

# Disentangling Syntax and Semantics in the Brain with Deep Networks



**Charlotte Caucheteux (INRIA/FAIR), Alexandre Gramfort (INRIA), Jean-Remi King (FAIR/ENS)**

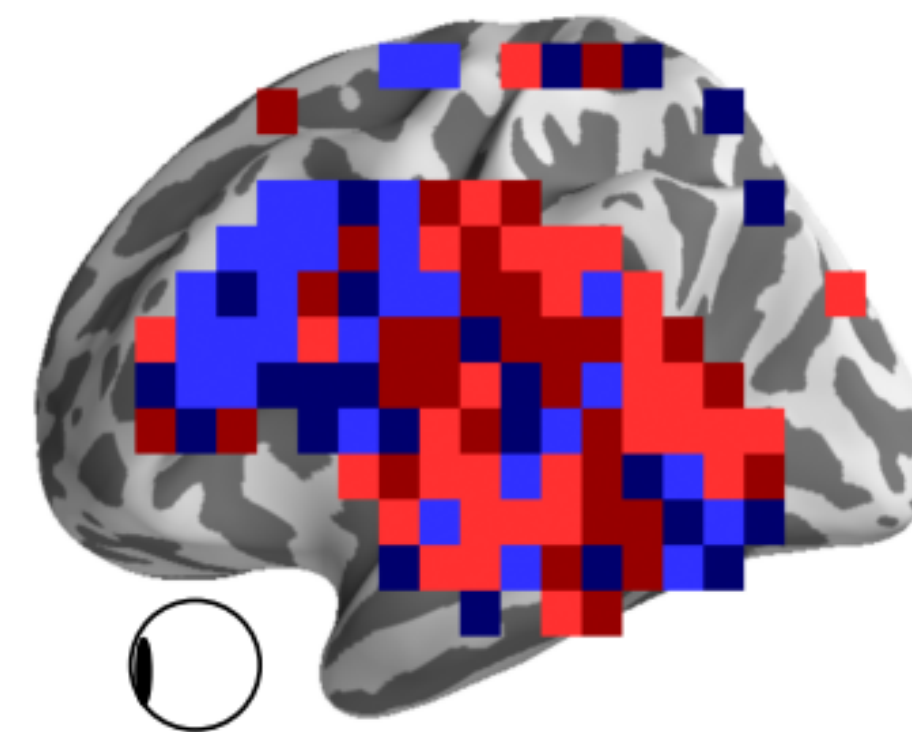
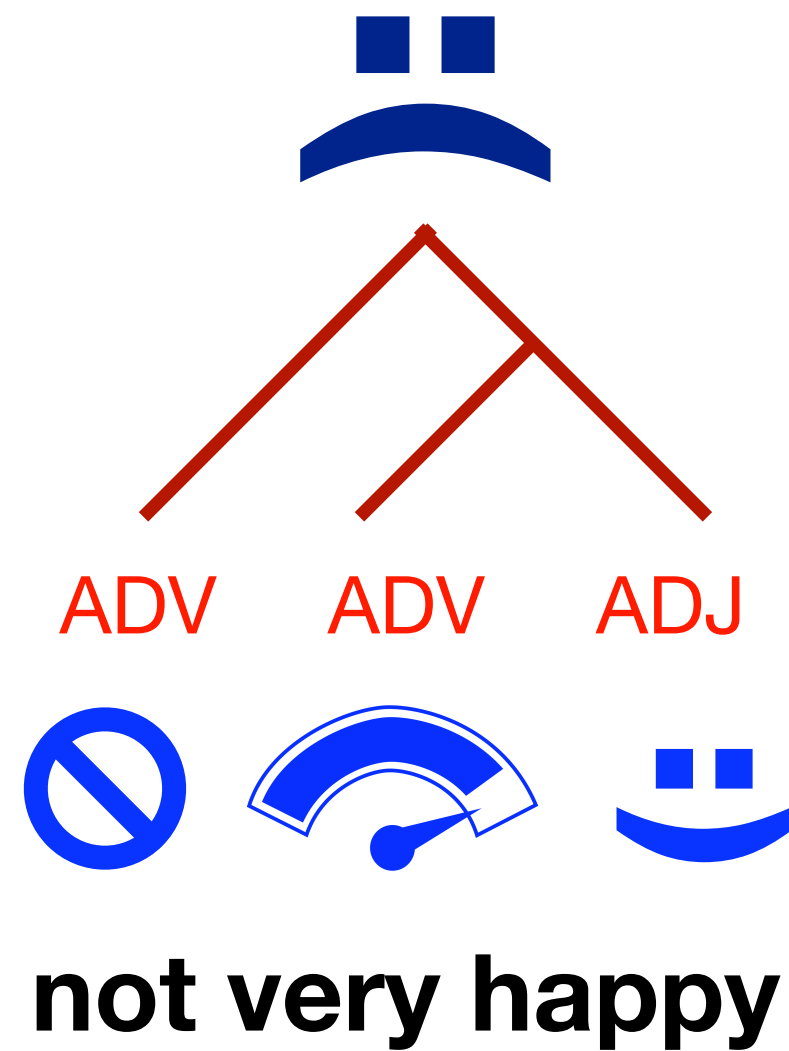
# What are the neural bases of syntax and semantics?

Compositional  
**meaning**

Compositional  
**syntax**

Lexical **syntax**

Lexical **meaning**



not very happy (🔊)

Semantic  
representations

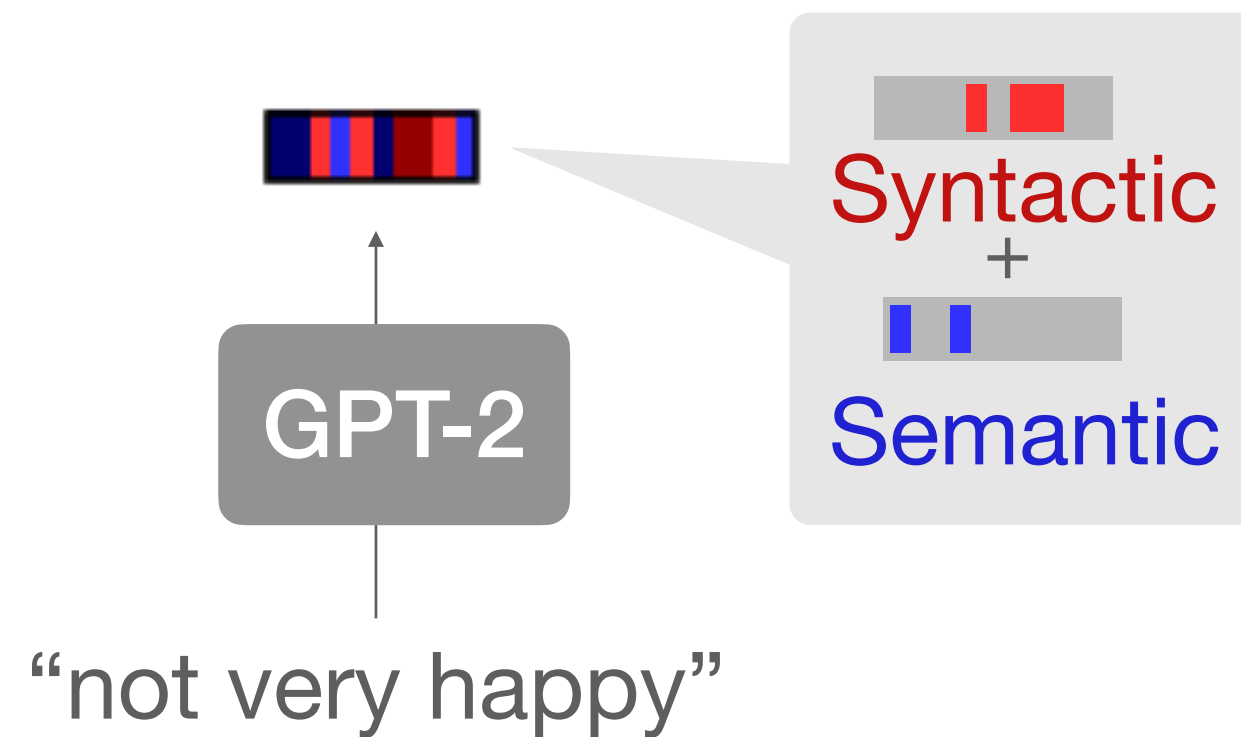


Syntactic  
representations



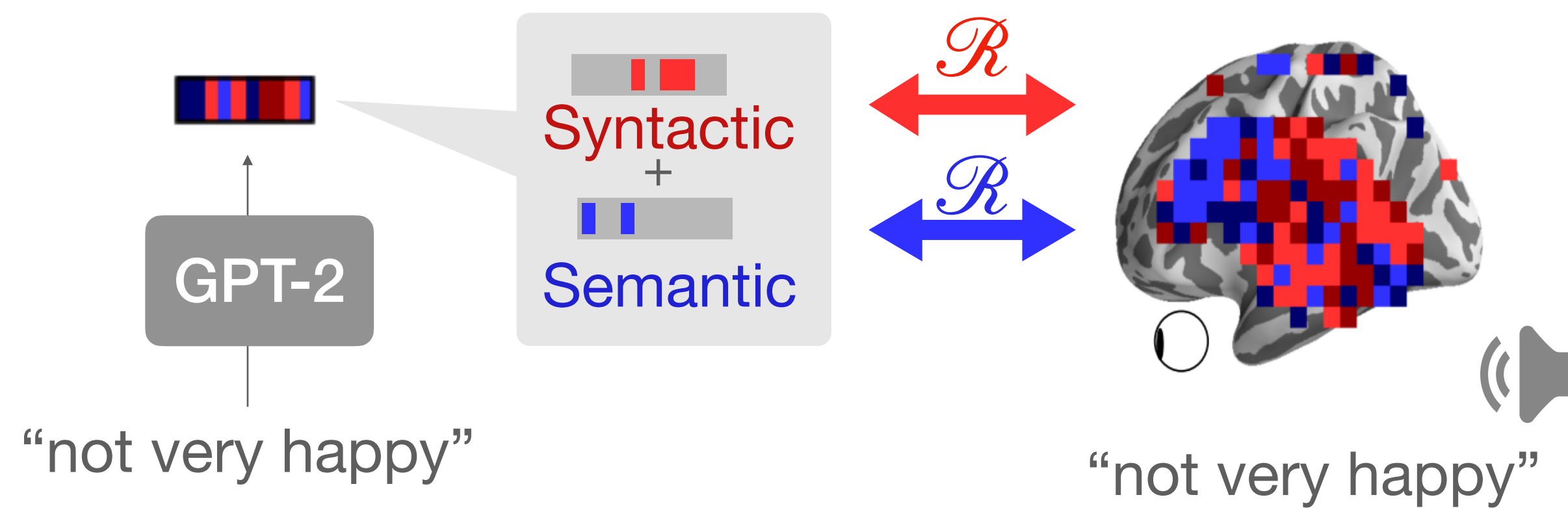
# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations
2. We use the disentangled activations to decompose language in the brain
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*



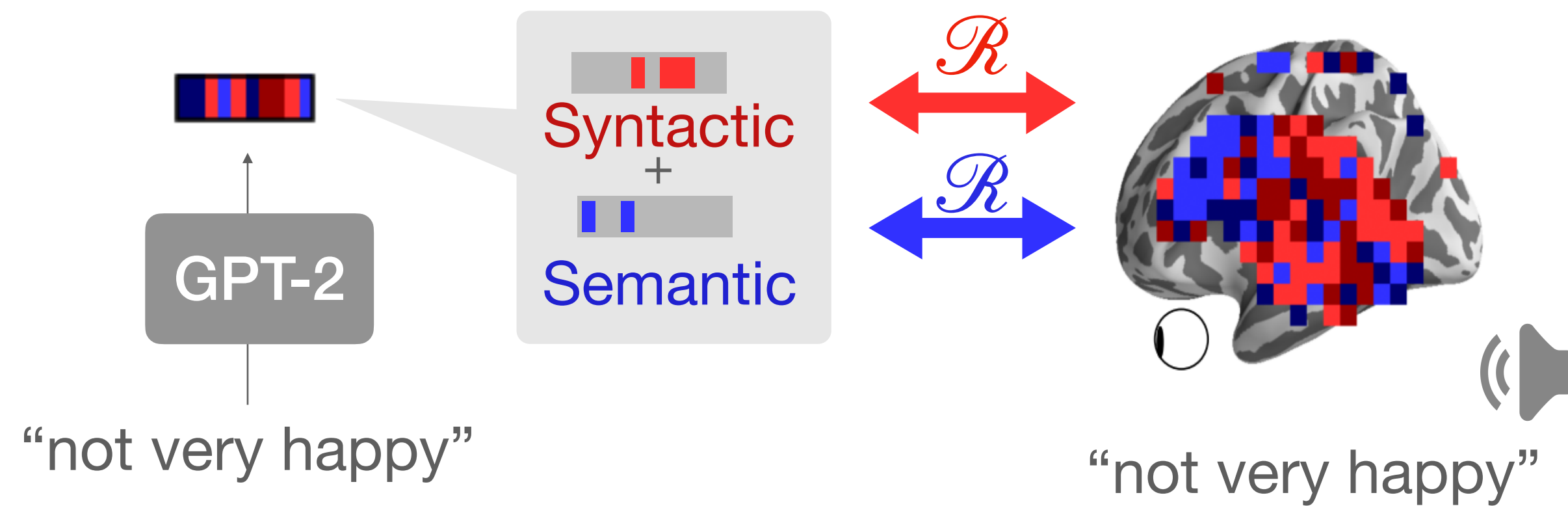
# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations
2. We use the disentangled activations to decompose language in the brain
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

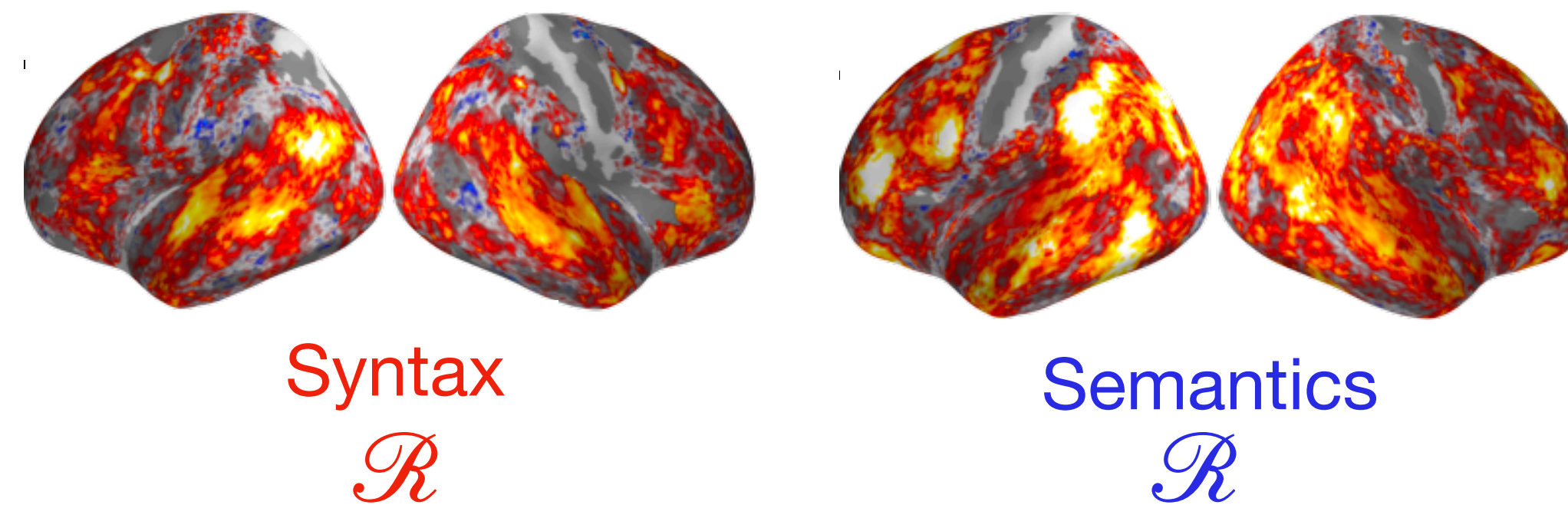


# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations



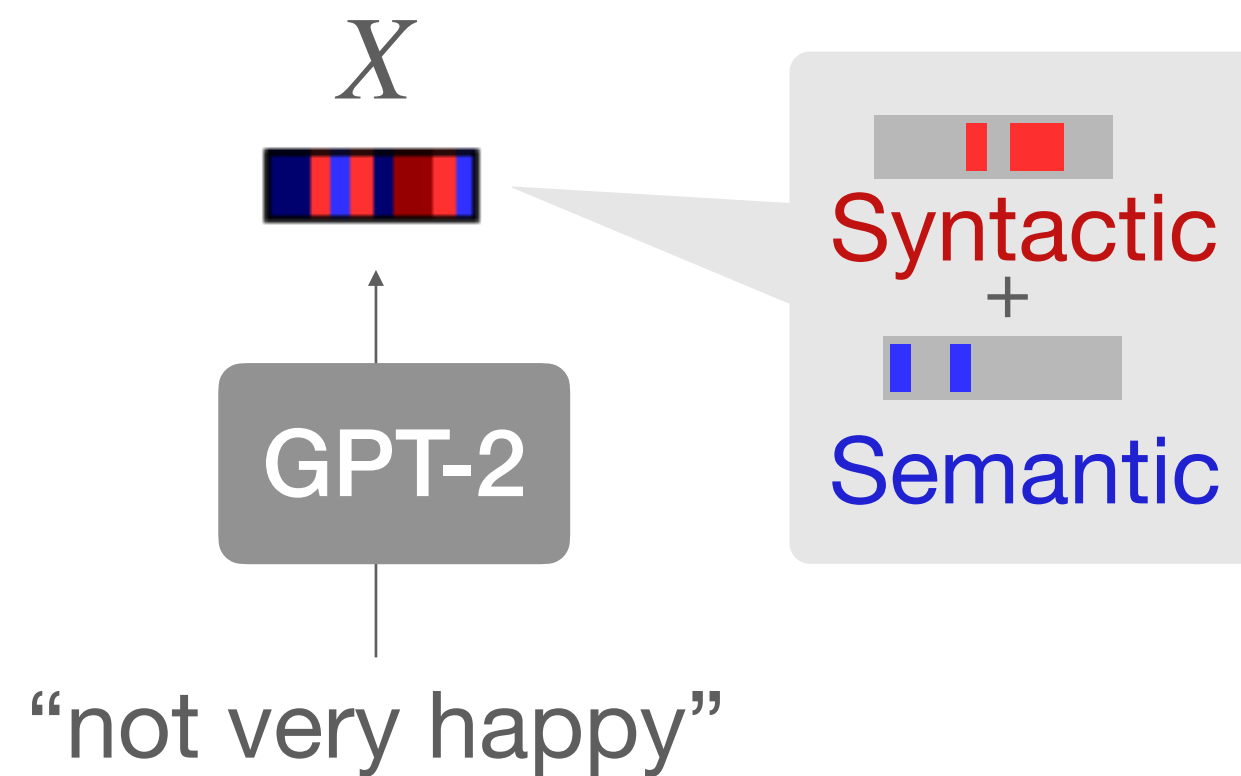
2. We use the disentangled activations to decompose language in the brain



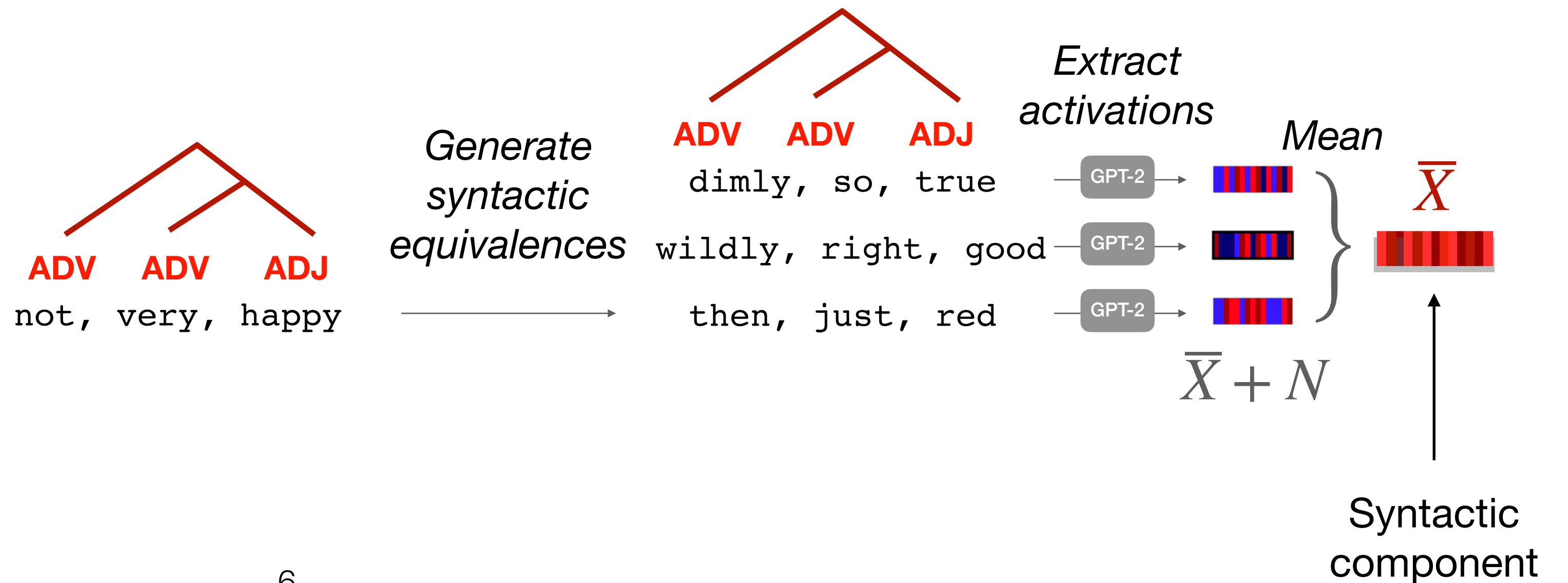
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle **syntax** and **semantics** in **deep nets' activations**

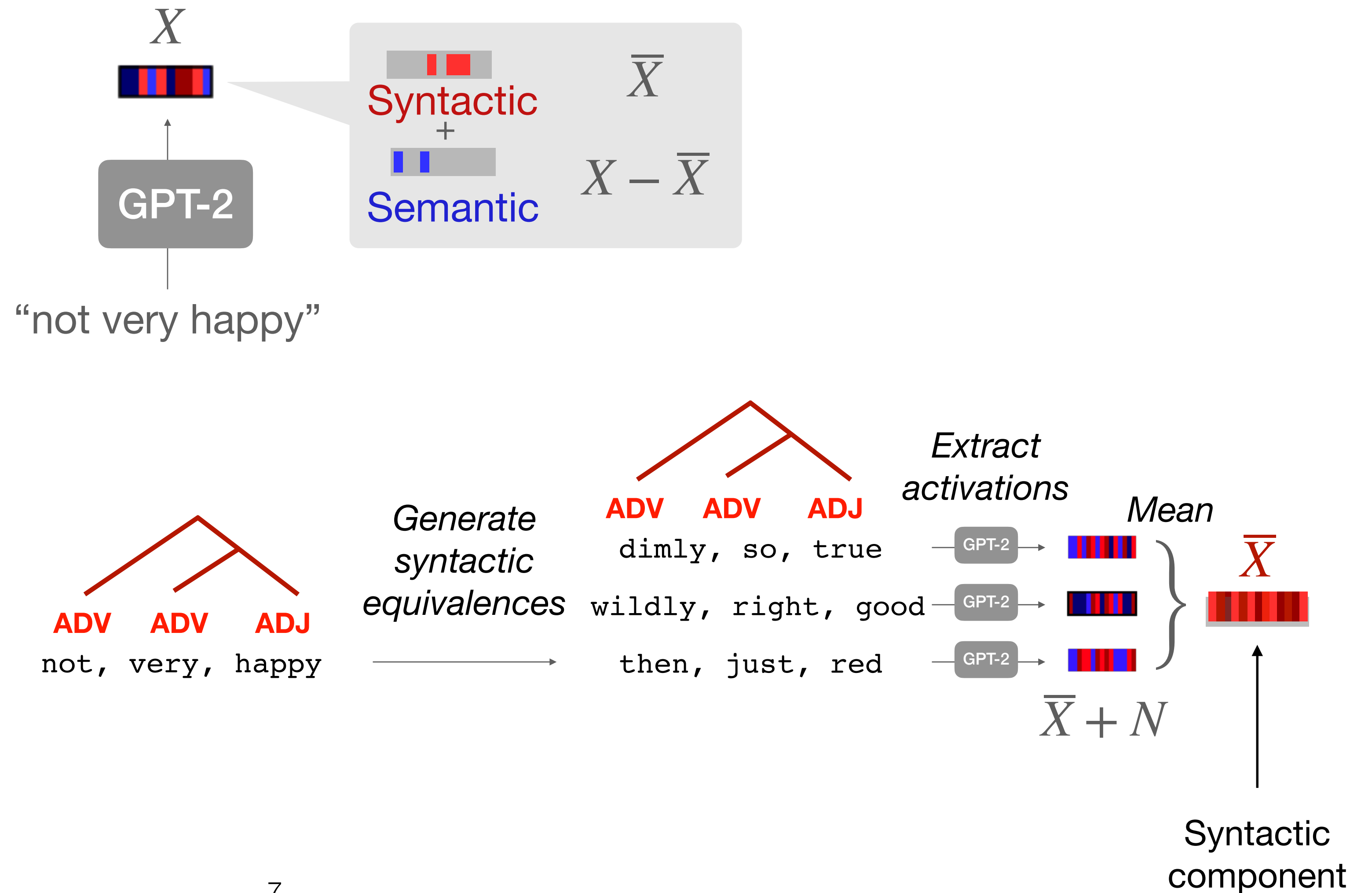
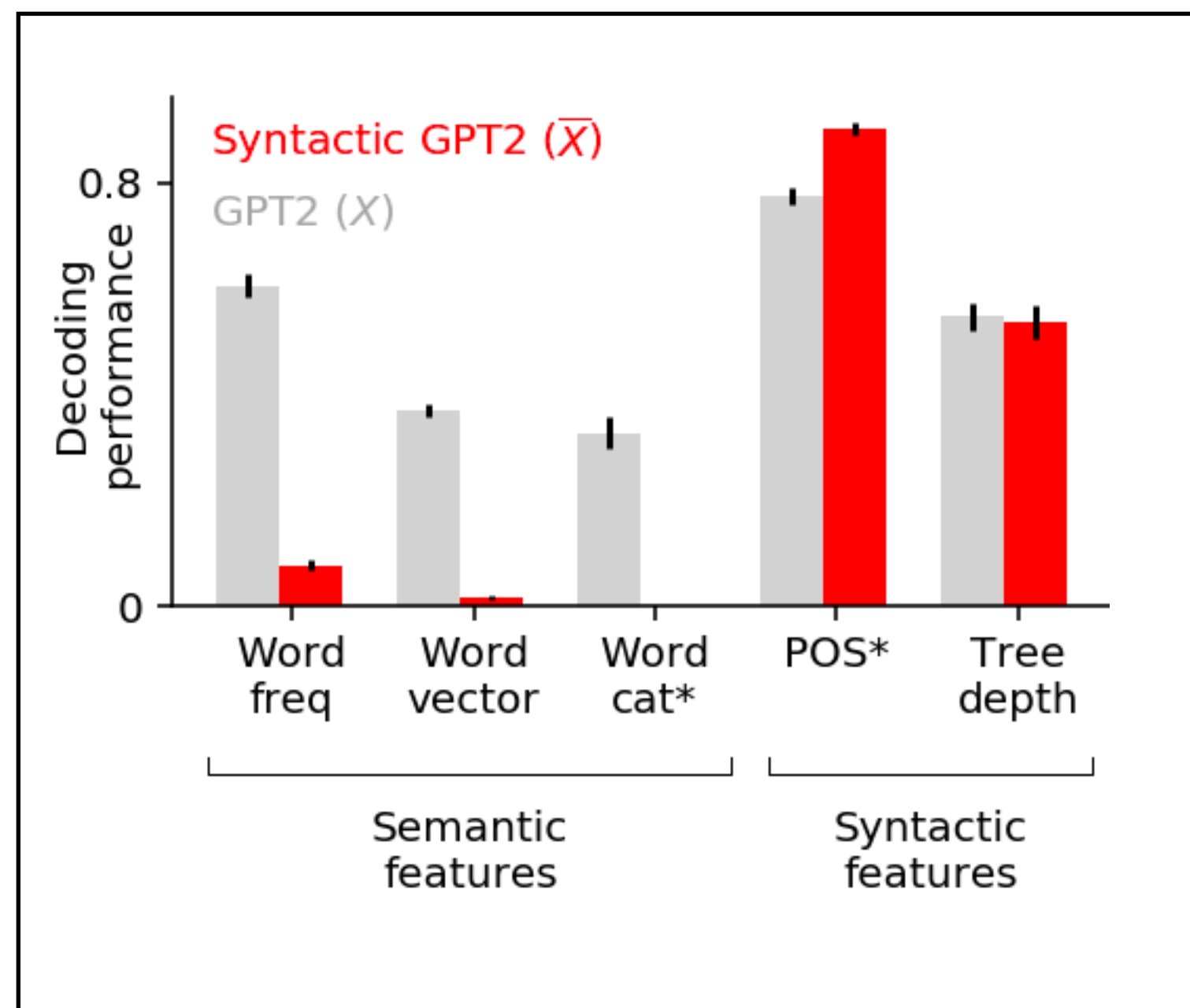


2. We use the disentangled activations to decompose language in the brain



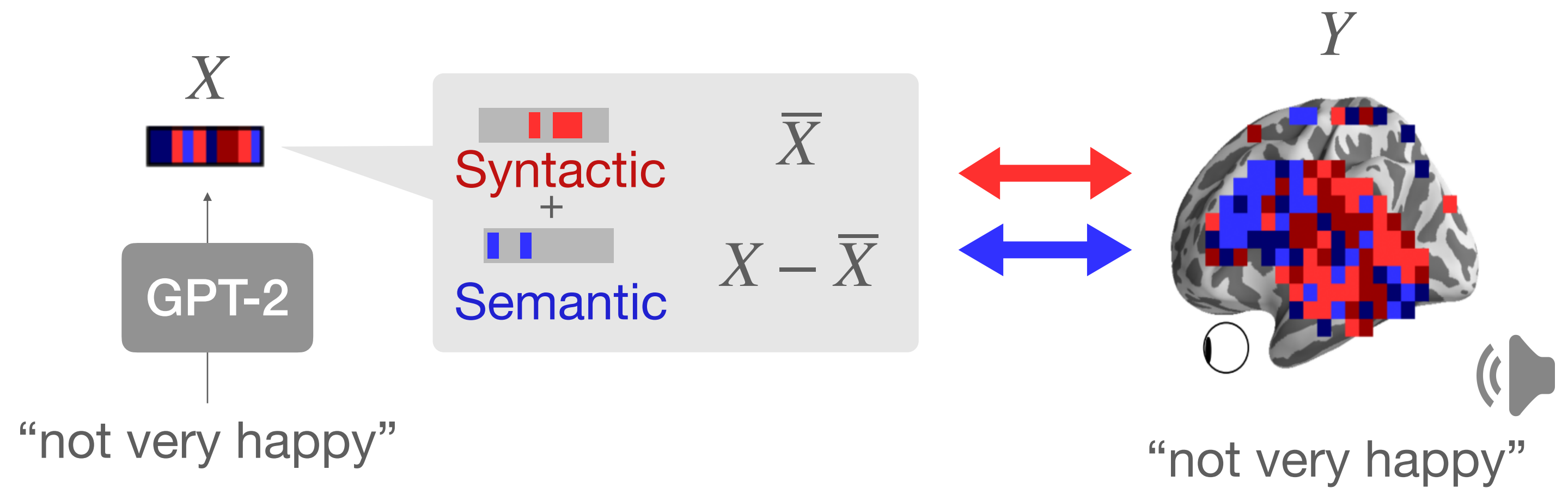
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

# What are the neural bases of syntax and semantics?



# What are the neural bases of syntax and semantics?

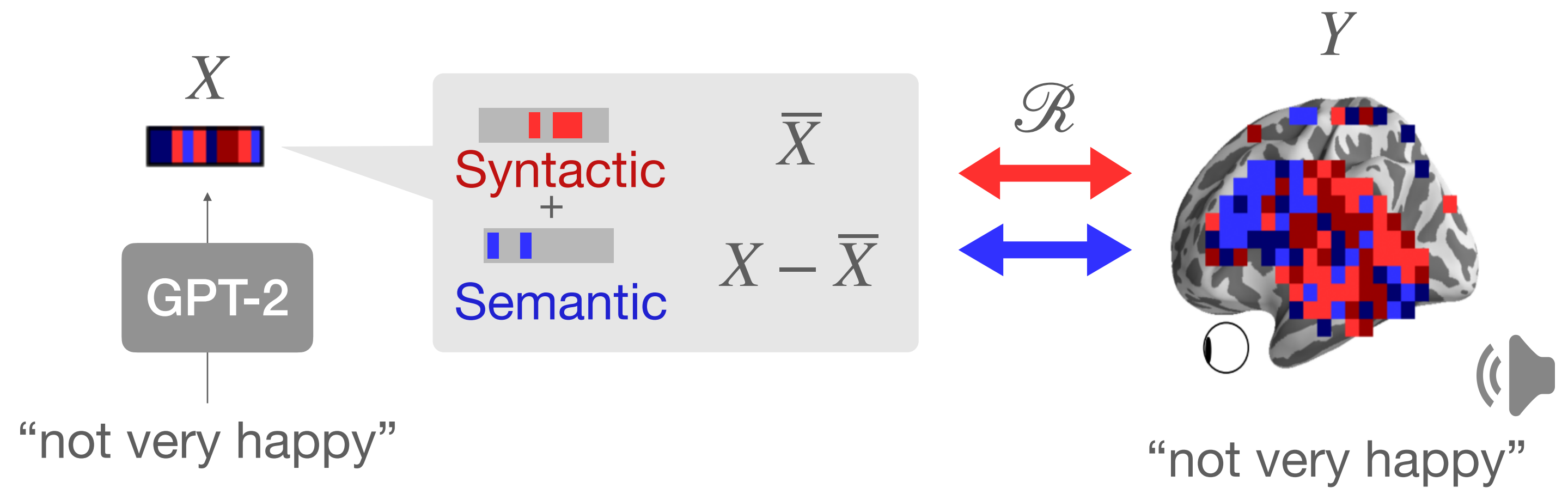
1. We introduce a method to disentangle syntax and semantics in deep nets' activations
2. We use the disentangled activations to **decompose language in the brain**
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*





# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations
2. We use the disentangled activations to **decompose language in the brain**
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

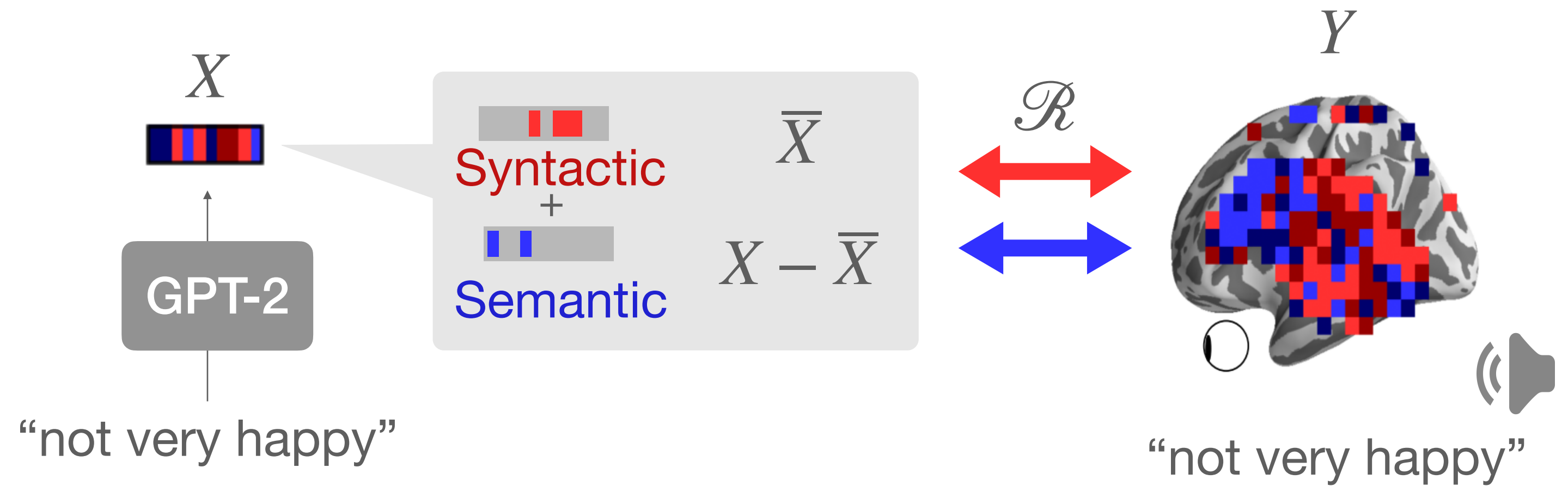


$$\text{Brain score} : \mathcal{R}(\bar{X}) = \text{Corr}(W^T \bar{X}, Y)$$

↑  
Estimated by Ridge regression on training data

# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations



2. We use the disentangled activations to **decompose language in the brain**

$$\text{Brain score} : \mathcal{R}(\bar{X}) = \text{Corr}(W^T \bar{X}, Y)$$

3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

## Measures of brain representations

Syntactic score  $\mathcal{R}(\bar{X})$

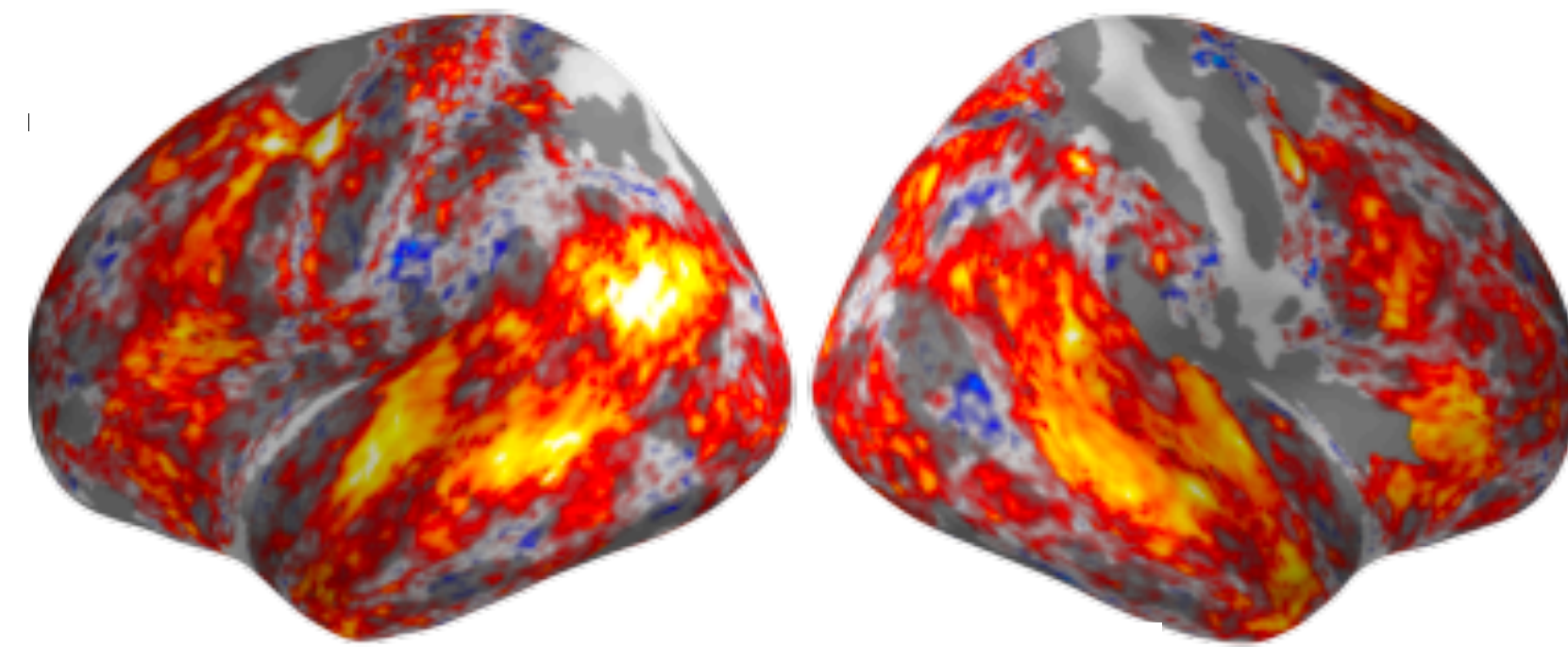
Semantic score  $\mathcal{R}(X) - \mathcal{R}(\bar{X})$

# What are the neural bases of syntax and semantics?

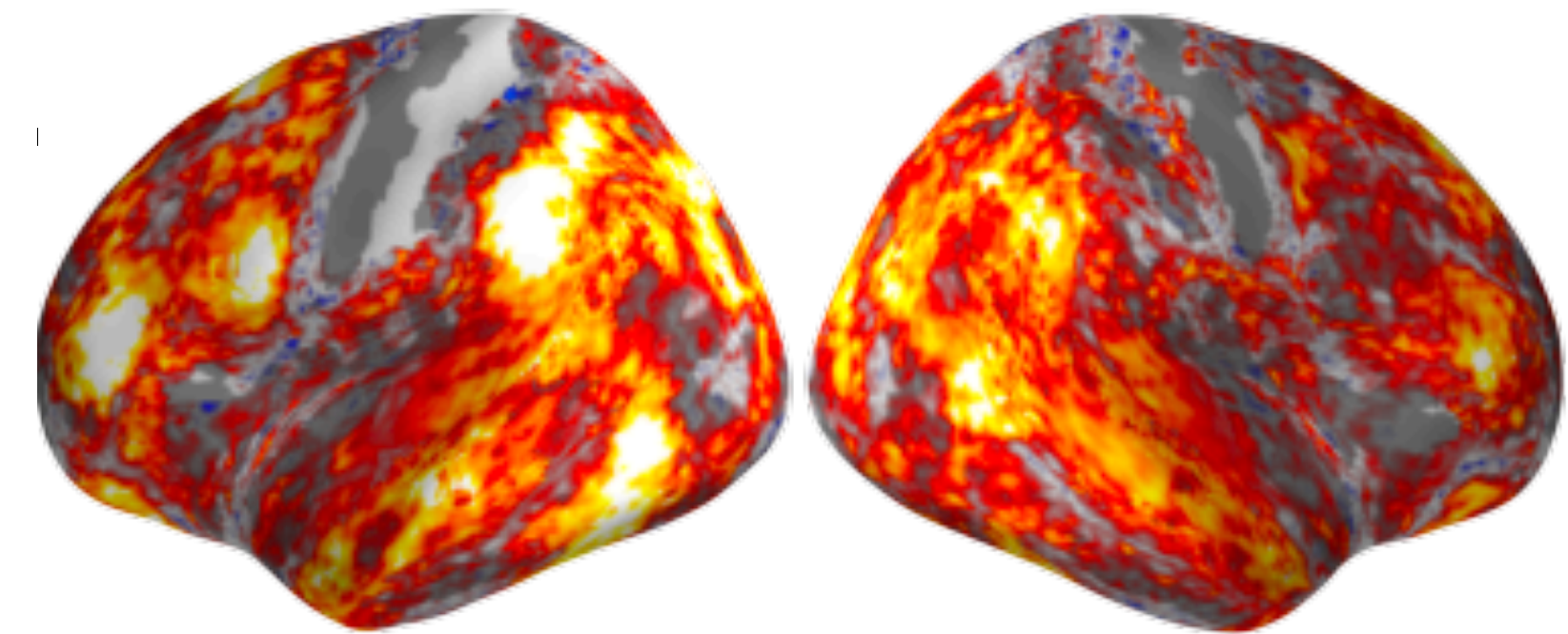
1. We introduce a method to disentangle syntax and semantics in deep nets' activations

2. We use the disentangled activations to decompose language in the brain

3. We apply our method to the **fMRI** recordings of **345 subjects** listening to stories\*



Syntax  
 $\mathcal{R}(\bar{X})$



Semantics  
 $\mathcal{R}(X) - \mathcal{R}(\bar{X})$

## Measures of brain representations

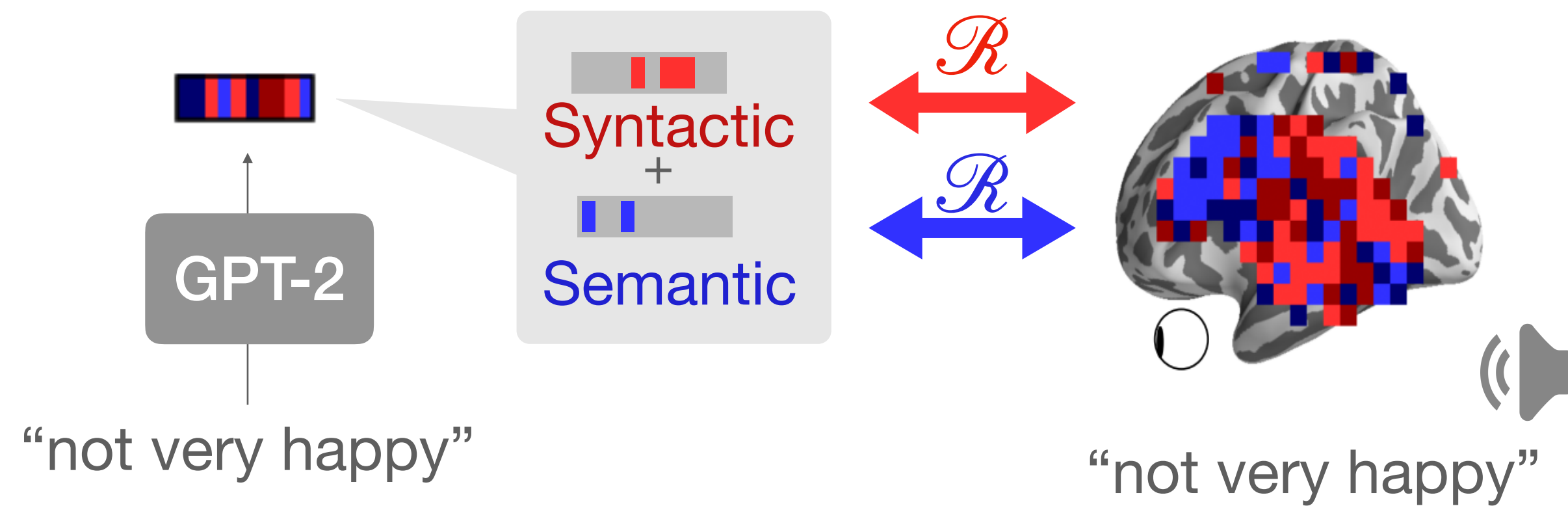
Syntactic score  $\mathcal{R}(\bar{X})$

Semantic score  $\mathcal{R}(X) - \mathcal{R}(\bar{X})$

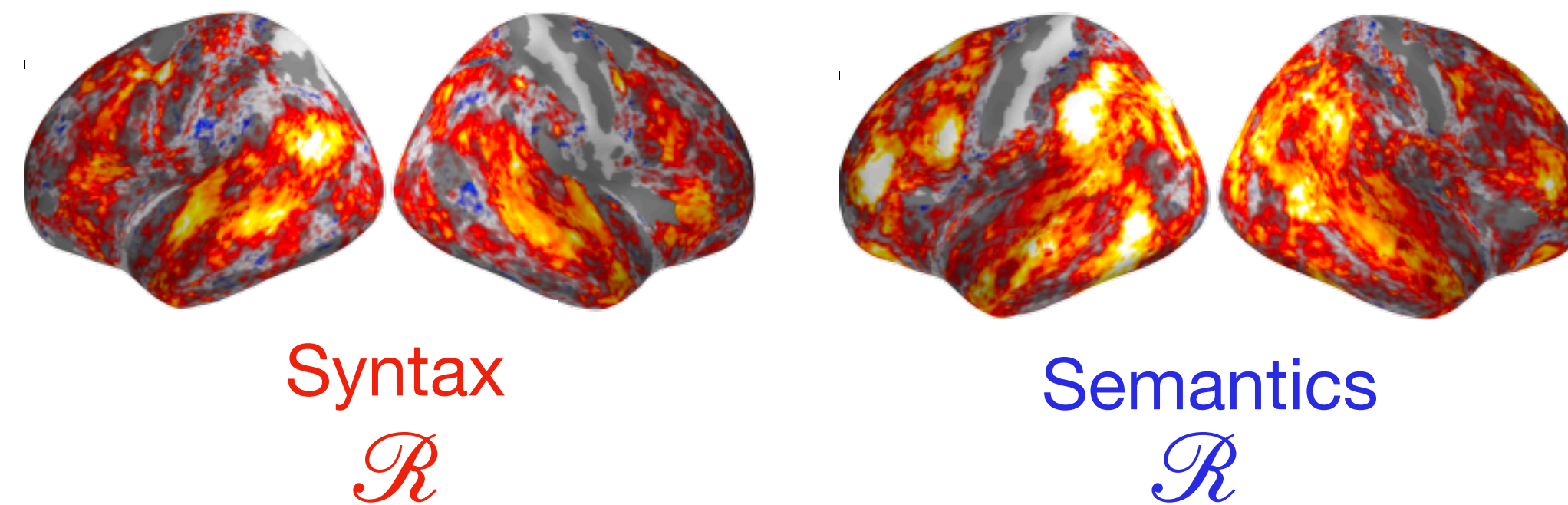
\*345 subjects listening to ~4 hours of unique audio from the “Narratives” dataset (*Nastase et al. 2020*)

# What are the neural bases of syntax and semantics?

1. We introduce a method to disentangle syntax and semantics in deep nets' activations



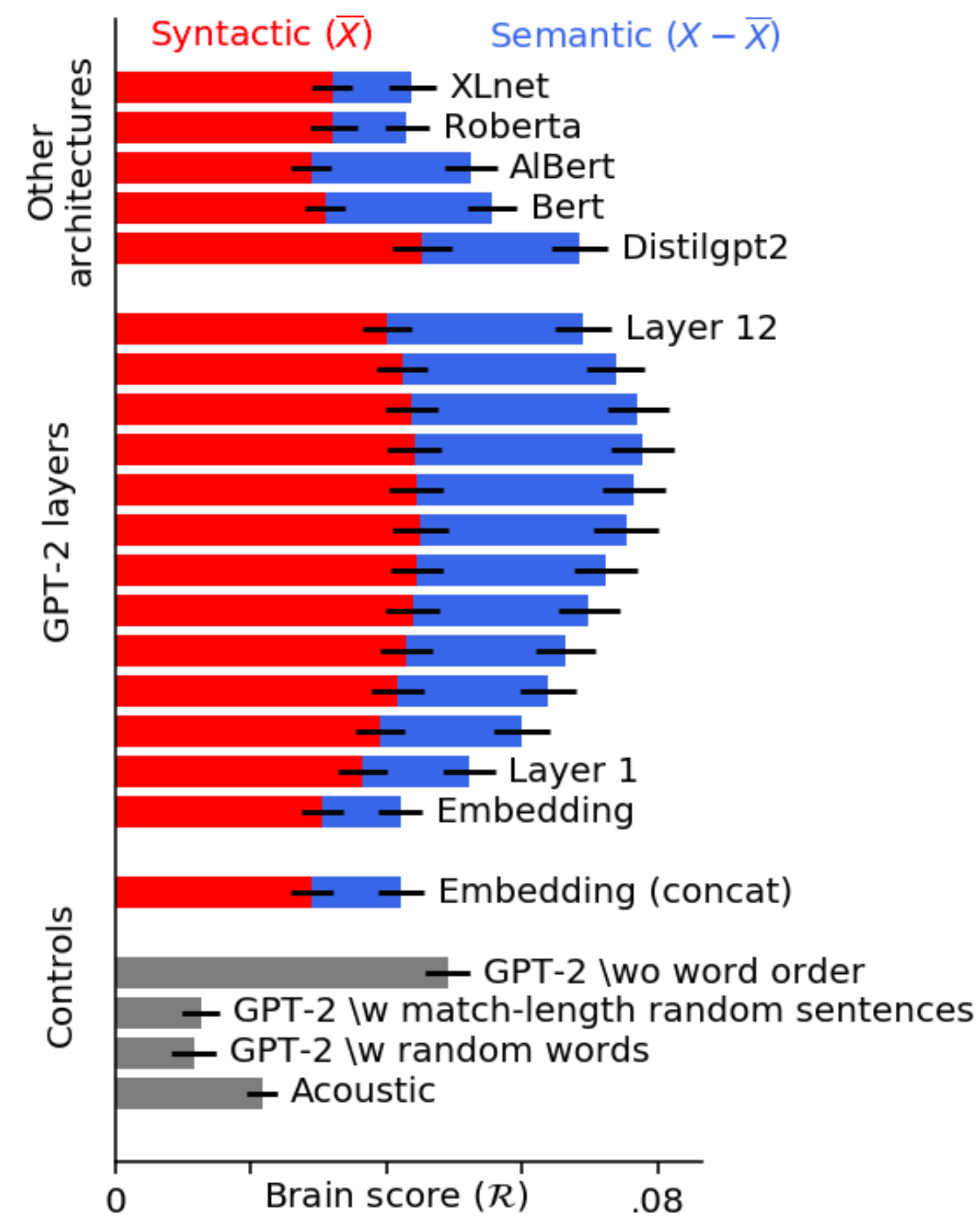
2. We use the disentangled activations to decompose language in the brain



3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

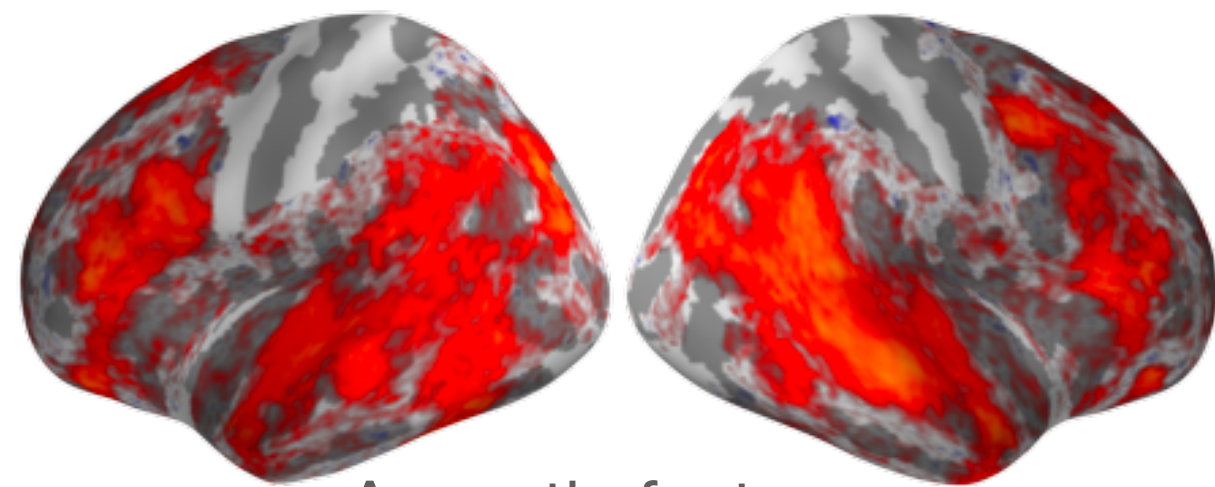
# Generalisation to other architectures and layers

1. We introduce a method to disentangle syntax and semantics in deep nets' activations
2. We use the disentangled activations to decompose language in the brain
3. We apply our method to the fMRI recordings of 345 subjects listening to stories\*

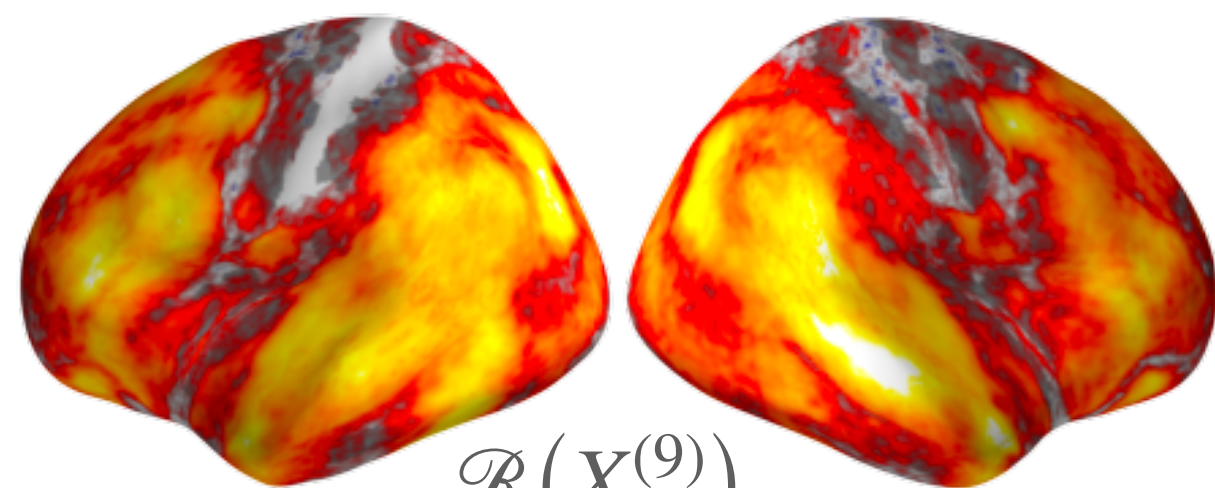
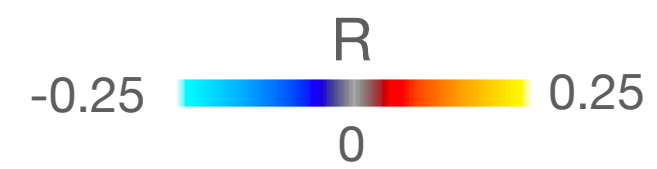


# Fine grained decomposition

A.



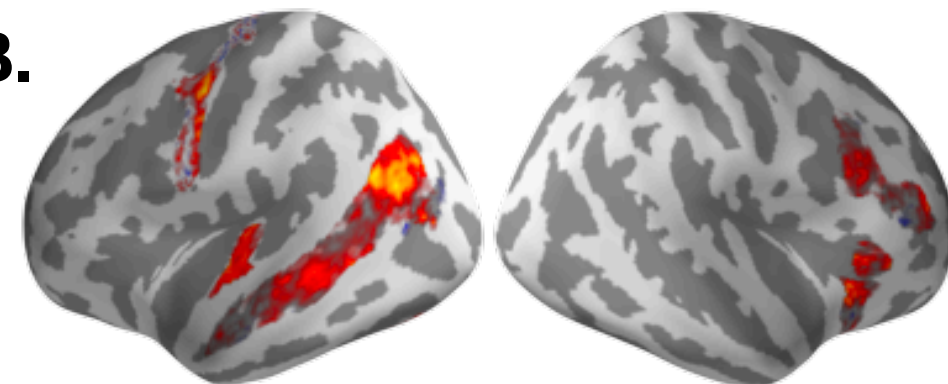
Acoustic features



$\mathcal{R}(X^{(9)})$   
GPT2 activations

Syntactic

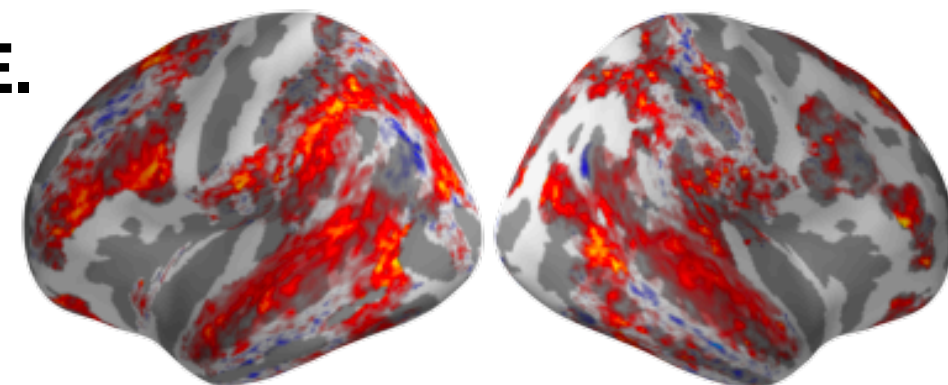
B.



$\mathcal{R}(\overline{X^{(0)}})$

Semantic

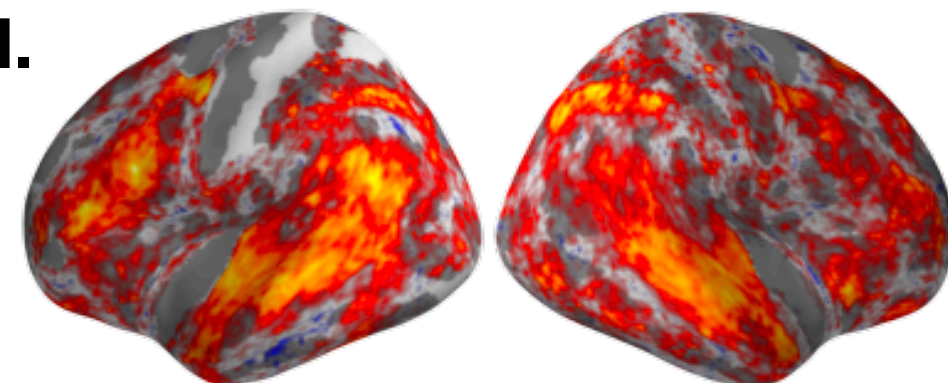
E.



$\mathcal{R}(X^{(0)}) - \mathcal{R}(\overline{X^{(0)}})$

Syntactic  
+ Semantic

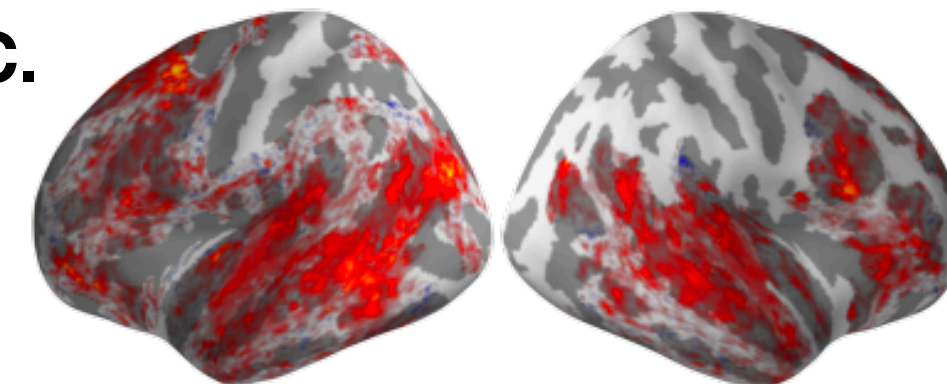
H.



$\mathcal{R}(X^{(0)})$

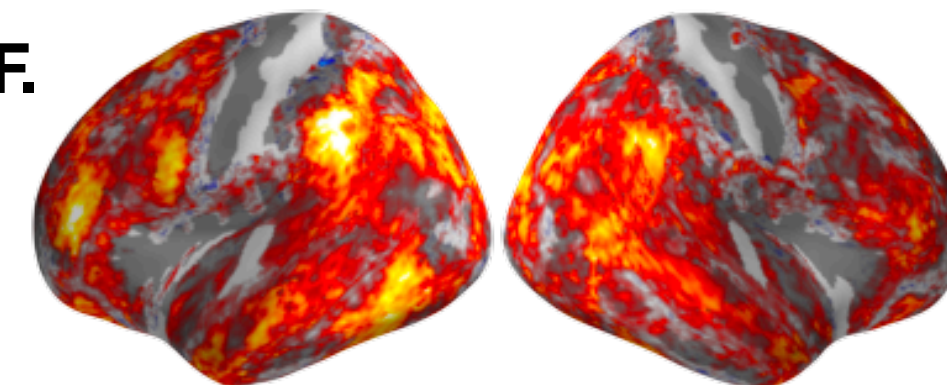
Lexical

C.



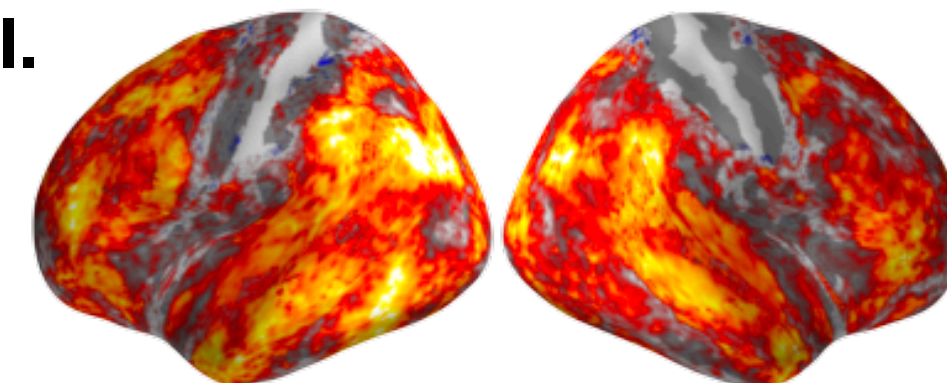
$\mathcal{R}(\overline{X^{(9)}}) - \mathcal{R}(\overline{X^{(0)}})$

F.



$\mathcal{R}(X^{(9)}) - \mathcal{R}(\overline{X^{(9)}} \oplus X^{(0)})$

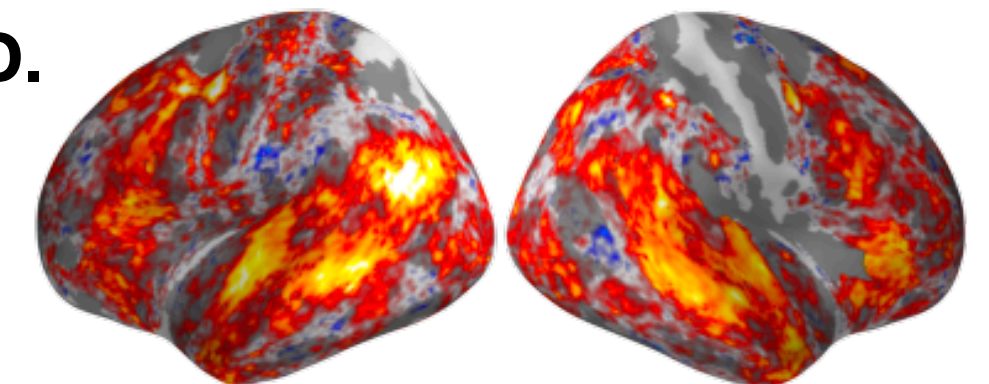
I.



$\mathcal{R}(X^{(9)}) - \mathcal{R}(X^{(0)})$

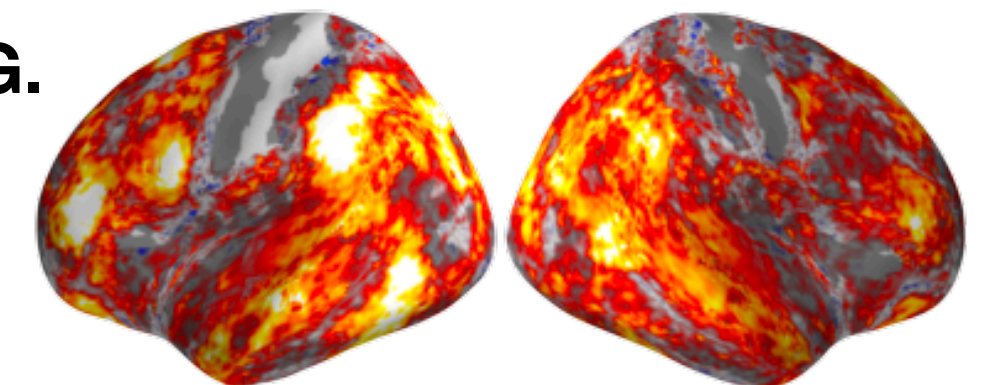
Compositional

D.



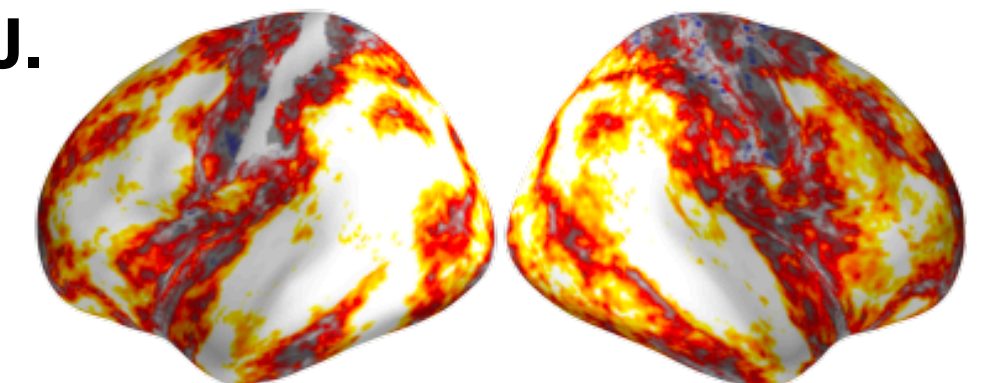
$\mathcal{R}(\overline{X^{(9)}})$

G.



$\mathcal{R}(X^{(9)}) - \mathcal{R}(\overline{X^{(9)}})$

J.



$\mathcal{R}(X^{(9)})$



# Thank you for your attention!



**facebook**  
Artificial Intelligence Research



**Charlotte Caucheteux (INRIA/FAIR), Alexandre Gramfort (INRIA), Jean-Remi King (FAIR/ENS)**