Observatory releases unprecedented survey of activity in the mouse visual cortex.

**Helen Shen**

13 July 2016

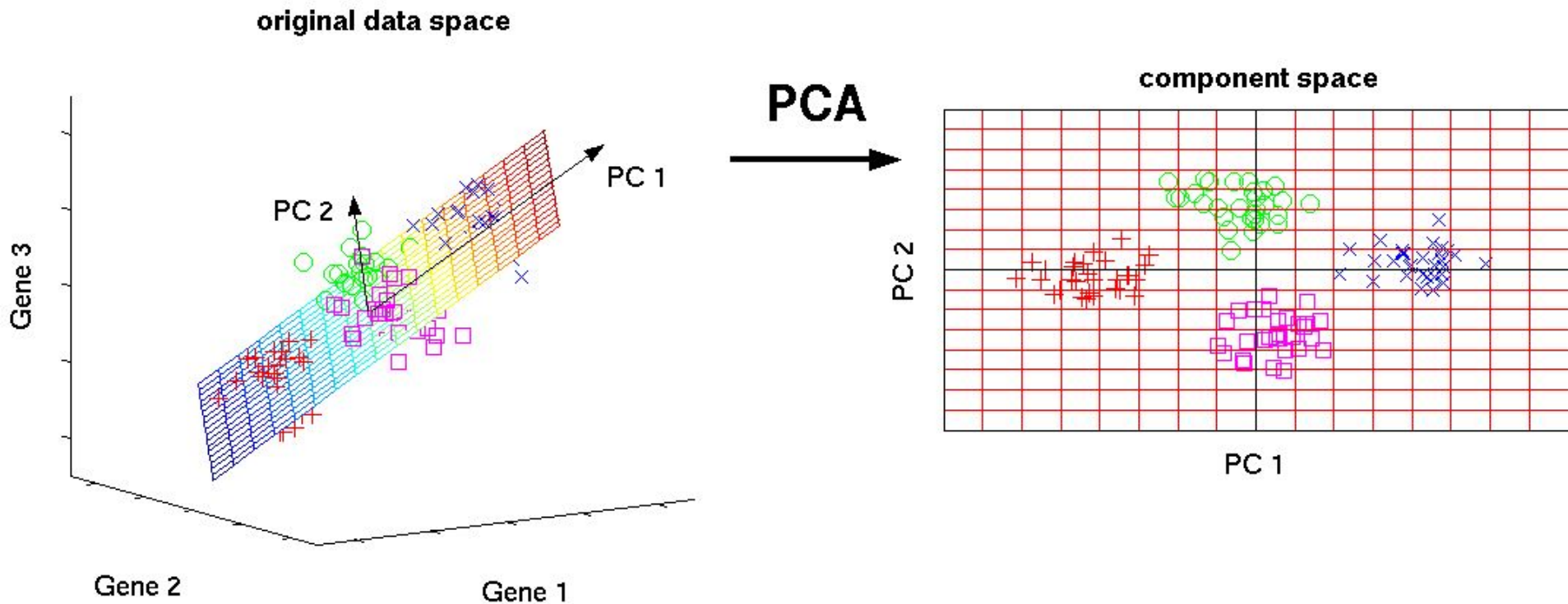
# Noises Across Neurons in the Network

- Average human brain: 100 billion neurons
- Proper dimensionality reduction techniques needed to analyze the variability across neurons
- Current algorithm: scalable
- Neuroscientists lack tools to analyze network covariability at this scale
- Solution: distributed dimensionality reduction algorithms

# Overview

- Scientific Motivation
- Dimensionality Reduction Models:
  - PCA / sPCA / FA
- Allen Institute Dataset and Data Pre-processing
  - Distributed Implementation
- Distributed Factor Analysis (FA) Algorithm
- Conclusion

# Dimensionality Reduction



# 1. Principal Component Analysis (PCA)

$$\begin{array}{c} D \times d \\ \downarrow \\ \vec{y} - \vec{\mu} = C\vec{x} \end{array}$$

- Closed-Form Solution: Singular Value Decomposition
- Complexity:  $O(ND^2)$
- Distributed Implementation:  
Spark: MLlib-PCA; R: RScaLAPACK



## 2. Probabilistic Principal Component Analysis (PPCA)

$$\begin{array}{ccc} & D \times d & N(0, \sigma I) \\ & \downarrow & \downarrow \\ \vec{y} - \vec{\mu} & = & C\vec{x} + \vec{\epsilon} \end{array}$$

- Expectation-Maximization (EM) Algorithm
- Complexity:  $O(ND)$ , better than PCA
- Distributed Implementation:

Spark/Hadoop: Stochastic Principal Component Analysis (sPCA)

### 3. Factor Analysis (FA)

$$\vec{y} - \vec{\mu} = \mathbf{C}\vec{x} + \vec{\epsilon}$$

$D \times d$                        $D \times D$   
↓                                      ↓  
 $N(0, \Psi)$

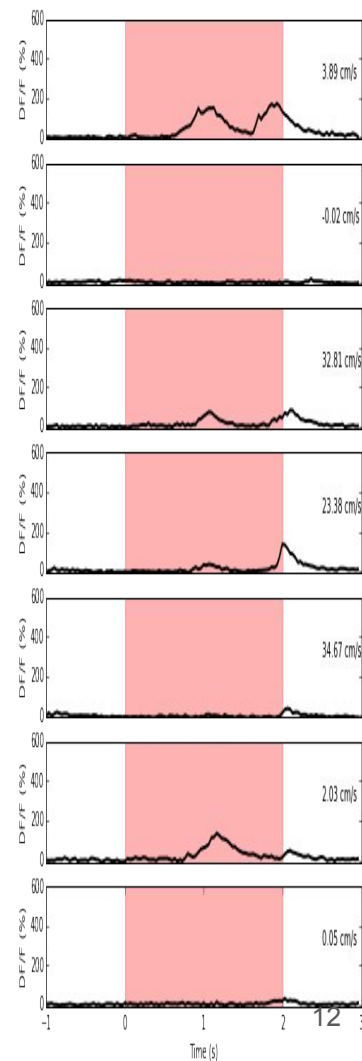
- Expectation-Maximization (EM) Algorithm
- Complexity:  $O(ND)$
- Distributed Implementation:  
    Our implementation!

# Overview

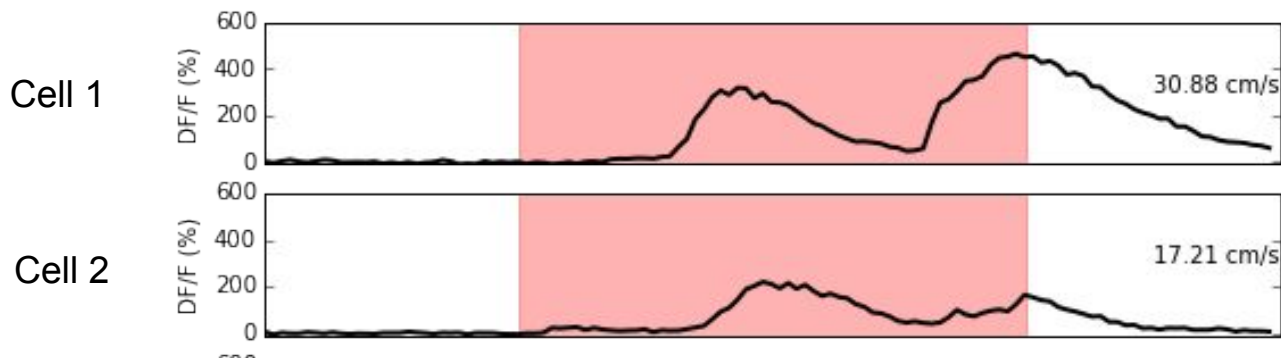
- Scientific Motivation
- Dimensionality Reduction Models:
  - PCA / sPCA / FA
- **Allen Institute Dataset and Data Pre-processing**
  - **Distributed Implementation**
- Distributed Factor Analysis (FA) Algorithm
- Conclusion

# Allen Institute Dataset

- AllenSDK:
  - For reading and processing data from Allen Institute
  - Contains an API to fetch data and experiment images
  - Also contains analytical functions for neuroscientists
- We focus on:
  - Raw data from Calcium Imaging
  - Cross-covariance between two cells' calcium response
  - Final output: covariance matrix



# Bottleneck: Calculating Covariance Matrix



- Between two cells: sum of cross-covariances up to lag  $p$
- Complexity:  $O(pTN^2) \approx O(TN^2)$
- On my laptop: 2hr for 1000 cells
- Solution: distributed matrix operations in Spark

# Overview

- Scientific Motivation
- Dimensionality Reduction Models:
  - PCA / sPCA / FA
- Allen Institute Dataset and Data Pre-processing
  - Distributed Implementation
- **Distributed Factor Analysis (FA) Algorithm**
- Conclusion

**EM algorithm for FA**  $\vec{y} - \vec{\mu} = C\vec{x} + \vec{\epsilon}$

$\uparrow$   
 $N(0, \Psi)$

- **E step:**

- Fix C and  $\Psi$
- Compute conditional likelihood  $L(X|Y)$ 
  - $\langle X|Y_c \rangle$
  - $\langle (X|Y_c)(X|Y_c)^T \rangle$

- **M step:**

- Fix conditional likelihood
- Compute new C and  $\Psi$

# EM Algorithm for FA

$$Y_c = Y - \mu$$

$$M = CC^T + \Psi$$

$$\langle X|Y_c \rangle = X_m = C^T M^{-1} Y_c$$

$$\langle (X|Y_c)(X|Y_c)^T \rangle = \Sigma_{X_m} = I - C^T M^{-1} C + X_m X_m^T$$

$$Y_{proj} = Y_c X_m^T$$

$$C_{new} = \Sigma_{X_m} Y_{proj}$$

$$\psi_{proj} = C_{new} X_m Y_m^T$$

$$\Psi_{new} = \frac{1}{N} \text{diag}(Y_c Y_c^T - \psi_{proj})$$



# sFA Optimizations: Distribute operations minimally

Driver program does most operations locally, launches only a few Spark jobs  
distribute computations where you have dimension  $N$

$$Y_c = Y - \mu$$

$$M = CC^T + \Psi$$

$$\langle X|Y_c \rangle = X_m = C^T M^{-1} Y_c \quad \text{MapReduce}$$

$$\langle (X|Y_c)(X|Y_c)^T \rangle = \Sigma_{X_m} = I - C^T M^{-1} C + X_m X_m^T \quad \text{MapReduce}$$

$$Y_{proj} = Y_c X_m^T \quad \text{MapReduce}$$

$$C_{new} = \Sigma_{X_m} Y_{proj}$$

$$\psi_{proj} = C_{new} X_m Y_m^T \quad \text{MapReduce}$$

$$\Psi_{new} = \frac{1}{N} \text{diag}(Y_c Y_c^T - \psi_{proj})$$

# sFA Optimizations: Distribute operations minimally

Use same MapReduce job for operations w/o dependencies

$$Y_c = Y - \mu$$

$$M = CC^T + \Psi$$

$$\langle X|Y_c \rangle = X_m = C^T M^{-1} Y_c \quad \text{MapReduce}$$

$$\langle (X|Y_c)(X|Y_c)^T \rangle = \Sigma_{X_m} = I - C^T M^{-1} C + X_m X_m^T \quad \text{MapReduce}$$

$$Y_{proj} = Y_c X_m^T$$

$$C_{new} = \Sigma_{X_m} Y_{proj}$$

$$\psi_{proj} = C_{new} X_m Y_m^T \quad \text{MapReduce}$$

$$\Psi_{new} = \frac{1}{N} \text{diag}(Y_c Y_c^T - \psi_{proj})$$

# sFA Optimizations: Minimize Intermediary Data

Recompute  $X$  and  $Y$  at each job rather than storing and exchanging

$$Y_c = Y - \mu$$

$$M = CC^T + \Psi$$

$$\langle X|Y_c \rangle = X_m = C^T M^{-1} Y_c$$

$$\langle (X|Y_c)(X|Y_c)^T \rangle = \Sigma_{X_m} = I - C^T M^{-1} C + X_m X_m^T$$

$$Y_{proj} = Y_c X_m^T$$

$$C_{new} = \Sigma_{X_m} Y_{proj}$$

$$\psi_{proj} = C_{new} X_m Y_m^T$$

$$\Psi_{new} = \frac{1}{N} \text{diag}(Y_c Y_c^T - \psi_{proj})$$

# sFA Optimizations: Minimize Intermediary Data

Recompute X and Y at each job rather than storing and exchanging

$$Y_c = Y - \mu$$

$$M = CC^T + \Psi$$

~~$$\langle X|Y_c \rangle = X_m = C^T M^{-1} Y_c$$~~

$$\langle (X|Y_c)(X|Y_c)^T \rangle = \Sigma_{X_m} = I - C^T M^{-1} C + \boxed{X_m} X_m^T \quad \text{MapReduce}$$

$$Y_{proj} = Y_c X_m^T$$

$$C_{new} = \Sigma_{X_m} Y_{proj}$$

$$\psi_{proj} = C_{new} \boxed{X_m} Y_m^T \quad \text{MapReduce}$$

$$\Psi_{new} = \frac{1}{N} \text{diag}(Y_c Y_c^T - \psi_{proj})$$

# sFA Optimizations: Leverage Sparsity

(Don't lose opportunity to do computations on 0 values)

Mean sparsity:  $M^{-1}Y_c = M^{-1}(Y - \mu) = M^{-1}\boxed{Y} - M^{-1}\mu$

**sparse**

Matrix Inversion Lemma:  $M^{-1} = (CC^T + \Psi)^{-1}$   
 $= \boxed{\Psi^{-1}} - \boxed{\Psi^{-1}C}(I + C^T\Psi^{-1}C)^{-1}\boxed{C^T\Psi^{-1}}$

**Diagonal matrix  
multiplication tricks**

# sFA Optimizations: Efficient Matrix Multiplication



**Optimizes computations when both matrices are large**

$$(A * B)_i = A_i * B$$

**Optimization when one matrix is small enough to fit in memory**

# Overview

- Scientific Motivation
- Dimensionality Reduction Models:
  - PCA / sPCA / FA
- Allen Institute Dataset and Data Pre-processing
  - Distributed Implementation
- Distributed Factor Analysis (FA) Algorithm
- **Conclusion**

# Conclusion: Format for Large Distributed Files

## H5Spark: Bridging the I/O Gap between Spark and Scientific Data Formats on HPC Systems

Jialin Liu<sup>1</sup>, Evan Racah<sup>1</sup>, Quincey Koziol<sup>1</sup>, Richard Shane Canon<sup>1</sup>,  
Alex Gittens<sup>2</sup>, Lisa Gerhardt<sup>1</sup>, Suren Byna<sup>1</sup>, Mike F. Ringenburt<sup>3</sup>, Prabhat<sup>1</sup>.

- HDF5:



- Hierarchical Data Format V
- Flexible and efficient I/O
- High volume and complex
- NWB from Allen Institute

- HFSpark



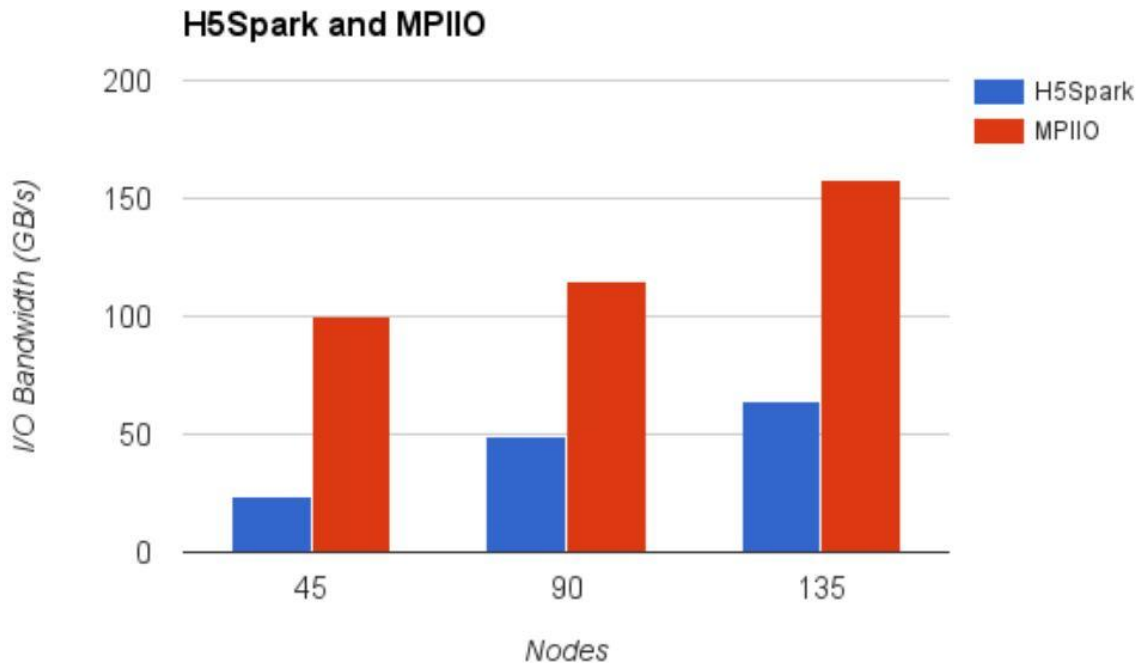
- Department of Energy
- From HDF5 to Spark RDD
- Implemented in JavaSpark
- Superior performance



# Conclusion: Format for Large Distributed Files

## MPI IO

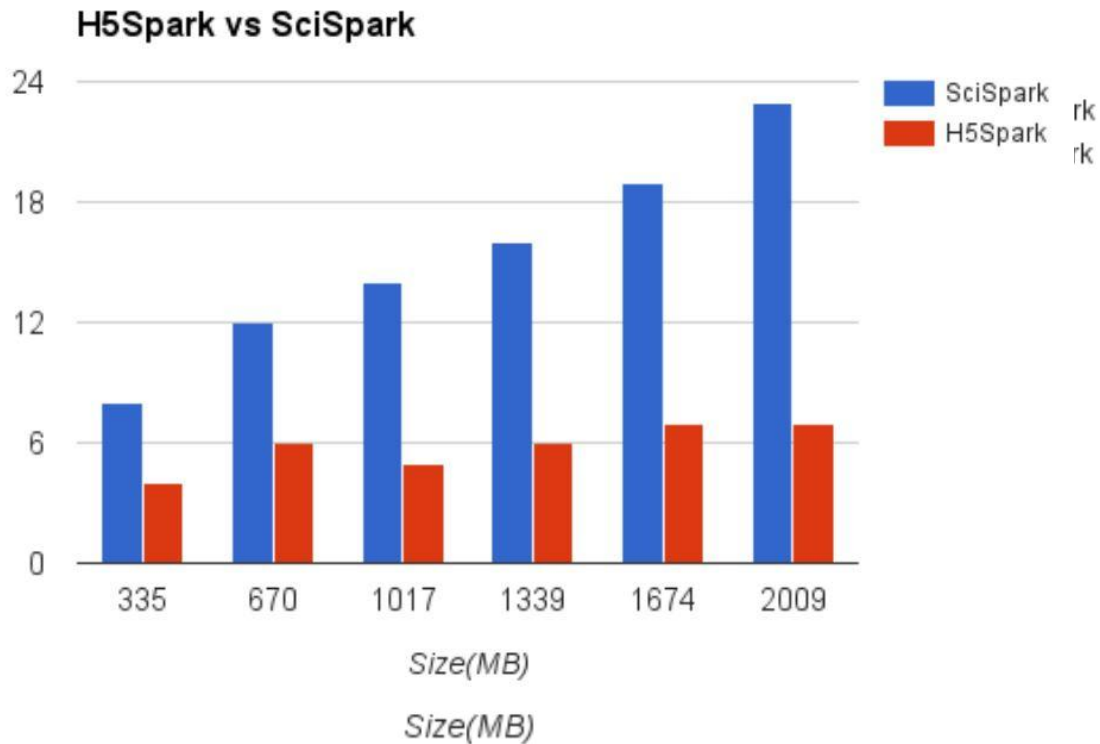
- IBM's Watson Laboratory
- Parallel I/O



# Conclusion: Format for Large Distributed Files

## SciSpark

- NASA
- scientific RDD (SRDD)



# Reference

1. Allen Institute. "Brain Observatory Trace Analysis." [http://alleninstitute.github.io/AllenSDK/\\_static/examples/nb/brain\\_observatory\\_analysis.html#Drifting-Gratings](http://alleninstitute.github.io/AllenSDK/_static/examples/nb/brain_observatory_analysis.html#Drifting-Gratings)
2. Tarek Elgamal, Maysam Yabandeh, Ashraf Aboulnaga, Waleed Mustafa, and Mohamed Hefeeda. 2015. sPCA: Scalable Principal Component Analysis for Big Data on Distributed Platforms. In Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15). ACM, New York, NY, USA, 79-91. DOI: <http://dx.doi.org/10.1145/2723372.2751520>
3. [http://johnsonhsieh.github.io/DSR\\_workshop/image/fig\\_pca\\_principal\\_component\\_analysis.png](http://johnsonhsieh.github.io/DSR_workshop/image/fig_pca_principal_component_analysis.png)
4. [https://cug.org/proceedings/cug2016\\_proceedings/includes/files/pap137-file2.pdf](https://cug.org/proceedings/cug2016_proceedings/includes/files/pap137-file2.pdf)
5. J.L. Liu, E. Racah, Q. Koziol, R. S. Canon, A. Gittens, L. Gerhardt, S. Byna, M. F. Ringenburg, Prabhat. "H5Spark: Bridging the I/O Gap between Spark and Scientific Data Formats on HPC Systems", Cray User Group, 2016