

*All models are wrong,
but some are useful.*

– George E. P. Box

Ilya Kuzovkin

Understanding Information Processing in Human Brain by Interpreting Machine Learning Models

Supervised by Raul Vicente



UNIVERSITY
OF TARTU

Modeling is a well-proven way of obtaining knowledge

MANUAL

AUTOMATED

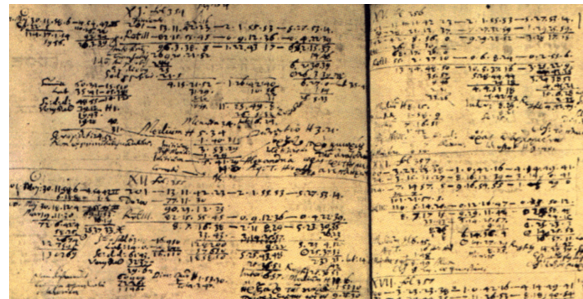
Modeling is a well-proven way of obtaining knowledge

DATA

ALGORITHM

MODEL

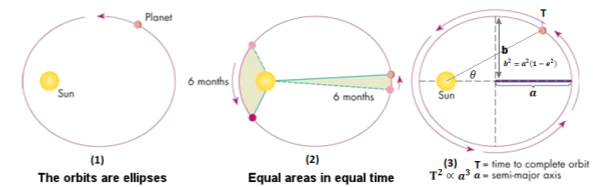
MANUAL



Tycho Brahe's observations



Johannes Kepler



Three Laws of Planetary Motion

AUTOMATED

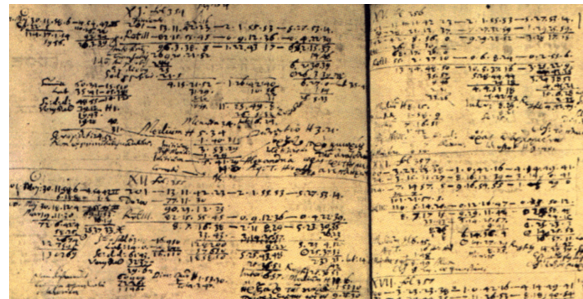
Modeling is a well-proven way of obtaining knowledge

DATA

ALGORITHM

MODEL

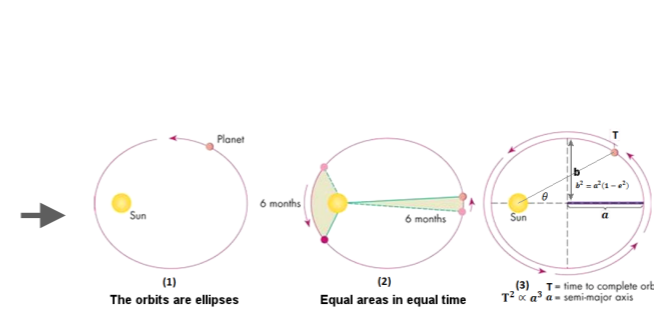
MANUAL



Tycho Brahe's observations



Johannes Kepler

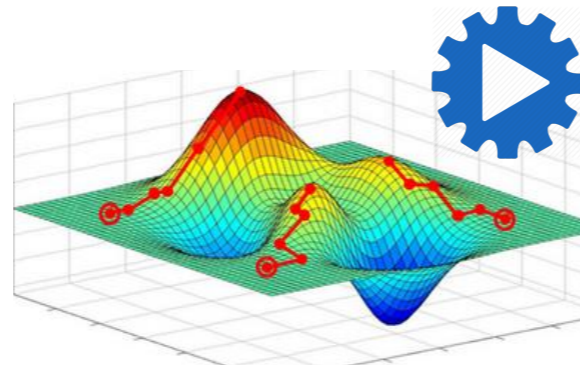


Three Laws of Planetary Motion

AUTOMATED



Some data



A machine learning algorithm



The model

Modeling is a well-proven way of obtaining **knowledge**

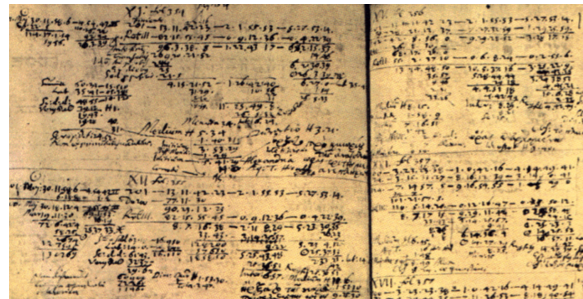
DATA

ALGORITHM

MODEL

KNOWLEDGE

MANUAL



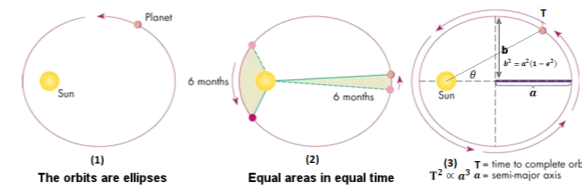
Tycho Brahe's observations



Johannes Kepler



27 years



Three Laws of Planetary Motion

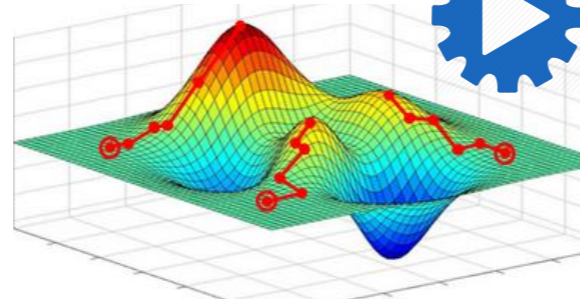
=

MODEL

AUTOMATED



Some data



A machine learning algorithm



The model

IT'S THERE
BUT
HIDDEN

Modeling is a well-proven way of obtaining knowledge

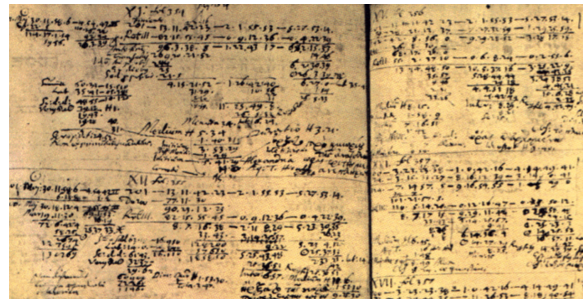
DATA

ALGORITHM

MODEL

KNOWLEDGE

MANUAL



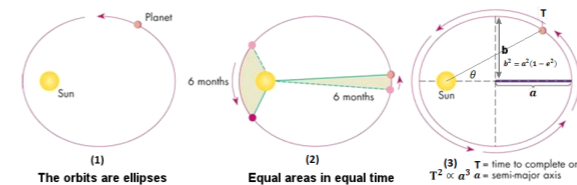
Tycho Brahe's observations



Johannes Kepler



27 years



Three Laws of Planetary Motion

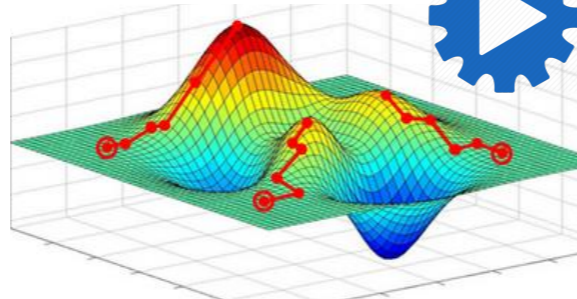
=

MODEL

AUTOMATED



Some data



A machine learning algorithm



The model

IT'S THERE
BUT
HIDDEN

Machine learning can
uncover knowledge

Model interpretation is
required to articulate it

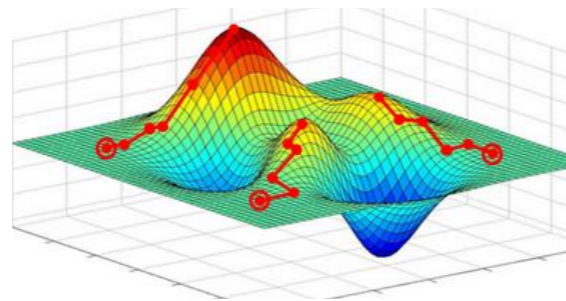
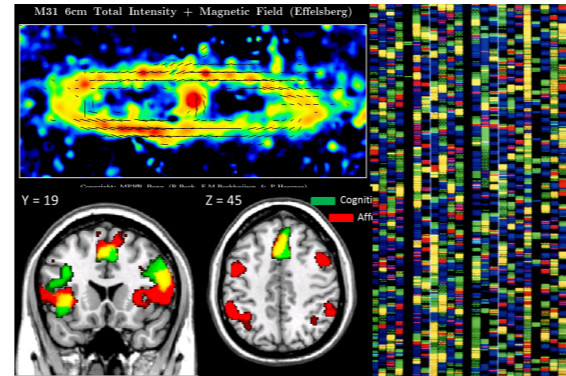
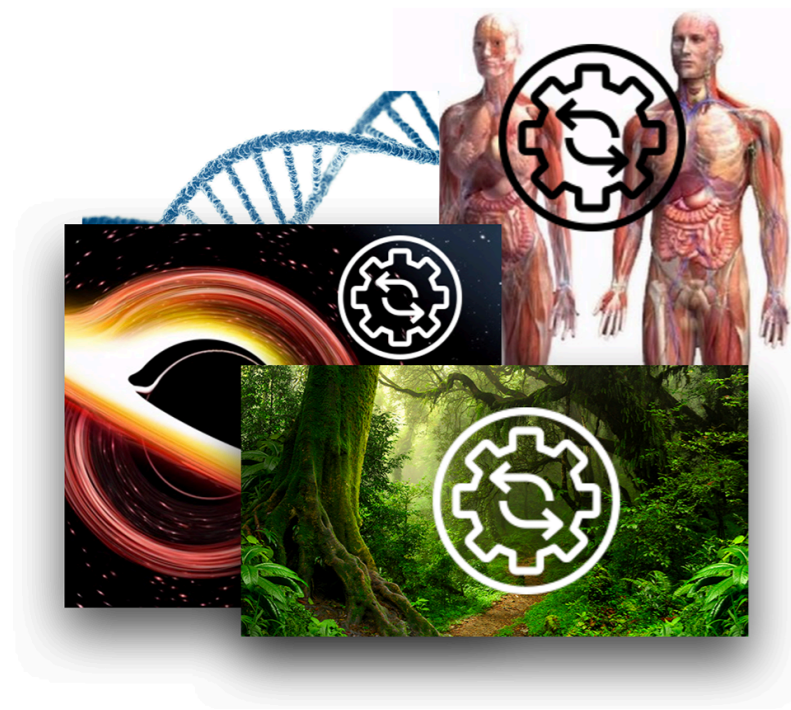
In life sciences model interpretation has a **special** significance

In life sciences model interpretation has a **special** significance

Interpretability in ML

- Trust in model's decision
- Legal transparency
- Debugging

In life sciences model interpretation has a **special** significance

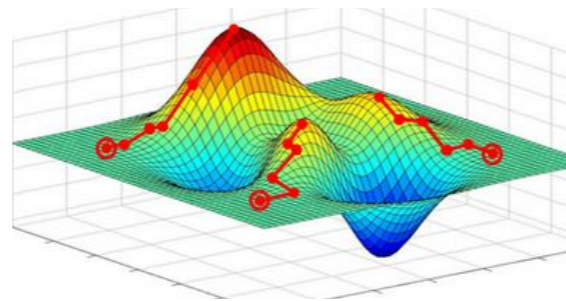
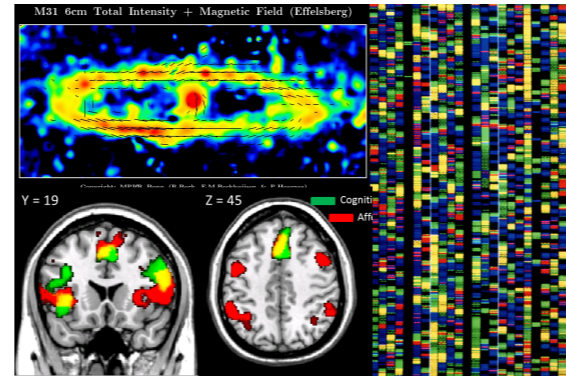
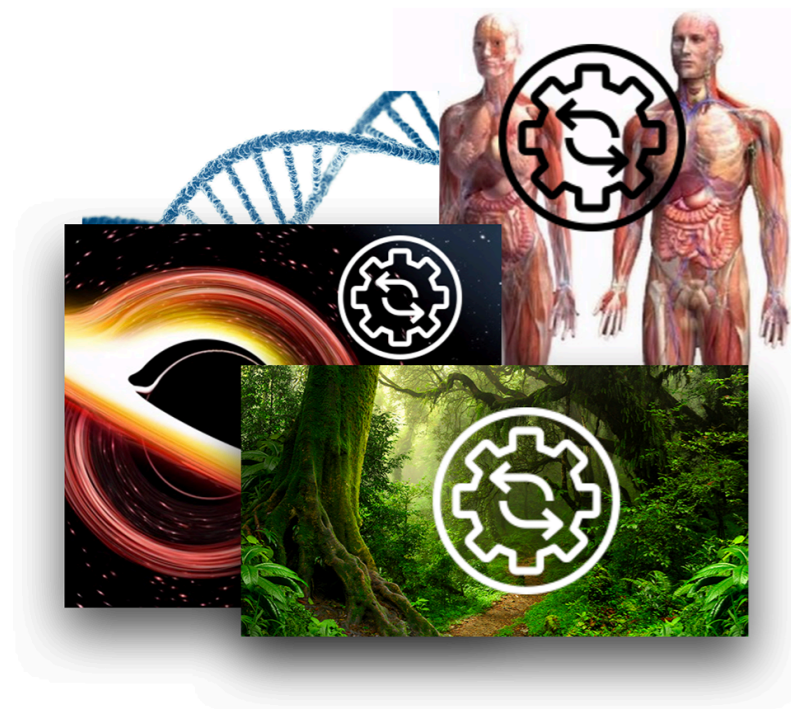


A model that can make correct predictions about reality

Interpretability in ML

- Trust in model's decision
- Legal transparency
- Debugging

In life sciences model interpretation has a **special** significance



A model that can make correct predictions about reality

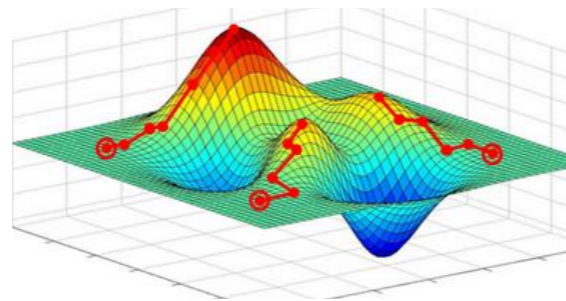
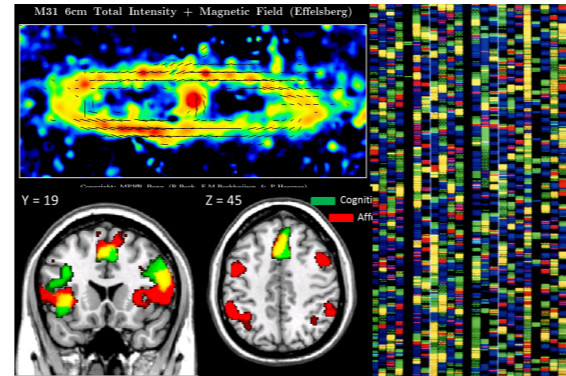
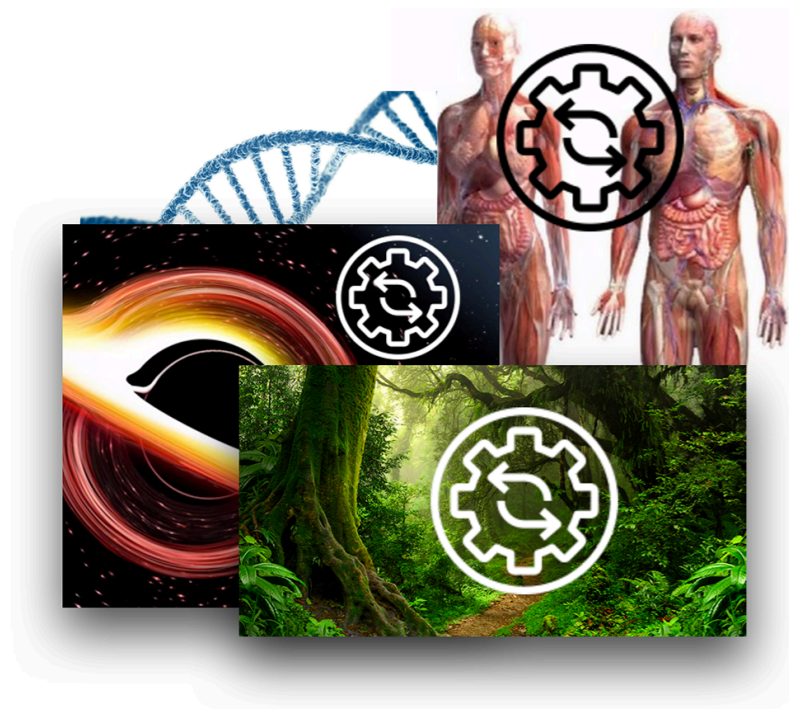


At this point the model knows something about the world that the scientist does not

Interpretability in ML

- Trust in model's decision
- Legal transparency
- Debugging

In life sciences model interpretation has a **special** significance



A model that can make correct predictions about reality



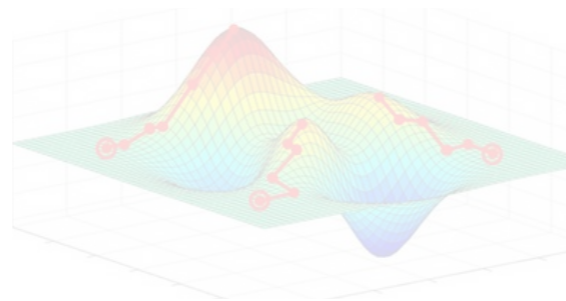
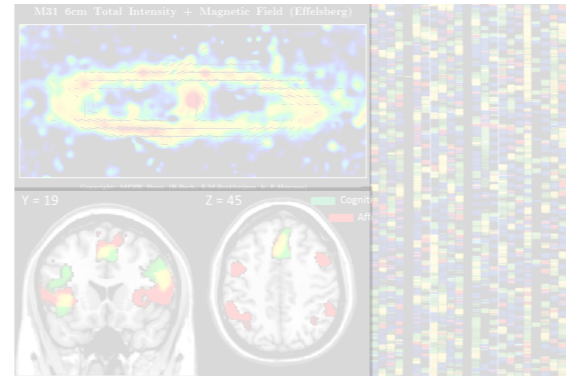
At this point the model knows something about the world that the scientist does not

Interpretability in ML

- Trust in model's decision
- Legal transparency
- Debugging
- Articulating the knowledge about the underlying process that the model has captured



In life sciences model interpretation has a **special** significance



A model that can make correct predictions about reality



At this point the model knows something about the world that the scientist does not

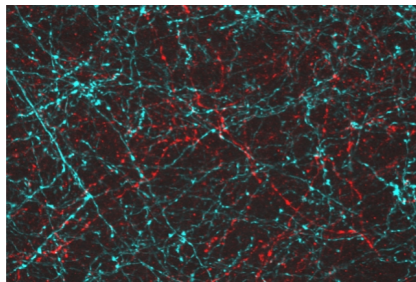
Interpretability in ML

- Trust in model's decision
- Legal transparency
- Debugging
- Articulating the knowledge about the underlying process that the model has captured

Machine-learned models as scientific theories

- Both capture and model observations
- Both make correct predictions on new observations
- Computers can generate and test theories faster than humans
- In the big data regime there is not enough scientists to sift through all potential explanations of the data
- Algorithms have different bias than humans, hence find different solutions

Applying this principle in Neuroscience



"Identifying task-relevant spectral signatures of perceptual categorization in the human cortex"

Scientific Reports, 2020



"Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex"

Communications Biology, 2018



"Mental state space visualization for interactive modeling of personalized BCI control strategies"

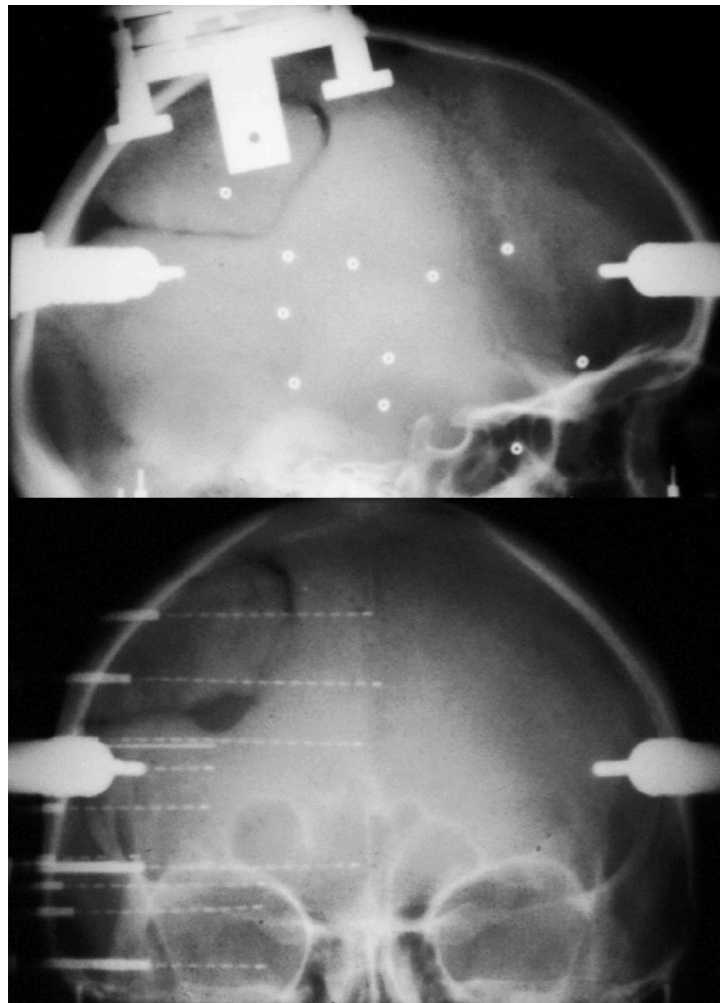
Journal of Neural Engineering, 2020

Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex



100 patients
12,000 electrodes

Centre de Recherche
Neurosciences Lyon

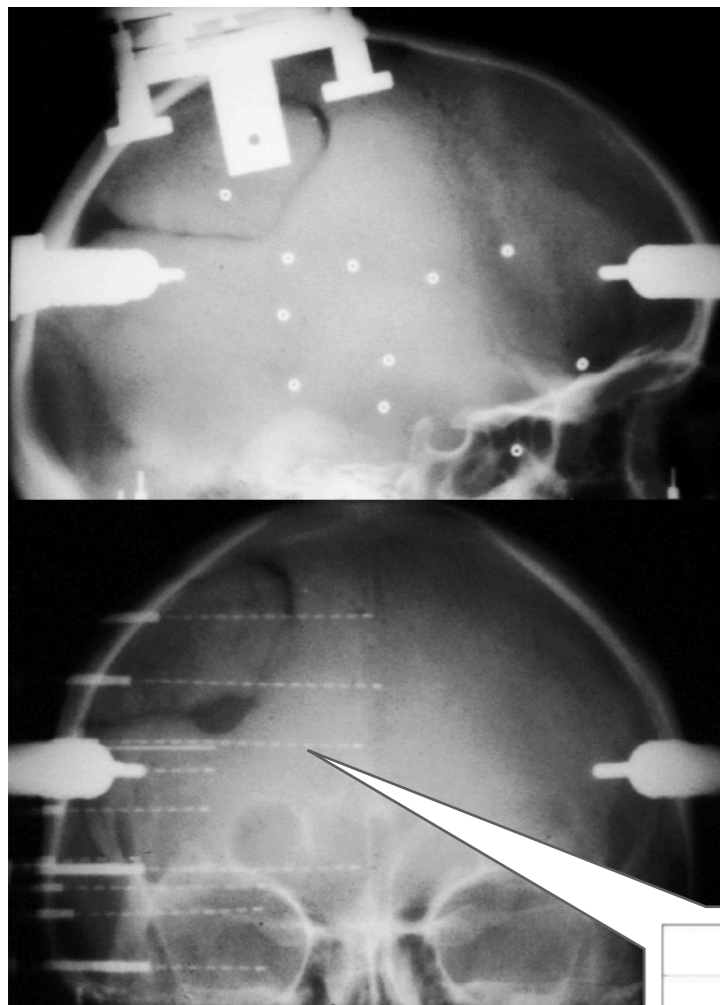


Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex

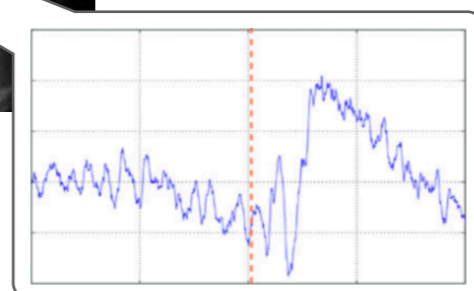


100 patients
12,000 electrodes

Centre de Recherche
Neurosciences Lyon



400 images from
8 categories



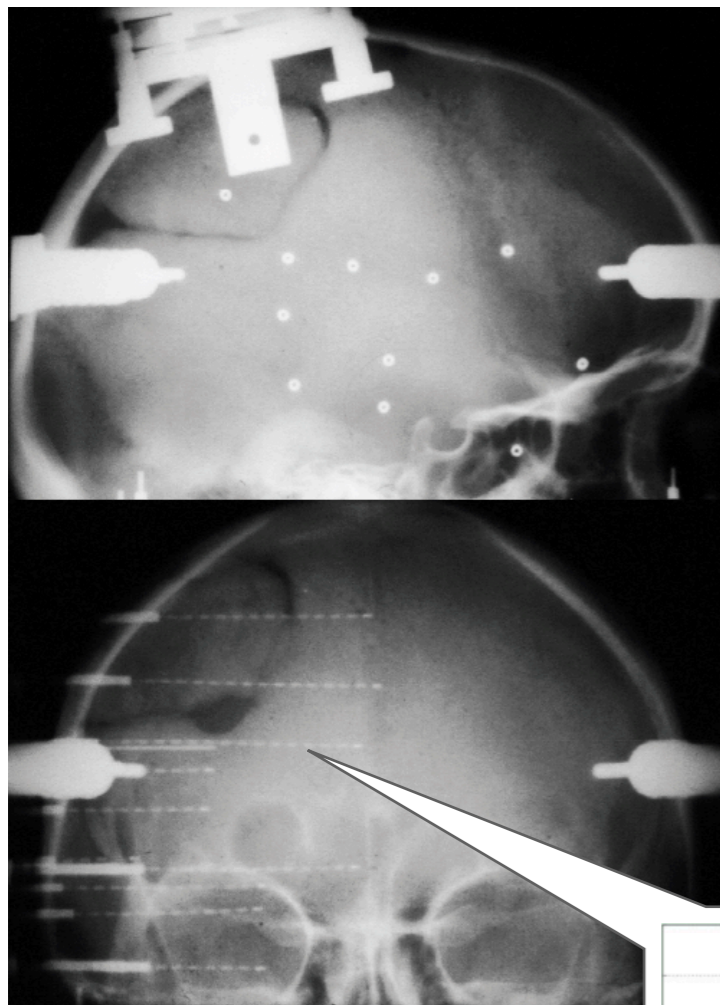
4,528,400 responses

Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex

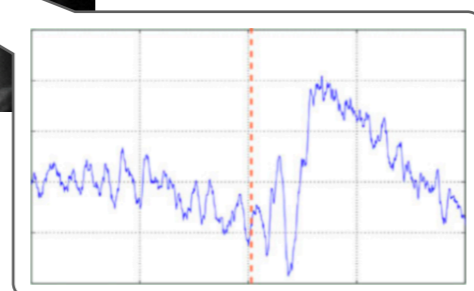


Centre de Recherche
Neurosciences Lyon

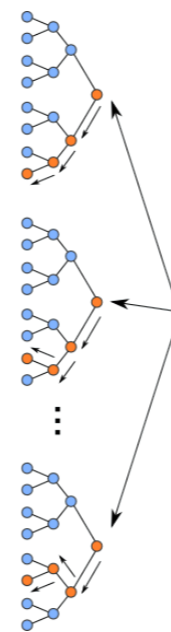
100 patients
12,000 electrodes



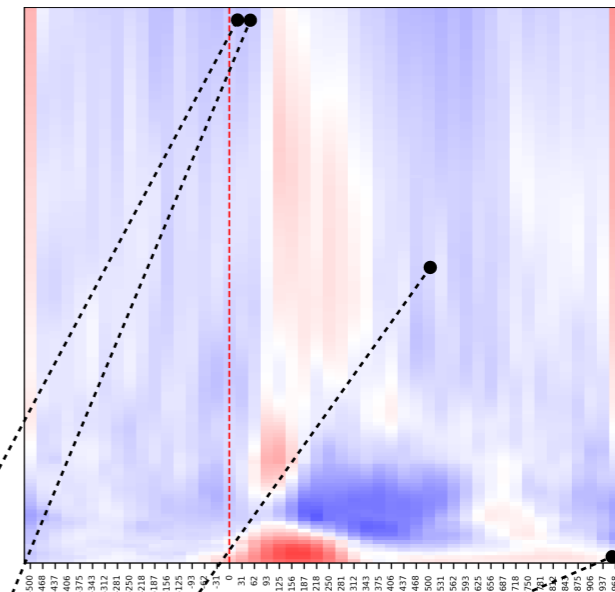
400 images from
8 categories



4,528,400 responses



$$\mathbf{x} = \{f_1, f_2, \dots, f_{3504}, \dots, f_{7008}\}$$

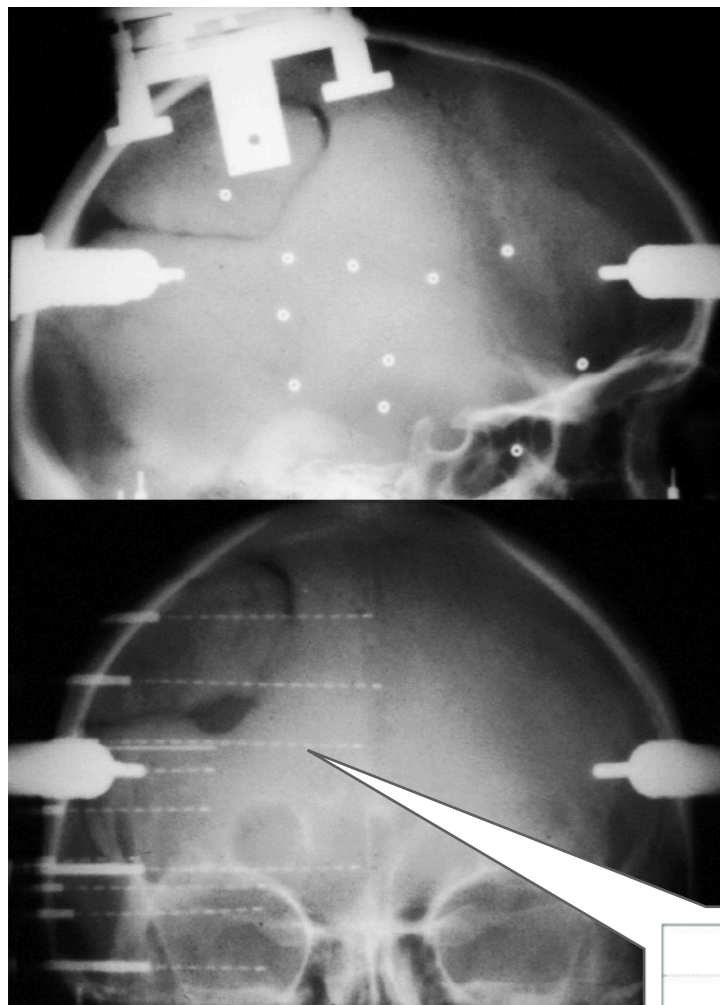


Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex

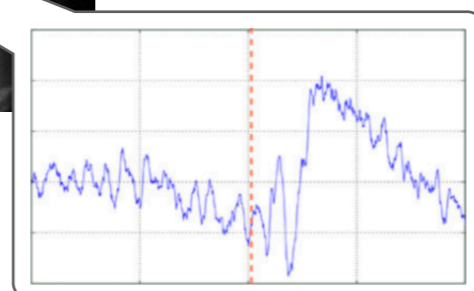


100 patients
12,000 electrodes

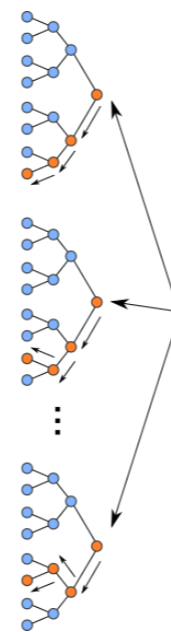
Centre de Recherche
Neurosciences Lyon



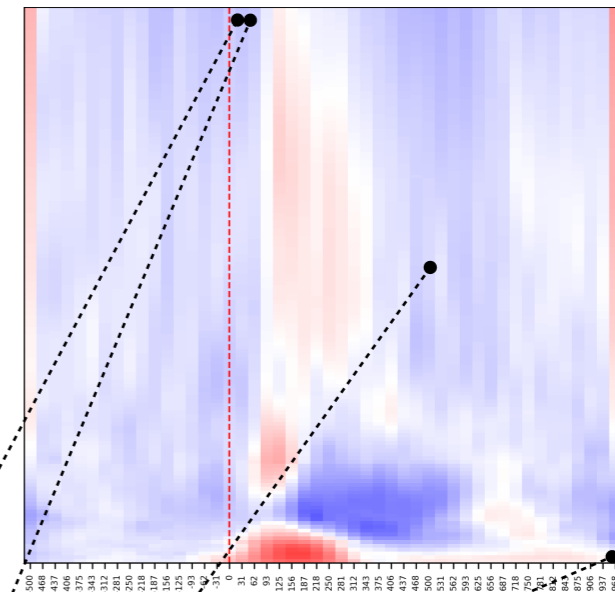
400 images from
8 categories



4,528,400 responses



$$\mathbf{x} = \{f_1, f_2, \dots, f_{3504}, \dots, f_{7008}\}$$



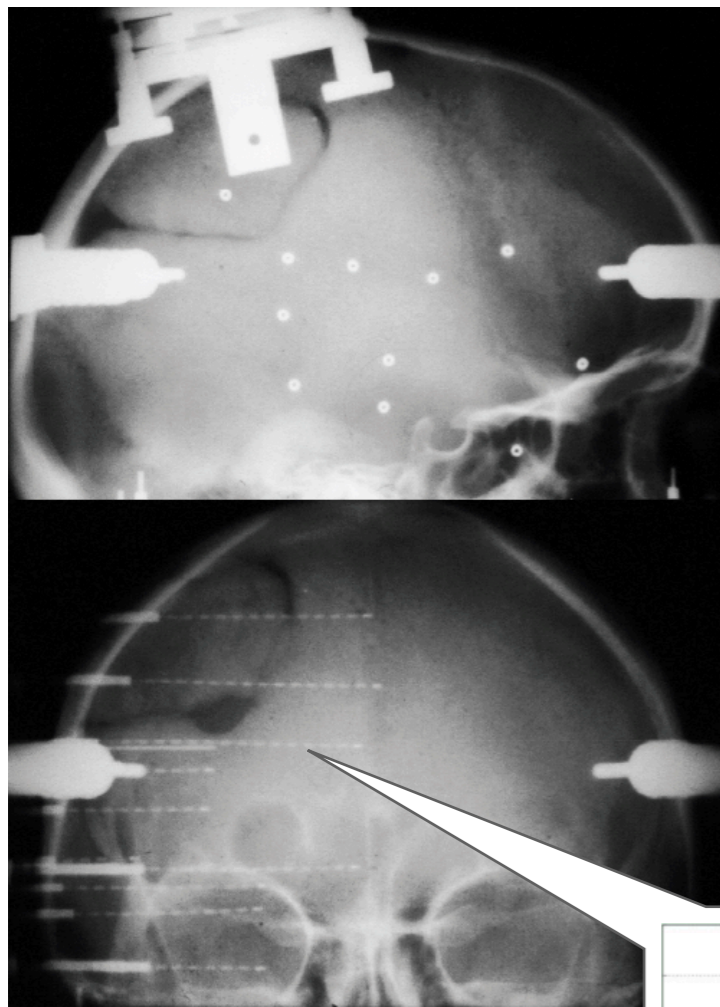
Which parts of this brain activity are generated by the underlying process of mental image categorization?

Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex

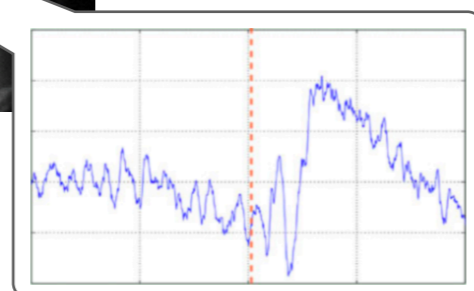


Centre de Recherche
Neurosciences Lyon

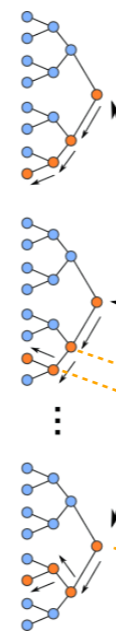
100 patients
12,000 electrodes



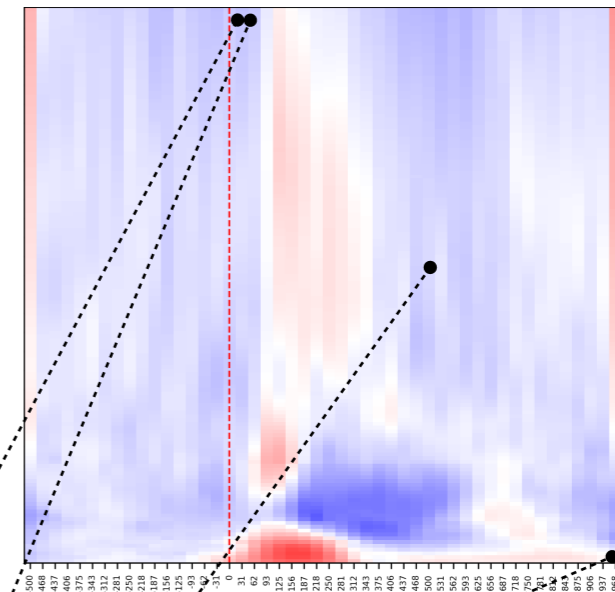
400 images from
8 categories



4,528,400 responses

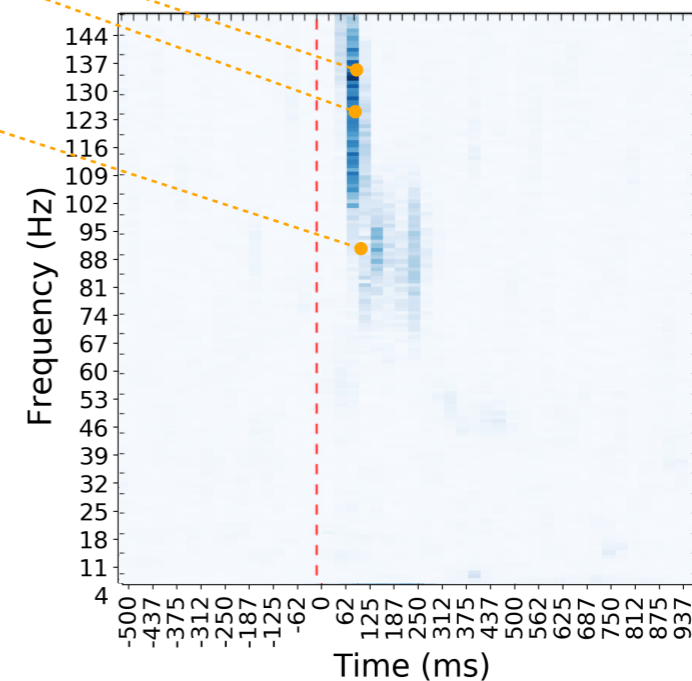


$$\mathbf{x} = \{f_1, f_2, \dots, f_{3504}, \dots, f_{7008}\}$$

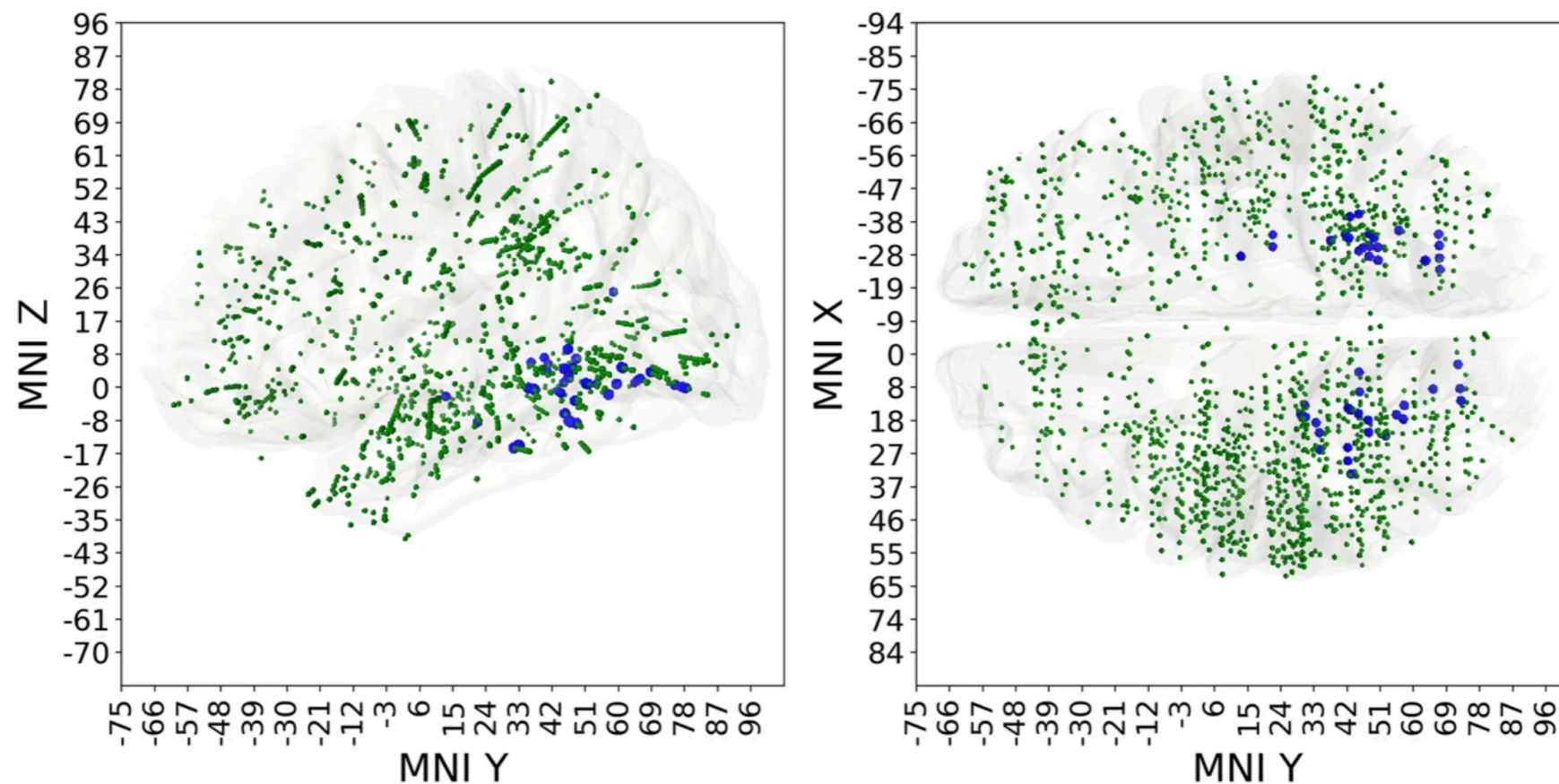


Which parts of this brain activity are generated by the underlying process of mental image categorization?

Feature importance map

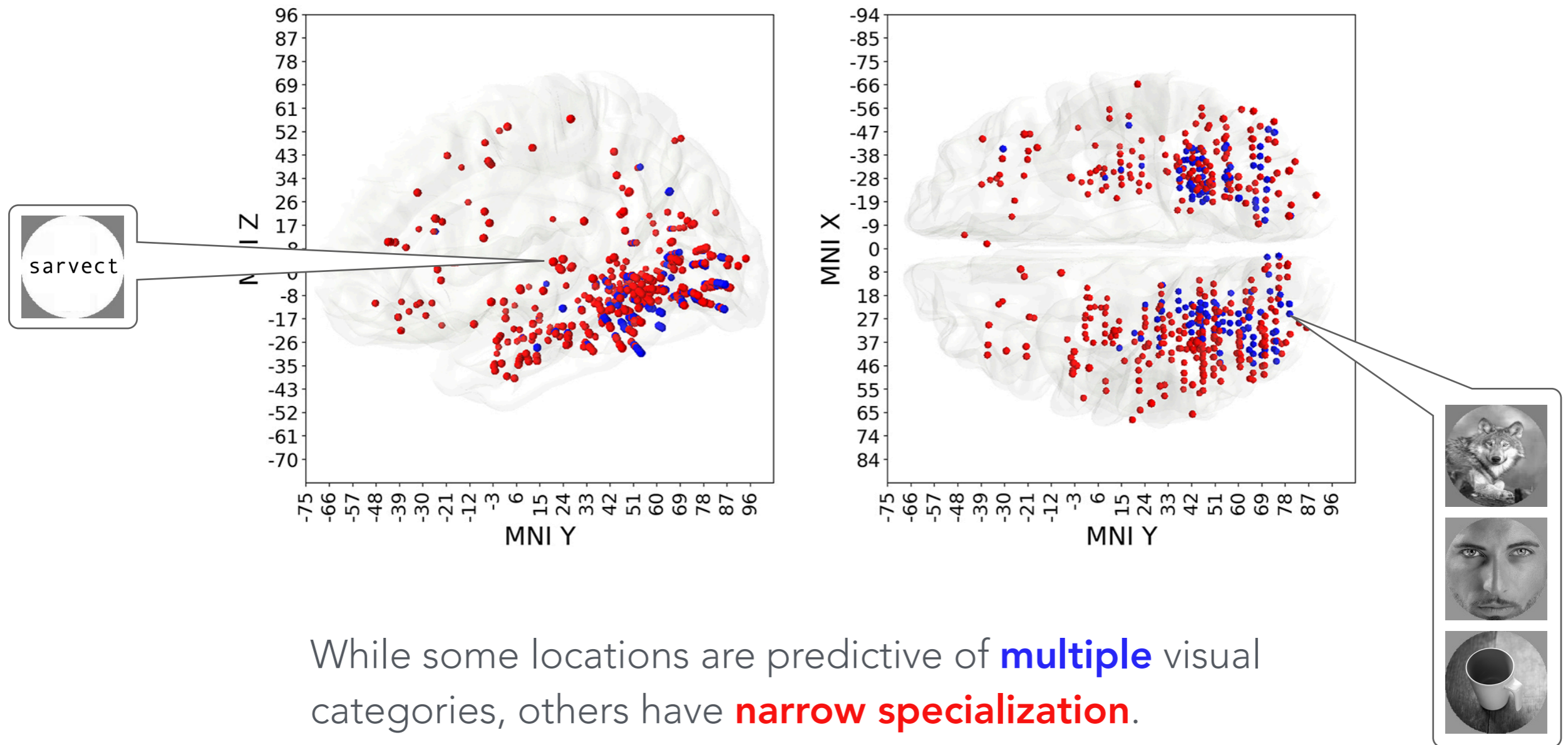


Identifying **task-relevant spectral signatures** of perceptual categorization in the **human visual cortex**

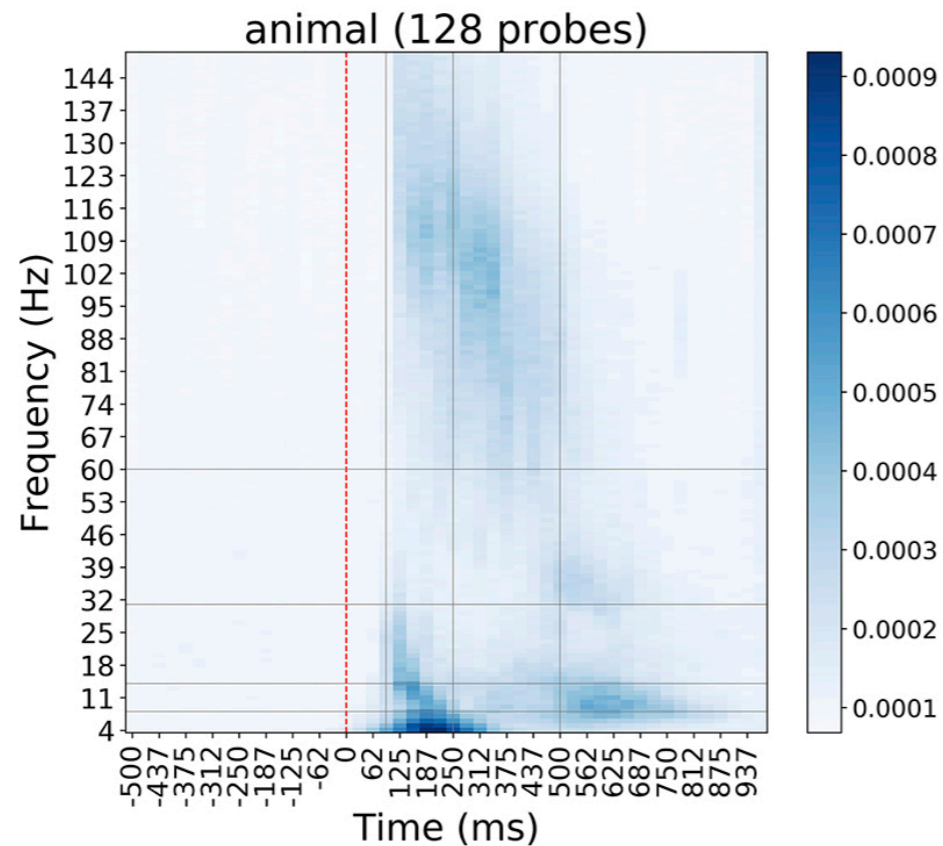


Multiple brain regions **respond** when a stimulus is shown, but only a small fraction (5%) are **predictive** of a category.

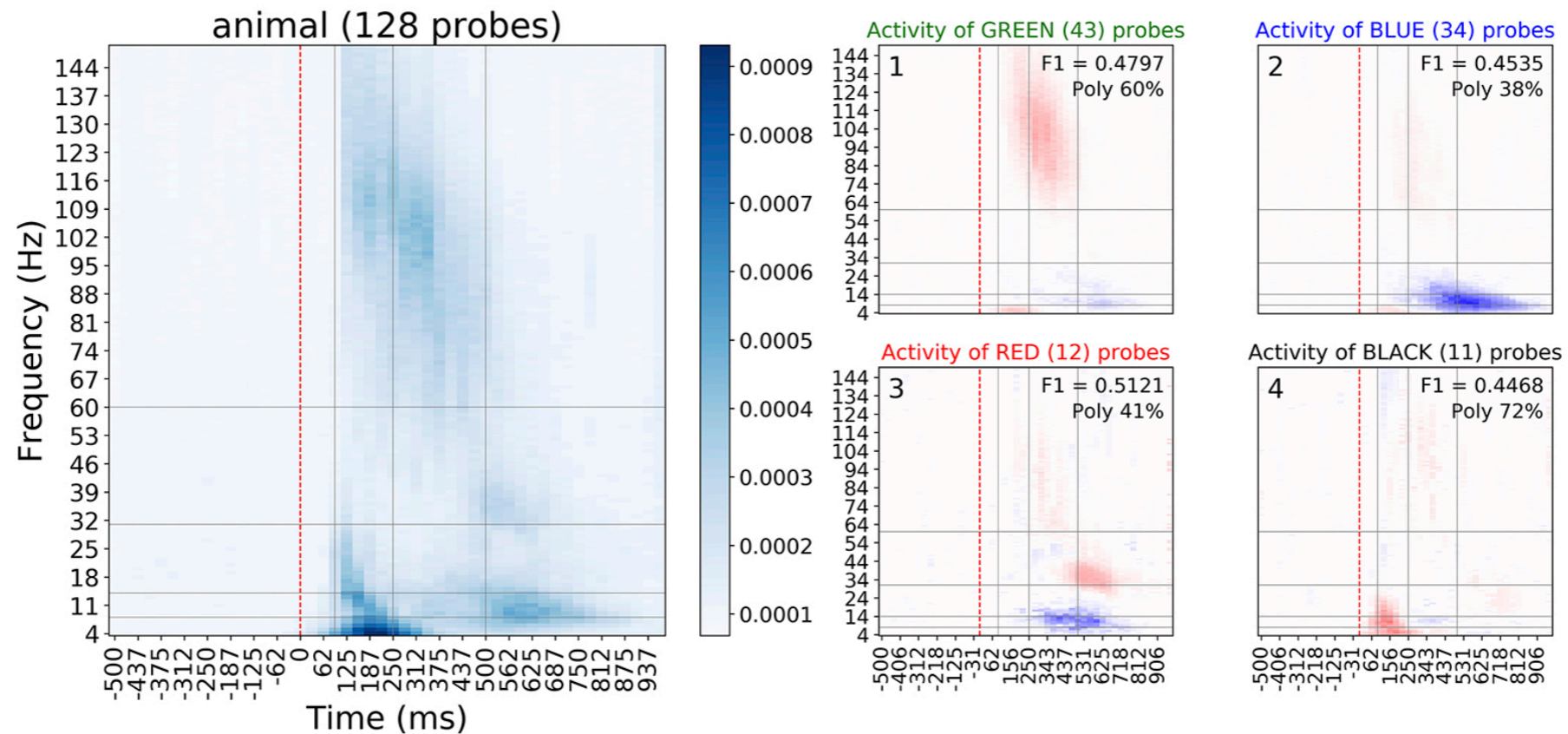
Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex



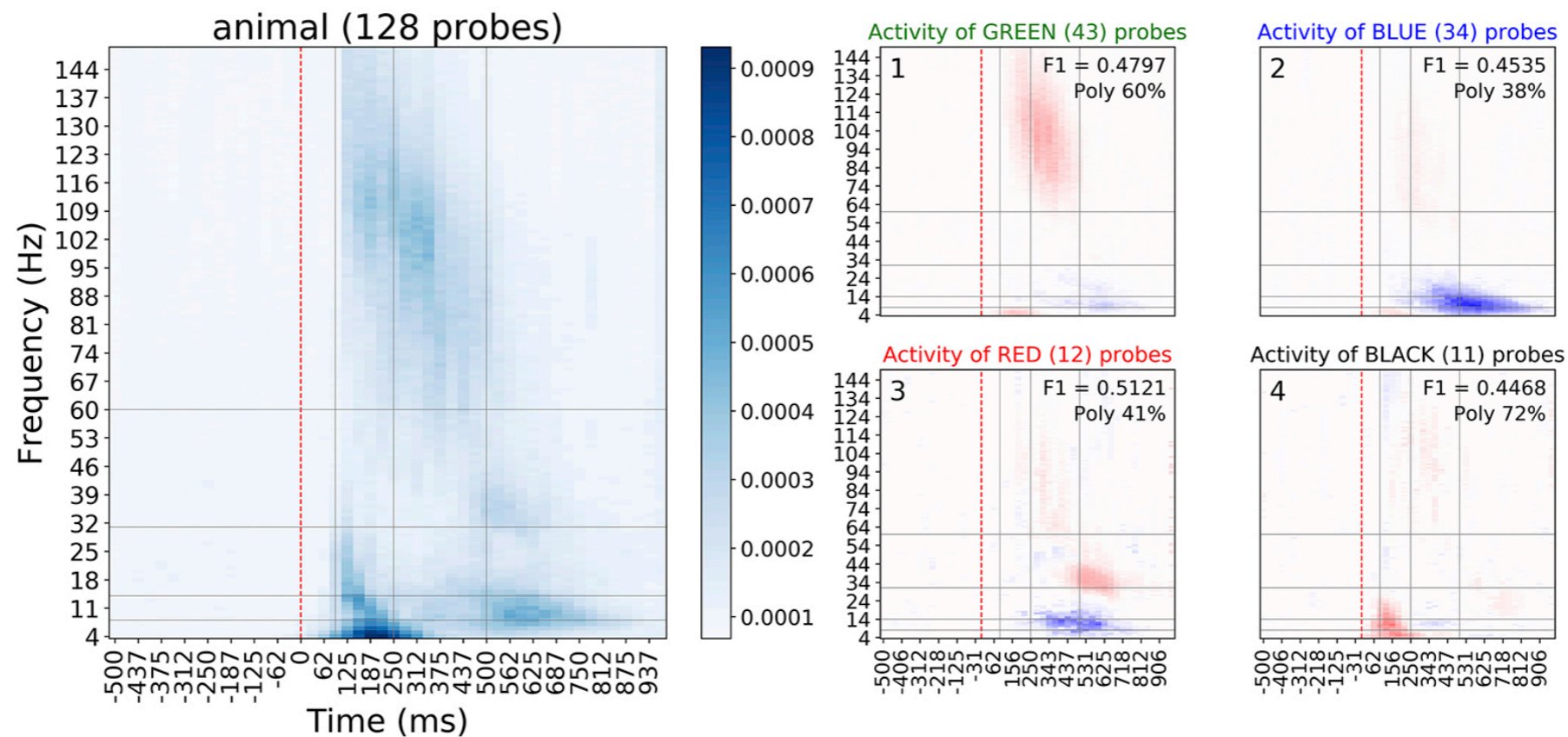
Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex



Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex



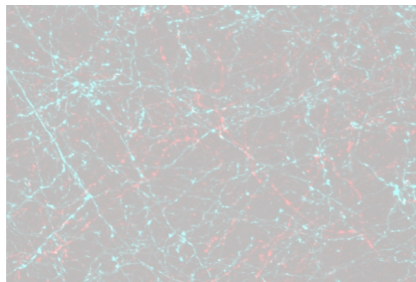
Identifying task-relevant spectral signatures of perceptual categorization in the human visual cortex



It is believed that high-frequency activity reflects the information processing during high cognitive tasks. We show that low-frequency activity is almost as important for the task at hand.

Across all categories the classifier relied on power decreases in different brain networks, not only on the increases to perform the classification

Applying this principle in Neuroscience



"Identifying task-relevant spectral signatures of perceptual categorization in the human cortex"

Scientific Reports, 2020



"Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex"

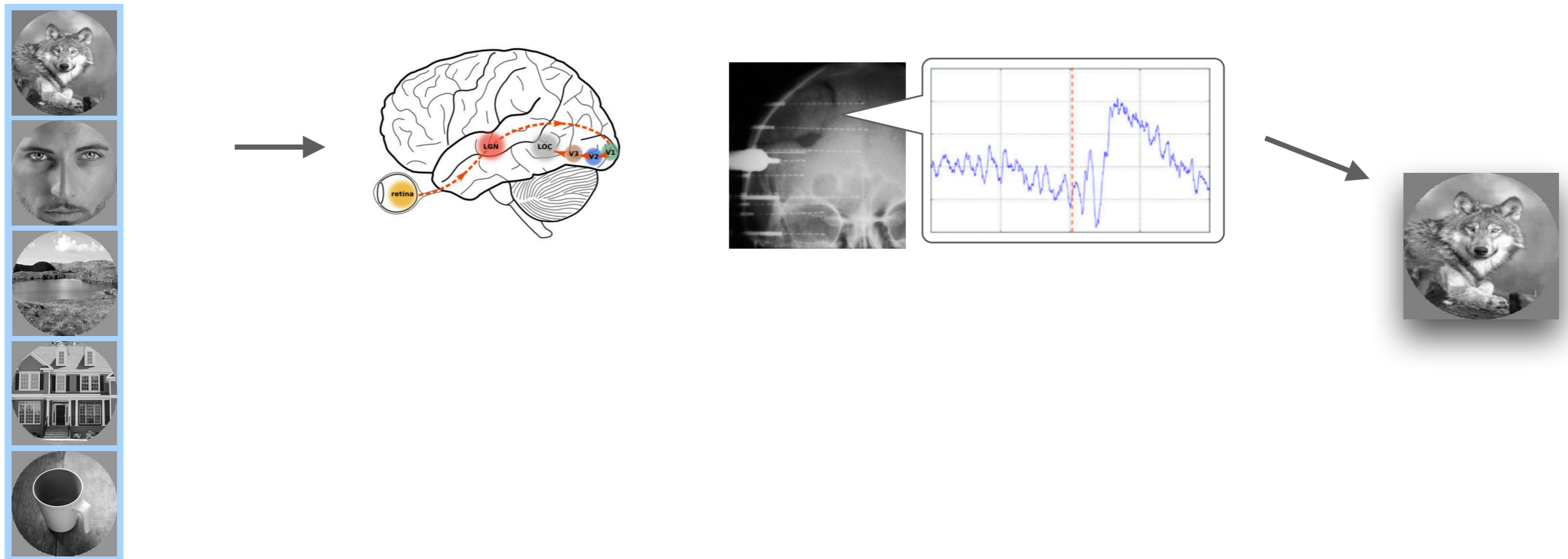
Communications Biology, 2018



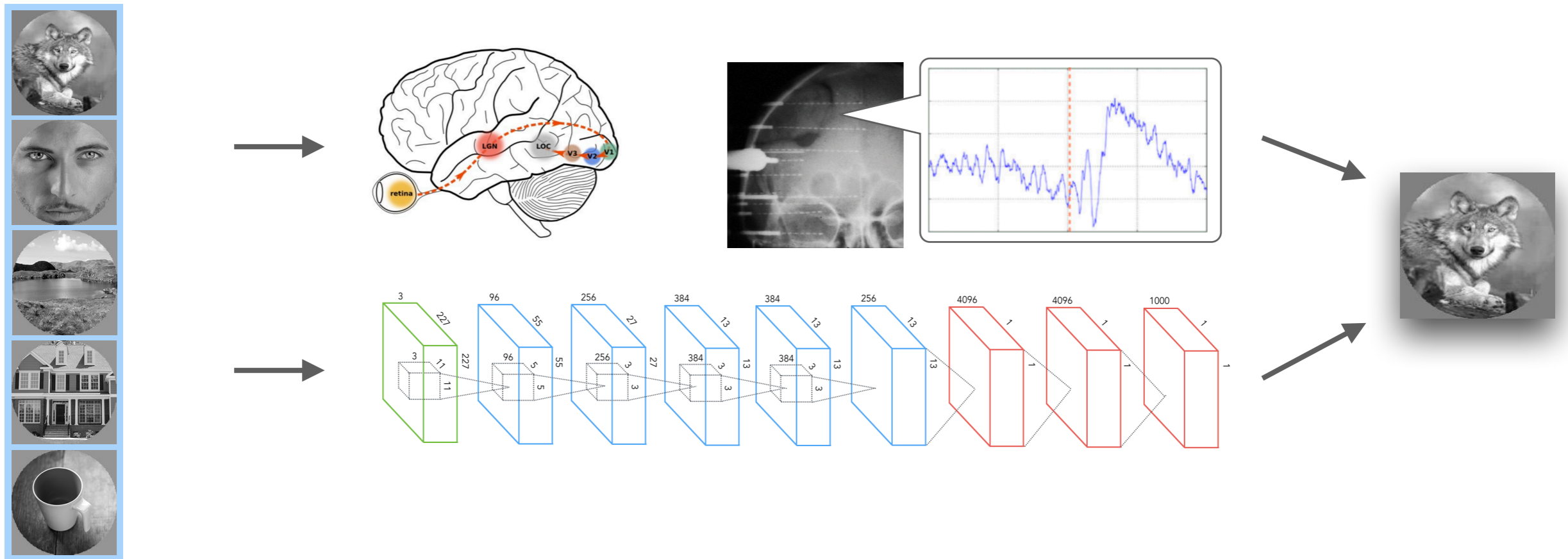
"Mental state space visualization for interactive modeling of personalized BCI control strategies"

Journal of Neural Engineering, 2020

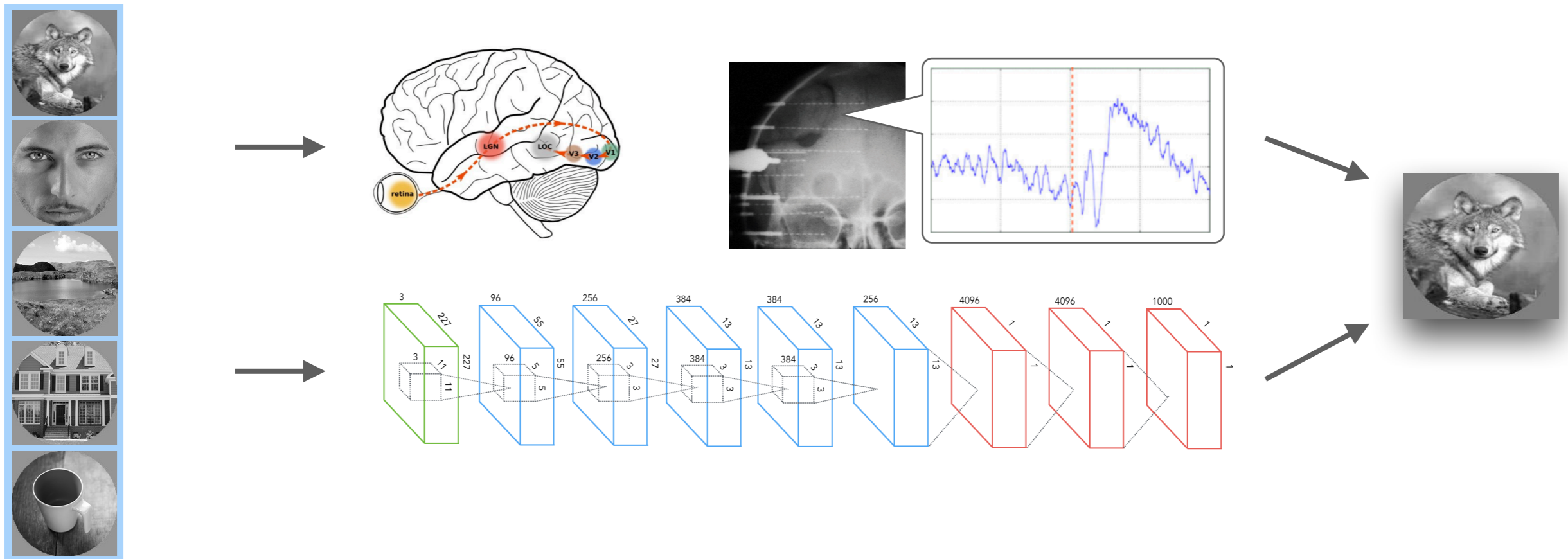
Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex

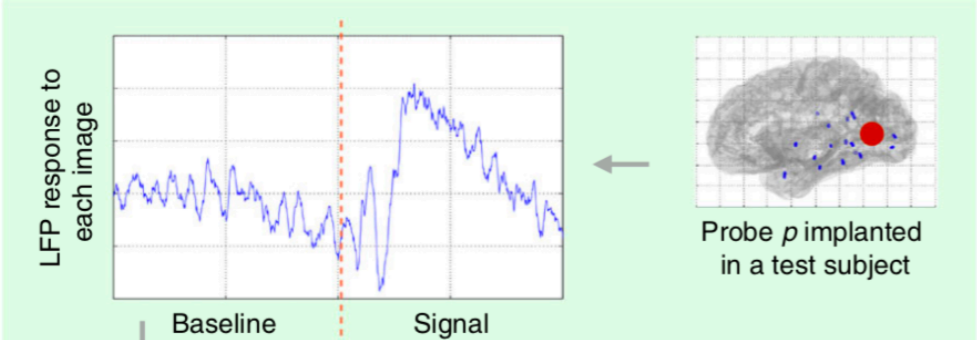


From previous fMRI research we knew that there is a mapping between the hierarchies

Intracranial electrophysiological data allowed us to learn **when** and **at which frequencies**

Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex

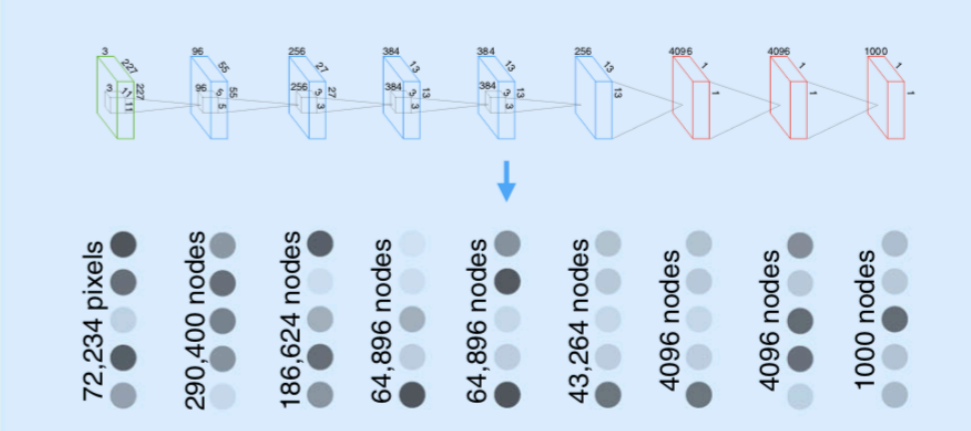
Step 1



250 images

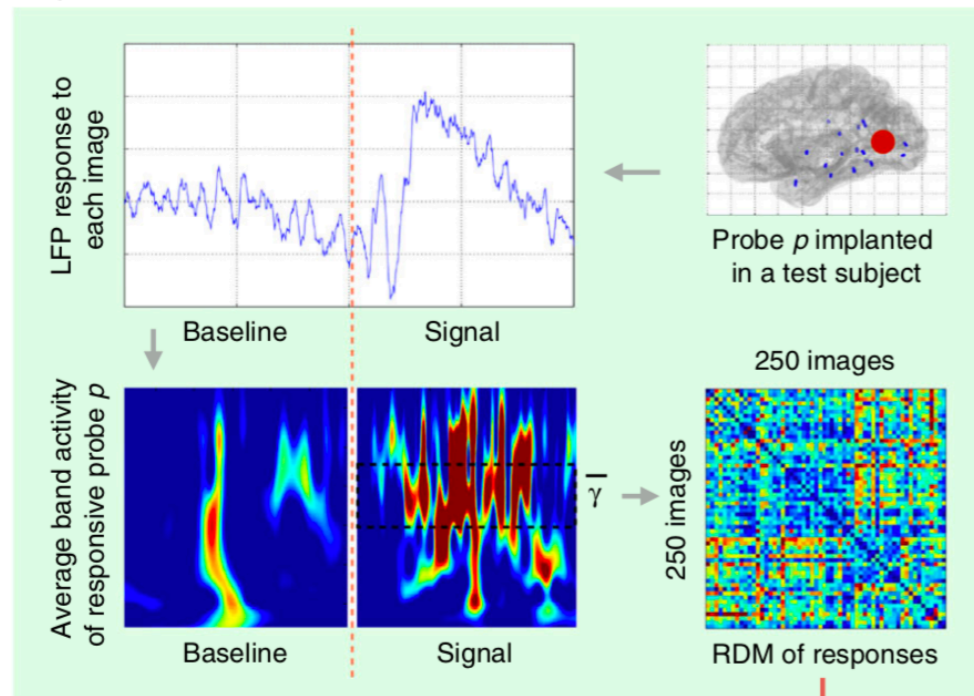


Step 2



Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex

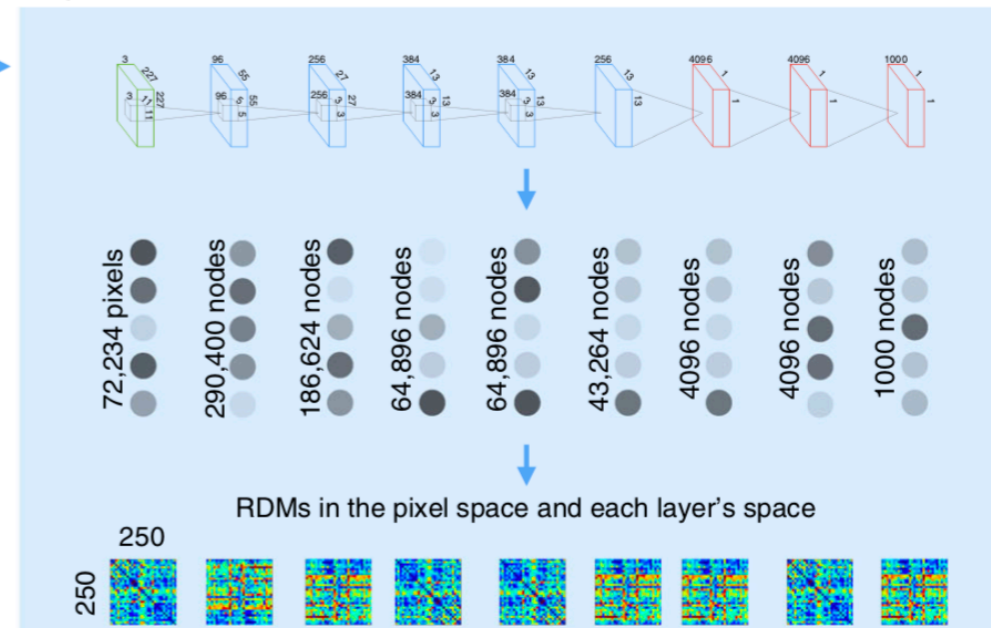
Step 1



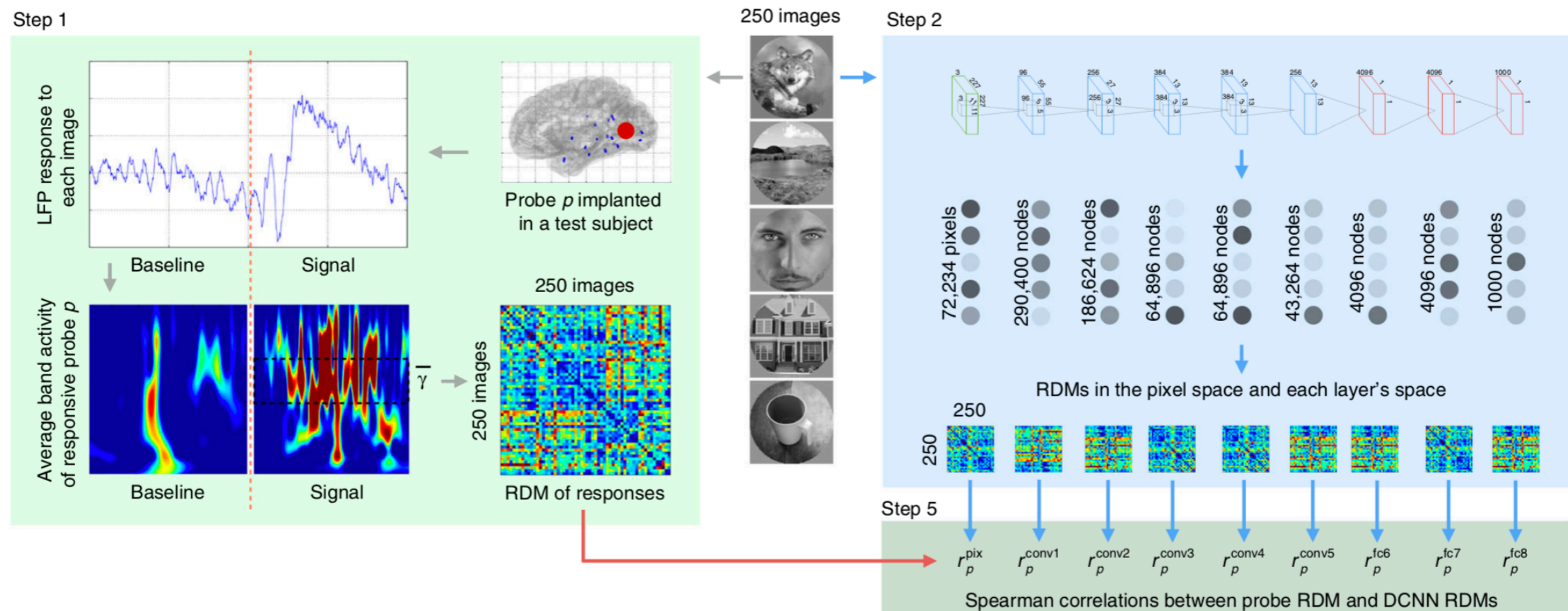
250 images



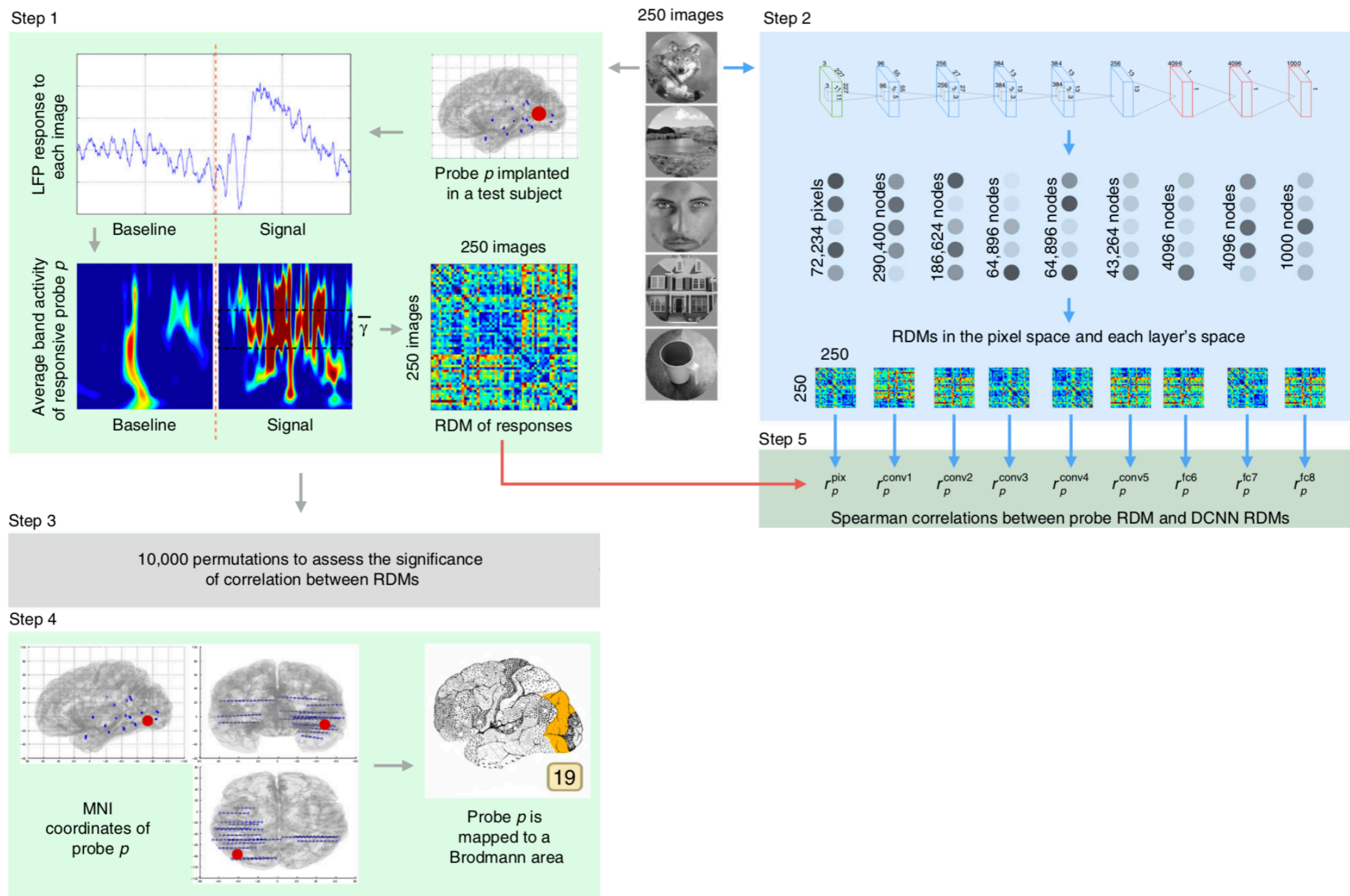
Step 2



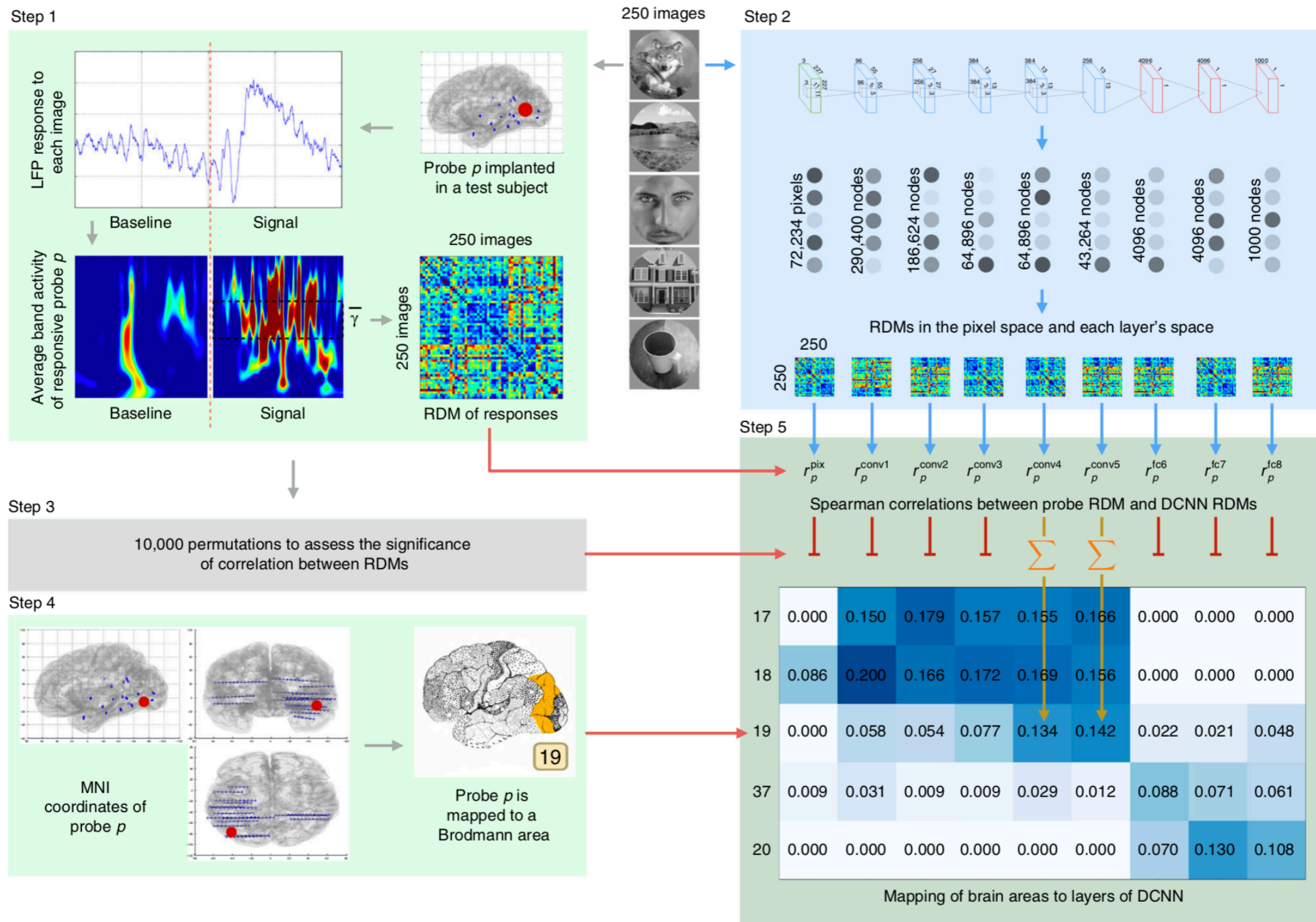
Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



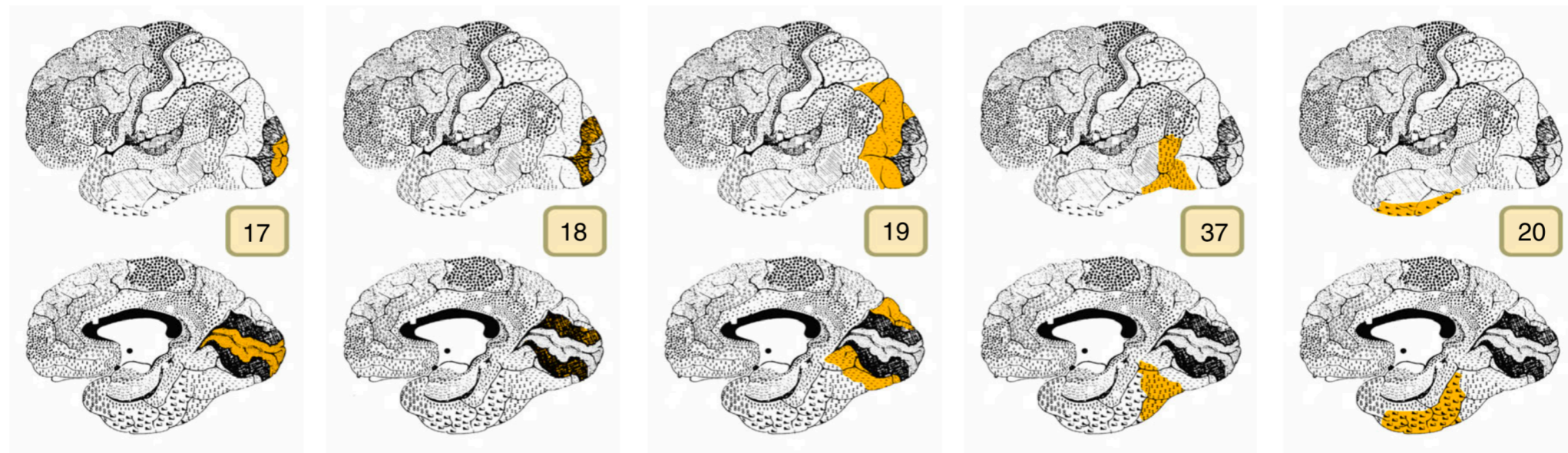
Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex

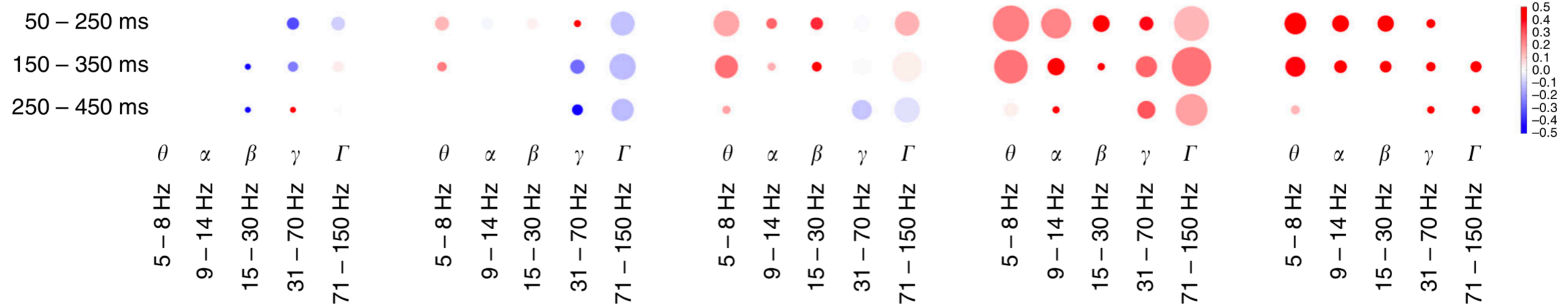
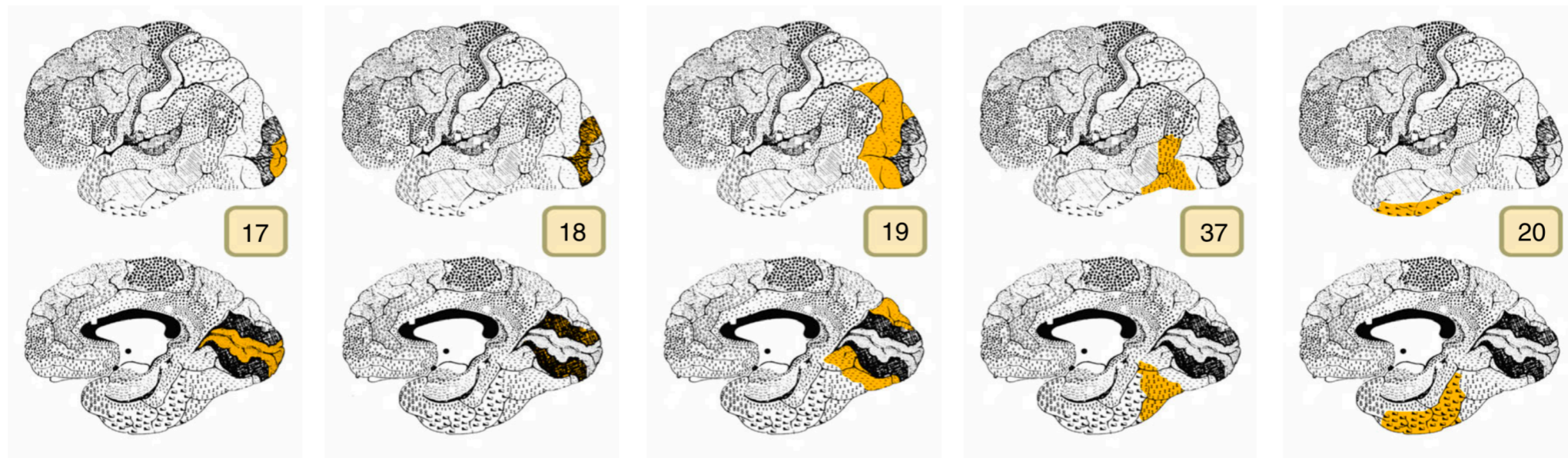


Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



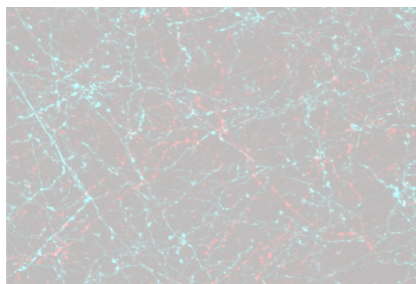
- - simple features (mapped to lower layers of DCNN)
- - complex features (higher layers)

Activations of deep convolutional neural networks are aligned with **gamma band activity** of human visual cortex



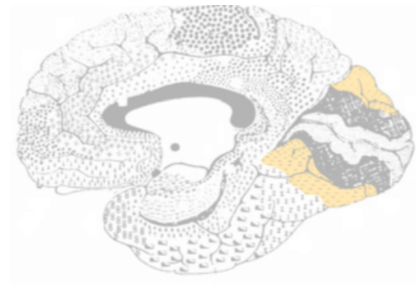
- - simple features (mapped to lower layers of DCNN)
- - complex features (higher layers)

Applying this principle in Neuroscience



"Identifying task-relevant spectral signatures of perceptual categorization in the human cortex"

Scientific Reports, 2020



"Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex"

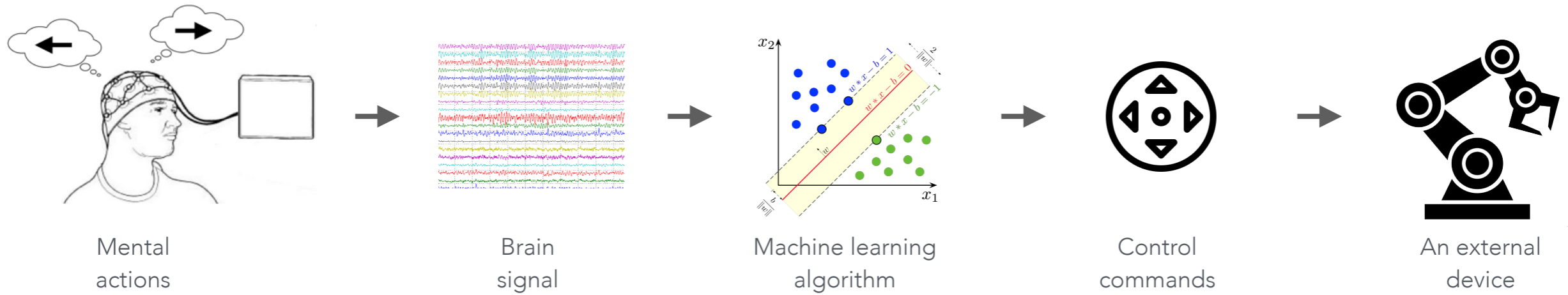
Communications Biology, 2018



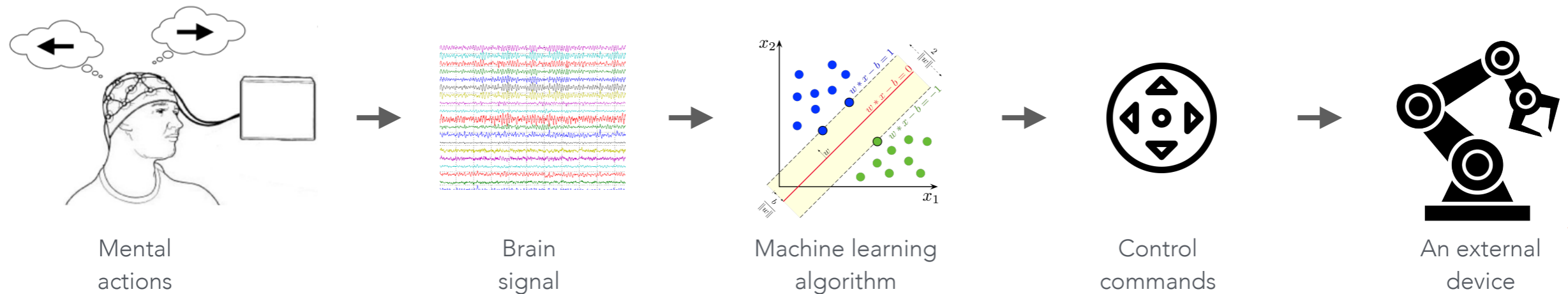
"Mental state space visualization for interactive modeling of personalized BCI control strategies"

Journal of Neural Engineering, 2020

Mental state space visualization for interactive modeling of personalized BCI control strategies



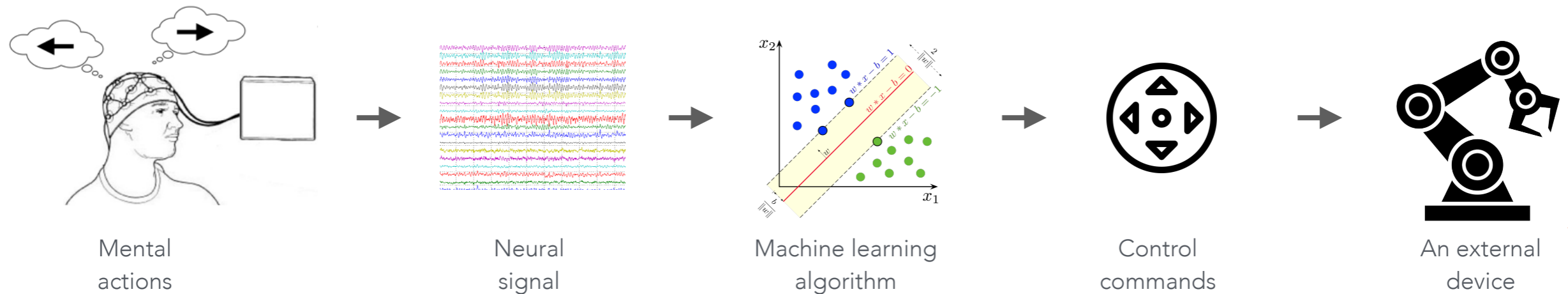
Mental state space visualization for interactive modeling of personalized BCI control strategies



The user must produce *mental actions* that are distinguishable by the machine and do that consistently

The algorithm must distinguish between different *mental actions*

Mental state space visualization for interactive modeling of personalized BCI control strategies

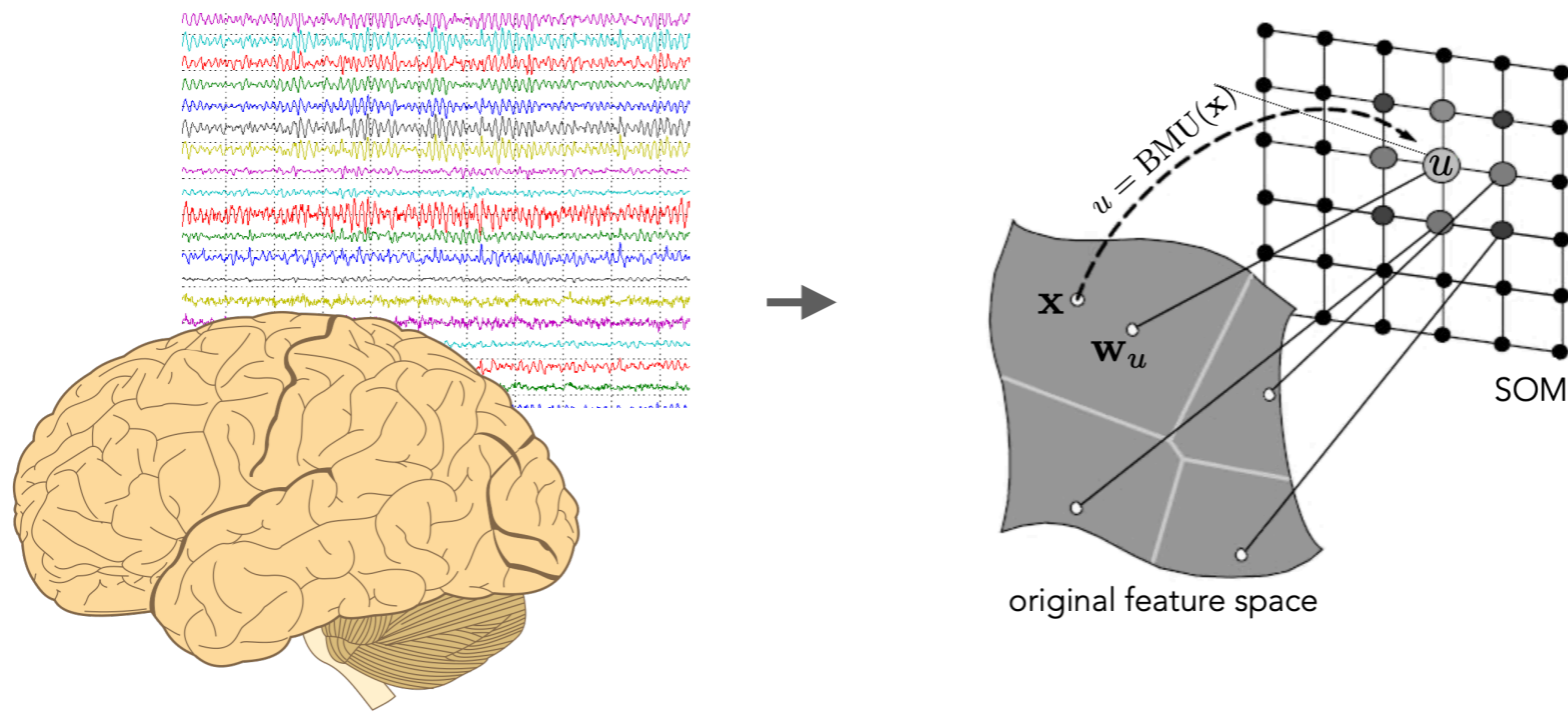


The user must produce *mental actions* that are distinguishable by the machine and do that consistently

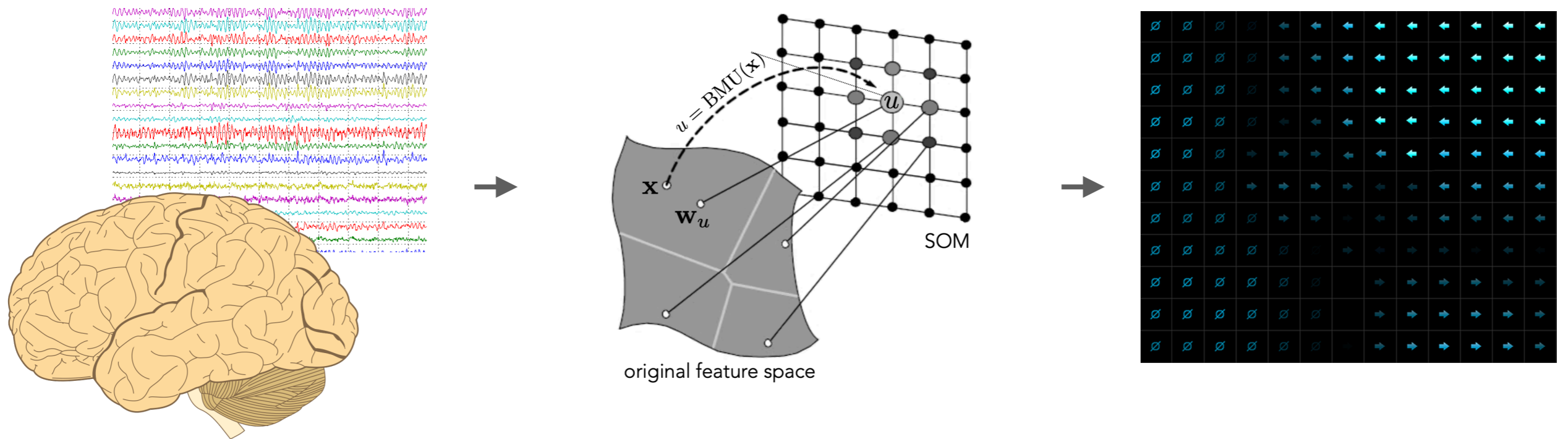
The algorithm must distinguish between different *mental actions*

If the user could see machine's representation of his mental actions he could find out which ones are suitable

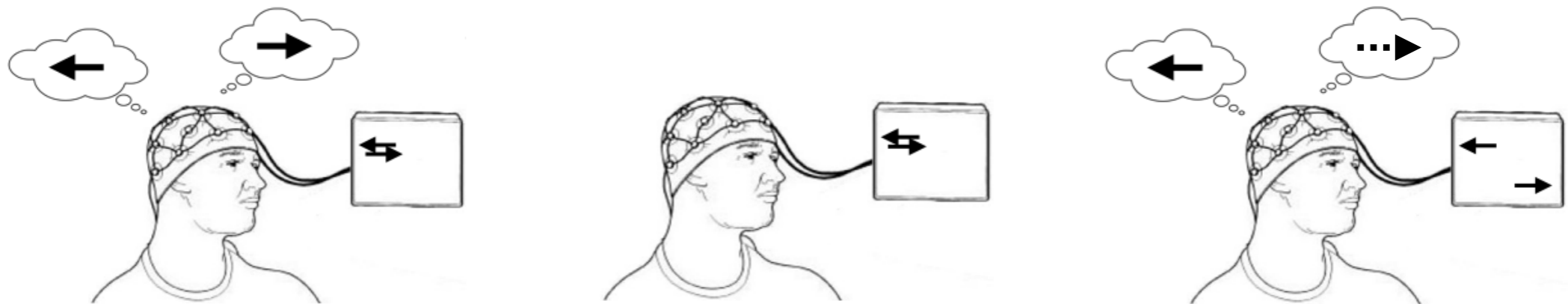
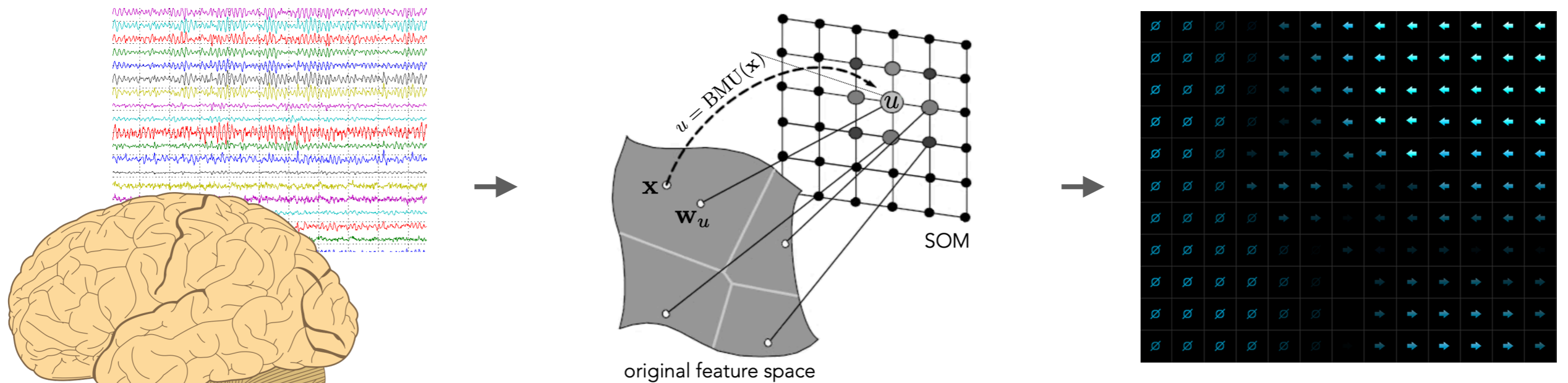
Mental state space visualization for interactive modeling of personalized BCI control strategies



Mental state space visualization for interactive modeling of personalized BCI control strategies



Mental state space visualization for interactive modeling of personalized BCI control strategies



I see that the signal for \leftarrow is very similar to the signal for \rightarrow

I should try to think \rightarrow differently

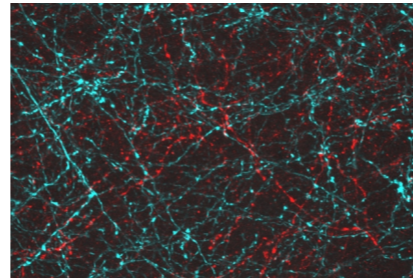
RESEARCH PROJECT

Kuzovkin et al.,
"Identifying task-relevant spectral signatures of perceptual categorization in the human cortex"
Scientific Reports, 2020 (in review)

Kuzovkin et al.,
"Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex"
Communications Biology, 2018

Kuzovkin et al.,
"Mental state space visualization for interactive modeling of personalized BCI control strategies"
Journal of Neural Engineering, 2020

LEVEL OF NEURAL ORGANIZATION



Local responses of a neural population

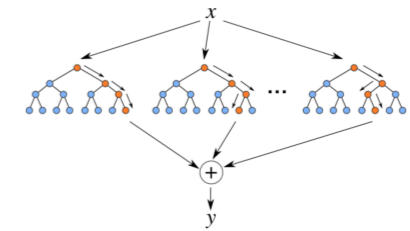


Hierarchy of visual areas

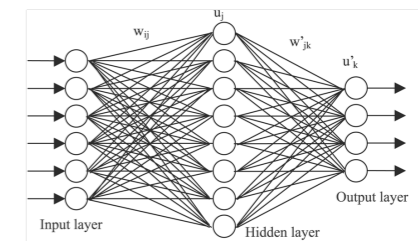


Correlates of mental states ("thoughts")

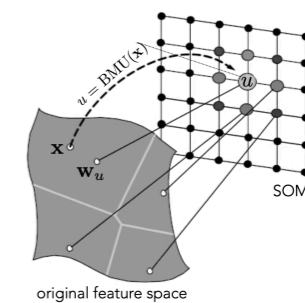
KNOWLEDGE REPRESENTATION IN THE MODEL



Random Forests: feature-based **rules**



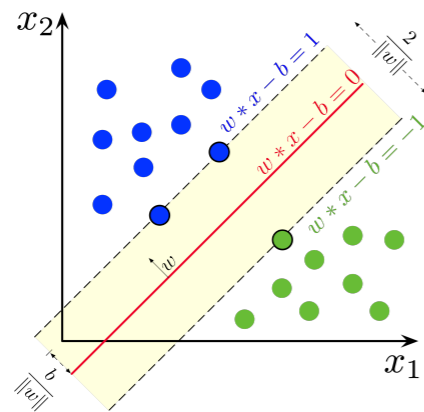
DNNs: **distributed representations** over features (input or latent)



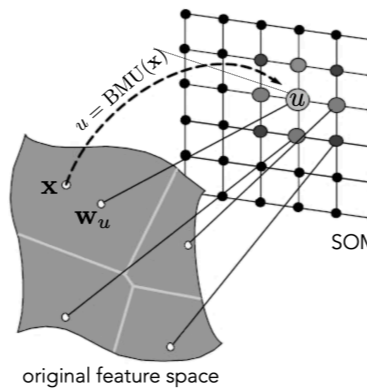
Self-Organizing Maps: **topology** of the samples

Interpretability adds a **new axis** for algorithm selection

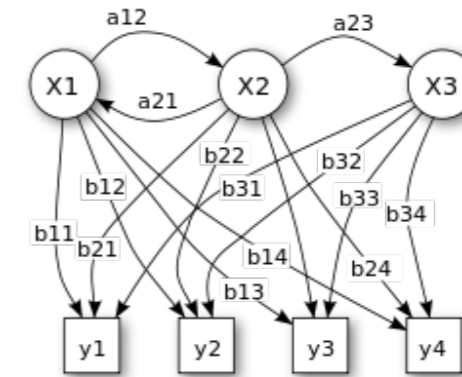
Different machine learning algorithms capture knowledge into **different representations**



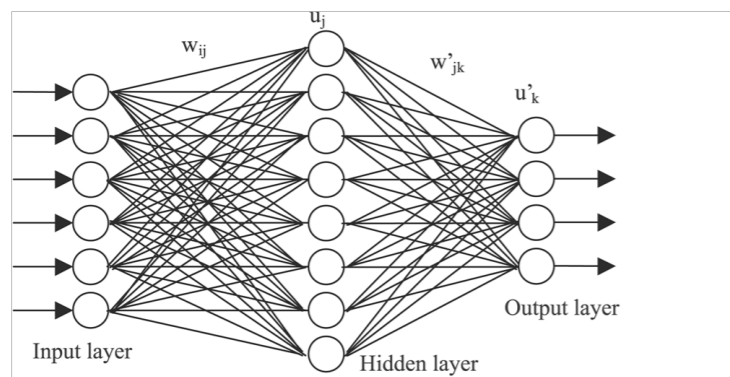
Support Vector Machines:
points in the feature space



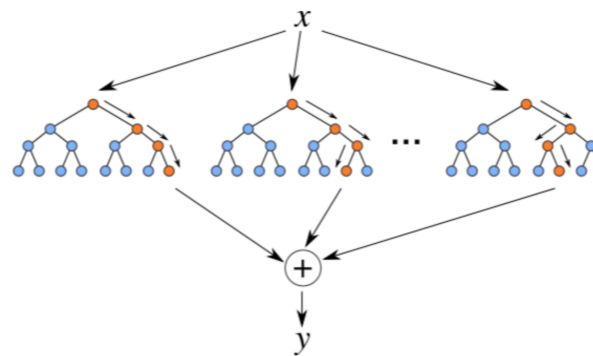
Self-Organizing Maps:
topology of the samples



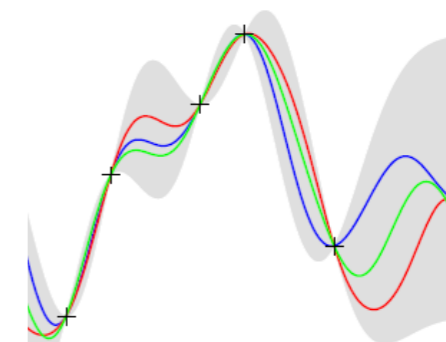
Hidden Markov Models:
states and transitions



DNNs: **distributed representations**
over features (input or latent)



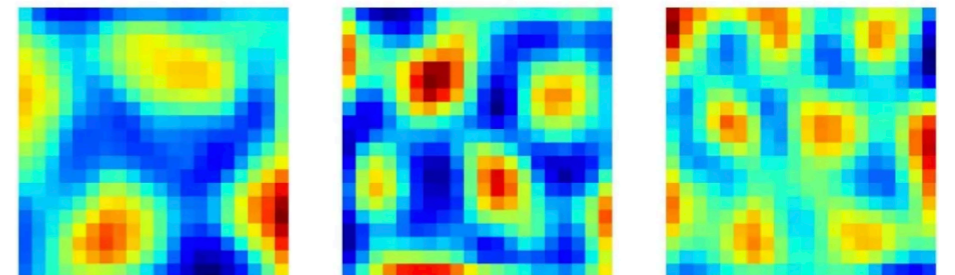
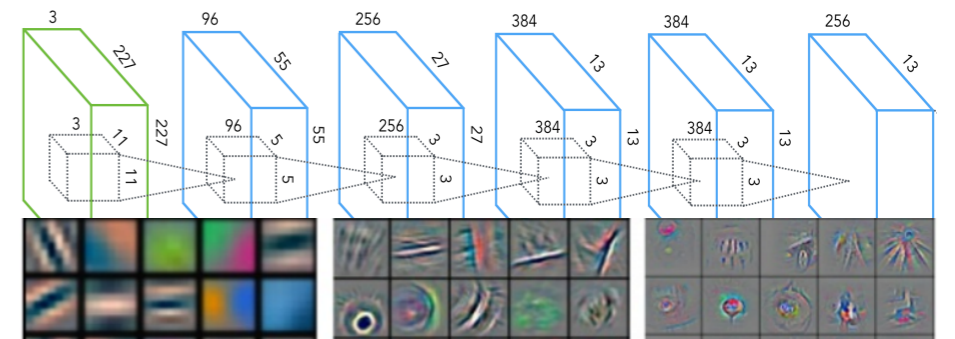
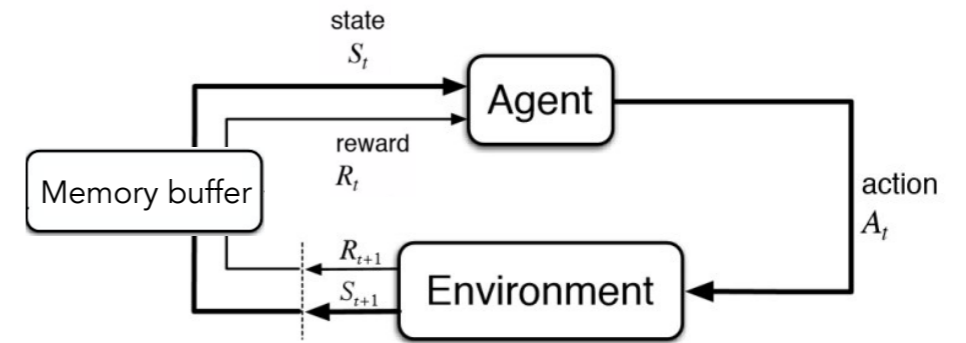
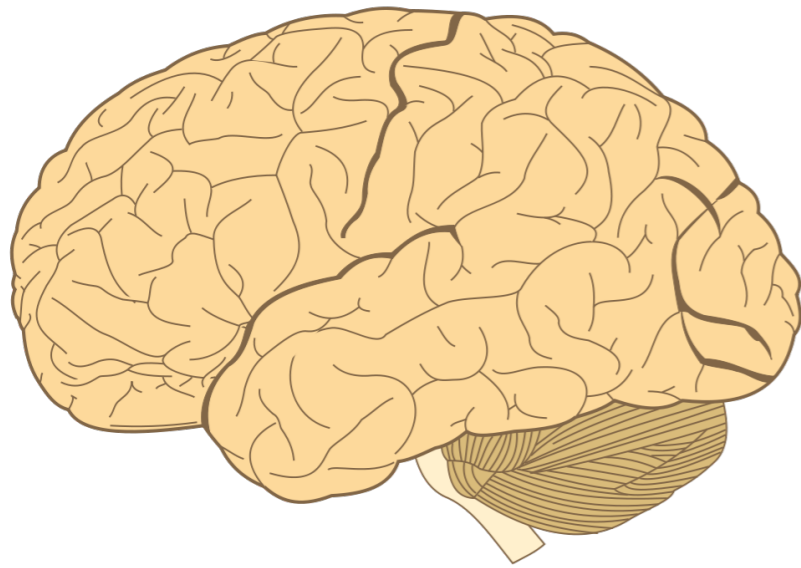
Random Forests:
feature-based **rules**



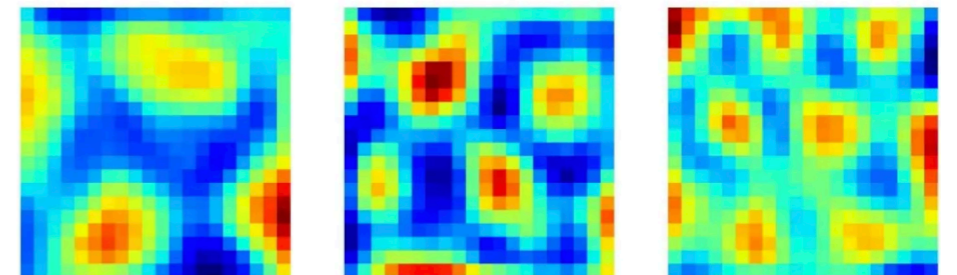
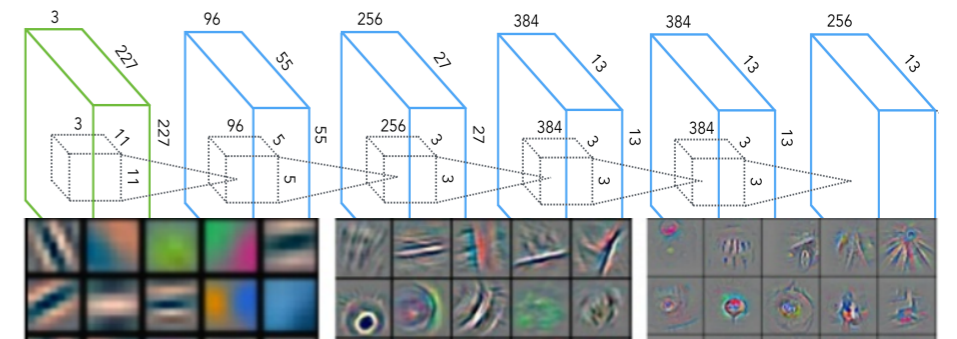
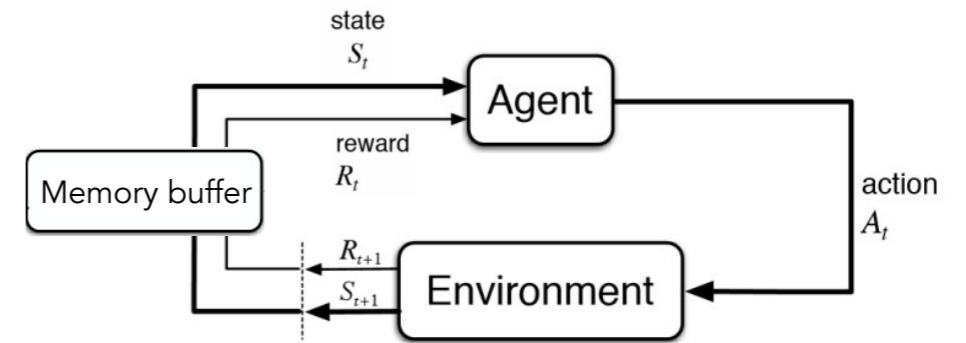
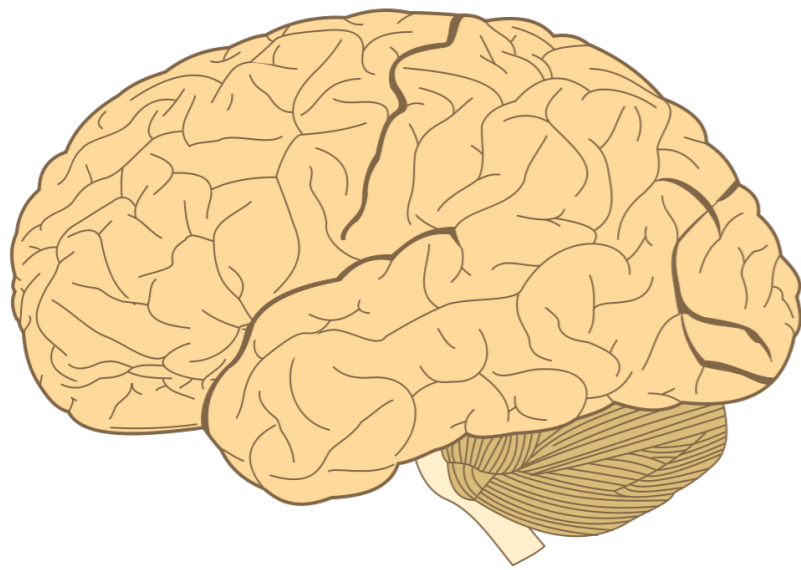
Gaussian Processes:
functions

Pick the one that will **reveal the knowledge you are after**,
not the one that just gives the best performance on a metric.

Curiously similar mechanisms in **biological** and **artificial** systems



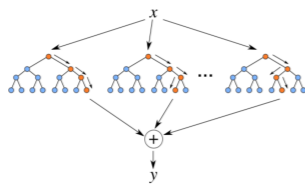
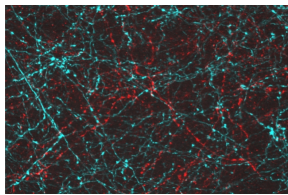
Curiously similar mechanisms in **biological** and **artificial** systems



Interpreting the mechanisms of machine learning models can shed light on the mechanisms of the brain

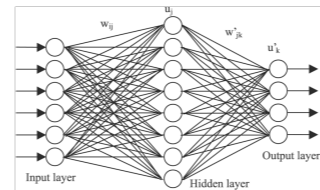
Discussion

- Modeling is a well-proven way of obtaining knowledge
- Machine-learned models do capture the knowledge, but an additional step of interpretation is required
- In life sciences model interpretation has a special significance
- Three examples of applying this principle in Neuroscience



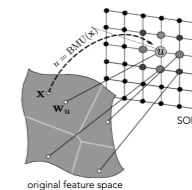
"Identifying task-relevant spectral signatures of perceptual categorization in the human cortex"

Scientific Reports, 2020



"Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex"

Communications Biology, 2018



"Mental state space visualization for interactive modeling of personalized BCI control strategies"

Journal of Neural Engineering, 2020

- Interpretability adds a new axis for algorithm selection
- Interpreting the mechanisms of machine learning models can shed light on the mechanisms of the brain

This would not be possible without

*the support and knowledge given by the **University of Tartu** and the **Institute of Computer Science**,*

***Raul Vicente** establishing the **Computational Neuroscience Lab** and supervising my PhD studies,*

***Konstantin Tretyakov** introducing me to the field of machine learning,*

***Sven Laur** further developing my understanding of machine learning,*

*the work done by **Juan R. Vidal** and the co-authors from **Lyon Neuroscience Research Center**,*

*the constant flow of ideas born at the seminars, lunch breaks, discussions with **Anna Leontjeva**,
Tambet Matiisen, **Jaan Aru**, **Ardi Tampuu**, **Kristjan Korjus** and all lab members and alumni, students,
and co-authors,*

*my parents and family, who showed me the value of knowledge and provided with the opportunity
to pursue it.*

Thank you!

Neuroscience

Aims to understand learning systems and intelligence by analyzing and reverse engineering the existing example



a special kind of synergy that leads to **curious similarities**

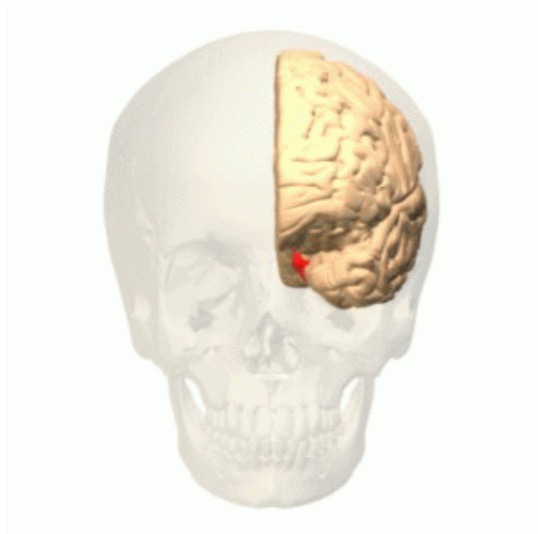


Aims to build an intelligent system ground-up by figuring out the building blocks and rules of interactions between them

Machine Learning & AI

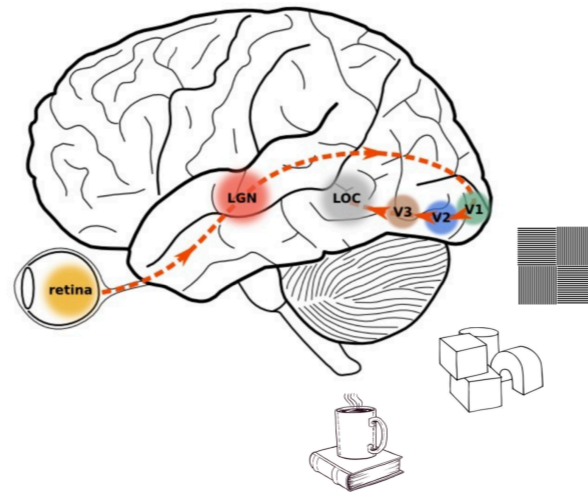
Curiously similar mechanisms in **biological** and **artificial** systems

Experience replay



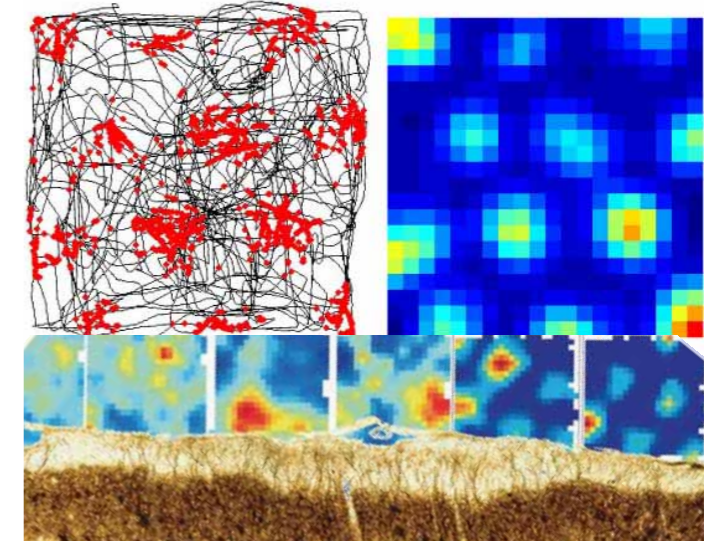
Hippocampus

Hierarchy of visual layers

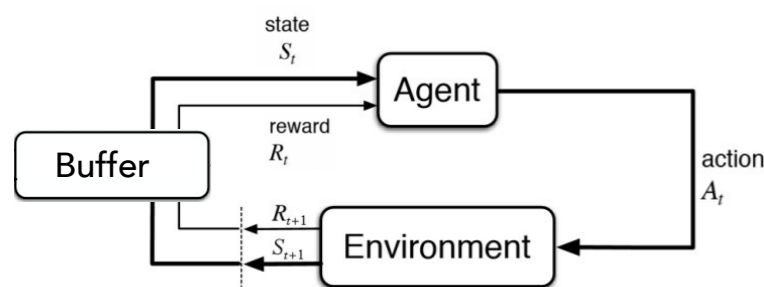


Layers of visual cortex

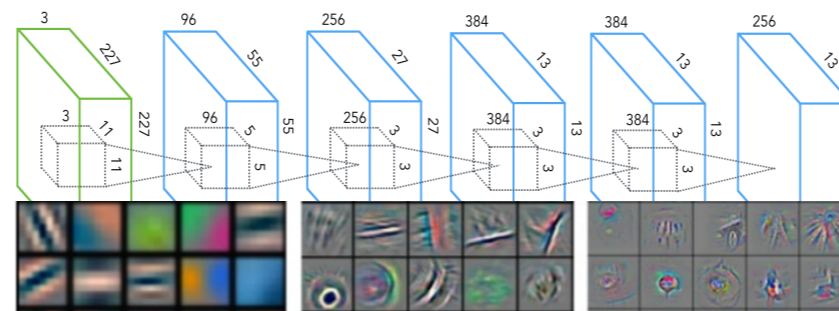
An efficient spacial code



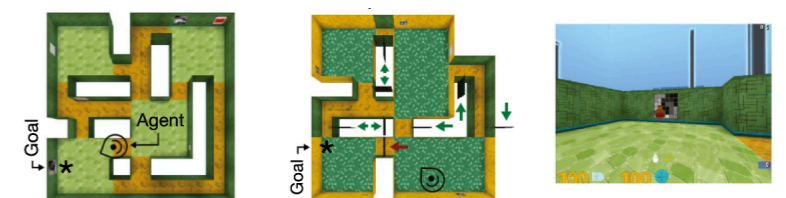
Grid cells in entorhinal cortex



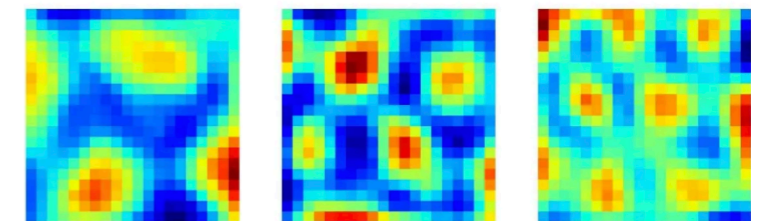
Deep Q-Learning



Deep convolutional neural network



Artificial (Agent)



Self-emergent spacial code

Memory consolidation Experience replay

- Layered structure
- Hierarchy of representational complexity
- Receptive fields convolutional filters

Grid-based code for location