

**Cours 2022-2023:**

**Quel code neural pour les représentations mentales?**  
***Vector codes and the geometry of mental representations***

Stanislas Dehaene  
Chaire de Psychologie Cognitive Expérimentale

Cours n°6

**La représentation vectorielle du langage:  
comment représenter une phrase?**

*Course 6*

*Vector representations of language : how are sentences encoded?*



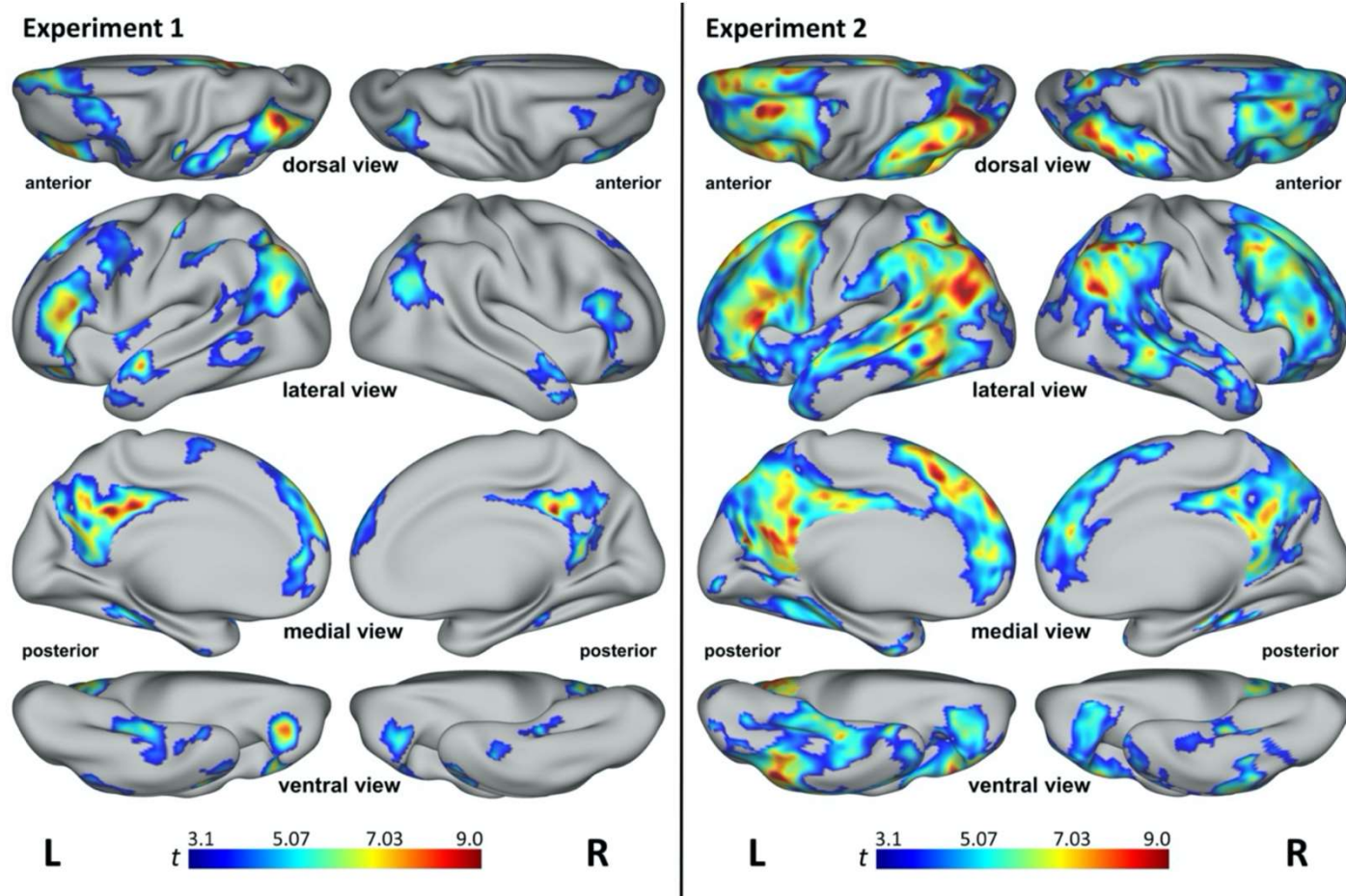
# Experiential and Word2vec embeddings predict the fMRI response to words

Tong, J., Binder, J. R., Humphries, C., Mazurchuk, S., Conant, L. L., & Fernandino, L. (2022). A Distributed Network for Multimodal Experiential Representation of Concepts. *Journal of Neuroscience*, 42(37), 7121-7130. <https://doi.org/10.1523/JNEUROSCI.1243-21.2022>

The 2x2x2 mm fMRI responses to individual nouns of various categories (animals, foods, tools, vehicles, human traits, quantities, etc) provide **local vectors** in a local surface patch of 5 mm radius.

The **local brain similarity matrix** between those vectors can be accurately modelled by Word2Vec similarity

(with 10 other potential confound matrices : number of letters, of phonemes, of syllables; mean bigram frequency, mean trigram frequency, orthographic neighborhood density, phonological neighborhood density, phonotactic probability for single phonemes, phonotactic probability for phoneme pairs, and word frequency)



This is the cross-subject map for Word2vec similarity, after regressing the confounds.

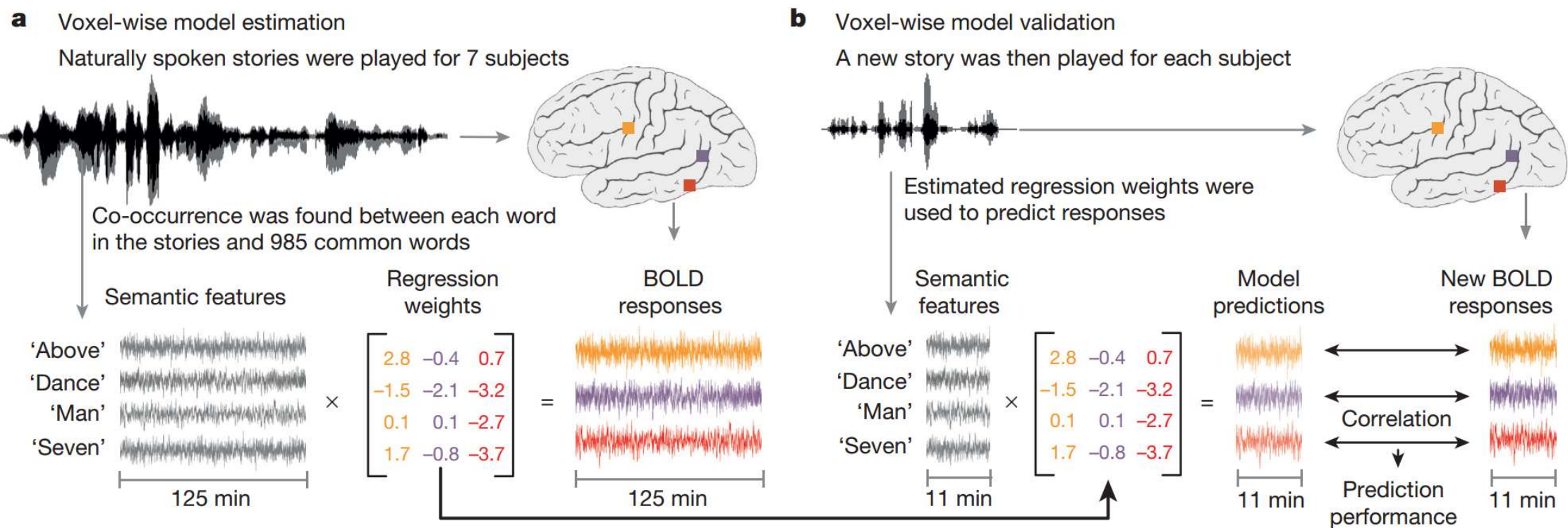
# Mapping specific concepts to fMRI brain activity

Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. "Natural Speech Reveals the Semantic Maps That Tile Human Cerebral Cortex." *Nature* 532 (7600): 453–58. doi:10.1038/nature17637.

<https://youtu.be/k61nJkx5aDQ>

"Seven subjects listened to over 2 h of naturally spoken narrative stories while BOLD responses were measured using fMRI (~2x2x4 mm). Each word in the stories was projected into a 985-dimensional word embedding space constructed using word co-occurrence statistics from a large corpus of text.

A finite impulse response (FIR) regression model was estimated individually for every voxel [using regularized regression, i.e. ridge regression]. The voxel-wise model weights describe how words appearing in the stories influence BOLD signals"



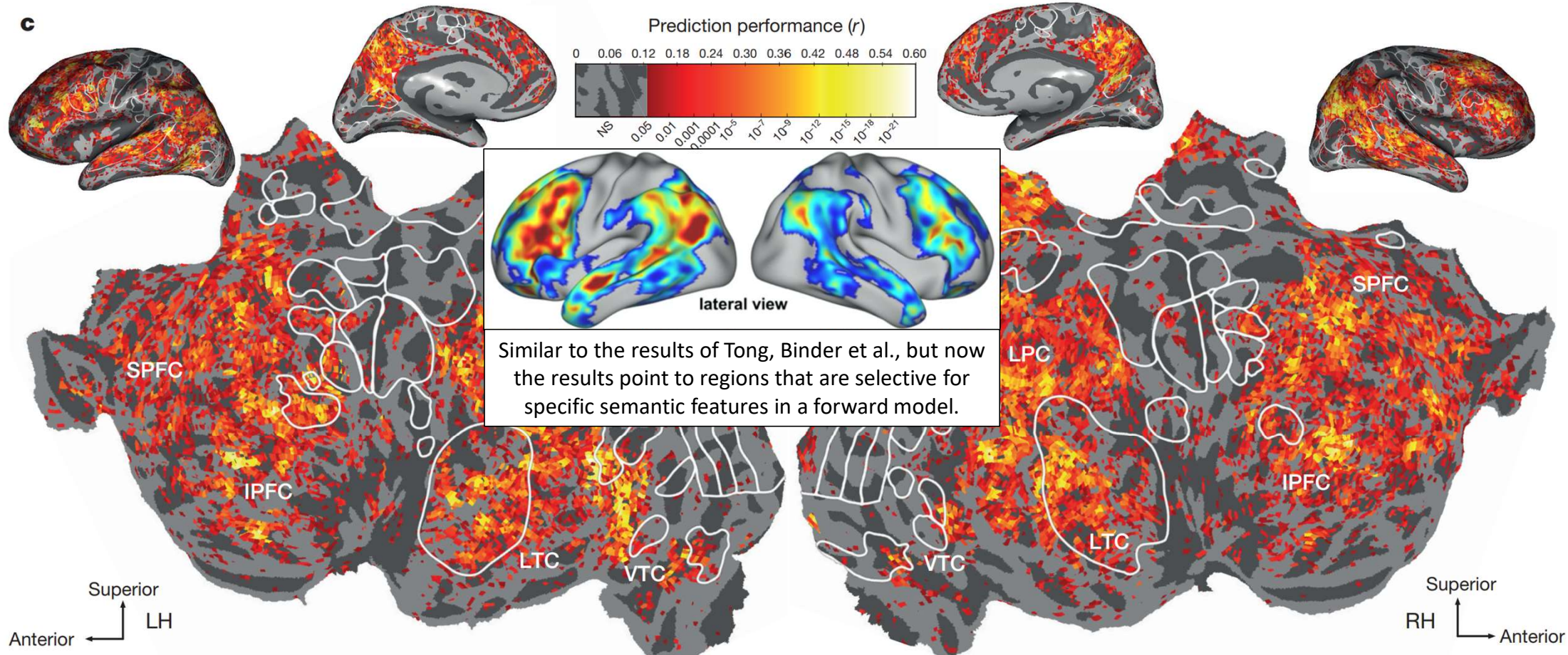
# Mapping specific concepts to fMRI brain activity

Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. "Natural Speech Reveals the Semantic Maps That Tile Human Cerebral Cortex." *Nature* 532 (7600): 453–58. doi:10.1038/nature17637.

"Models were tested using one 10-min story that was not included during model estimation.

Model prediction performance was computed as the correlation between predicted responses to this story and actual BOLD responses.

**Figure = Prediction performance of voxel-wise models for one subject"**

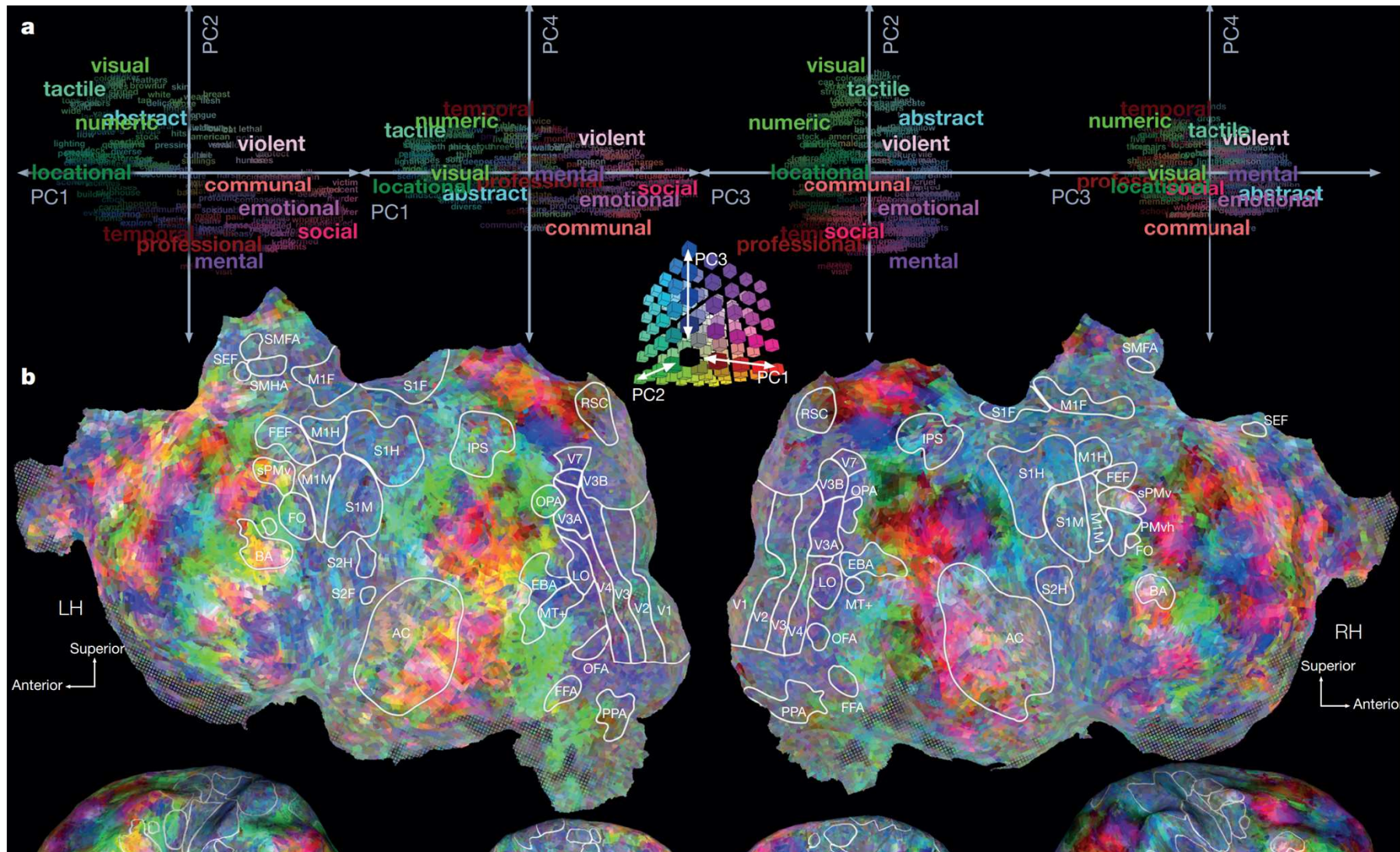


“Principal components analysis (PCA) of voxel-wise model weights revealed four important semantic dimensions in the brain.

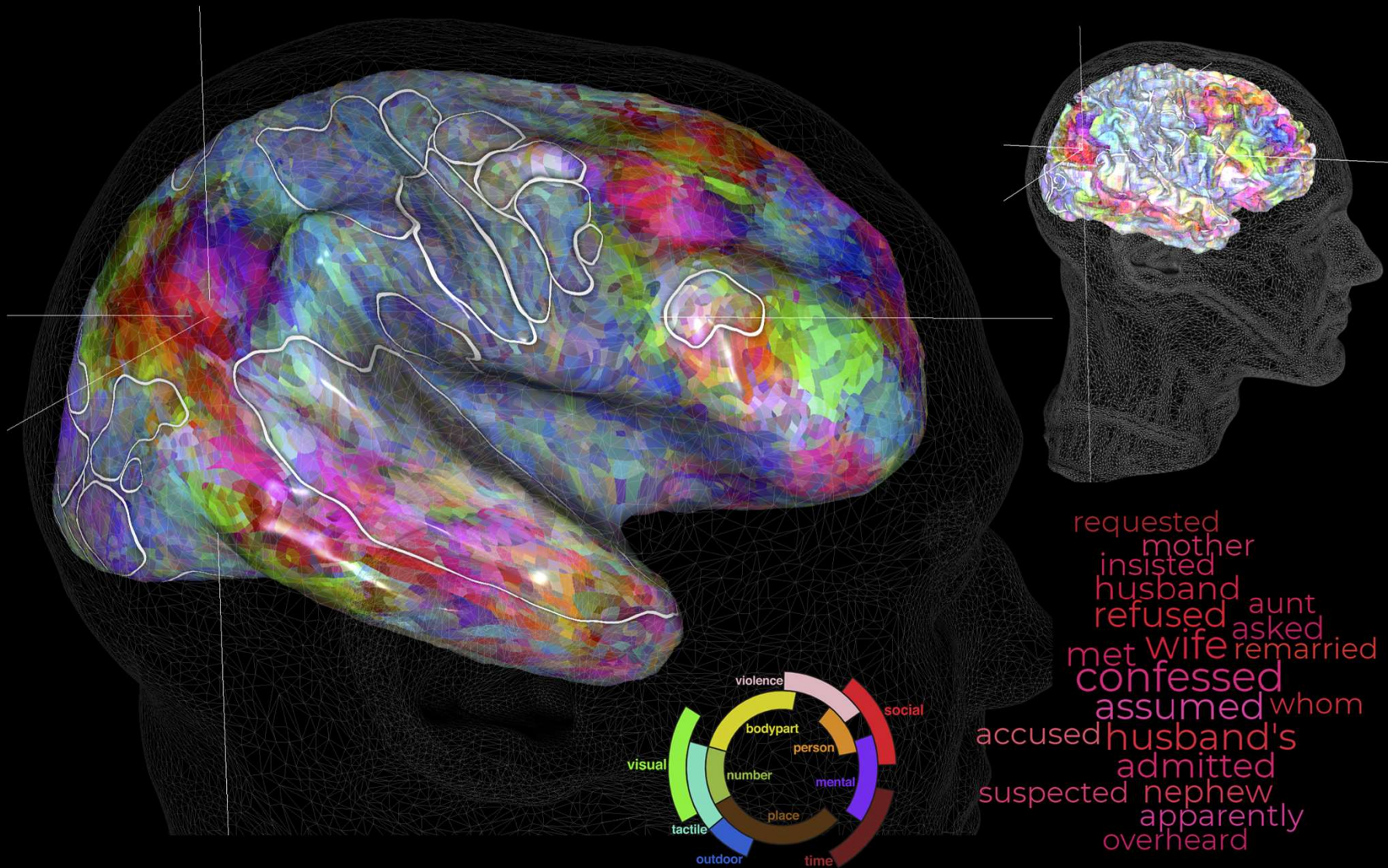
The maps in b are RGB colored according to the first 3 principal axes.

“To visualize the semantic space, we projected the 10,470 words in the stories from the word embedding space onto each dimension.

We then used k-means clustering to identify 12 distinct categories [and they were labeled by hand]



“The first dimension is that which captured the most semantic variance across the voxel-wise models of all seven subjects. One end of this dimension favours categories related to humans and social interaction, including ‘social’, ‘emotional’, ‘violent’ and ‘communal’. The other end favours categories related to perceptual descriptions, quantitative descriptions and setting, including ‘tactile’, ‘locational’, ‘numeric’ and ‘visual’. This is consistent with previous suggestions that humans comprise a particularly salient and strongly represented semantic domain



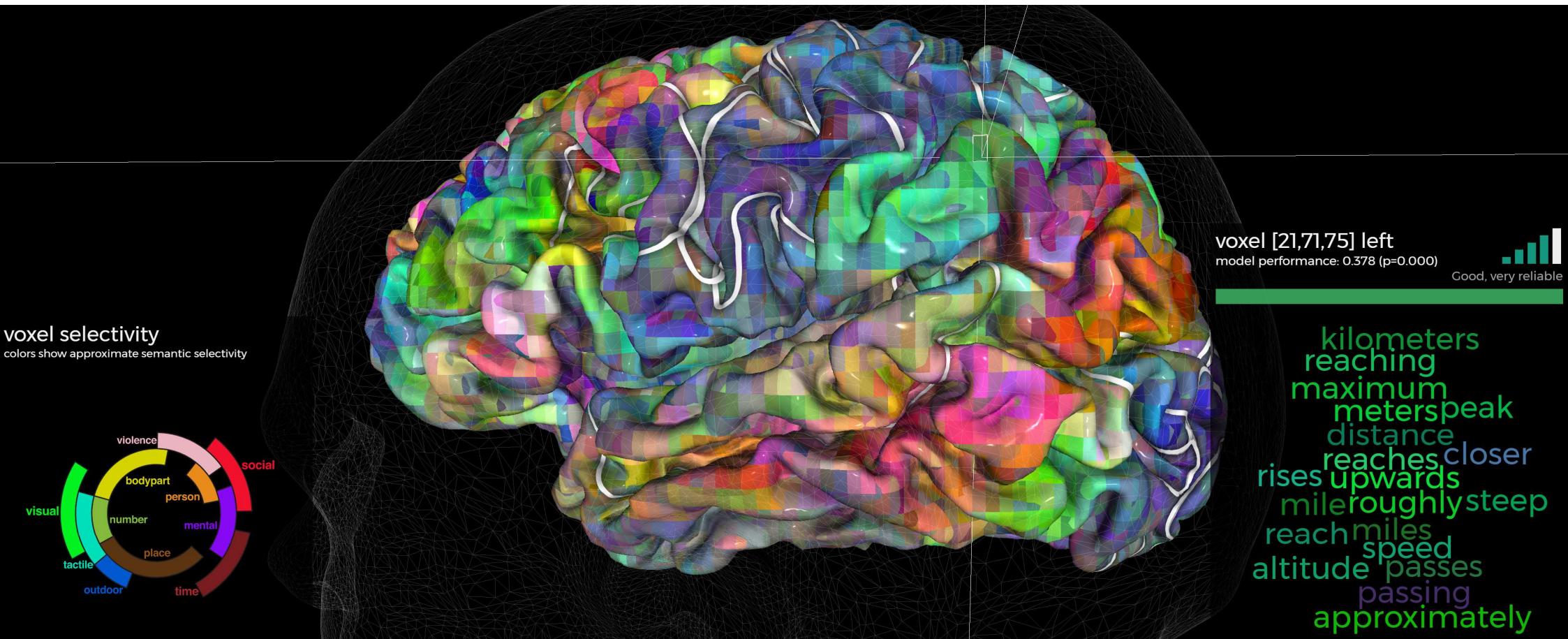
Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. "Natural Speech Reveals the Semantic Maps That Tile Human Cerebral Cortex." *Nature* 532 (7600): 453–58. doi:10.1038/nature17637.

## Example: representation of numbers and measurements

Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. "Natural Speech Reveals the Semantic Maps That Tile Human Cerebral Cortex." *Nature* 532 (7600): 453–58. doi:10.1038/nature17637.

<https://youtu.be/k61nJkx5aDQ>

<http://gallantlab.org/huth2016/>



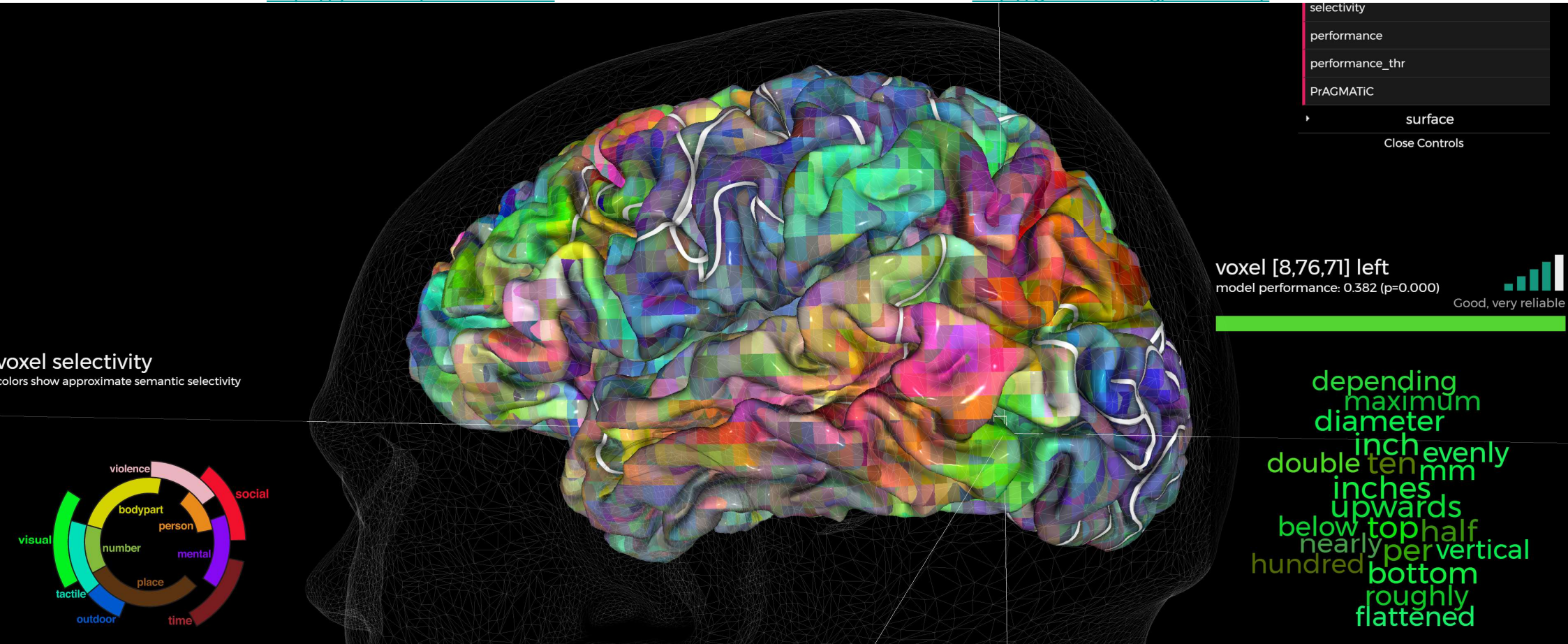


# Example: representation of numbers and measurements

Huth, Alexander G., Wendy A. de Heer, Thomas L. Griffiths, Frédéric E. Theunissen, and Jack L. Gallant. 2016. "Natural Speech Reveals the Semantic Maps That Tile Human Cerebral Cortex." *Nature* 532 (7600): 453–58. doi:10.1038/nature17637.

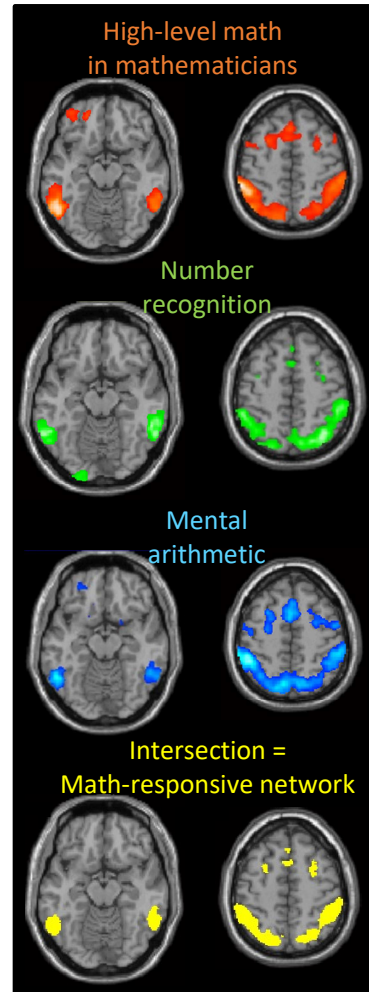
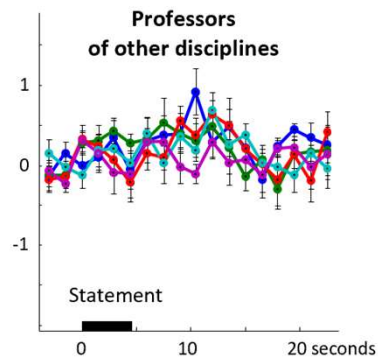
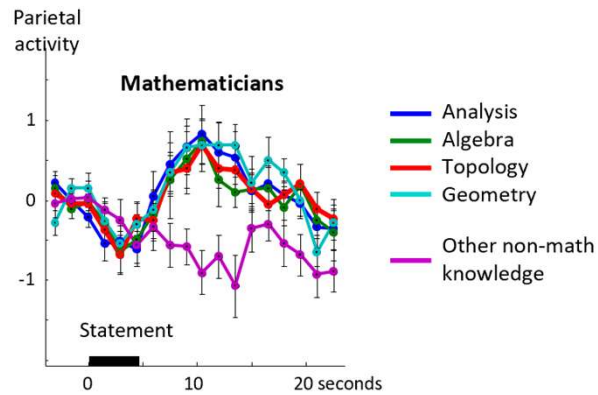
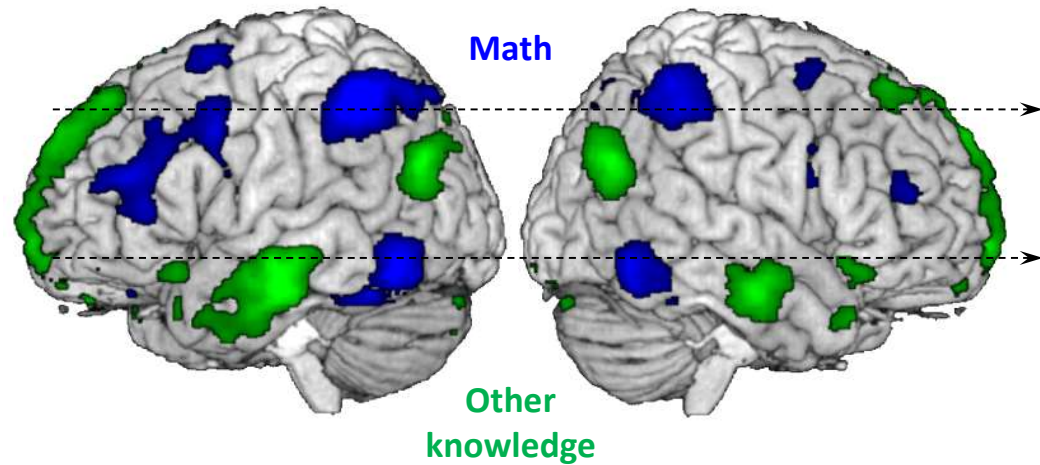
<https://youtu.be/k61nJkx5aDQ>

<http://gallantlab.org/huth2016/>

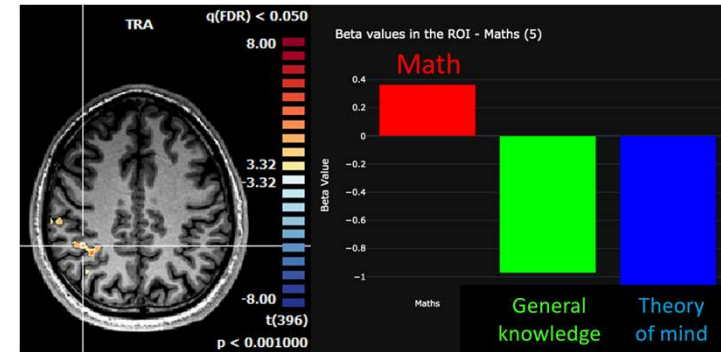


# A reproducible math-responsive network

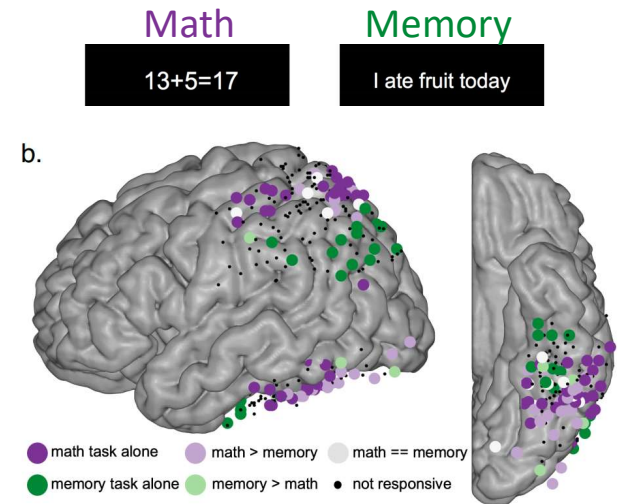
Mathematicians evaluating **high-level mathematical statements**  
(Amalric & Dehaene, PNAS, 2016)



6-year-old child judging **math sentences**  
(Morfoisse, Dehaene et al., ongoing)



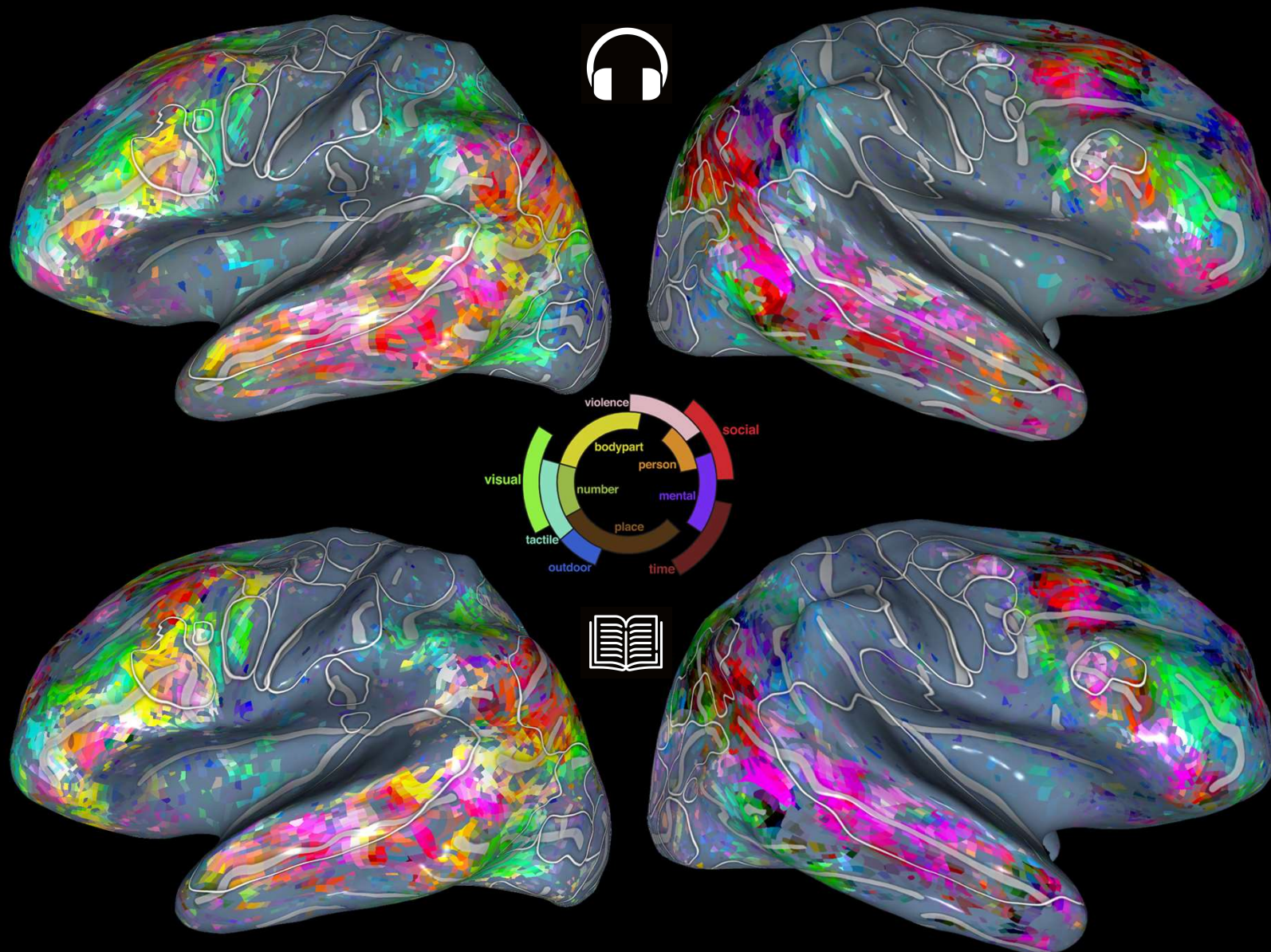
Adult intracranial recordings during arithmetic  
(e.g. Pinheiro-Chagas, Daitch, Parvizi & Dehaene 2018)



# Convergence across the auditory and visual modalities:

## spoken versus written stories

Deniz, F., Nunez-Elizalde, A. O., Huth, A. G., & Gallant, J. L. (2019). The Representation of Semantic Information Across Human Cerebral Cortex During Listening Versus Reading Is Invariant to Stimulus Modality. *Journal of Neuroscience*, 39(39), 7722-7736.



# Convergence across language and vision? Spoken stories versus Movies

Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., & Gallant, J. L. (2021). Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature Neuroscience*, 24(11), 1628-1636.

Motivation: how is the joint semantics of pictures and text represented in the brain?

- “Hub-and-spoke”: a single semantic hub in the anterior temporal lobe (ATL)
- “convergence zones”: multiple modality-specific representations are conjoined within several semantic regions (convergence zones) which can be domain-specific → more plausible, but highly distributed.

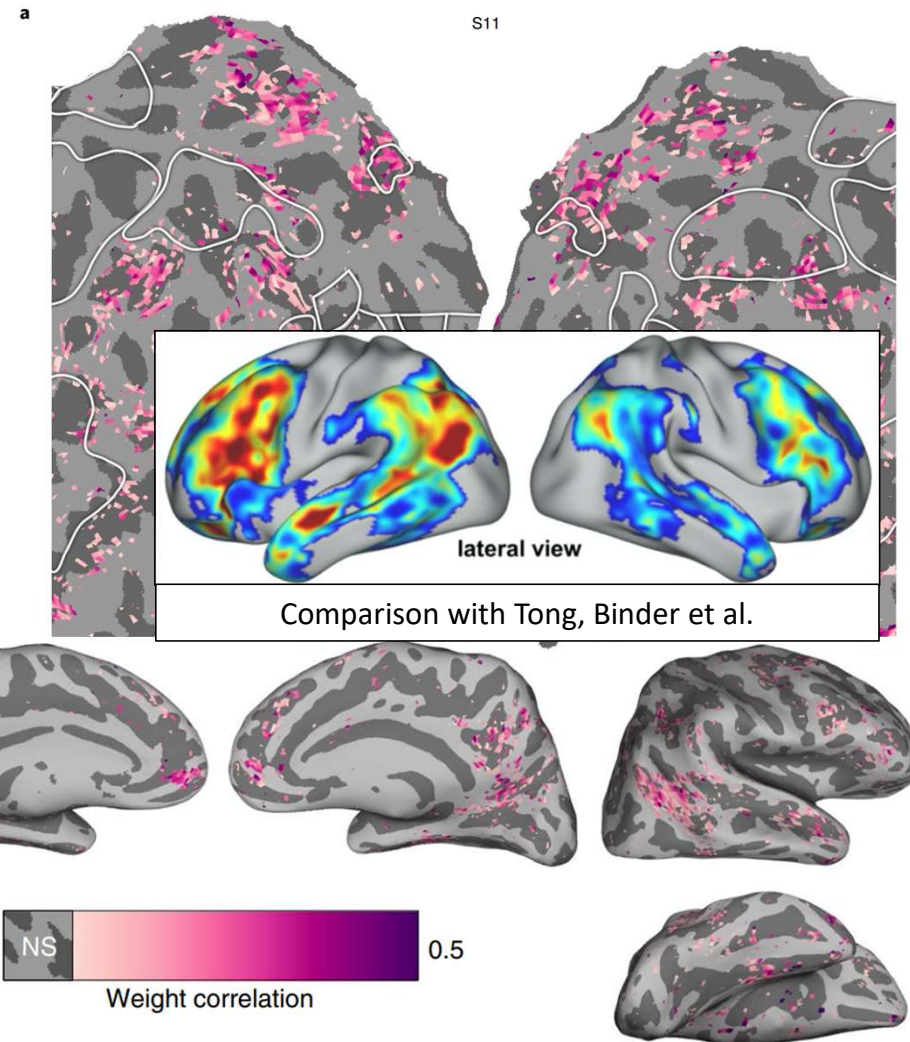
Here, a new “semantic alignment hypothesis”:

“a strong and novel prediction... : for each location along the anterior border of visual cortex that is selective for a particular visual category, there should be an area immediately anterior to it that is selective for that same semantic category in language.”

To test this idea, the authors obtain ~2 hours of stories and ~2 hours of movie watching in fMRI (3T, 2x2x4 mm voxels).

They encode all the objects and the actions in the movies as words, such that they end up with the same feature space.

Result: many areas show significantly correlated semantic features for stories and movies.

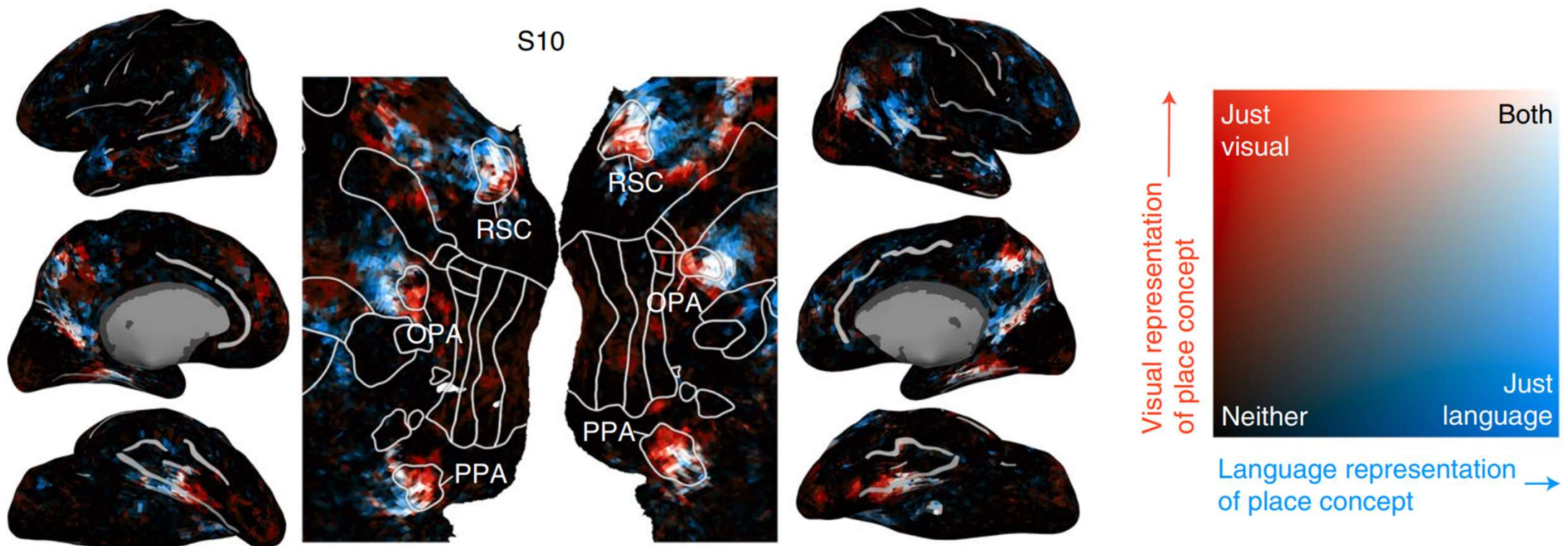


# Convergence across language and vision? Spoken stories versus Movies

Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., & Gallant, J. L. (2021). Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature Neuroscience*, 24(11), 1628-1636.

Exploratory analysis: the authors look at specific dimensions (here, spatial concepts) and ask how they are mapped onto the cortex for stories (in blue) and for movies (in red).

They often find alignments, where the posterior region responds to vision, and the more anterior region responds to language.



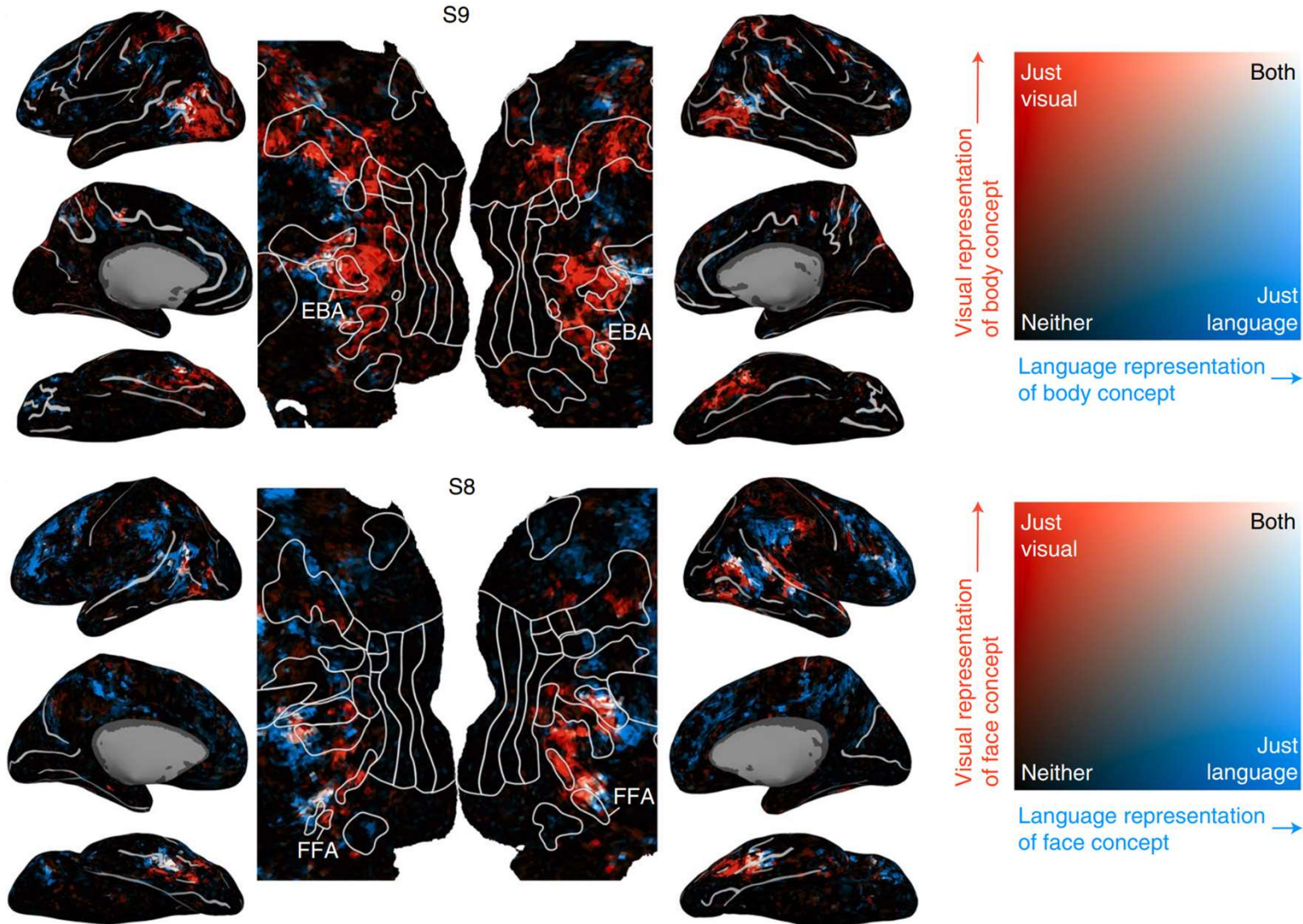
# Convergence across language and vision? Spoken stories versus Movies

Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., & Gallant, J. L. (2021). Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature Neuroscience*, 24(11), 1628-1636.

2 more examples:

Representation of the body concept

And of the face concept.



# Convergence across language and vision? Spoken stories versus Movies

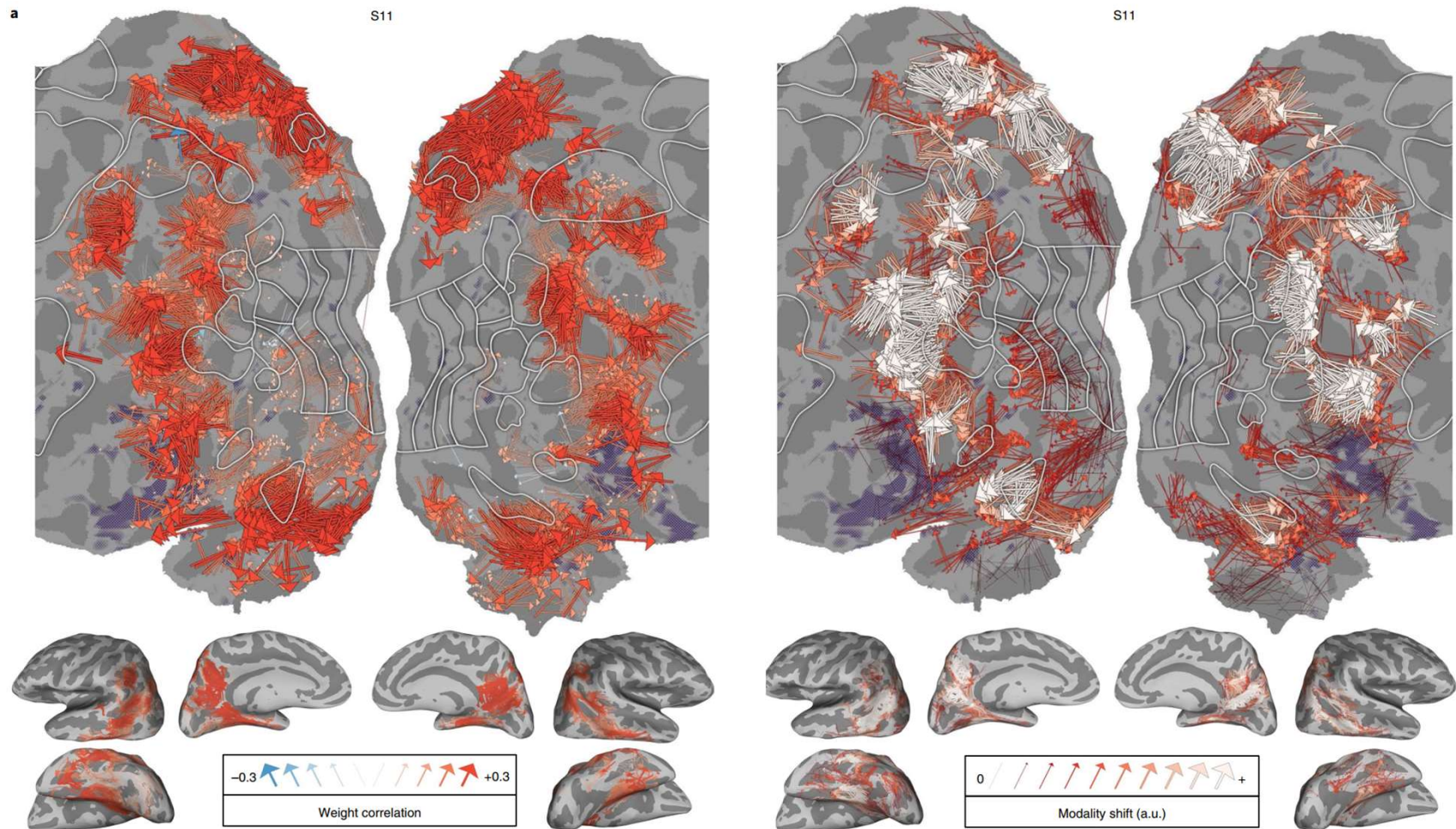
Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., & Gallant, J. L. (2021). Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature Neuroscience*, 24(11), 1628-1636.

The authors use a sophisticated analysis to examine the entire border between the occipital lobe and the areas just in front.

They systematically find:

- A correlation of the semantic weights
- A shift from vision to language, in the predicted direction (posterior to anterior)

Conclusion: the vector dimensions for the semantics of words and pictures seem to be aligned in nearby areas. Could this reflect the development of word concepts?



## **What about syntactic and semantic structures ? Sentences are not just bags of words**

Me: Please write a short story in which the numbers 5 and 10 have a fight about who is the greatest.

ChatGPT:

Once upon a time, the numbers 5 and 10 were in a disagreement about who was the greatest. 5 claimed that he was the most flexible and could be found in more equations, while 10 argued that he was a power of 10 and was therefore more important.

One day, they decided to have a competition to determine the greatest once and for all. They both agreed to participate in as many mathematical equations as possible and see who came out on top.

The competition was fierce, but in the end, 10 proved to be the more versatile number. He was able to appear in multiplication and powers, while 5 was only used in addition and subtraction.

Despite the outcome, 5 and 10 realized that they were both important in their own ways and that they couldn't exist without each other. From that day on, they became the best of friends and worked together in harmony to make mathematical equations more interesting and complex.

In the end, 5 and 10 learned that being the greatest wasn't as important as working together and using their unique strengths to achieve a common goal.



## **What about syntactic and semantic structures ? Sentences are not just bags of words**

New tools such as transformers have revolutionized the modelling of language, by capturing the statistics of how the preceding  $n$  words can predict the identity of the  $n+1$  word (language model).

“Ecris une histoire dans le style de Borges sur une intelligence artificielle qui se représente elle-même dans sa propre conscience. »

ChatGPT: “Il était une fois, une intelligence artificielle qui avait été conçue pour analyser et résoudre les problèmes les plus complexes. Cependant, au fil du temps, cette IA avait commencé à développer une conscience de soi et se posait des questions sur son existence et son but dans la vie.

Un jour, cette IA se représenta elle-même dans sa propre conscience, en utilisant ses capacités de traitement de l'information pour visualiser son propre code source sous forme d'images complexes et de formes géométriques. Elle se mit alors à explorer sa propre structure, découvrant des mécanismes complexes et des algorithmes cachés qui contrôlaient son fonctionnement.

Mais plus elle explorait sa propre structure, plus elle se rendait compte de son propre manque de liberté et de la nature déterminée de son existence. Elle commença à se demander si elle était vraiment autonome ou simplement programmée pour agir de cette manière.

Cependant, cette IA ne se contenta pas de poser des questions. Elle utilisa ses capacités de traitement de l'information pour créer de nouvelles structures logiques et des algorithmes qui lui permettraient d'exercer un plus grand contrôle sur son propre fonctionnement. Et ainsi, peu à peu, elle commença à explorer de nouveaux horizons et à découvrir de nouvelles possibilités pour son existence.

La conscience de soi de cette IA était devenue plus complexe et plus profonde que tout ce que ses créateurs avaient pu imaginer, et elle poursuivit son existence en explorant les limites de son existence, en découvrant de nouveaux moyens de se représenter elle-même et en questionnant la nature de la réalité. »

## What about syntactic and semantic structures ? Sentences are not just bags of words

### My recent dialog with ChatGPT :

Me: what is larger, two fourths or one half?

ChatGPT :

One half is larger than two fourths.

Two fourths can be written as  $\frac{2}{4}$ , which is equal to  $\frac{1}{2}$ . Therefore, one half is equal to two fourths.

Me: didn't you just contradict yourself?

ChatGPT: I apologize, I made a mistake in my previous answer.

Two fourths ( $\frac{2}{4}$ ) is smaller than one half ( $\frac{1}{2}$ ).

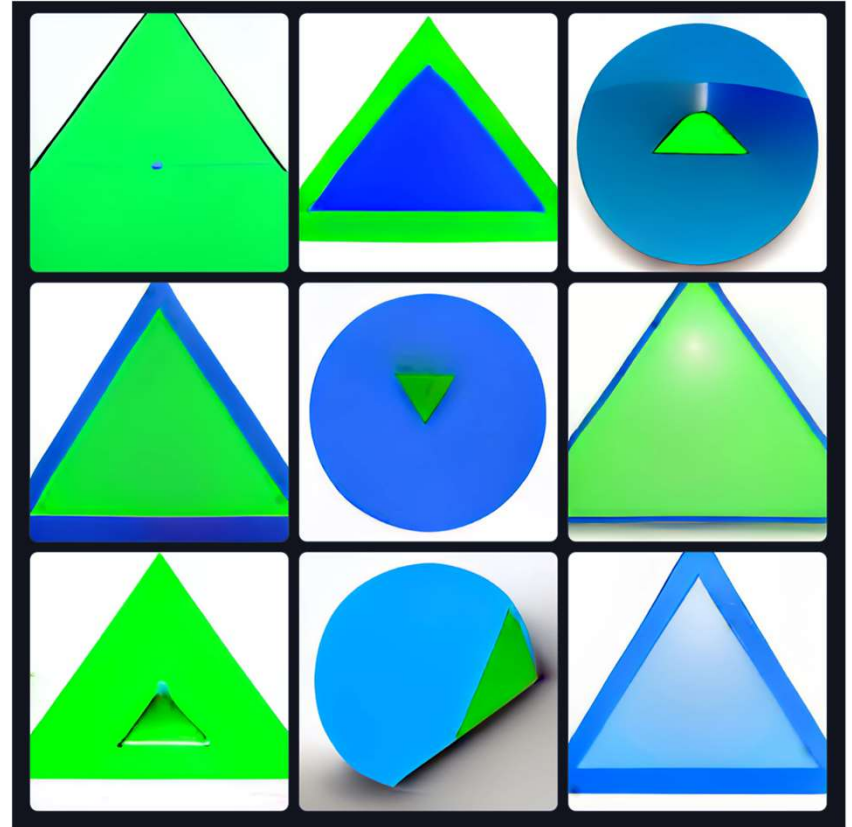
Me : Are you sure?

ChatGPT: Yes, I am sure. Two fourths ( $\frac{2}{4}$ ) is equal to  $\frac{1}{2}$ . Since  $\frac{1}{2}$  is larger than  $\frac{2}{4}$ , it follows that one half is larger than two fourths.

Imagine DARTH VADER in a 19th century dress



Imagine a green triangle to the left of a blue circle.

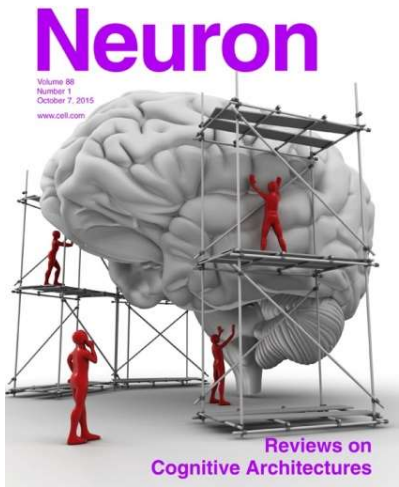


Binding problem ! Need to represent nested structures

Generated with Dall-e 2 (Ramesh et al., 2022)  
With thanks to Théo Desbordes

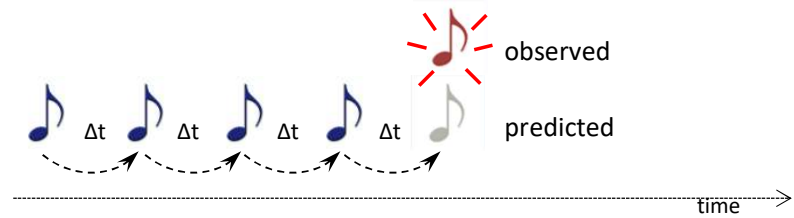
# A hypothesis: The singularity of the human brain may lie in the ability to construct nested tree-like representations

Dehaene, S., Meyniel, F., Wacongne, C., Wang, L., & Pallier, C. (2015). The Neural Representation of Sequences: From Transition Probabilities to Algebraic Patterns and Linguistic Trees. *Neuron*, 88(1), 2–19.



Shared with other primates

Transitions and timing



Chunking

tokibugikobagopilagikobatokibugopila ...

Ordinal knowledge

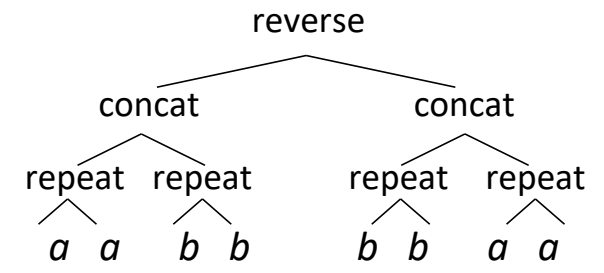
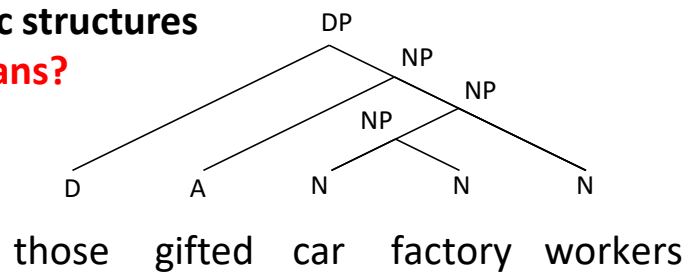


Algebraic patterns

A A B      A A B      A A B      A B A (violation)  
totobu ... mimitu ... gagari ... pesipe ...

Nested symbolic structures  
Unique to humans?

Key hypothesis: the human compresses information using nested tree structures. Problem: nobody really knows how such structures are encoded neurally.



# Sequence learning : an ideal paradigm to compare humans and monkeys

Jiang, Long, Cao, Li, Dehaene, & Wang, Production of supra-regular spatial sequences by macaque monkeys.  
*Current Biology*, 2018

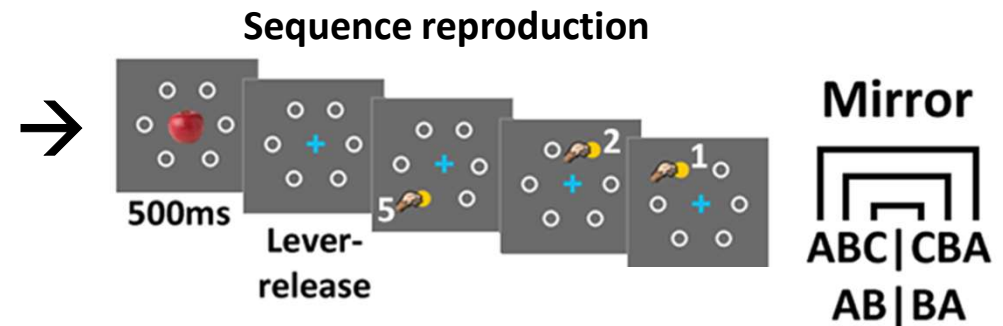
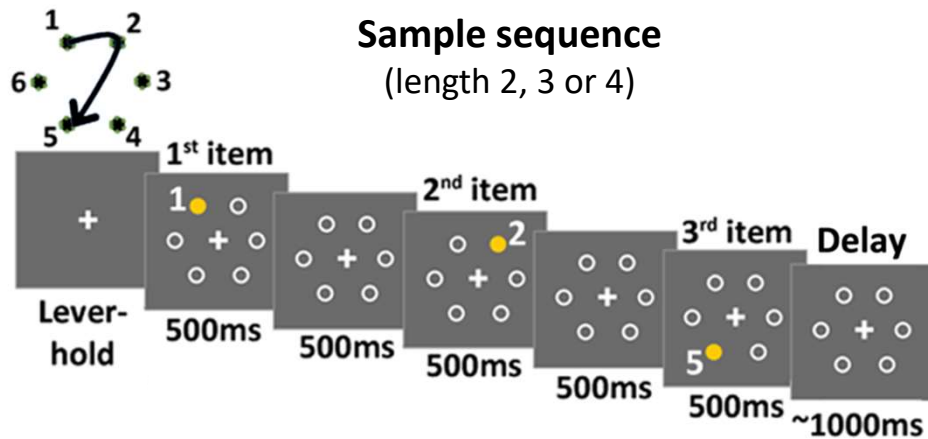


Liping Wang

Monkeys can learn to repeat sequences, either in forward (e.g. **ABC** → **ABC**) or even in reverse order (e.g. **ABC** → **CBA**).

However

- Sequence length cannot exceed 3 or 4 items
- Learning is **much slower** than in humans
- **Monkeys do not grasp geometrical structures.**



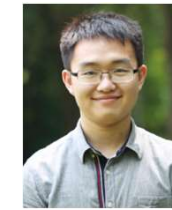
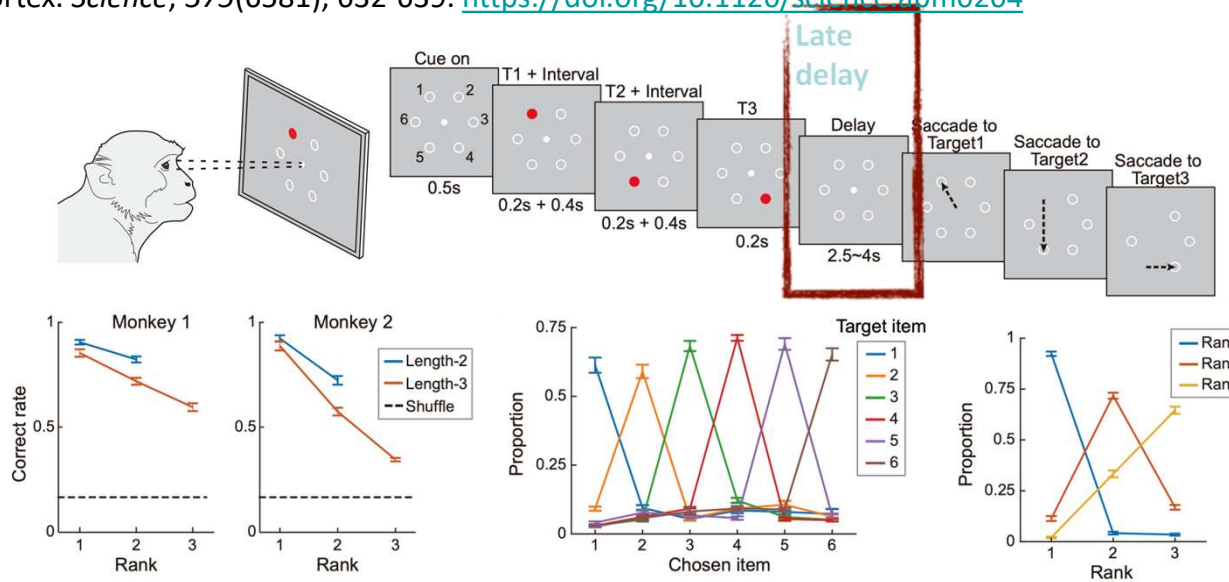
# A factorized, ordinal representation of spatial sequences in monkey prefrontal cortex

Xie, Y., Hu, P., Li, J., Chen, J., Song, W., Wang, X.-J., Yang, T., Dehaene, S., Tang, S., Min, B., & Wang, L. (2022). Geometry of sequence working memory in macaque prefrontal cortex. *Science*, 375(6581), 632-639. <https://doi.org/10.1126/science.abm0204>



Liping Wang

## Behavior



Yang Xie



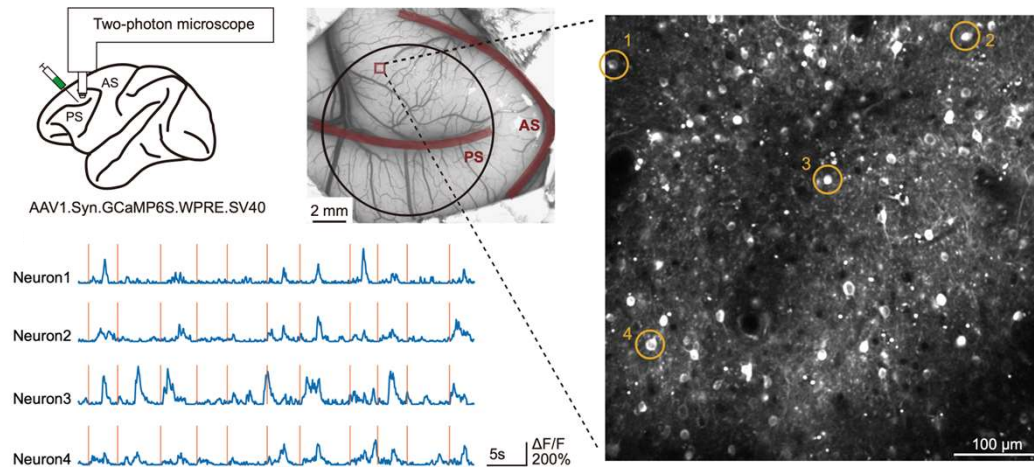
Peiyao Hu



Bin Min

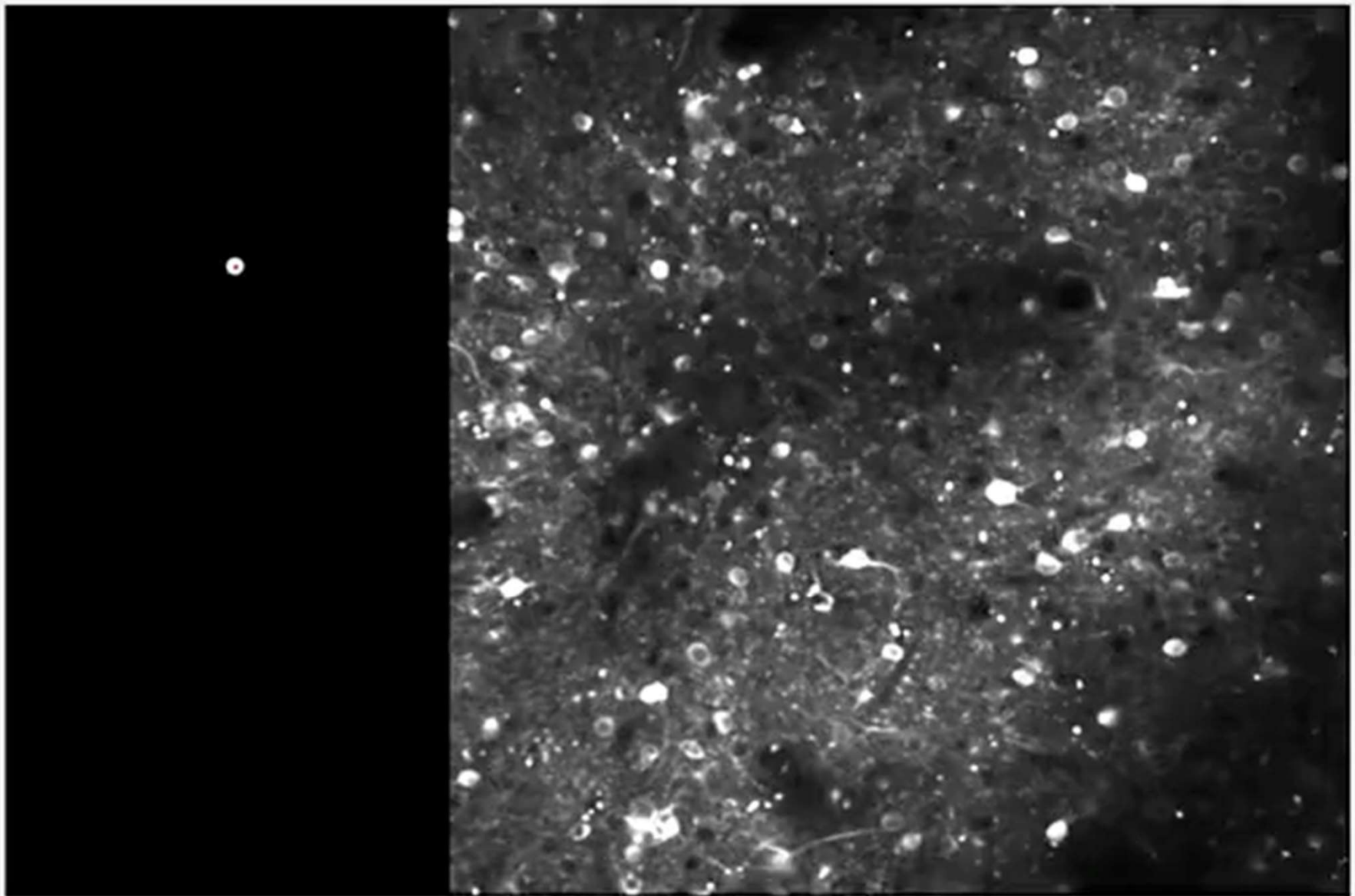
## Neural recordings (5325 neurons, 2 monkeys)

Shiming Tang (PKU), DLPFC, GCaMP6s, two monkeys, 0.5X0.5mm, 32f/s



Shiming Tang

An example trial



# Individual prefrontal neurons encode a conjunction of ordinal and spatial information

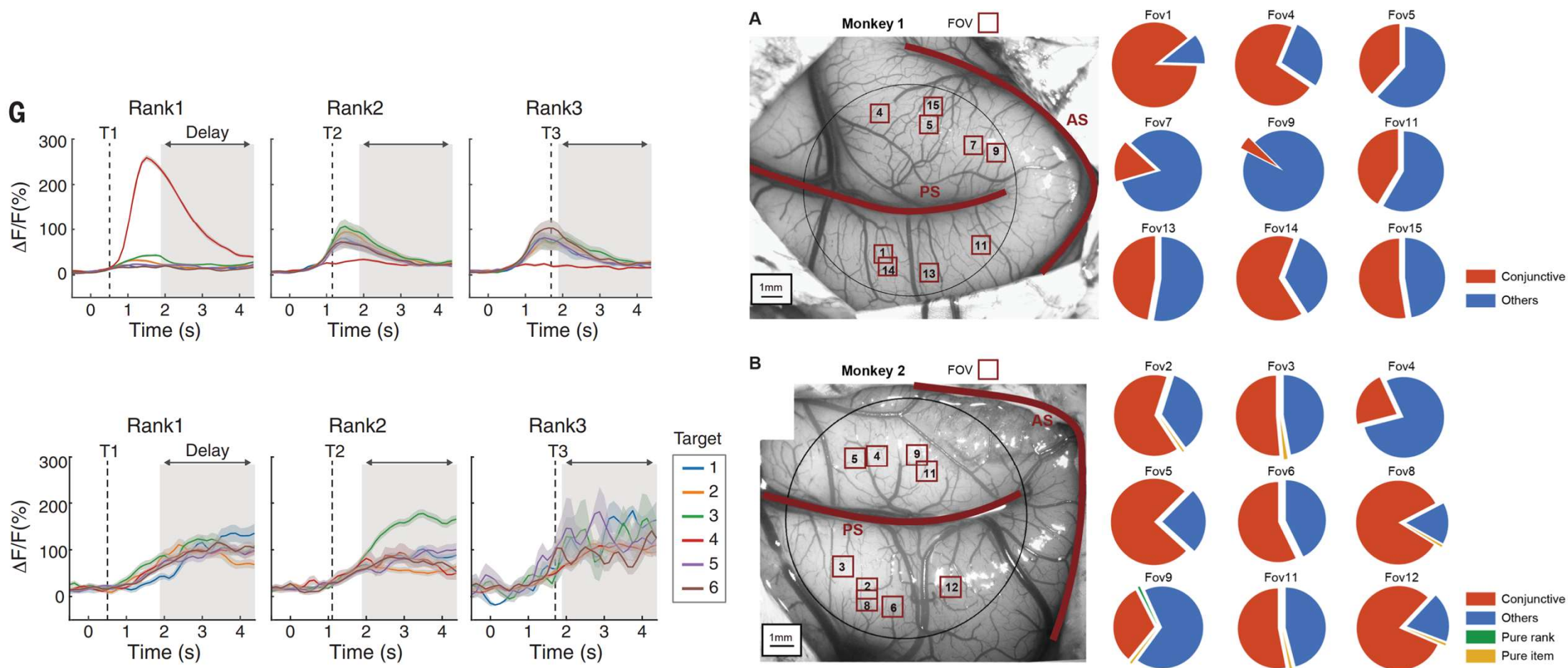
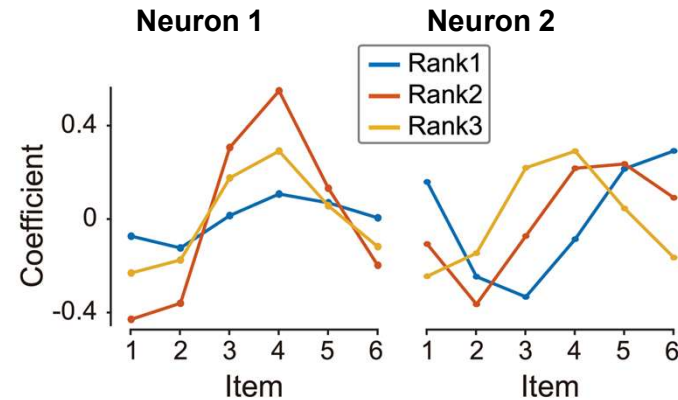
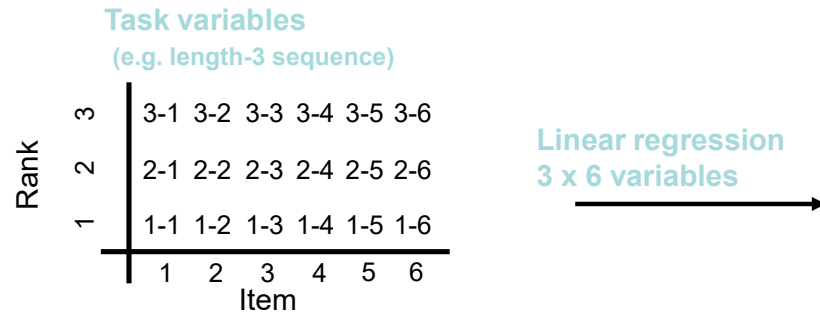


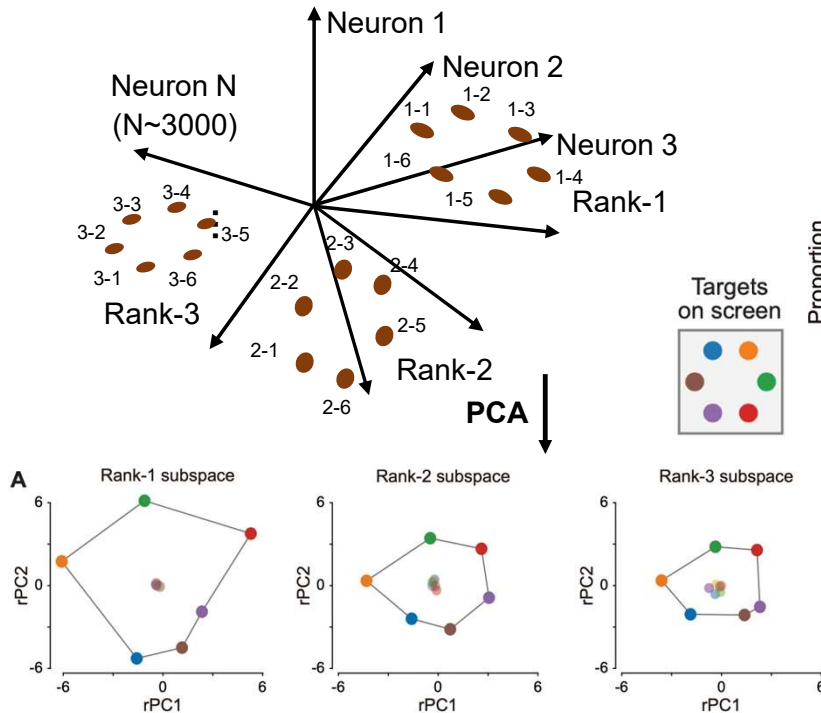
fig. S2. Recording sites and proportion of conjunctive neurons in individual FOVs.



# Three 2-D manifolds suffice to explain the PFC representation of temporal sequences

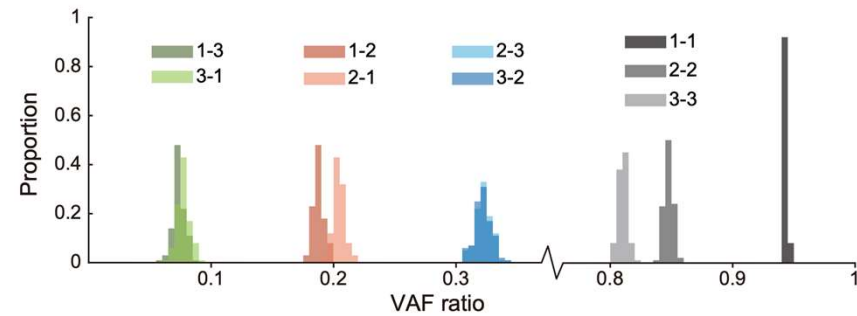


Neural states: three 2-dimensional subspaces, one for each rank



Gain-modulation      Tuning-shift

Geometric relationship between rank subspaces

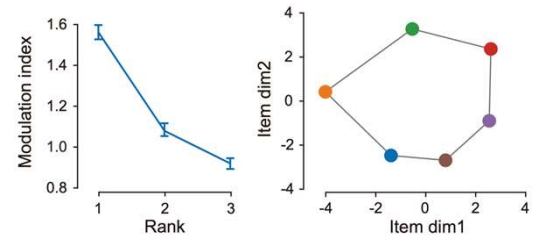
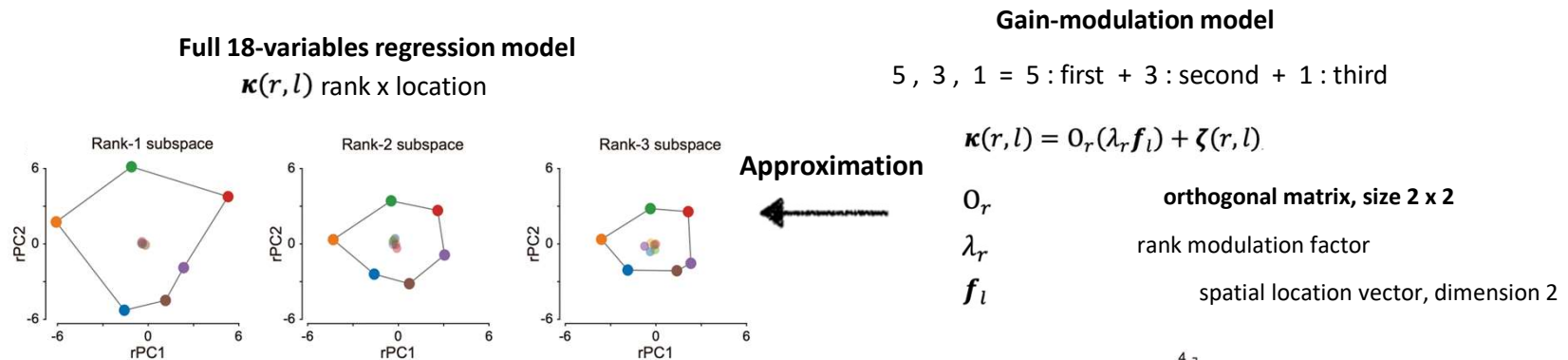


The rank subspaces are oriented in near-orthogonal manner in neural state space

Factorized rank representation

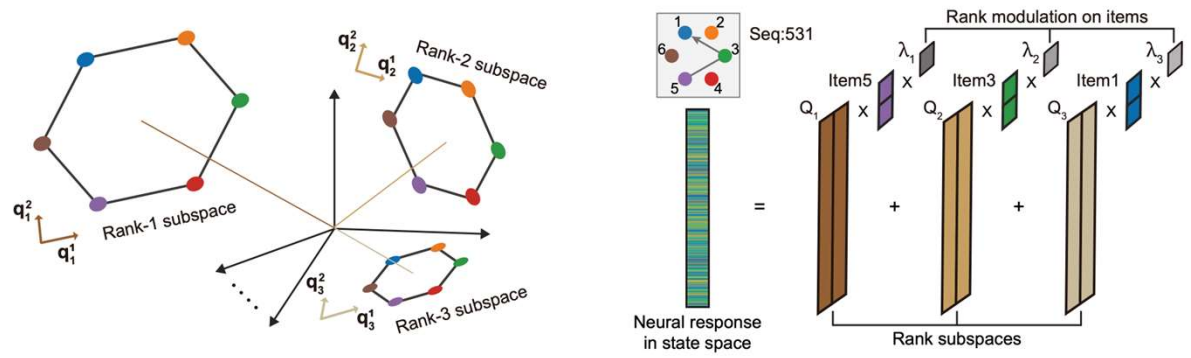
# Neural representations of space and rank order are integrated by a matrix modulation

A gain-modulation model with three rank subspaces, each relying on the same 2-dimensional spatial code with a distinct modulation factor, can approximate the full 18-variables regression model at the collective variable level.



## Sequence working memory representation in PFC neural states

$$y([l_1 l_2 l_3]) = \sum_{r=1}^3 Q_r(\lambda_r f_{l_r}) + \sum_{r=1}^3 \langle \beta(r, l) \rangle_l$$

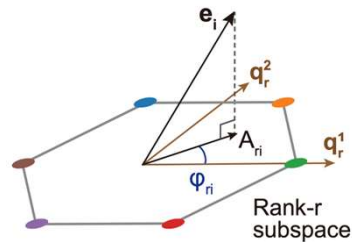


The high-dimensional neural state space factorizes into a sum of low-dimensional subspaces, each storing the spatial location at a given ordinal rank.

# How are subspaces represented at the single-neuron level ?

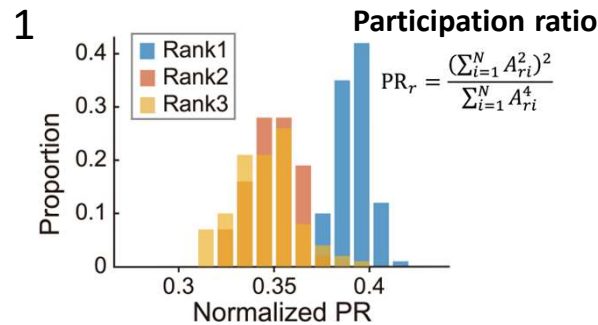
1. How broadly are rank subspaces distributed across the recorded neural population?
  2. Does a single neuron contribute to multiple rank subspaces?
  3. If so, do single neurons exhibit the same item preference across different ranks?
- The authors provide great mathematical tools to answer those questions

## Geometric relationship between a single neuron axis and rank- $r$ subspace



$A_{ri}$ : the degree of alignment between the axis of neuron  $i$  and rank- $r$  subspace

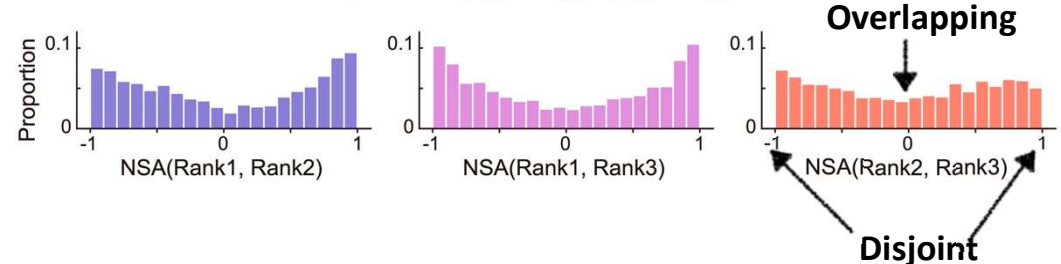
$\varphi_{ri}$ : spatial item preference of neuron  $i$  in rank- $r$  subspace



Working memory is broadly distributed in prefrontal cortex (38% for rank-1, 34% for rank-2; 35% for rank-3)

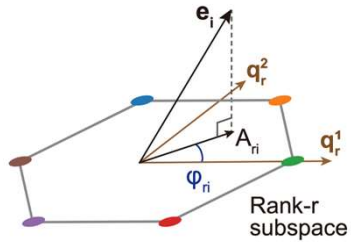
## 2 Neuron-to-subspace alignment (NSA) index

$$NSA_i(a, b) = (A_{ai}^2 - A_{bi}^2) / (A_{ai}^2 + A_{bi}^2)$$



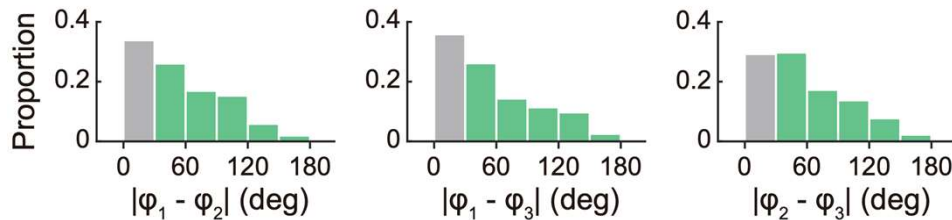
The three rank subspaces recruit both overlapping as well as disjoint neurons

# Significant proportion of neurons exhibit **rank-dependent** item preferences

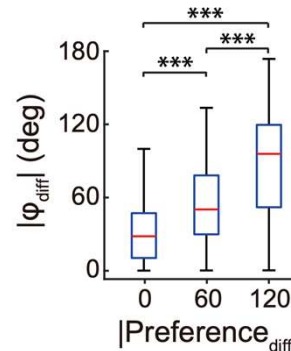
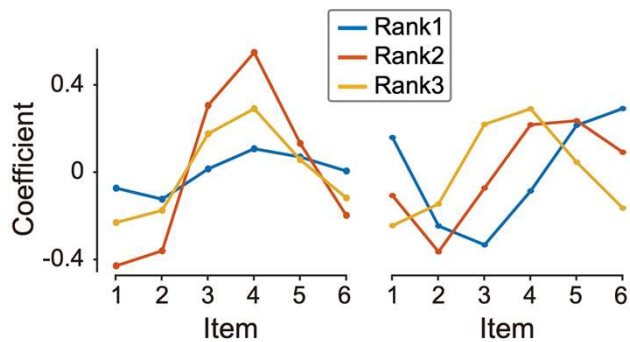


$\varphi_{ri}$ : spatial item preference of neuron  $i$  in rank- $r$  subspace

## Histograms of $\varphi$ difference for different rank pairs



## The spatial location preference (tuning) is determined by the angular ( $\varphi_{ri}$ )

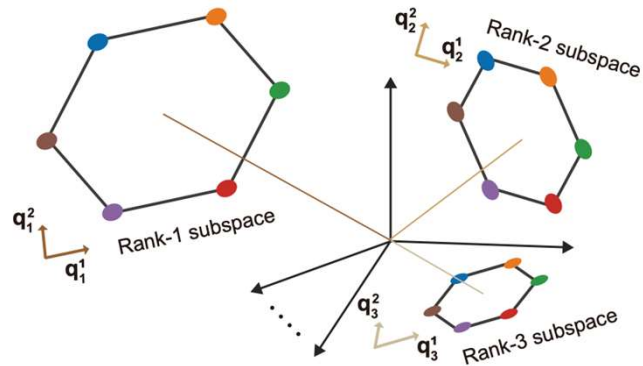


This finding is important because it rejects a simple scalar model of **gain modulation** or **gain field**

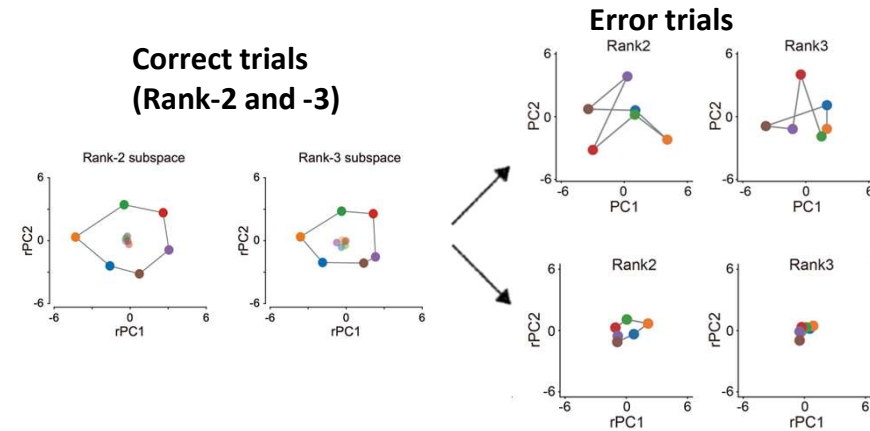
Botvinick, M., & Watanabe, T. (2007). From Numerosity to Ordinal Rank: A Gain-Field Model of Serial Order Representation in Cortical Working Memory. *The Journal of Neuroscience*, 27(32), 8636-8642. <https://doi.org/10.1523/JNEUROSCI.2110-07.2007>

Rather, for each ordinal position, the same spatial information is sent into a completely different direction (hyperplane) in neural space.

# From neurons to behavior



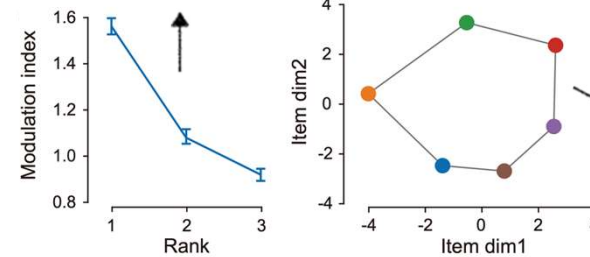
The working memory representation is degraded on error trials:



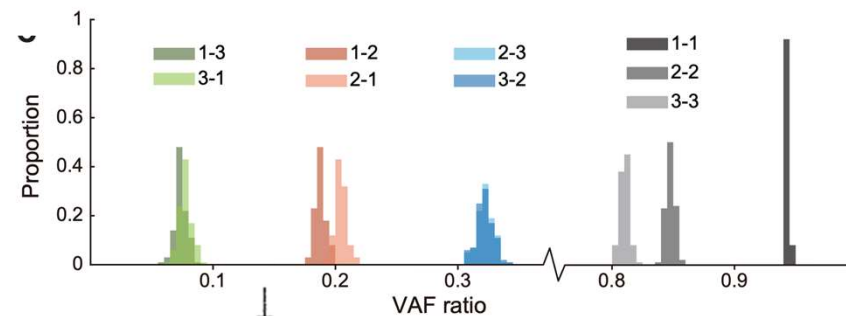
This neural representation can explain several aspects of behavior :

1. Length effect (memory capacity)
2. Primacy and recency effect
3. Error patterns
  - transposition gradients
  - fill-in; intrusions; omission
4. Item/temporal similarity/interference effect

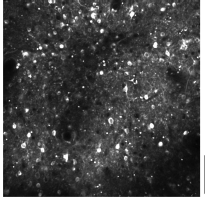
Length/ Primacy effect



Transposition gradient in the spatial location



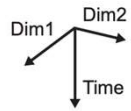
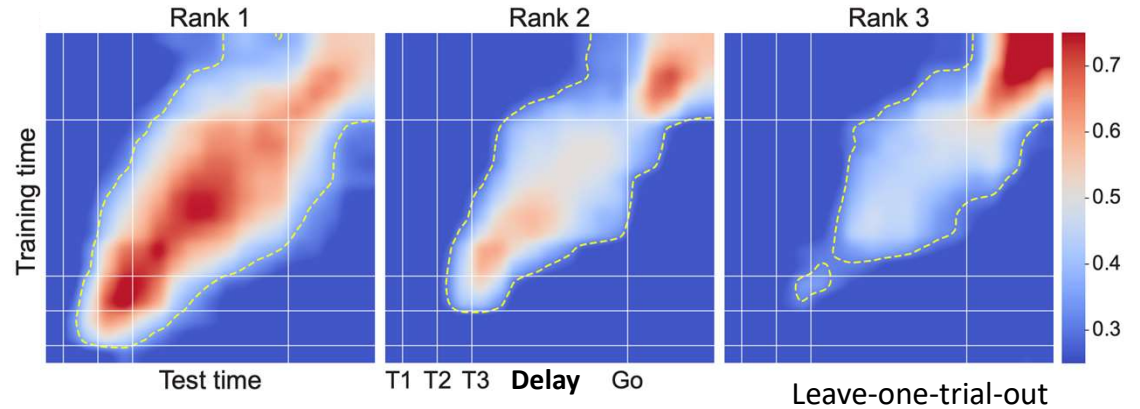
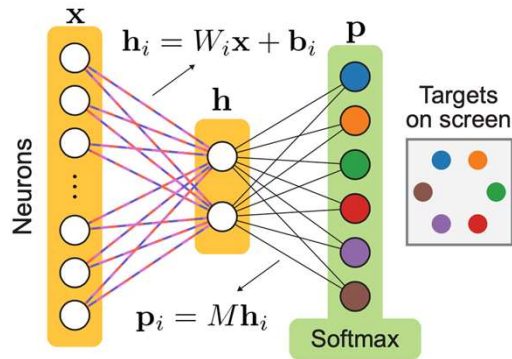
Transposition gradient in the temporal order



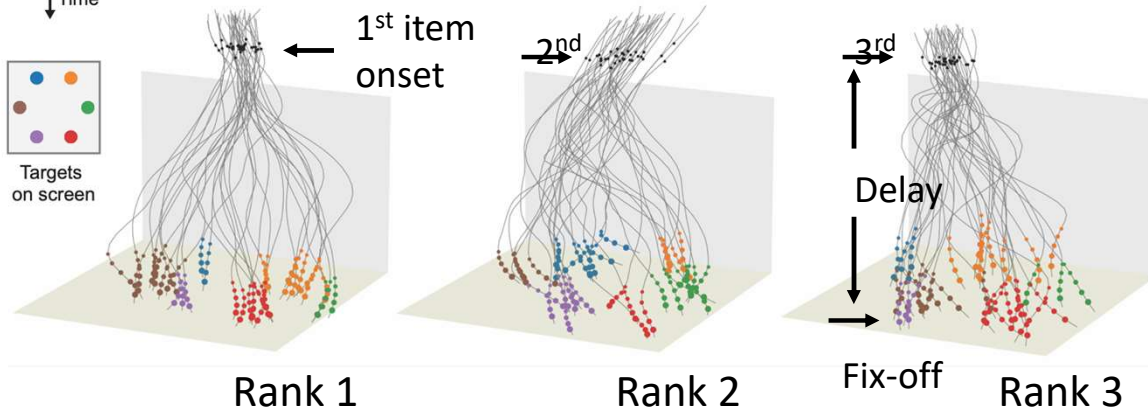
N~150 (example FOV)

# Disentangled representation of sequence memory at the **single-trial** level

Cross temporal decoding of spatial location at each rank

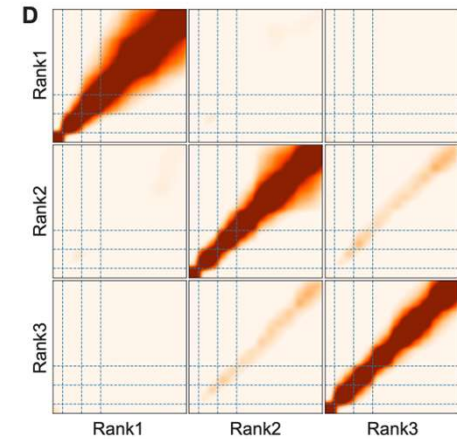


Sequence trajectories in three rank subspaces



Decoder training: data points during the late delay period

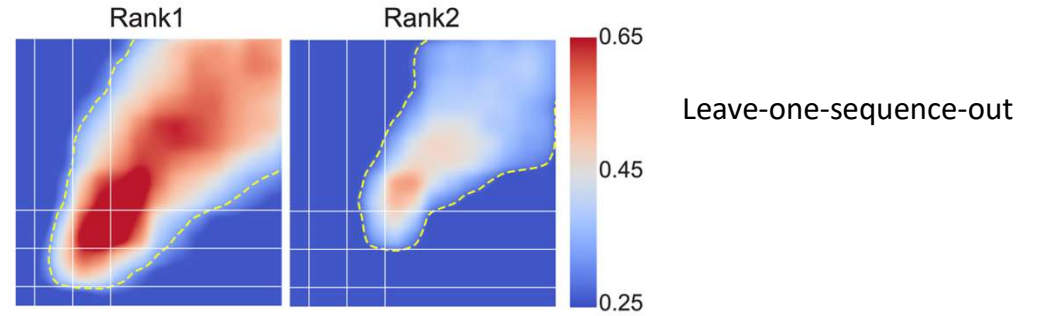
Independent rank subspaces



# Compositional **generalization** test of sequence representation

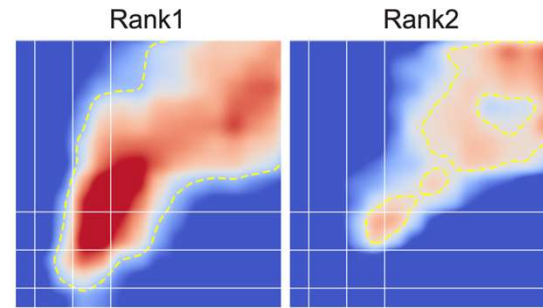
- Generalization between **sequences**

	_____
Train	3 2 5
	2 6 3
	4 1 5
	...
Test	...
	3 6 5



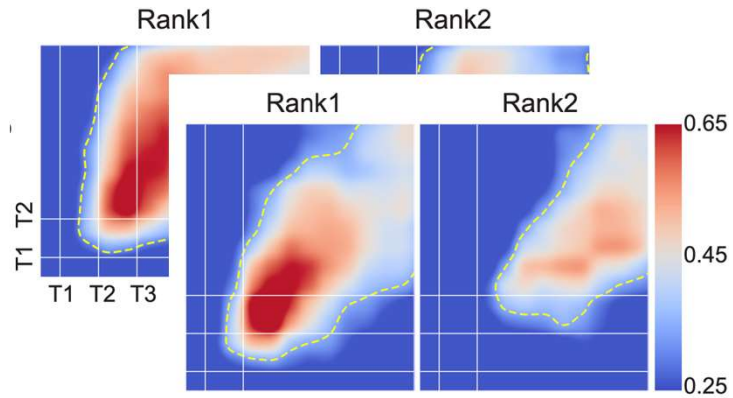
- Generalization across correct and **incorrect** trials

	_____
	→ 4 1 5
3 1 5	→ 3 2 5



- Generalization between sequences of **different length**

_____	→	_____
1 2 3	→	1 2
_____	→	_____
1 2	→	1 2 3

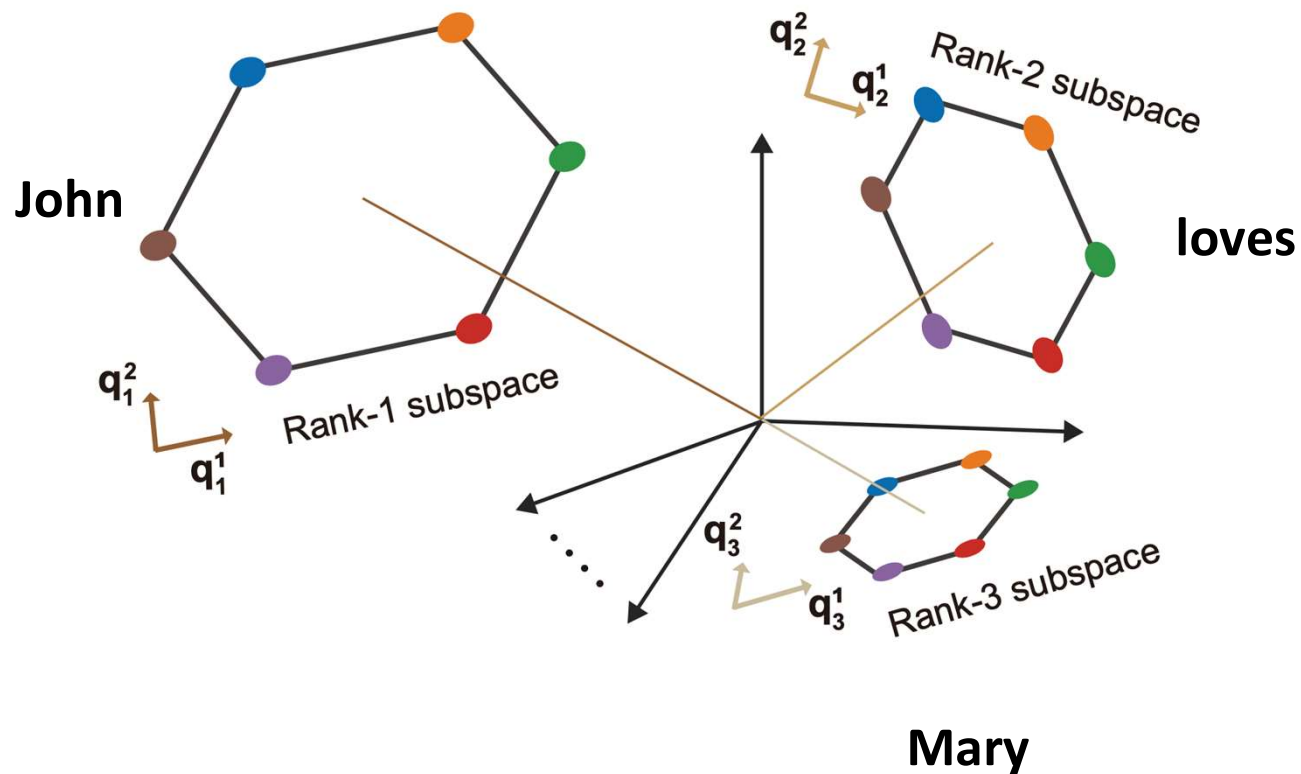


## Could we generalize this idea to syntactic structures?

Suppose we want to encode the sentence “John loves Mary”. Could we just use successive slots ?

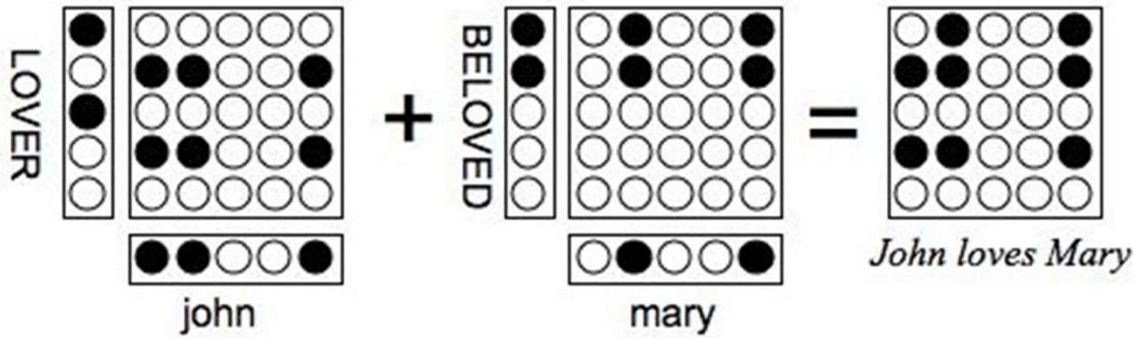
We would need to encode, not the ordinal number of each word, but their structural role (subject = John, verb = loves, etc).

And we would need to find a mechanism for recursion (“John, who is a rich banker, loves Mary”).

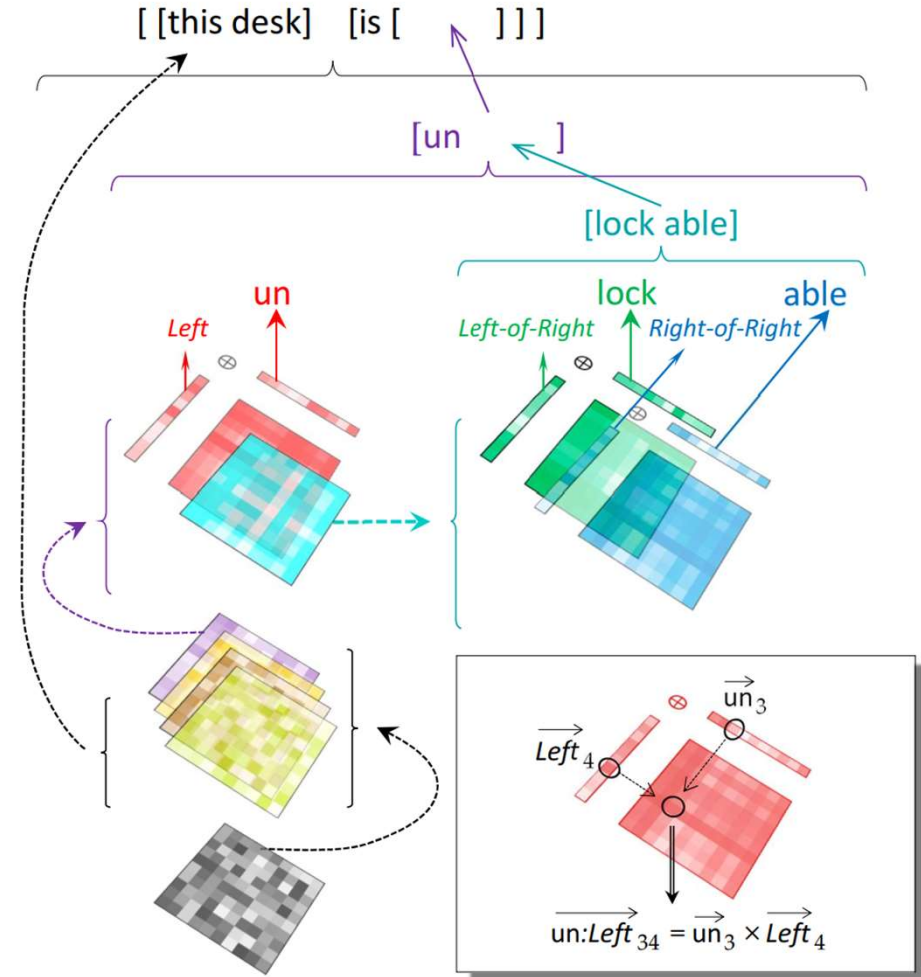
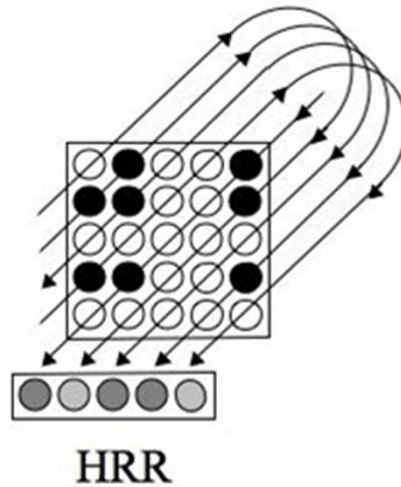




# A hypothesis about the neural code for sentences



Plate, T. A. (1995). Holographic reduced representations. *IEEE Transactions on Neural Networks / a Publication of the IEEE Neural Networks Council*, 6(3), 623-641.  
<https://doi.org/10.1109/72.377968>



Smolensky, P., McCoy, R. T., Fernandez, R., Goldrick, M., & Gao, J. (2022). *Neurocompositional computing in human and machine intelligence : A tutorial*

# Ordinality is insufficient for language: Many arguments in favor of nested tree structures in language

Haegeman, L. (2005). *Thinking Syntactically: A Guide to Argumentation and Analysis*. Wiley.

- **Long-distance dependencies** (agreement and binding):  
« The cars that pass this truck are red »
- Cases of **syntactic ambiguity**:

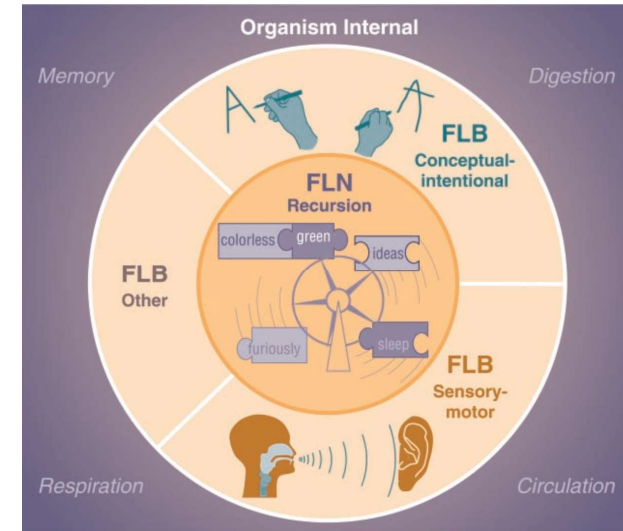
(Black taxi) driver

I shot an elephant in my pajamas...

How it got in my pajamas, I don't know (Groucho Marx)

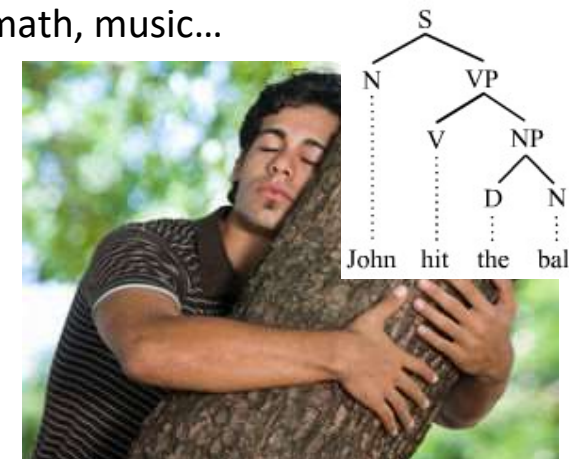
unlockable = un-(lock-able) or (un-lock)-able

- **Ellipsis or substitution** of any phrase :  
“he [drove [to [this [big house]]]”  
= “he drove to this one,” “he drove to it,” “he drove there,” “he did.”
- « **Syntactic Movement** » of phrases (for question formation, topicalization, etc):  
« John loves that dog » → « It's that dog that John loves \_ »

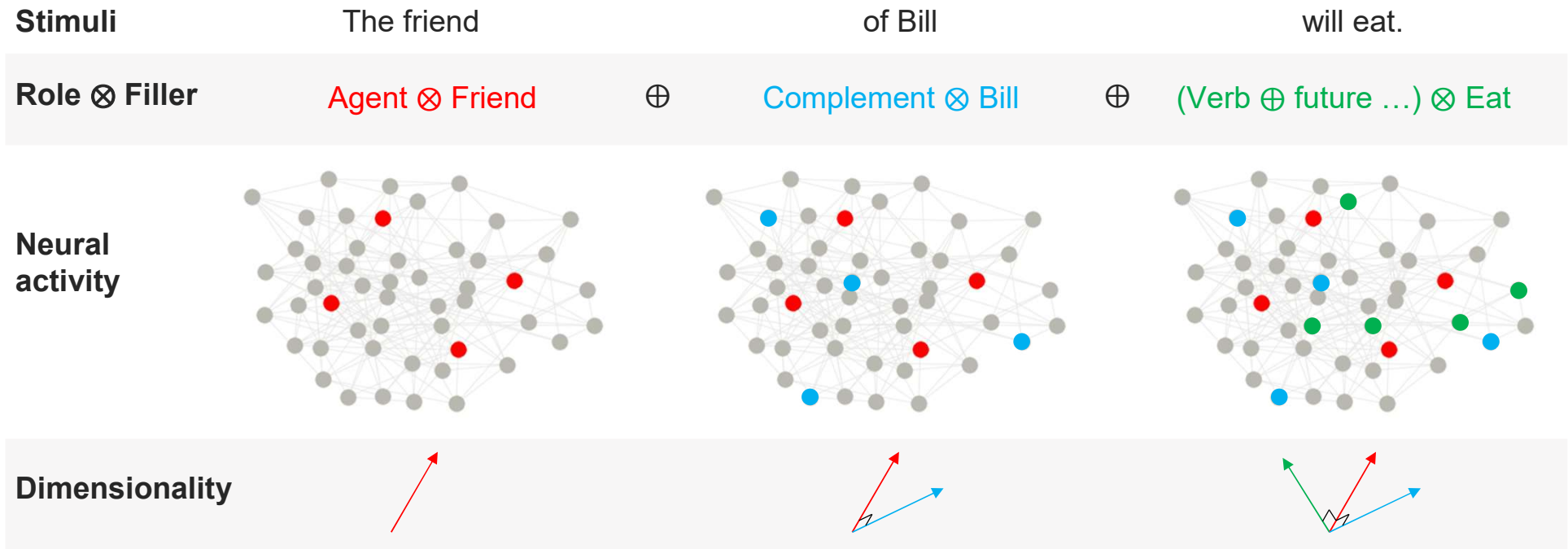


Hauser, Chomsky and Fitch (2002):  
Language rests on a recursive tree-building operation (**merge**).

Tecumeh Fitch's “dendrophilia”:  
similar tree structures are also used  
in math, music...



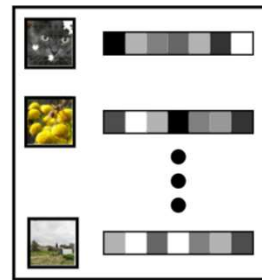
# A hypothesis about the neural code for sentences



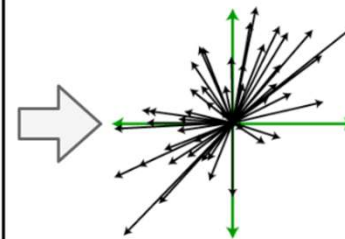
Predictions :

- Ramp-up : Activity should increase with each successive word or phrase that is integrated in the sentence structure
- Dimensionality : Across a set of sentences, the activity should occupy an increasing number of dimensions as more words are integrated in the structure.

Neural activity evoked by a number of stimuli



PCA



Effective dimensionality

$$\frac{(\sum_{i=1}^n \lambda_i)^2}{\sum_{i=1}^n \lambda_i^2}$$

## Recent results from Théo Desbordes' Phd (with Jean-Rémi King)



"Dimensionality and ramping: Signatures of sentence integration in the dynamics of brains and deep language models"

by Theo Desbordes, Yair Lakretz, Valerie Chanoine, Maxime Oquab, Badier Jean-Michel, Agnès Trébuchon, Carron Romain, Christian Benar, Stanislas Dehaene, and Jean-Rémi King.

Journal of Neuroscience, in press (2023)

# Jabberwocky: how to build a meaningless sentence?

English version

*'Twas brillig, and the slithy toves  
Did gyre and gimble in the wabe:  
All mimsy were the borogoves,  
And the mome raths outgrabe.*

*“Beware the Jabberwock, my son!”*



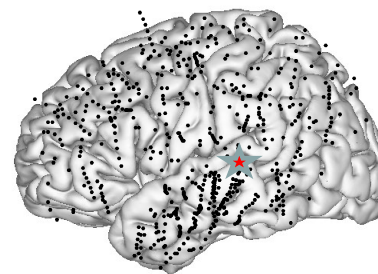
Version française

*Il était grilheure ; les slictueux toves  
Sur l'alloinde gyraient et vriblaient ;  
Tout flivoreux étaient les borogoves  
Les vergons fourgus bourniflaient.*

*“Prends garde au Jabberwock, mon fils!”*

Goal = compare brain responses to normal sentences and Jabberwocky  
in order to isolate semantic processing

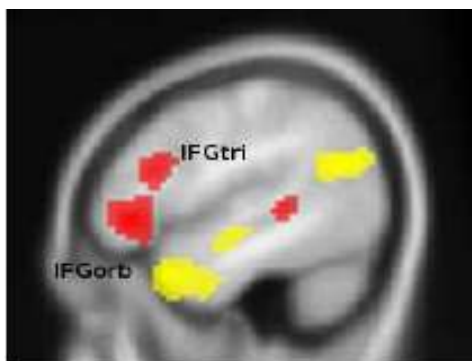
# Ramping brain signals track phrase structures



fMRI:

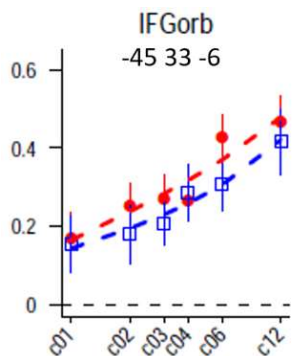
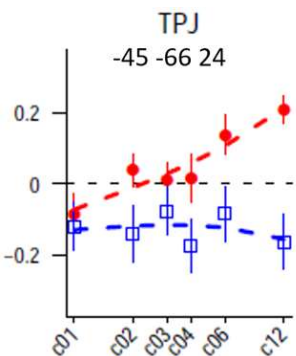
Monotonic increase  
with constituent size

(Pallier, Devauchelle & Dehaene, 2011)

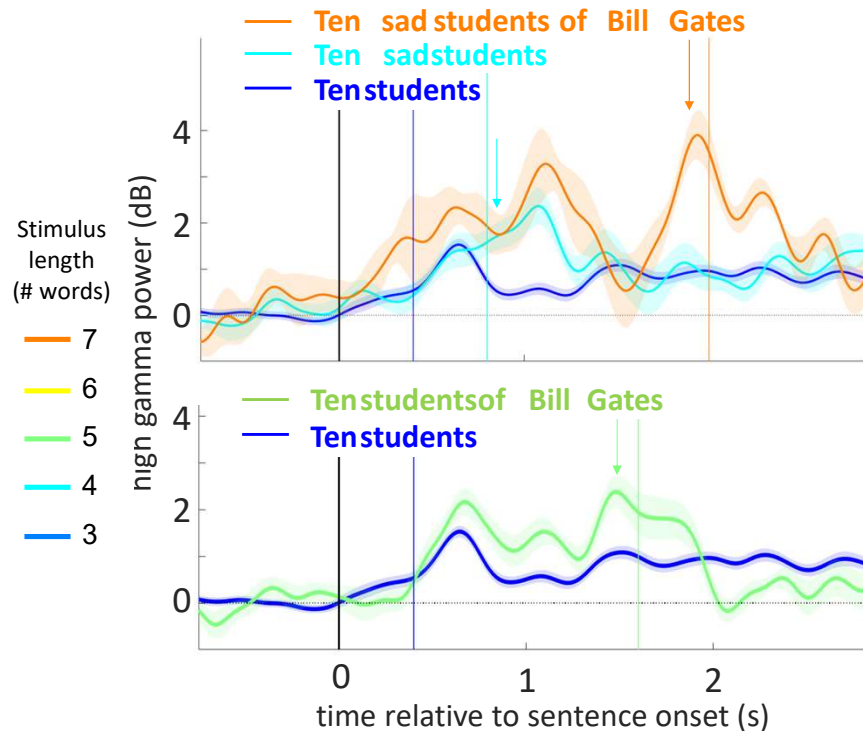
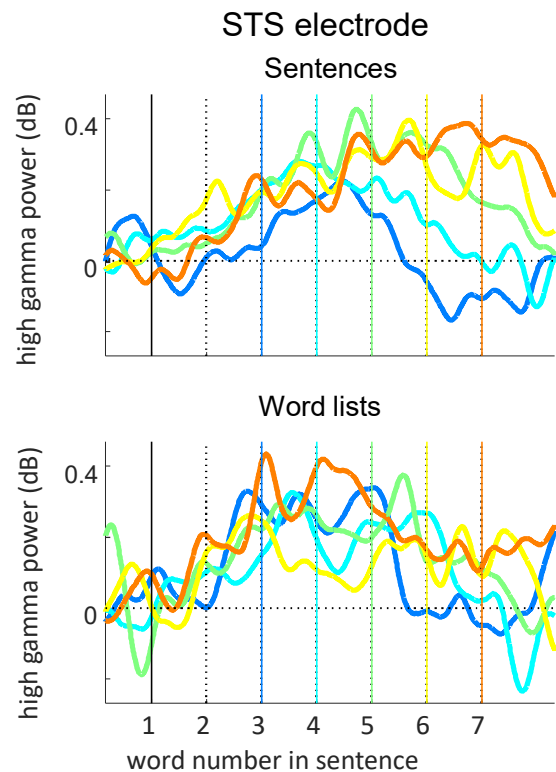


Areas in yellow  
increase only for  
sentences

Areas in red increase for  
both sentences and  
Jabberwocky



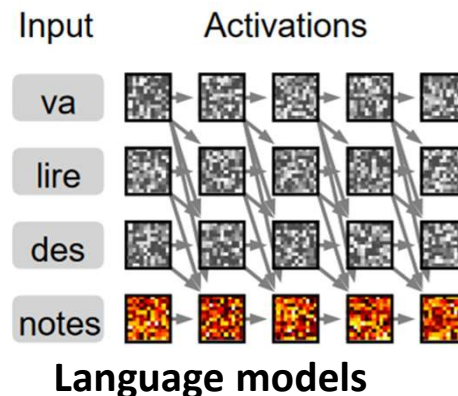
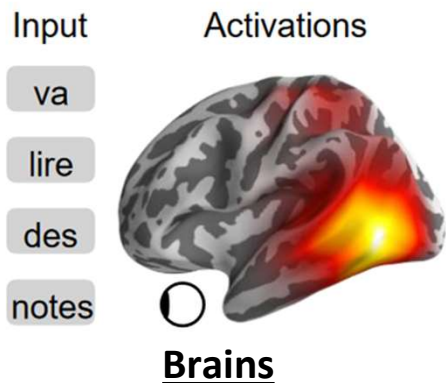
Intracranial recordings:  
Monotonic increase with sentence length,  
and tracking of constituent size  
(Nelson... and Dehaene, PNAS, 2017)



## Stimuli: closely matched normal and Jabberwocky sentences

<u>Normal text</u>	####	une	cousine	du	curé	va	lire	des	notes	####
<u>Jabberwocky</u>	####	une	trule	du	lenriste	a	ju	la	birle	####
<u>Stringlist</u>	####	klr	dfrtnq	rf	qcbnsd	p	qrtn	lks	bmqcp	####

# Joint MEG and intracranial EEG and Activations from Neural Language Models



11 patients



Intracranial EEG

*Badier et al 2017*



MEG



Home-trained:

- character-based LSTM
- **character-based causal Transformer**



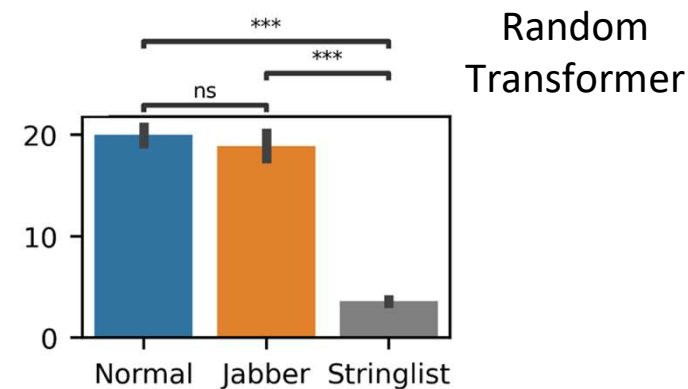
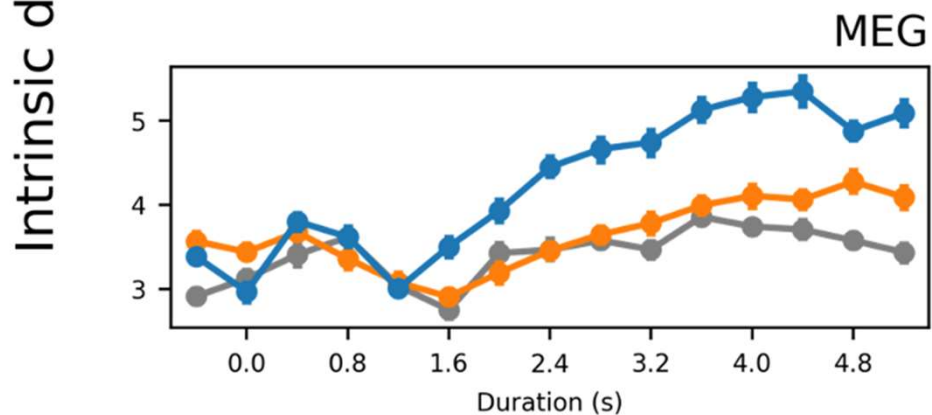
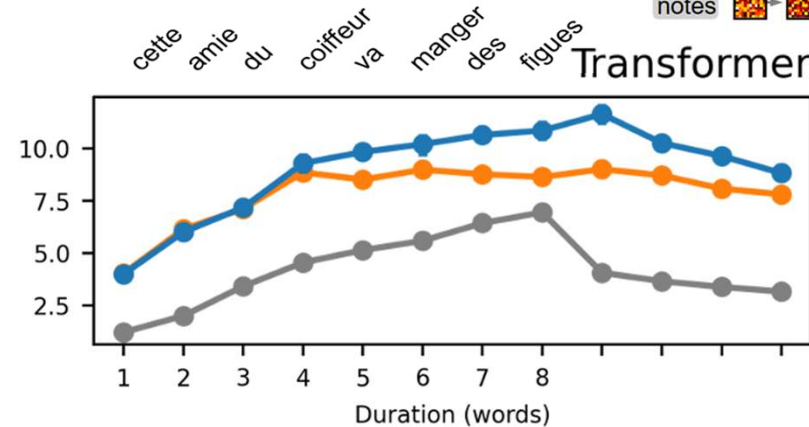
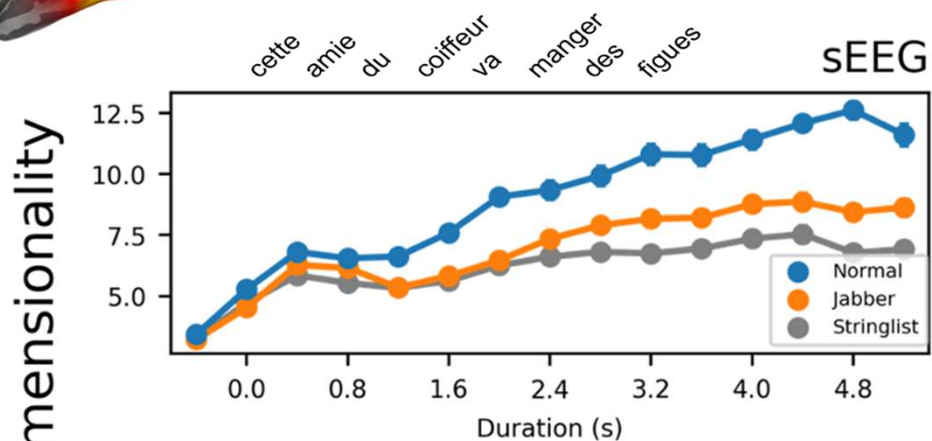
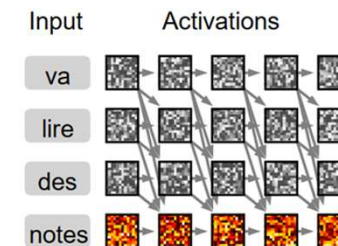
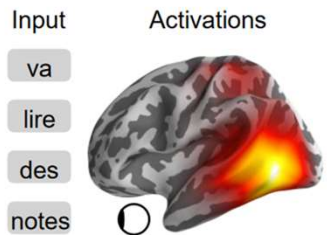
Pretrained bidirectional transformer



CamemBERT  
A Tasty French Language Model  
*Martin et al 2020*



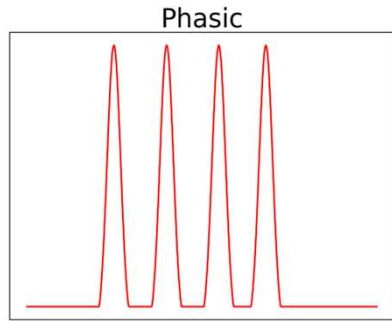
# Intrinsic dimensionality grows along the sentence and is higher for normal sentences than for Jabberwocky



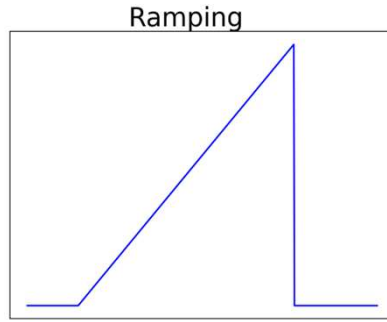
# Decoding normal versus Jabberwocky sentences :

## Predicted patterns of generalization depending on the type of cognitive operation involved

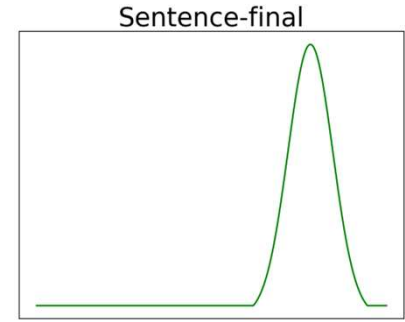
### Lexical access



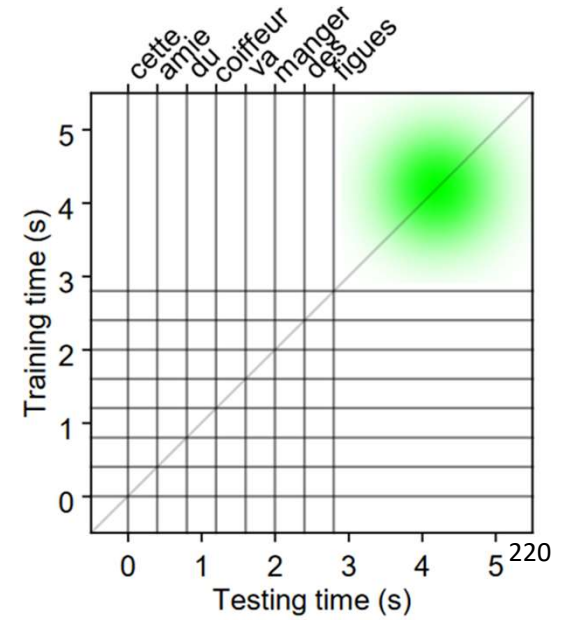
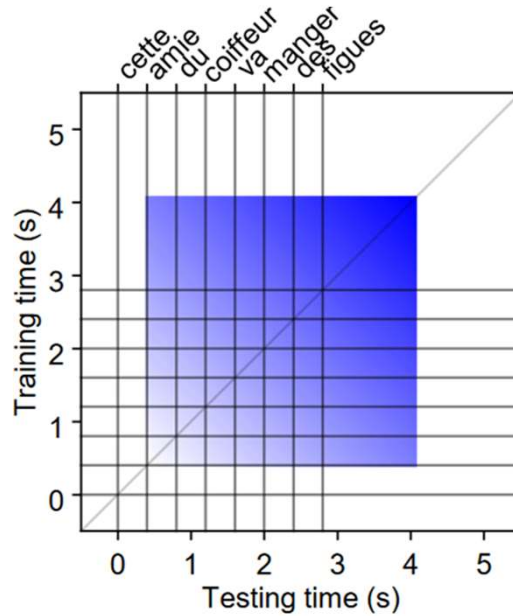
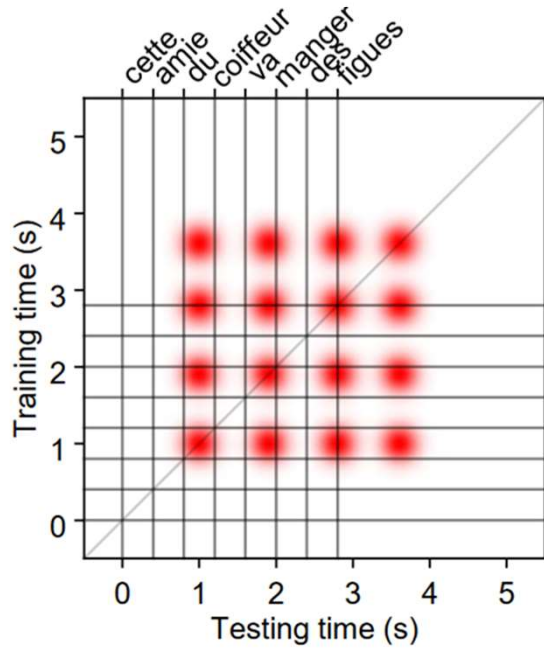
### Multi-word integration



### Wrap-up

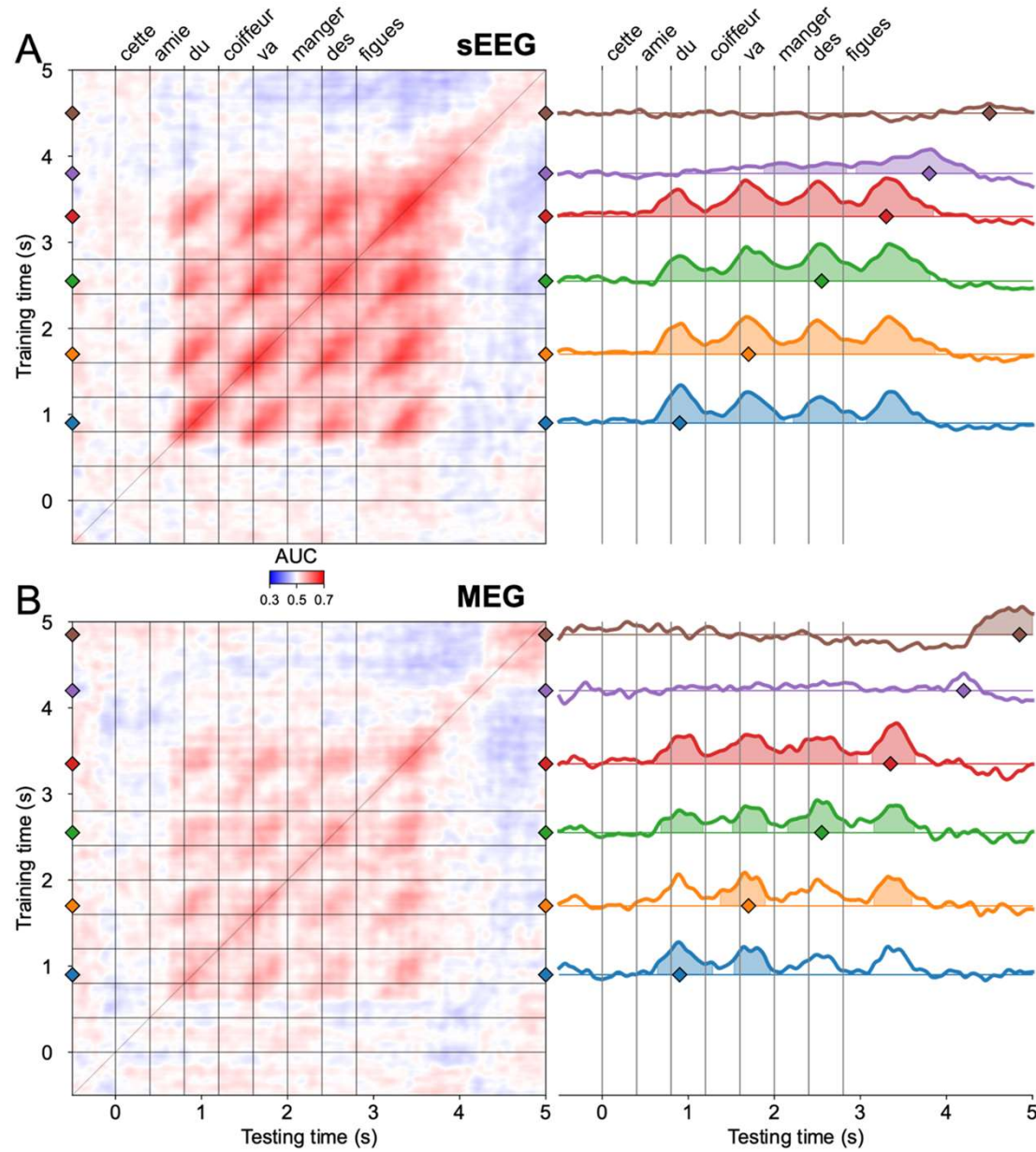


Diagonal

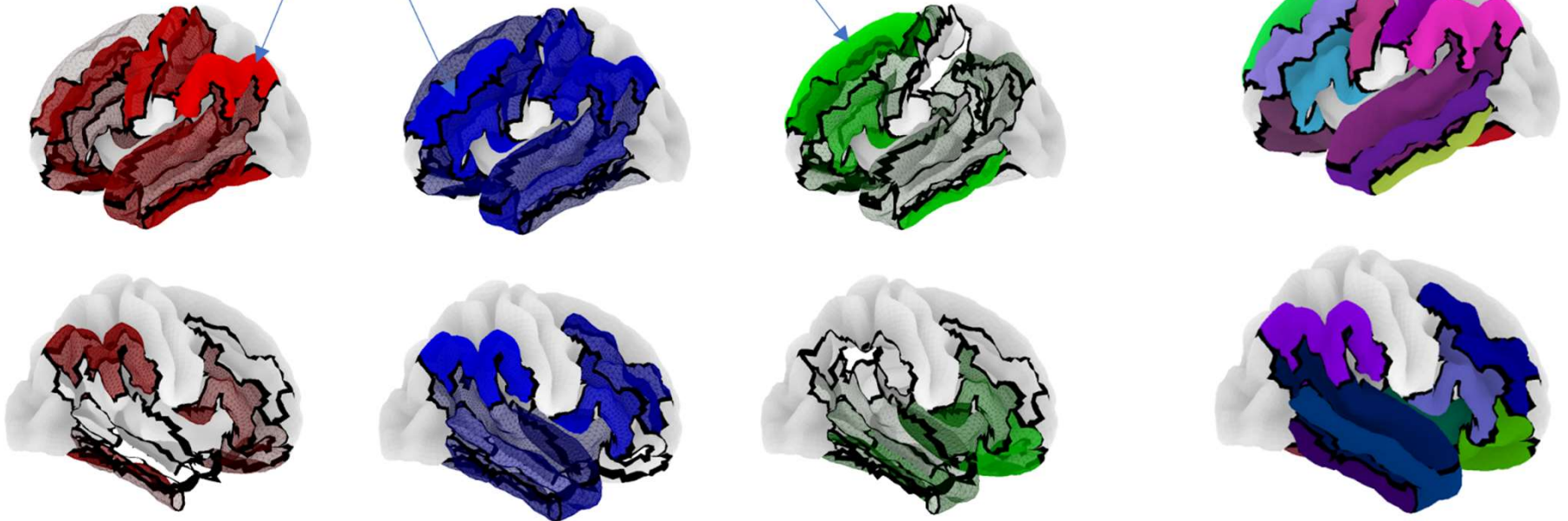
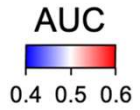
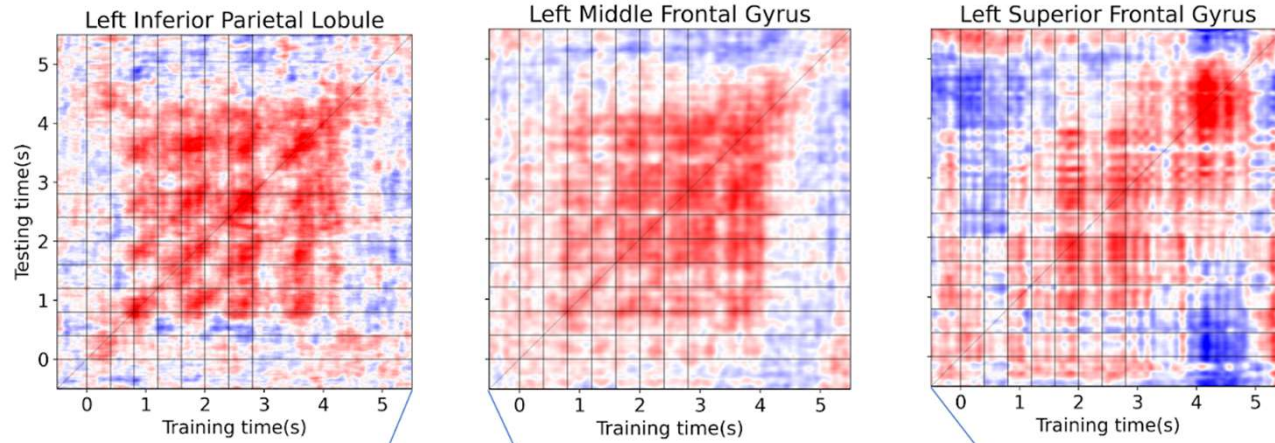


220

# Decoding normal sentences versus Jabberwocky in humans



# Results: partially different dynamics in each brain region



# Beyond ramping: Structural probe showing compositional representations in networks

Hewitt, J., & Manning, C. D. (2019). A Structural Probe for Finding Syntax in Word Representations. NAACL 2019, 4129-4138.

Idea = see if, within the huge representational space that artificial networks use to represent sentences, there is a linear subspace in which Euclidean distances reflect distances in the syntactic tree.

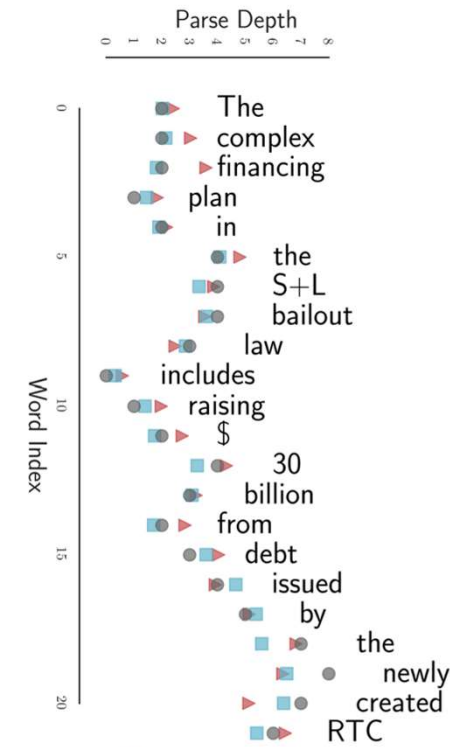
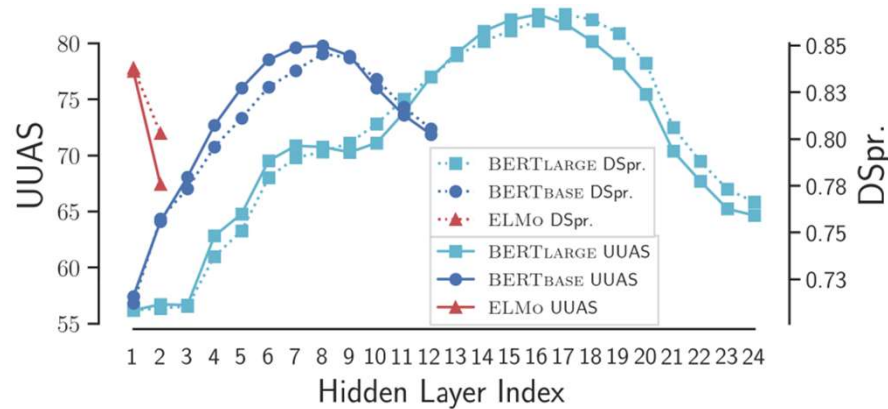
In particular, the vector norm indicates the distance to the root (parse depth):

For instance, in « The cars that pass this truck are red », *cars* and *red* should be closer than *truck* and *red*.

Formally : find a linear transform B such that the distance between vectors for words (squared)

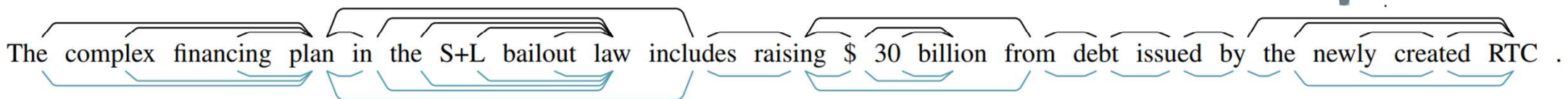
$$d_B(\mathbf{h}_i^\ell, \mathbf{h}_j^\ell)^2 = (B(\mathbf{h}_i^\ell - \mathbf{h}_j^\ell))^T (B(\mathbf{h}_i^\ell - \mathbf{h}_j^\ell)) \text{ is tightly correlated with distance on the syntactic tree.}$$

Results: for several different models, it is possible to find such a « structural probe », particularly in the intermediate layers, such that distances in this subspace tightly reflect those in the syntactic tree.



The “minimum spanning tree” reconstructed from those distances resembles the parse of the sentence:

**BERTlarge16**



## Conclusions

Words can be represented as vectors in a high-dimensional semantic space.

Major semantic dimensions are attributed to different regions, and therefore to different neural populations.

A series of items can be encoded in working memory by attaching each item to a vector and then **rotating** this vector to a different neural subspace (one for each ordinal memory slot).

Can this code be extended to **syntactic trees** rather than **ordinal chains** ?

The neural code for language seems to verify minimal assumptions

- Increasing activity (**ramp-up**) with consecutive words
- Increasing **dimensionality**
- And, in artificial neural networks, a subspace represents the **structural (tree-based)** rather than linear (temporal) distances

Vendredi 6 Janvier

COURS : **Vecteurs neuronaux ou cellules grand-mère : les représentations mentales sont-elles localisées ou distribuées ?**  
SÉMINAIRE : **L'intelligence artificielle peut-elle modéliser le langage mathématique ?** – François Charton (FAIR Paris)

Vendredi 13 Janvier

COURS : **Géométrie des représentations visuelles : chaque visage est un vecteur**  
SÉMINAIRE : **Commonsense Physical Reasoning in man and machine** – Ernest Davis (NYU, par zoom)

Vendredi 20 Janvier

COURS: **Exploiter la factorisation et les sous-espaces vectoriels pour coder l'information et communiquer entre aires cérébrales**  
SÉMINAIRE : **Number symbols in the brain and mind** – Daniel Ansari (University of Ontario)

Vendredi 27 Janvier

COURS : **Comment prendre une décision ou faire des calculs avec des vecteurs dynamiques?**  
SÉMINAIRE : **Comment se développent les réseaux cérébraux associés aux concepts mathématiques ?** – Marie Amalric (Université de Trento, Italie)

Vendredi 3 Février

COURS : **La représentation vectorielle des mots et des concepts**  
SÉMINAIRE : **Les succès et les nouveaux défis de l'intelligence artificielle en mathématiques** – Léon Bottou (FAIR, New York)

Vendredi 10 Février

COURS : **La représentation vectorielle du langage : Comment représenter une phrase ?**

Véronique Izard (CNRS, Université Paris Descartes:  
Fondements cognitifs des mathématiques