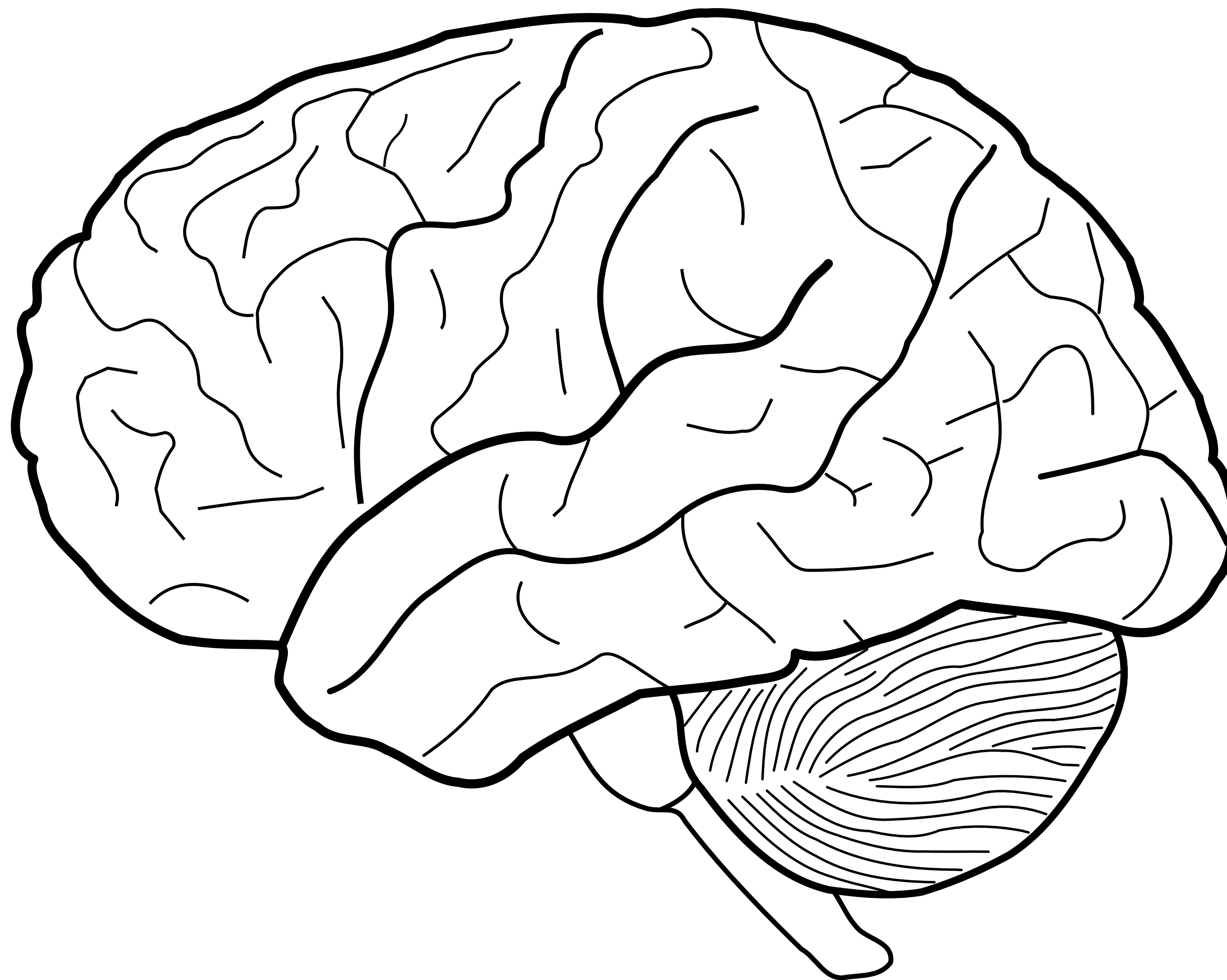# Neuroscience Introduction
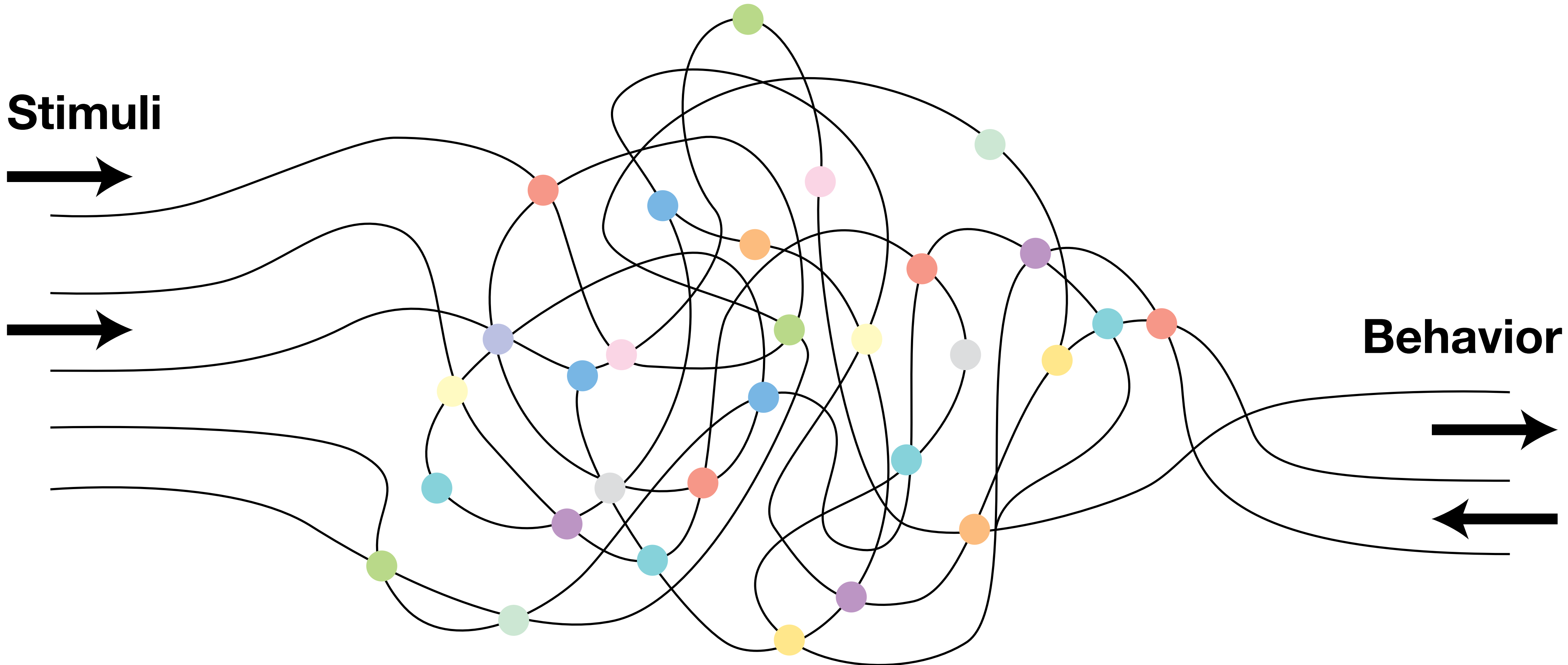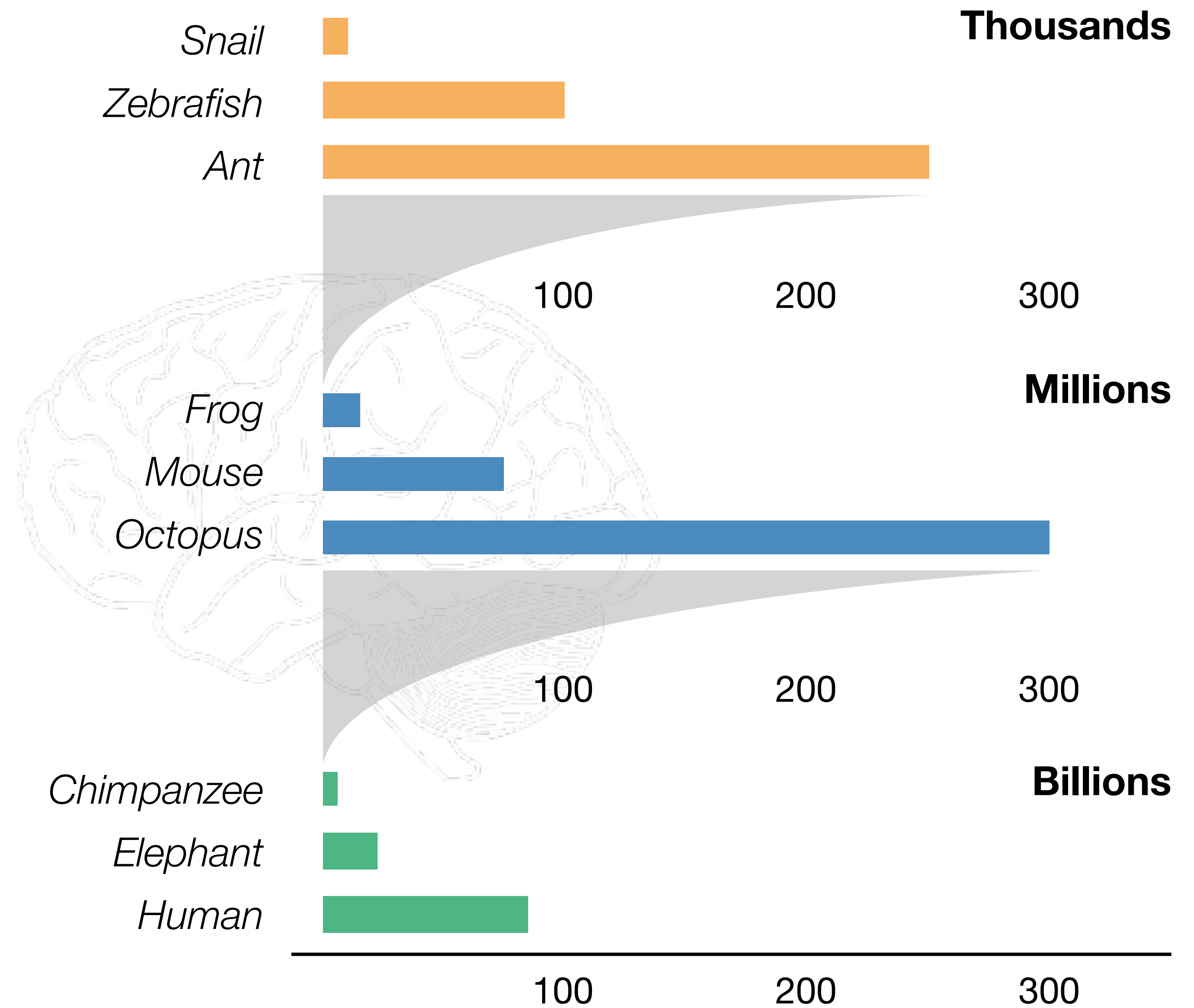
# The brain

"As humans, we can identify galaxies light years away, we can study particles smaller than an atom. **But we still haven't unlocked the mystery of the three pounds of matter that sits between our ears**."

*President Obama*
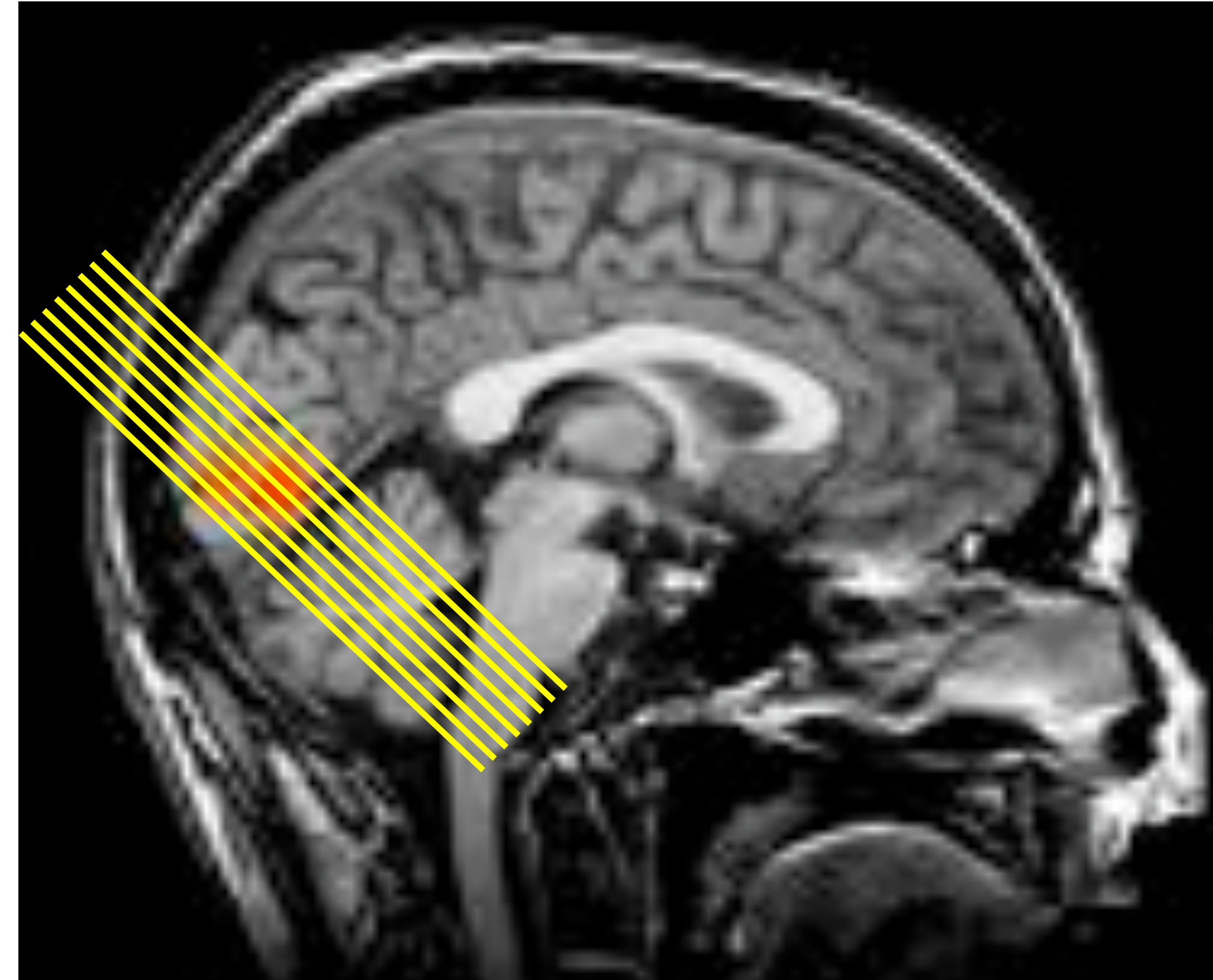
# The brain

**Stimuli**

**Behavior**

# Numbers of neurons



**Thousands**

*Snail*
*Zebrafish*
*Ant*

100   200   300

**Millions**

*Frog*
*Mouse*
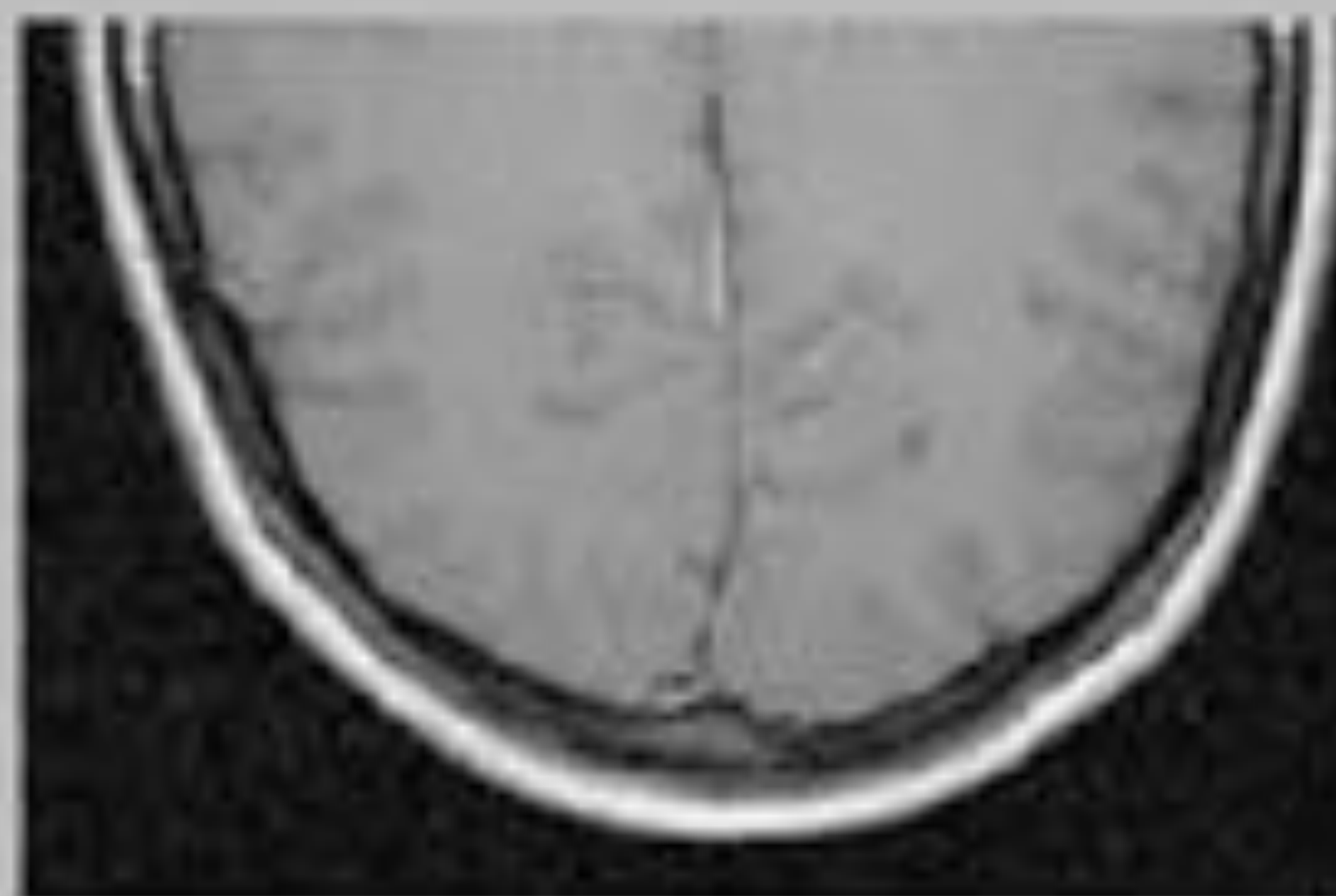*Octopus*

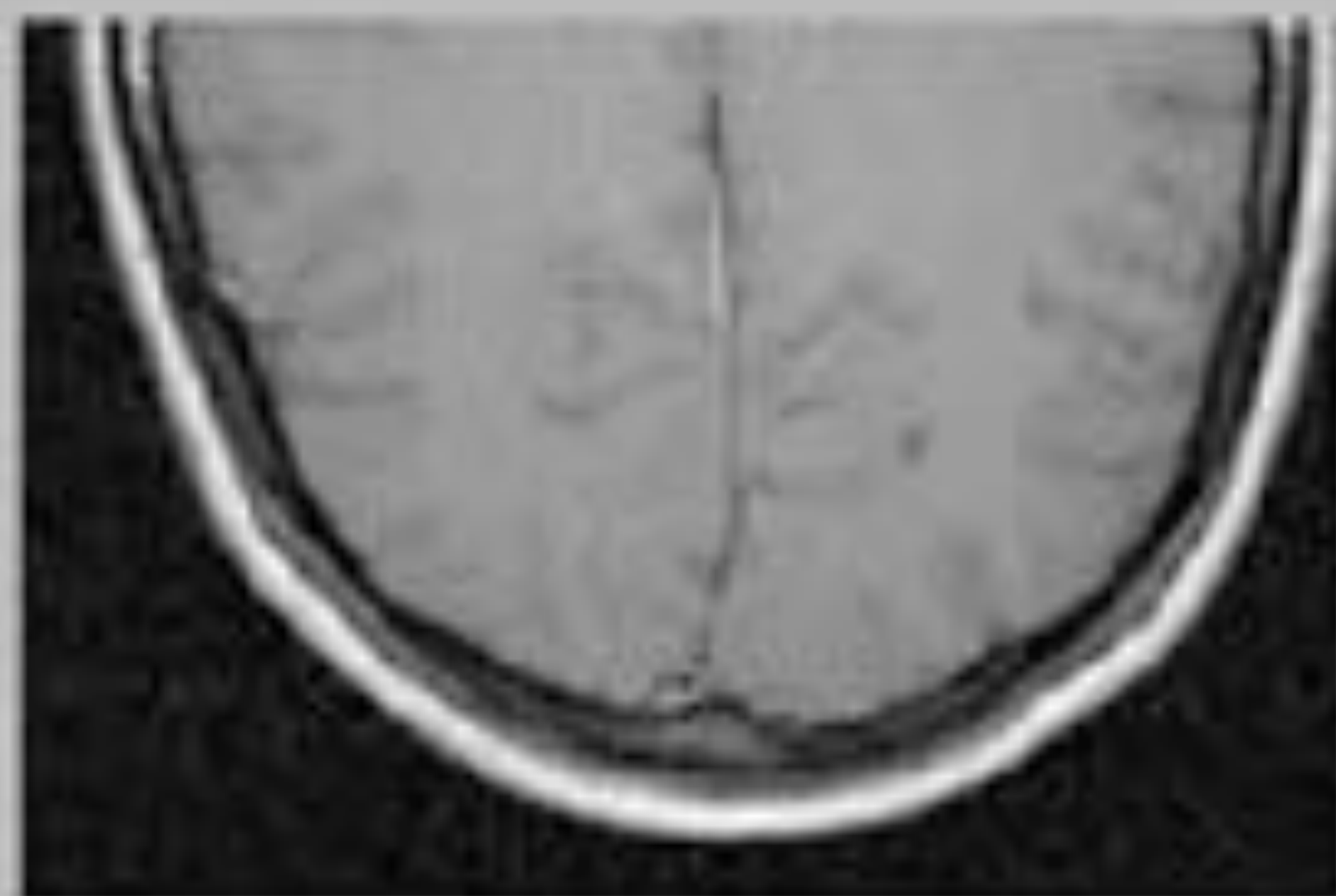100   200   300

**Billions**

*Chimpanzee*
*Elephant*
*Human*

100   200   300

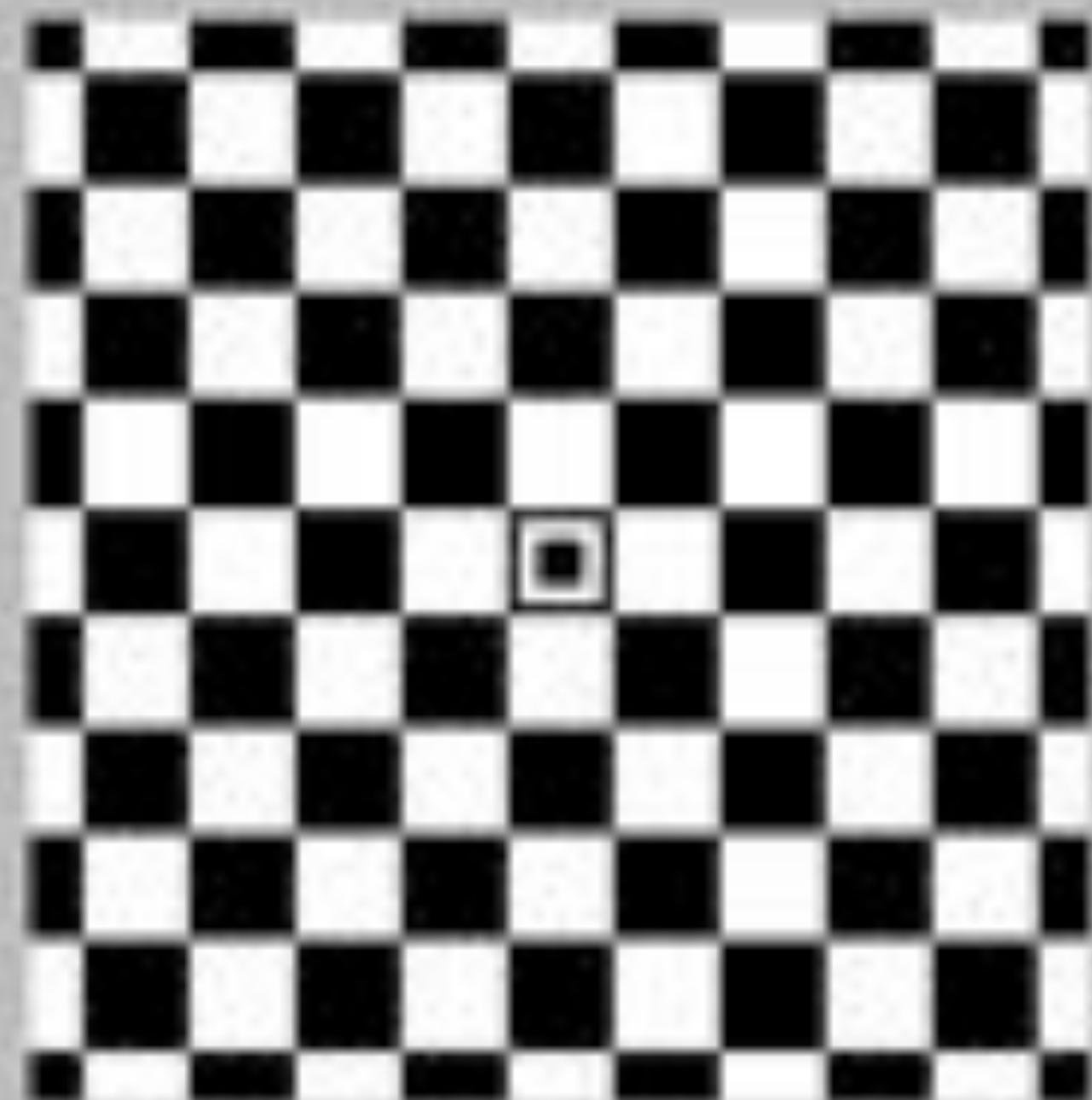# Studying the brain in humans



fMRI scanner



human brain

~50,000 neurons per cubic millimeter
-> need higher resolution!

*multielectrode*
*10-100*

*two-photon*
*100-1000*

*light-sheet*
*100000*

PHOTON DETECTOR

CAMERA

Vladimirov, et al., 2014

Sofroniew, et al., 2014

# relating neuronal responses to properties of an animal and its environment



"place cell"

"grid cell"

● position of a mouse in maze

Moser et al., 2008

fine-scale
sensory
tuning

Hubel & Weisel, 1959

Ohki et al., 2006

Mouse, somatosensory cortex
~1,000 neurons

| 0.1 TB / experiment

Larval zebrafish, whole-brain
~100,000 neurons

■ 1 TB

* Entire mouse brain
~80,000,000 neurons

▬▬▬▬▬▬ // ■ >100 TB

* hypothetical

# Exploratory
# Data Analysis

This is really big

Raw data

This is complex

Extracted signals

Analysis

Visualization

Sharing

Exploring

Interactive feedback

## Supervised methods

$$\mathbf{y} = f(\mathbf{X})$$

predict our data

as a function of other data

## Unsupervised methods

$$f(\mathbf{X})$$

find structure in the data on its own

# Clustering for preprocessing

- Raw data is complex and high-dimensional

- Clustering finds collections of inputs that are similar to one another

- These groups of clusters may be the more meaningful "unit" of measurement

Raw data
Clustered data

# Clustering to find waveforms associated with individual neurons based on their traces across multiple electrodes

# Dimensionality reduction for insight

- Raw data is complex and high-dimensional

- Dimensionality reduction describes the data using a simpler, more compact representation

- This representation may make interesting patterns in the data more clear or easier to see

"Shoe store" space



"Size" space

$$\underbrace{\begin{array}{cc} \bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet \\ \rule{4cm}{0.4pt} \\ 6 \qquad 12 \\ \text{American sizes} \\[2mm] \bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet \\ \rule{4cm}{0.4pt} \\ 39 \qquad 45 \\ \text{European sizes} \end{array}} = \overset{\text{weights}}{\begin{bmatrix} 6s \\[4mm] 6s + 33 \end{bmatrix}} \star \begin{bmatrix} s \end{bmatrix} \longleftrightarrow \begin{array}{c} \bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet\!-\!\bullet \\ 1 \; \dfrac{\rule{3cm}{0.4pt}}{s} \; 2 \end{array}$$

**Dimensionality reduction**

Neurons

$r_1$

$t$

$r_2$

$t$

$r_3$

$t$

Weights

$$\begin{bmatrix} 0.8 & -0.6 \\ -0.5 & -0.6 \\ 0.3 & 0.5 \end{bmatrix} * \begin{bmatrix} s_1(t) \\ s_2(t) \end{bmatrix}$$

Population space

$r_3$

$s_2$

Time

$s_1$

$r_2$

$r_1$

**Dimensionality reduction**

Yu and Cunningham, 2014

Swim Trials
Crawl Trials

Time (s)

Briggman et al., 2005

When the
leech changes
its mind!

When the leech crawls

When the
leech swims

Briggman et al., 2005

# Principal Component Analysis (PCA) Overview

# Raw data can be Complex, High-dimensional

To understand a phenomenon we measure various related quantities

If we knew what to measure or how to represent our measurements we might find simple relationships

But in practice we often *measure redundant signals*, e.g., US and European shoe sizes

We also *represent data via the method by which it was gathered*, e.g., pixel representation of brain imaging data

# Dimensionality Reduction

**Issues**

- *Measure redundant signals*
- *Represent data via the method by which it was gathered*

**Goal**: Find a 'better' representation for data

- To visualize and discover hidden patterns
- Preprocessing for supervised task, e.g., feature hashing

How do we define 'better'?

# E.g., Shoe Size

We take noisy measurements on European and American scale

- Modulo noise, we expect perfect correlation

How can we do 'better', i.e., find a simpler, compact representation?

- Pick a direction and project onto this direction

# E.g., Shoe Size

We take noisy measurements on European and American scale

- Modulo noise, we expect perfect correlation

How can we do 'better', i.e., find a simpler, compact representation?

- Pick a direction and project onto this direction

# Goal: Minimize Reconstruction Error

Minimize Euclidean distances between original points and their projections

PCA solution solves this problem!

**Linear Regression** — predict $y$ from $x$. Evaluate accuracy of predictions (represented by blue line) by **vertical** distances between points and the line

**PCA** — reconstruct 2D data via 2D data with single degree of freedom. Evaluate reconstructions (represented by blue line) by **Euclidean** distances

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

Can we do 'better', i.e., find a compact representation that captures variation?

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

Can we do 'better', i.e., find a compact representation that captures variation?

# Another Goal: Maximize Variance

To identify patterns we want to study variation across observations

Can we do 'better', i.e., find a compact representation that captures variation?

PCA solution finds directions of maximal variance!

# PCA Assumptions and Solution

# PCA Formulation

PCA: find lower-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance constraints

Linearity assumption ( $\mathbf{Z} = \mathbf{XP}$ ) simplifies problem

$$\mathbf{Z} = \mathbf{X} \, \mathbf{P}$$

Given $n$ training points with $d$ features:

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix storing points
- $x_j^{(i)}$ : $j$th feature for $i$th point
- $\mu_j$ : mean of $j$th feature

Variance of 1st feature $\quad \sigma_1^2 = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( x_1^{(i)} - \mu_1 \right)^2$

Variance of 1st feature (assuming zero mean) $\quad \sigma_1^2 = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( x_1^{(i)} \right)^2$

Given $n$ training points with $d$ features:

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : matrix storing points

- $x_j^{(i)}$ : $j$th feature for $i$th point

- $\mu_j$ : mean of $j$th feature

Covariance of 1st and 2nd features (assuming zero mean)

$$\sigma_{12} = \frac{1}{n} \sum_{i=1}^{n} x_1^{(i)} x_2^{(i)}$$

- Symmetric: $\sigma_{12} = \sigma_{21}$

- Zero $\rightarrow$ uncorrelated

- Large magnitude $\rightarrow$ (anti) correlated / redundant

- $\sigma_{12} = \sigma_1^2 = \sigma_2^2$ $\rightarrow$ features are the same

# Covariance Matrix

Covariance matrix generalizes this idea for many features

$d \times d$ covariance matrix with zero mean features

$$\mathbf{C_X} = \frac{1}{n}\mathbf{X}^\top\mathbf{X}$$

- $i$th diagonal entry equals variance of $i$th feature
- $ij$th entry is covariance between $i$th and $j$th features
- Symmetric (makes sense given definition of covariance)

Variance: $\sigma_1^2 = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( x_1^{(i)} \right)^2$

$$\begin{bmatrix} 2 & -1 & -1 \\ 3 & 2 & -5 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ -1 & 2 \\ -1 & -5 \end{bmatrix} = \begin{bmatrix} 6 & \\ & \end{bmatrix}$$

$\mathbf{X}^\top \qquad\qquad \mathbf{X} \qquad\qquad \mathbf{X}^\top\mathbf{X}$

Variance: $\sigma_1^2 = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( x_1^{(i)} \right)^2$

Covariance: $\sigma_{12} = \dfrac{1}{n} \sum\limits_{i=1}^{n} x_1^{(i)} x_2^{(i)}$

$$\begin{bmatrix} 2 & -1 & -1 \\ 3 & 2 & -5 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ -1 & 2 \\ -1 & -5 \end{bmatrix} = \begin{bmatrix} 6 & 9 \\ 9 & 38 \end{bmatrix}$$

$\mathbf{X}^\top \qquad\qquad \mathbf{X} \qquad\qquad \mathbf{X}^\top \mathbf{X}$

Dividing by $n$ yields
covariance matrix

# PCA Formulation

PCA: find lower-dimensional representation of raw data
- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{XP}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance / Covariance constraints

What constraints make sense in reduced representation?
- No feature correlation, i.e., all off-diagonals in $\mathbf{C_Z}$ are zero
- Rank-ordered features by variance, i.e., sorted diagonals of $\mathbf{C_Z}$

# PCA Formulation

PCA: find lower-dimensional representation of raw data

- $\mathbf{X}$ is $n \times d$ (raw data)
- $\mathbf{Z} = \mathbf{X}\mathbf{P}$ is $n \times k$ (reduced representation, PCA 'scores')
- $\mathbf{P}$ is $d \times k$ (columns are $k$ principal components)
- Variance / Covariance constraints

$\mathbf{P}$ equals the top $k$ eigenvectors of $\mathbf{C_X}$

$$\begin{bmatrix} \mathbf{Z} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \end{bmatrix}\begin{bmatrix} \mathbf{P} \end{bmatrix}$$

# PCA Solution

All covariance matrices have an eigendecomposition

- $\mathbf{C_X} = \mathbf{U \Lambda U}^\top$ (eigendecomposition)
- $\mathbf{U}$ is $d \times d$ (column are eigenvectors, sorted by their eigenvalues)
- $\mathbf{\Lambda}$ is $d \times d$ (diagonals are eigenvalues, off-diagonals are zero)

The $d$ eigenvectors are orthonormal directions of max variance

- Associated eigenvalues equal variance in these directions
- 1st eigenvector is direction of max variance (variance is $\lambda_1$)

In lab, we'll use the `eigh` function from `numpy.linalg`

# Choosing $k$

How should we pick the dimension of the new representation?

**Visualization**: Pick top 2 or 3 dimensions for plotting purposes

**Other analyses**: Capture 'most' of the variance in the data
- Recall that eigenvalues are variances in the directions specified by eigenvectors, and that eigenvalues are sorted

- Fraction of retained variance: $\dfrac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{d} \lambda_i}$

Can choose $k$ such that we retain some fraction of the variance, e.g., 95%

# Other Practical Tips

PCA assumptions (linearity, orthogonality) not always appropriate
- Various extensions to PCA with different underlying assumptions, e.g., manifold learning, Kernel PCA, ICA

Centering is crucial, i.e., we must preprocess data so that all features have zero mean before applying PCA

PCA results dependent on scaling of data
- Data is sometimes rescaled in practice before applying PCA

# PCA Algorithm

# Orthogonal and Orthonormal Vectors

*Orthogonal* vectors are **perpendicular** to each other
- Equivalently, their dot product equals zero
- $\mathbf{a}^\top \mathbf{b} = 0$ and $\mathbf{d}^\top \mathbf{b} = 0$, but $\mathbf{c}$ isn't orthogonal to others



$$\mathbf{a} = \begin{bmatrix} 1 & 0 \end{bmatrix}^\top \qquad \mathbf{b} = \begin{bmatrix} 0 & 1 \end{bmatrix}^\top \qquad \mathbf{c} = \begin{bmatrix} 1 & 1 \end{bmatrix}^\top \qquad \mathbf{d} = \begin{bmatrix} 2 & 0 \end{bmatrix}^\top$$

*Orthonormal* vectors are orthogonal and have unit norm
- $\mathbf{a}$ are $\mathbf{b}$ are orthonormal, but $\mathbf{b}$ are $\mathbf{d}$ are not orthonormal

# PCA Iterative Algorithm

$k = 1$: Find direction of max variance, project onto this direction

- Locations along this direction are the new 1D representation

More generally, for $i$ in $\{1, ..., k\}$:

- Find direction of max variance that is *orthonormal* to previously selected directions, project onto this direction
- Locations along this direction are the $i$th feature in new representation