

Causal Interpretation Rules for Encoding and Decoding Models in Neuroimaging

Moritz Grosse-Wentrup

Max Planck Institute for Intelligent Systems
Department Empirical Inference
Tübingen, Germany

June 14, 2015



MAX-PLANCK-GESELLSCHAFT



Causal Terminology in Neuroimaging

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

*We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory.*

*(Anonymous authors, *Journal of Neuroscience*, 2014)*

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

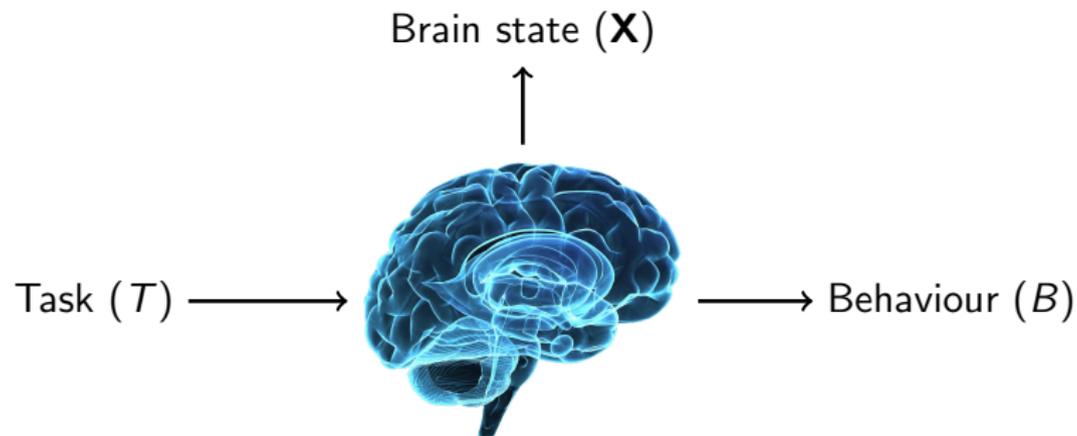
*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

*We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory.*

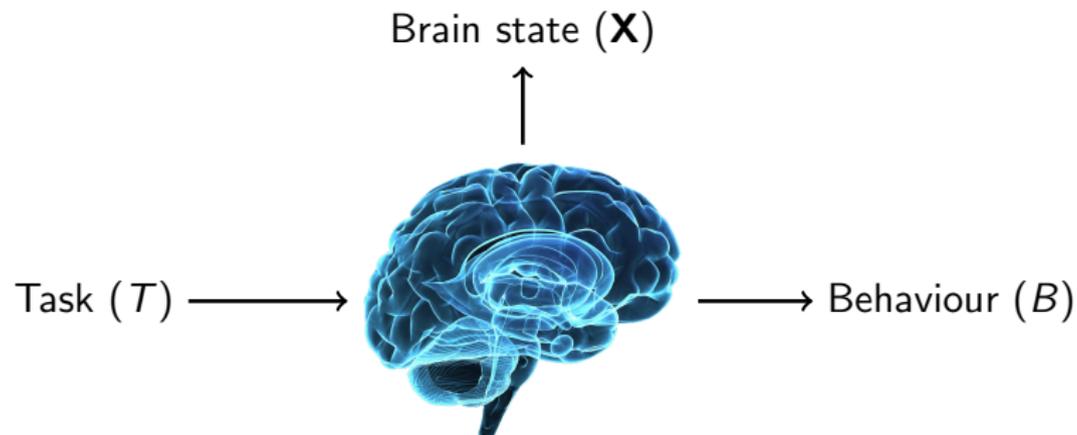
*(Anonymous authors, *Journal of Neuroscience*, 2014)*

- Which causal statements are warranted and which ones are not supported by empirical evidence?

Causal Modelling in Neuroimaging

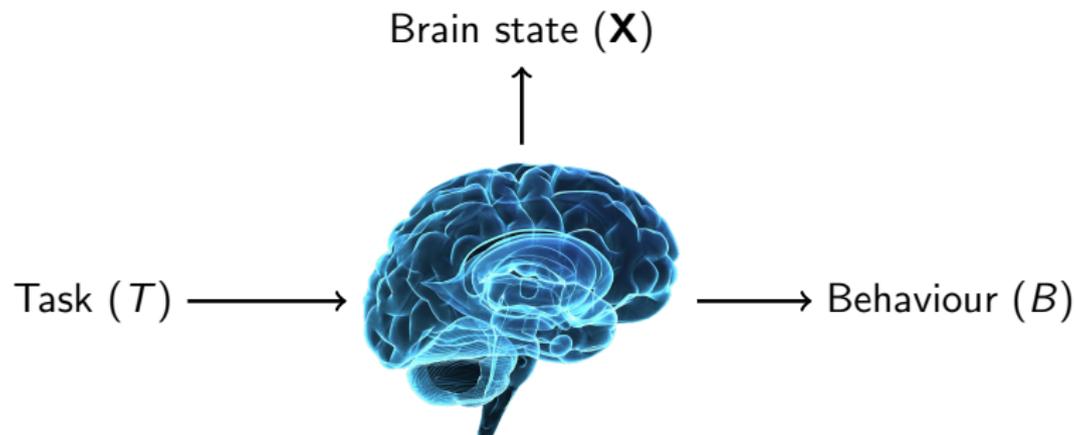


Causal Modelling in Neuroimaging



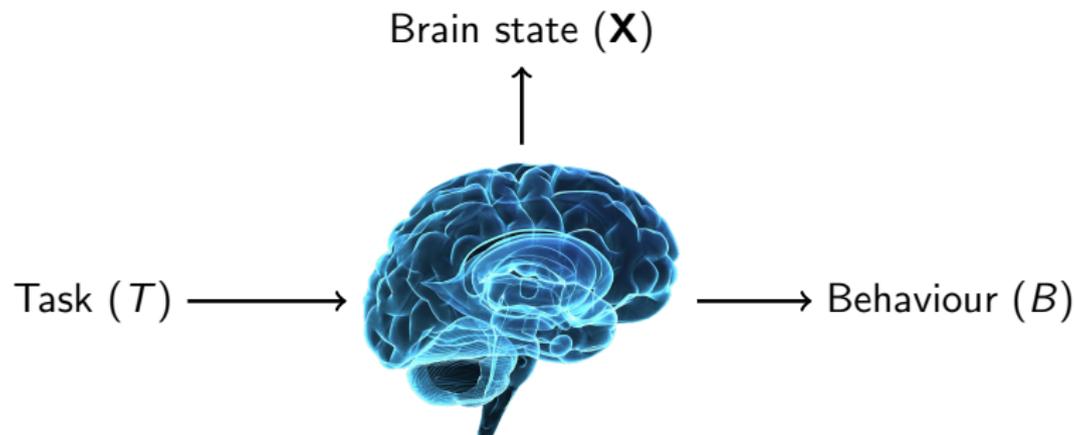
Encoding models: $p(\mathbf{X}|C)$

Causal Modelling in Neuroimaging



Encoding models: $p(\mathbf{X}|C)$

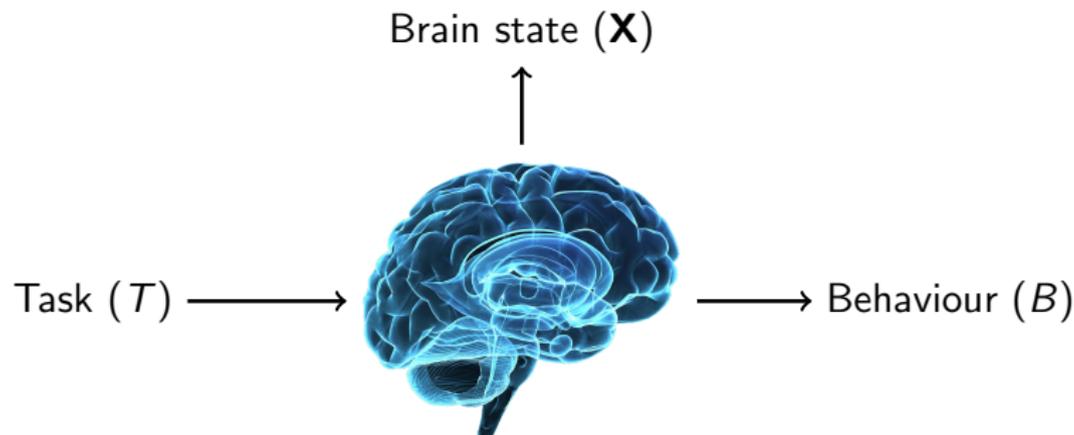
- Task-based: $p(\mathbf{X}|T) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction



Encoding models: $p(\mathbf{X}|C)$

- Task-based: $p(\mathbf{X}|T) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction
- Behaviour-based: $p(\mathbf{X}|B) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction

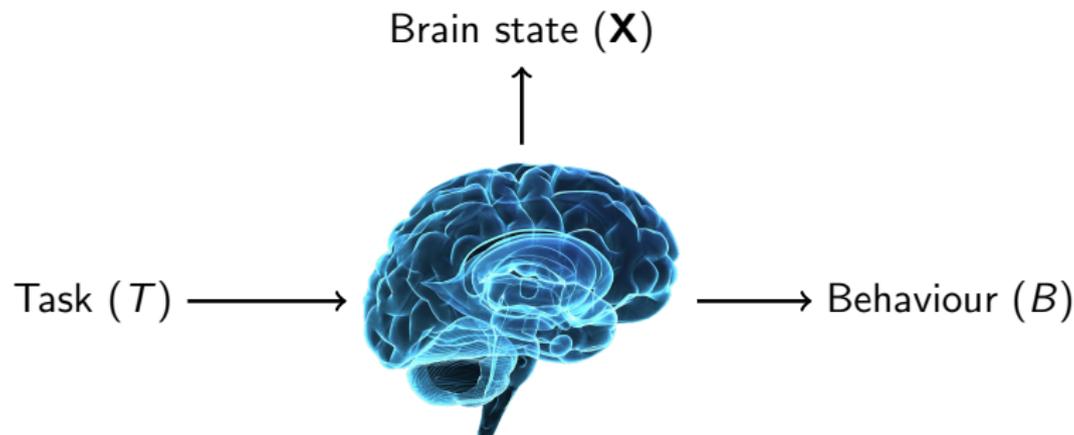
Causal Modelling in Neuroimaging



Encoding models: $p(\mathbf{X}|C)$

- Task-based: $p(\mathbf{X}|T) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction
- Behaviour-based: $p(\mathbf{X}|B) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction

Decoding models: $p(C|\mathbf{X})$

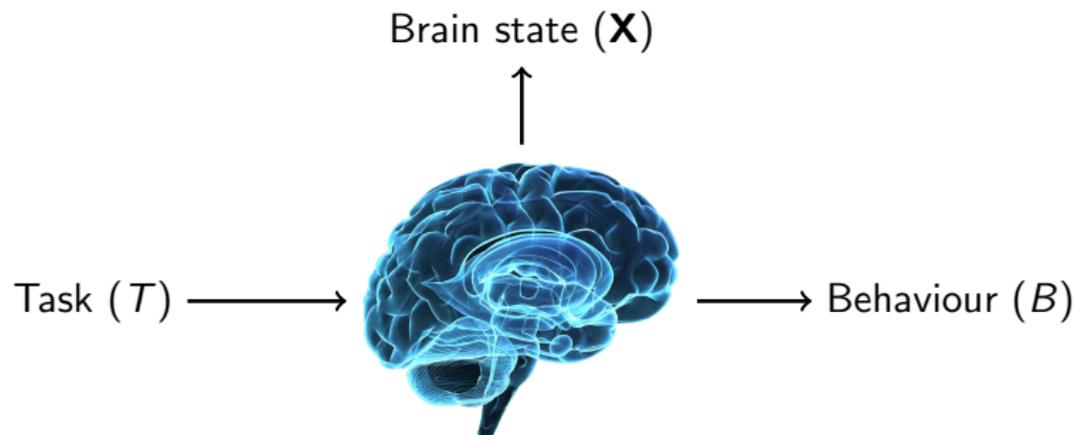


Encoding models: $p(\mathbf{X}|C)$

- Task-based: $p(\mathbf{X}|T) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction
- Behaviour-based: $p(\mathbf{X}|B) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction

Decoding models: $p(C|\mathbf{X})$

- Task-based: $p(T|\mathbf{X}) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction



Encoding models: $p(\mathbf{X}|C)$

- Task-based: $p(\mathbf{X}|T) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction
- Behaviour-based: $p(\mathbf{X}|B) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction

Decoding models: $p(C|\mathbf{X})$

- Task-based: $p(T|\mathbf{X}) \equiv p(\text{cause}|\text{effect}) \Rightarrow$ anti-causal direction
- Behaviour-based: $p(B|\mathbf{X}) \equiv p(\text{effect}|\text{cause}) \Rightarrow$ causal direction

- Experimental condition C

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons
 - ▶ ...

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons
 - ▶ ...
- Experimental data $\{(c, \mathbf{x})_1, \dots, (c, \mathbf{x})_N\}$, drawn i.i.d. from $p(C, \mathbf{X})$

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons
 - ▶ ...
- Experimental data $\{(c, \mathbf{x})_1, \dots, (c, \mathbf{x})_N\}$, drawn i.i.d. from $p(C, \mathbf{X})$
- Assumption: Oracle for properties of $p(C, \mathbf{X})$ available

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons
 - ▶ ...
- Experimental data $\{(c, \mathbf{x})_1, \dots, (c, \mathbf{x})_N\}$, drawn i.i.d. from $p(C, \mathbf{X})$
- Assumption: Oracle for properties of $p(C, \mathbf{X})$ available
 - ▶ Statistical independence:
 $C \perp\!\!\!\perp X_i \Leftrightarrow p(C, X_i) = p(C)p(X_i)$

- Experimental condition C , such as
 - ▶ type of stimulus displayed to a subject: $C = T \in \{-1, +1\}$
 - ▶ delay of a behavioural response: $C = B \in \mathbb{R}_+$
 - ▶ ...
- Brain state features $\mathbf{X} \in \mathbb{R}^M$, such as
 - ▶ BOLD response at M voxels measured by fMRI
 - ▶ bandpower at M EEG channels
 - ▶ mean spike-rate of M neurons
 - ▶ ...
- Experimental data $\{(c, \mathbf{x})_1, \dots, (c, \mathbf{x})_N\}$, drawn i.i.d. from $p(C, \mathbf{X})$
- Assumption: Oracle for properties of $p(C, \mathbf{X})$ available
 - ▶ Statistical independence:
 $C \perp\!\!\!\perp X_i \Leftrightarrow p(C, X_i) = p(C)p(X_i)$
 - ▶ Conditional statistical independence:
 $C \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i \Leftrightarrow p(C, X_i | \mathbf{X} \setminus X_i) = p(C | \mathbf{X} \setminus X_i)p(X_i | \mathbf{X} \setminus X_i)$

Causal Bayesian Networks (Pearl, 2000; Spirtes et al., 2000)

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain

$$X_1 \rightarrow X_2 \rightarrow X_3$$

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain

$$X_1 \rightarrow X_2 \rightarrow X_3$$

The fork

$$X_1 \leftarrow X_2 \rightarrow X_3$$

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j | \text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j | \text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$

$$X_1 \not\perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 | X_2$$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$
 $X_1 \perp\!\!\!\perp X_3 | X_2$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain

$$X_1 \rightarrow X_2 \rightarrow X_3$$
$$X_1 \not\perp\!\!\!\perp X_3$$
$$X_1 \perp\!\!\!\perp X_3 | X_2$$

The fork

$$X_1 \leftarrow X_2 \rightarrow X_3$$
$$X_1 \not\perp\!\!\!\perp X_3$$
$$X_1 \perp\!\!\!\perp X_3 | X_2$$

The collider

$$X_1 \rightarrow X_2 \leftarrow X_3$$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j | \text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$
 $X_1 \perp\!\!\!\perp X_3 | X_2$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$
 $X_1 \perp\!\!\!\perp X_3 | X_2$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$
 $X_1 \perp\!\!\!\perp X_3$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j|\text{do}\{x_i\}) \neq p(x_j)$.

The chain
 $X_1 \rightarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$
 $X_1 \perp\!\!\!\perp X_3 | X_2$

The fork
 $X_1 \leftarrow X_2 \rightarrow X_3$
 $X_1 \not\perp\!\!\!\perp X_3$
 $X_1 \perp\!\!\!\perp X_3 | X_2$

The collider
 $X_1 \rightarrow X_2 \leftarrow X_3$
 $X_1 \perp\!\!\!\perp X_3$
 $X_1 \not\perp\!\!\!\perp X_3 | X_2$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.

X_i is a cause of X_j ($X_i \rightarrow X_j$), iff there exist values of X_i and X_j such that $p(x_j | \text{do}\{x_i\}) \neq p(x_j)$.

The chain	The fork	The collider
$X_1 \rightarrow X_2 \rightarrow X_3$	$X_1 \leftarrow X_2 \rightarrow X_3$	$X_1 \rightarrow X_2 \leftarrow X_3$
$X_1 \not\perp\!\!\!\perp X_3$	$X_1 \not\perp\!\!\!\perp X_3$	$X_1 \perp\!\!\!\perp X_3$
$X_1 \perp\!\!\!\perp X_3 X_2$	$X_1 \perp\!\!\!\perp X_3 X_2$	$X_1 \not\perp\!\!\!\perp X_3 X_2$

- Causal Markov condition: Independence relations implied by a directed acyclic graph (DAG) are encoded in every $p(\mathbf{X})$ generated by this DAG.
- Faithfulness: $p(\mathbf{X})$ contains no additional independence relations beyond those implied by its generating DAG.

Causal Terminology Revisited

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

HC ~~↔~~ AM ⇒ AM → HC → EM

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

*(Anonymous author, *Trends in Cognitive Sciences*, 2001)*

$HC \not\perp AM \Rightarrow AM \rightarrow HC \rightarrow EM$

*We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory.*

*(Anonymous authors, *Journal of Neuroscience*, 2014)*

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus.*

(Anonymous author, *Trends in Cognitive Sciences*, 2001)

$$HC \not\perp AM \Rightarrow AM \rightarrow HC \rightarrow EM$$

*We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory.*

(Anonymous authors, *Journal of Neuroscience*, 2014)

$$\alpha \not\perp WM \Rightarrow \alpha \rightarrow WM$$

Causal Inference in Encoding/Decoding Models

Encoding models: $p(\mathbf{X}|C)$

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?

$$H_0 : C \perp\!\!\!\perp X_i$$

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Decoding models: $p(C|\mathbf{X})$

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Decoding models: $p(C|\mathbf{X})$

- Does a brain-state feature help in decoding the experimental condition?

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Decoding models: $p(C|\mathbf{X})$

- Does a brain-state feature help in decoding the experimental condition?

$$H_0 : C \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$$

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Decoding models: $p(C|\mathbf{X})$

- Does a brain-state feature help in decoding the experimental condition?
- Reject $H_0 : C \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i \Rightarrow X_i$ is a relevant feature

Encoding models: $p(\mathbf{X}|C)$

- Does a brain-state feature X_i change across experimental conditions?
- Reject $H_0 : C \perp\!\!\!\perp X_i \Rightarrow X_i$ is a relevant feature

Decoding models: $p(C|\mathbf{X})$

- Does a brain-state feature help in decoding the experimental condition?
- Reject $H_0 : C \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i \Rightarrow X_i$ is a relevant feature

Feature relevance $\Rightarrow C \not\perp\!\!\!\perp X_i / C \not\perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i \Rightarrow$ causal structure

- 1 Causal Interpretation Rules for Individual Encoding/Decoding Models
- 2 Causal Interpretation Rules for Joint Encoding/Decoding Models

- 1 Causal Interpretation Rules for Individual Encoding/Decoding Models
- 2 Causal Interpretation Rules for Joint Encoding/Decoding Models

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	x		
		✓		
			x	
			✓	
	Behaviour	x		
		✓		
			x	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i$
		✓		
			×	
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		
			×	
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		√		
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
			√	
	Behaviour	×		
		√		
			×	
			√	

$T \perp\!\!\!\perp X_i$ vs. $T \rightarrow X_j \rightarrow X_i$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

$T \perp\!\!\!\perp X_i$ vs. $T \rightarrow X_j \rightarrow X_i$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i$
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
	Behaviour	×		
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
	Behaviour	×		
		✓		
			×	
			✓	

$$T \rightarrow X_i \text{ or } T \rightarrow X_j \leftarrow X_i$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		
		✓		
			×	
			✓	

$$T \rightarrow X_i \text{ or } T \rightarrow X_j \leftarrow X_i$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i$
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		
			×	
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
			✓	

X_i B vs. $X_i \rightarrow X_j \rightarrow B$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	

X_i B vs. $X_i \rightarrow X_j \rightarrow B$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i$
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i$
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	

$$X_i \rightarrow B \text{ or } X_i \leftarrow X_j \rightarrow B$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	

$$X_i \rightarrow B \text{ or } X_i \leftarrow X_j \rightarrow B$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$

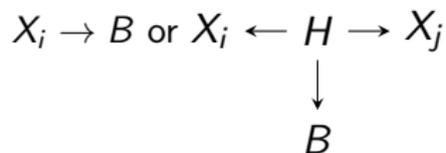
$$X_i \rightarrow B \text{ or } X_i \leftarrow H \rightarrow X_j$$

$$\downarrow$$

$$B$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive



Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive

We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory. (Anonymous authors, *Journal of Neuroscience*, 2014)

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			×	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	×		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			×	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive

We tested [...] whether pre-stimulus alpha oscillations **measured** with electroencephalography (EEG) **influence** the encoding of items into working memory. (Anonymous authors, *Journal of Neuroscience*, 2014)

$$\alpha \not\perp\!\!\!\perp \text{WM} \Rightarrow \alpha \rightarrow \text{WM}$$

Causal Interpretation Rules for Individual Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	x		$T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of T
		✓		$T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of T
			x	$T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
	Behaviour	x		$B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of B
		✓		$B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive
			x	$B \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive
			✓	$B \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i \Rightarrow$ inconclusive

- 1 Causal Interpretation Rules for Individual Encoding/Decoding Models
- 2 Causal Interpretation Rules for Joint Encoding/Decoding Models

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	
		×	✓	
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	$T \not\perp X_i$ & $T \perp X_i \mathbf{X} \setminus X_i$
		×	✓	
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	$T \not\perp\!\!\!\perp X_i$ & $T \perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
		×	✓	
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

$$T \rightarrow X_j \rightarrow X_i$$

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	X_i is an indirect effect of T
		×	✓	
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

$$T \rightarrow X_j \rightarrow X_i$$

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	X_i is an indirect effect of T
		×	✓	$T \perp\!\!\!\perp X_i$ & $T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	X_i is an indirect effect of T
		×	✓	$T \perp\!\!\!\perp X_i$ & $T \not\perp\!\!\!\perp X_i \mathbf{X} \setminus X_i$
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

$$T \rightarrow X_j \leftarrow X_i \text{ or } T \rightarrow X_j \leftarrow H \rightarrow X_i$$

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	
	Behaviour	×	×	
		✓	×	
		×	✓	
		✓	✓	

$$T \rightarrow X_j \leftarrow X_i \text{ or } T \rightarrow X_j \leftarrow H \rightarrow X_i$$

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	X_i is no effect of T
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	X_i is an effect of T
	Behaviour	×	×	X_i is no cause of B
		✓	×	X_i is no direct cause of B
		×	✓	X_i provides context
		✓	✓	inconclusive

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	X_i is no effect of T
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	X_i is an effect of T
	Behaviour	×	×	X_i is no cause of B
		✓	×	X_i is no direct cause of B
		×	✓	X_i provides context
		✓	✓	inconclusive

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	X_i is no effect of T
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	X_i is an effect of T
	Behaviour	×	×	X_i is no cause of B
		✓	×	X_i is no direct cause of B
		×	✓	X_i provides context
		✓	✓	inconclusive

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus. (Anonymous author, Trends in Cognitive Sciences, 2001)*

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	X_i is no effect of T
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	X_i is an effect of T
	Behaviour	×	×	X_i is no cause of B
		✓	×	X_i is no direct cause of B
		×	✓	X_i provides context
		✓	✓	inconclusive

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus. (Anonymous author, Trends in Cognitive Sciences, 2001)*

$$HC \not\perp AM \Rightarrow AM \rightarrow HC \rightarrow EM$$

Causal Interpretation Rules for Joint Models

		Feature X_i relevant?		Causal interpretation
		Encoding	Decoding	
Experimental setting	Task	×	×	X_i is no effect of T
		✓	×	X_i is an indirect effect of T
		×	✓	X_i provides context
		✓	✓	X_i is an effect of T
	Behaviour	×	×	X_i is no cause of B
		✓	×	X_i is no direct cause of B
		×	✓	X_i provides context
		✓	✓	inconclusive

*Hippocampal activity in this study was **correlated** with amygdala activity, supporting the view that the amygdala **enhances** explicit memory by **modulating** activity in the hippocampus. (Anonymous author, Trends in Cognitive Sciences, 2001)*

$$AM \not\perp EM \ \& \ AM \perp EM|HC \Leftarrow AM \rightarrow HC \rightarrow EM$$

Wrapping Up

Wrapping Up

Caveats:

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Take home message:

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Take home message:

- If you don't like causal inference, don't use causal terminology.

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Take home message:

- If you don't like causal inference, don't use causal terminology.
- If you use causal terminology, make sure that

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Take home message:

- If you don't like causal inference, don't use causal terminology.
- If you use causal terminology, make sure that
 - ▶ your conclusions are supported by empirical data

Caveats:

- The probability of type II errors can not be quantified in the same manner as the probability of type I errors.
- The assumption of faithfulness is presently untestable.
- Permutation-based relevance tests in decoding models are biased towards dependence (Strobl et al., *BMC Bioinformatics*, 2008)
- Unbiased conditional independence tests are hard (Zhang et al., *UAI*, 2011)

Take home message:

- If you don't like causal inference, don't use causal terminology.
- If you use causal terminology, make sure that
 - ▶ your conclusions are supported by empirical data
 - ▶ you are explicit about inherent assumptions

Acknowledgements

People:

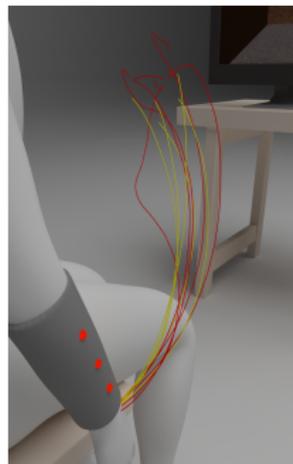
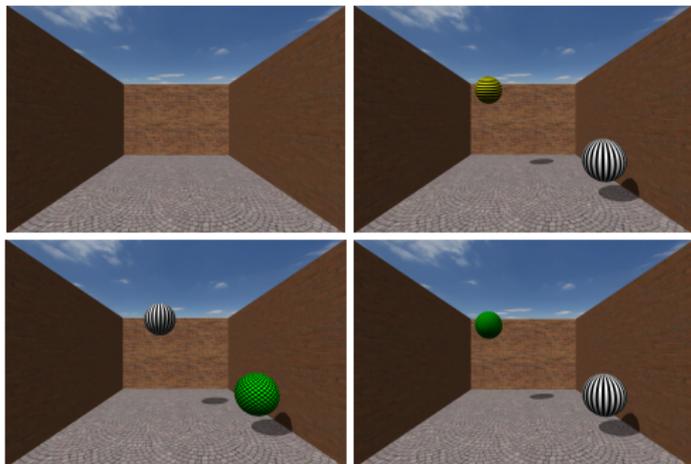
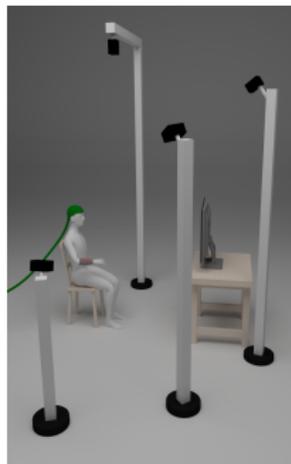
- Sebastian Weichwald (UCL)
- Timm Meyer (MPI-IS)
- Ozan Özdenizci (Sabanci University)
- Bernhard Schölkopf (MPI-IS)
- Tonio Ball (University of Freiburg)

Publications:

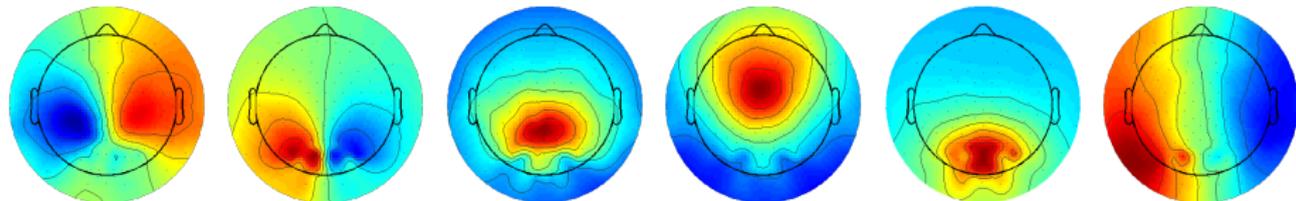
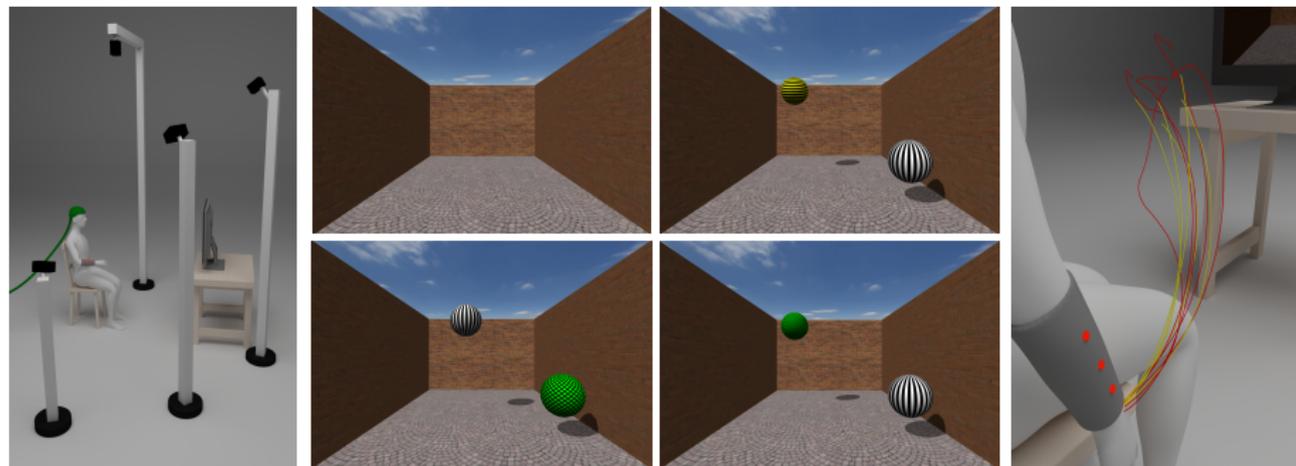
- Weichwald et al., Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage*, 2015.
- Weichwald et al., Causal and anti-causal learning in pattern recognition for neuroimaging. *PRNI*, 2014.

<http://brain-computer-interfaces.net>

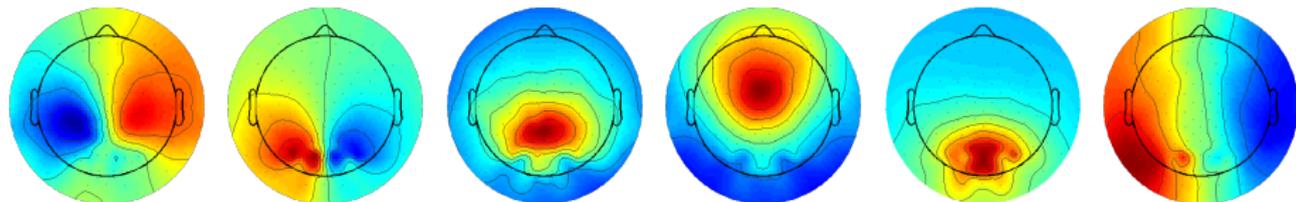
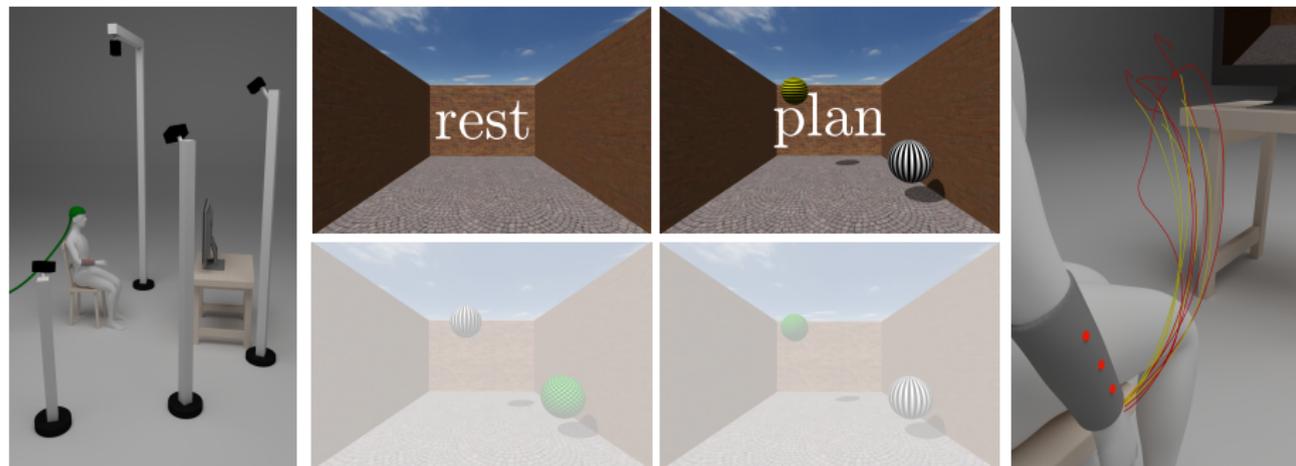
Experimental Setup



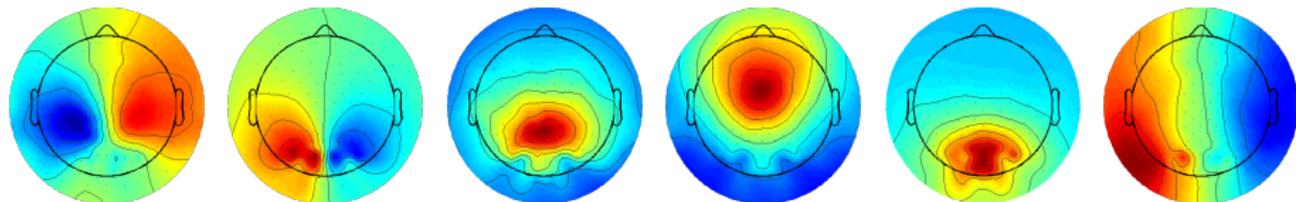
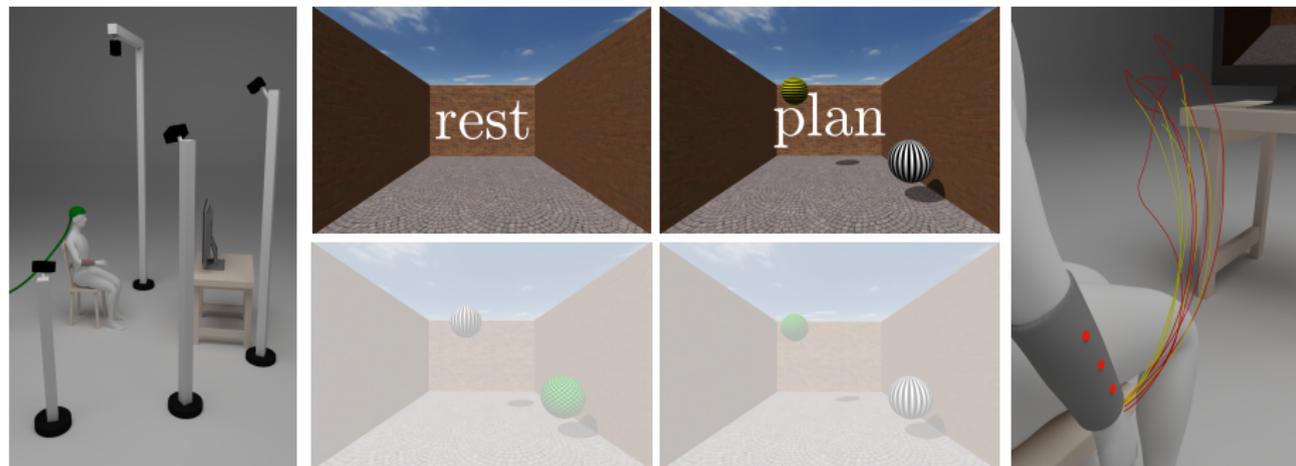
Experimental Setup



Experimental Setup



Experimental Setup



$$\{S, X_1, \dots, X_6\} = \{\text{rest/plan}, |\alpha_{IC_1}|, \dots, |\alpha_{IC_6}|\}$$

- Experimental data: 17 subjects with 444 - 498 trials each

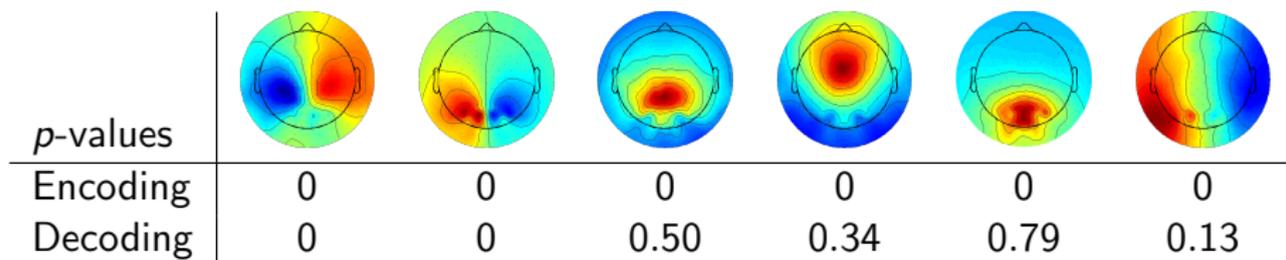
Data Analysis & Experimental Results

- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]

- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]

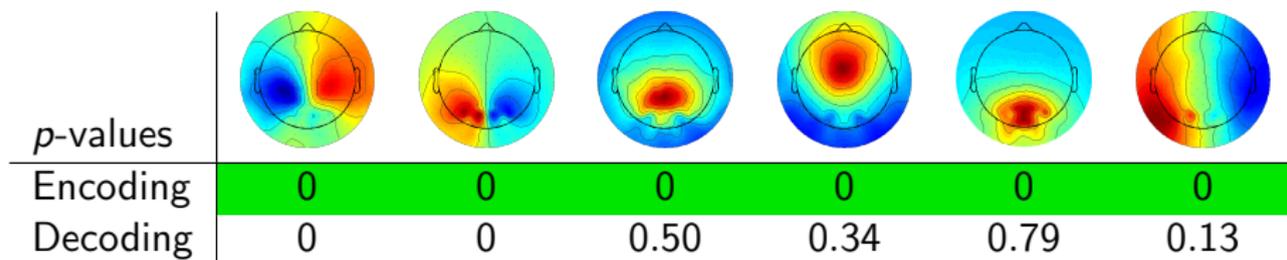
Data Analysis & Experimental Results

- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]



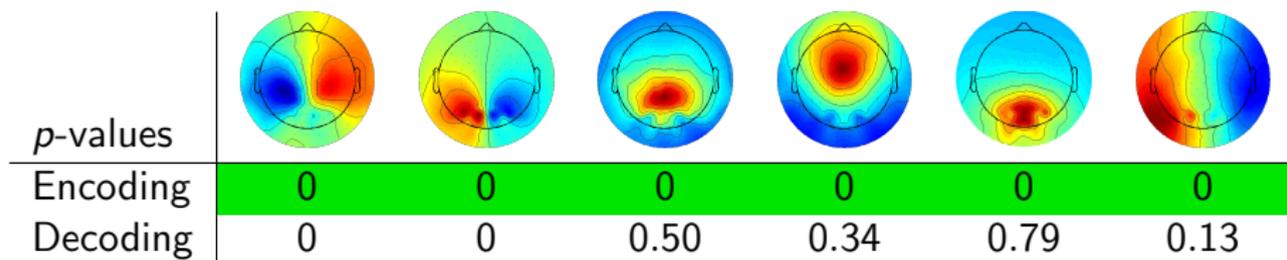
Data Analysis & Experimental Results

- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]



Data Analysis & Experimental Results

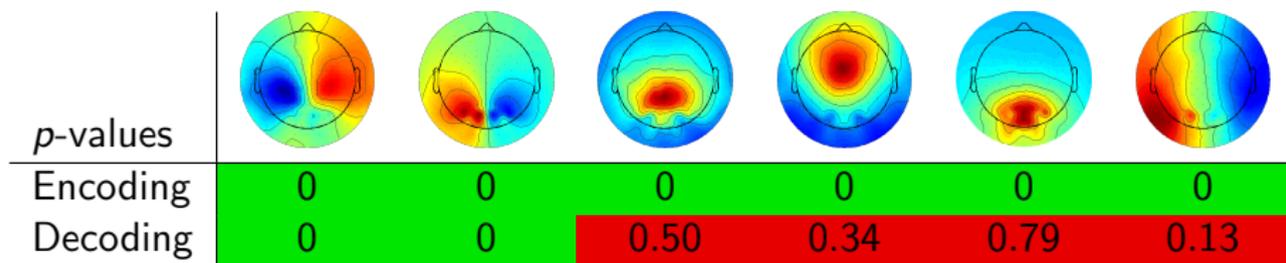
- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]



→ The instruction to plan a reaching movement causes modulation of α -power at every IC: $S \rightarrow \{|\alpha_{IC_i}|\}, i = 1, \dots, 6$

Data Analysis & Experimental Results

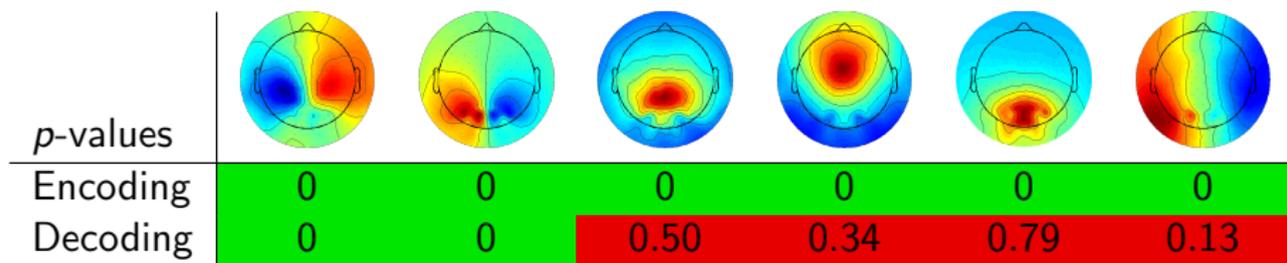
- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]



→ The instruction to plan a reaching movement causes modulation of α -power at every IC: $S \rightarrow \{|\alpha_{IC_i}|\}, i = 1, \dots, 6$

Data Analysis & Experimental Results

- Experimental data: 17 subjects with 444 - 498 trials each
- Encoding model: $H_0 : S \perp\!\!\!\perp X_i$ [HSIC (Gretton et al., *NIPS*, 2008)]
- Decoding model: $H_0 : S \perp\!\!\!\perp X_i | \mathbf{X} \setminus X_i$ [Random forest (Breiman, *Machine Learning*, 2001)]



- The instruction to plan a reaching movement causes modulation of α -power at every IC: $S \rightarrow \{|\alpha_{IC_i}|\}, i = 1, \dots, 6$
- Modulation of α -power at ICs 3–6 is only an indirect effect relative to ICs 1 & 2: $S \rightarrow \{|\alpha_{IC_1}|, |\alpha_{IC_2}|\} \rightarrow \{|\alpha_{IC_3}|, |\alpha_{IC_4}|, |\alpha_{IC_5}|, |\alpha_{IC_6}|\}$