MAX-PLANCK-GESELLSCHAFT

# An introduction to the different causal frameworks in neuroimaging

Sebastian Weichwald

Max Planck Institute for Intelligent Systems,
Max Planck ETH Center for Learning Systems

sweichwald.de/prni2017          neural.engineering

# Why causality?

*To paraphrase a old joke, there are two types of statisticians: those who do causal inference and those who lie about it.*

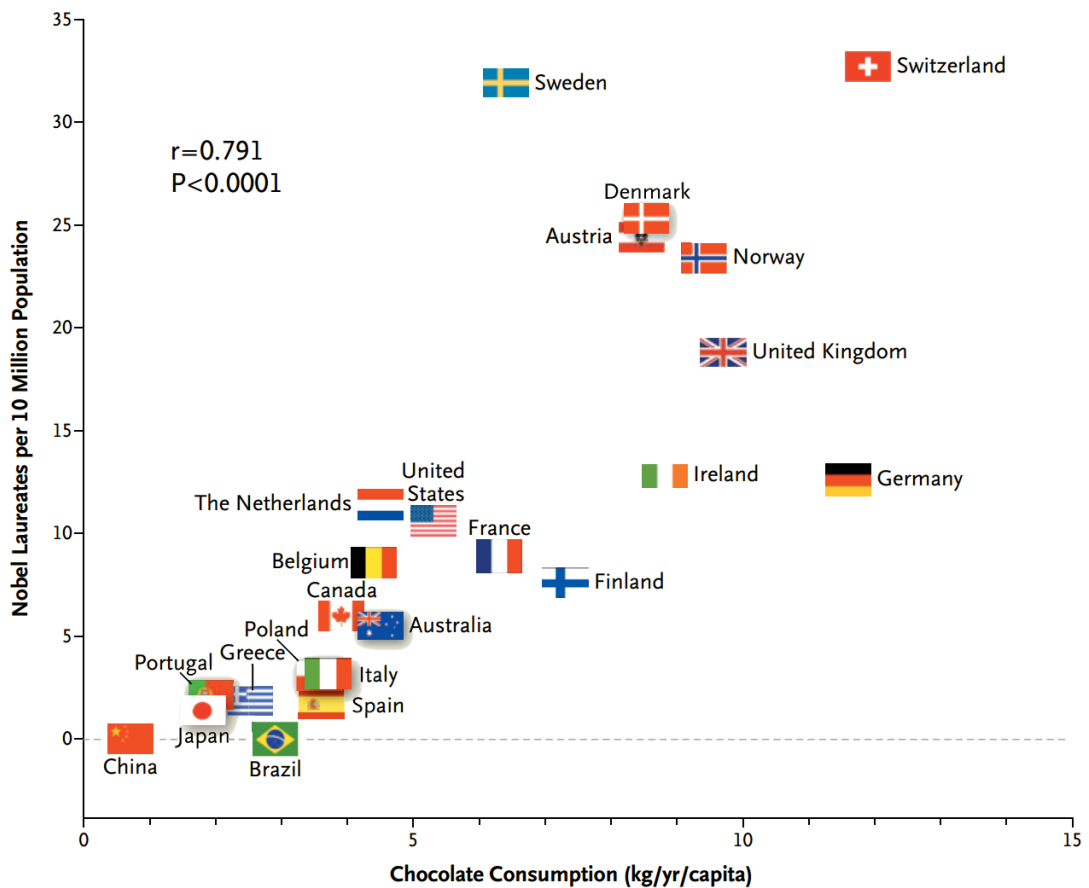(L Wasserman, *Journal of the American Statistical Association,* 1999)

**Figure 1.** Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

(FH Messerli, Chocolate Consumption, Cognitive Function, and Nobel Laureates. *New England Journal of Medicine,* 2012)

A scientific theory should

- ‣ Explain already observed data
- ‣ Predict future observations
  - ○ of a *passively* observed system
  - ○ of a system that is *actively* intervened upon

We want to predict the effect of interventions!

3

amygdala → hippocampus → explicit memory

*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 Authors, *Trends in Cognitive Sciences,* 2001)

# Common causal frameworks

‣ Potential Outcomes Framework

‣ Granger Causality

‣ Dynamic Causal Modelling

‣ Causal Bayesian Networks and Structural Equation Models

# Potential Outcomes Framework

Ingredients:

- ▸ Population $\mathcal{U}$ of units $u \in \mathcal{U}$,

    e. g. a patient group

- ▸ Treatment variable $S : \mathcal{U} \to \{\mathrm{t}, \mathrm{c}\}$,

    e. g. assignment to treatment/control

- ▸ Potential outcomes $Y : \mathcal{U} \times \{\mathrm{t}, \mathrm{c}\} \to \mathbb{R}$,

    e. g. survival times $Y_{\mathrm{t}}(u)$ and $Y_{\mathrm{c}}(u)$ of patient $u$

*Fundamental problem of causal inference:*
For each unit $u$ we get to observe *either $Y_t(u)$ or $Y_c(u)$* and hence the treatment effect $Y_t(u) - Y_c(u)$ cannot be computed.
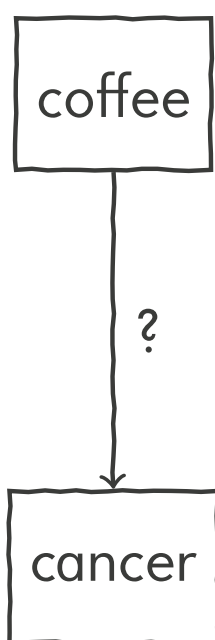
*Possible remedy assumptions:*

- ▸ Unit homogeneity: $Y_t(u_1) = Y_t(u_2)$ and $Y_c(u_1) = Y_c(u_2)$
- ▸ Causal transience: can measure $Y_t(u)$ and $Y_c(u)$ sequentially

*"Statistical solution":* Average Treatment Effect $\mathbb{E}[Y_t] - \mathbb{E}[Y_c]$

- ▸ Can observe $\mathbb{E}[Y_t | S = t]$ and $\mathbb{E}[Y_c | S = c]$
- ▸ which, when randomly assigning treatments, i. e. $(Y_t, Y_c) \perp\!\!\!\perp S$,
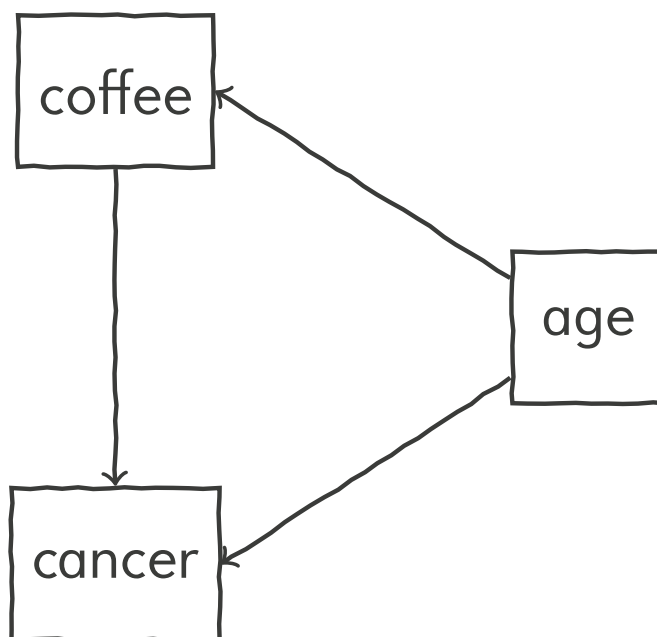- ▸ is equal to $\mathbb{E}[Y_t]$ and $\mathbb{E}[Y_c]$.

(PW Holland, Statistics and Causal Inference. *Journal of the American Statistical Association,* 1986)

---

- ‣ Split population $\mathcal{U}$ into

    - ○ 'consumed little': $S(u) = \square$
    - ○ 'consumed lots': $S(u) = \blacksquare$

- ‣ Observe whether they suffer from cancer or not, $Y \in \{0, 1\}$

- ‣ Assume older units have higher cumulative coffee consumption as well as an increased risk of cancer

- ▸ Split population $\mathcal{U}$ into
    - ○ 'consumed little': $S(u) = \square$
    - ○ 'consumed lots': $S(u) = \blacksquare$
- ▸ Observe whether they suffer from cancer or not, $Y \in \{0, 1\}$
- ▸ Assume older units have higher cumulative coffee consumption as well as an increased risk of cancer
    - ○ $(Y_\square, Y_\blacksquare) \not\perp S$
    - ○ $\mathbb{E}[Y_\square | S = \square] < \mathbb{E}[Y_\square]$

$\implies$ $\mathbb{E}[Y_\blacksquare] - \mathbb{E}[Y_\square]$ systematically overestimates the effect of cumulative coffee consumption on cancer

---

Common causal frameworks

---

- ▸ Potential Outcomes Framework

    may work under certain (untestable) assumptions

- ▸ Granger Causality

- ▸ Dynamic Causal Modelling

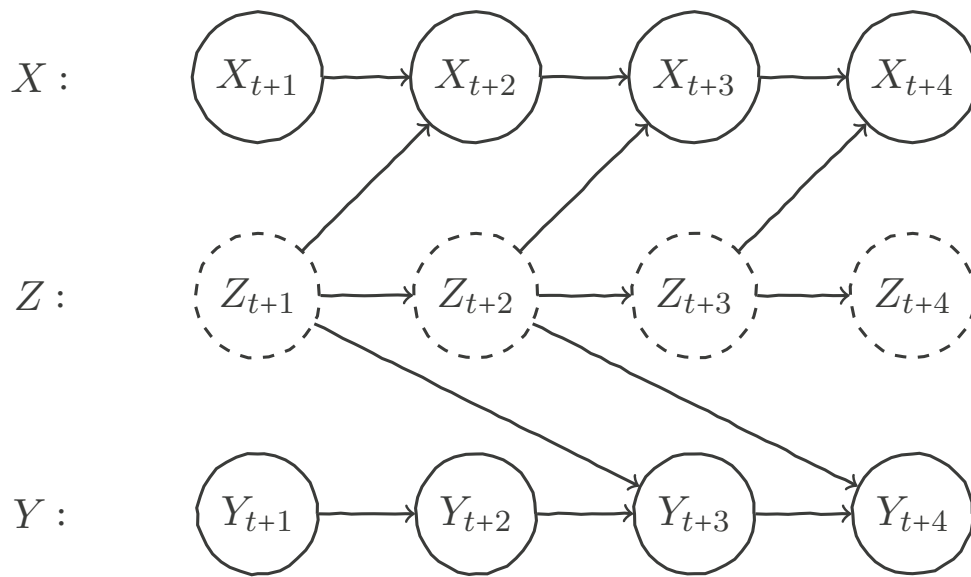- ▸ Causal Bayesian Networks and Structural Equation Models

# Granger Causality

Simplified Definition: One stochastic process $X$ is causal to a second $Y$ if the autoregressive predictability of the second process at a given time point is improved by *including* measurements from the past of the first, i. e. if

$$\mathrm{PredAcc}[Y_t|Y_{<t}] < \mathrm{PredAcc}[Y_t|Y_{<t}, X_{<t}]$$

(*not* by C Granger)

$$\mathrm{PredAcc}[Y_t|Y_{<t}] < \mathrm{PredAcc}[Y_t|Y_{<t}, X_{<t}]$$

Granger causality erroneously infers causal influence from $X$ to $Y$!

(J Peters et al. Causal discovery on time series using restricted structural equation models. *NIPS*, 2013)

Simplified Definition: One stochastic process $X$ is causal to a second $Y$ if the autoregressive predictability of the second process at a given time point is improved by *including* measurements from the past of the first, i. e. if
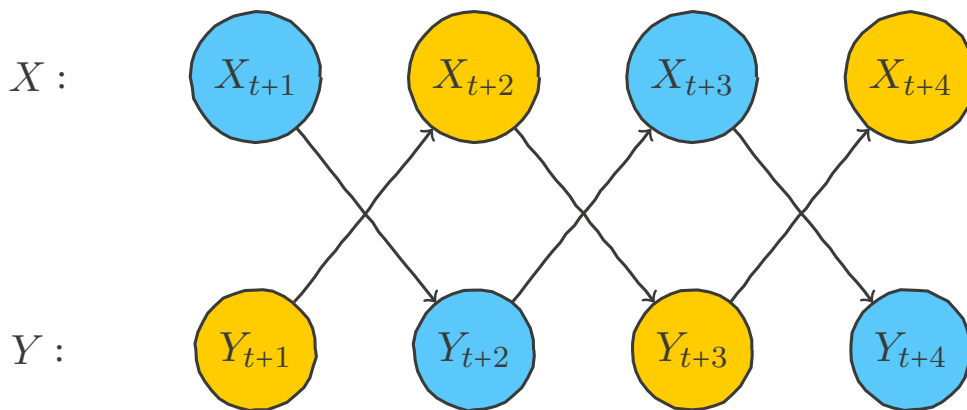
$$\mathrm{PredAcc}[Y_t|Y_{<t}] < \mathrm{PredAcc}[Y_t|Y_{<t}, X_{<t}]$$

(*not* by C Granger)

Granger's Definition: One stochastic process $X$ is causal to a second $Y$ if the predictability of the second process at a given time point is worsened by *removing* past measurements of the first from the universe's past, i. e. if

$$\mathrm{PredAcc}[Y_t|🌍_{<t}] > \mathrm{PredAcc}[Y_t|🌍_{<t} \smallsetminus X_{<t}]$$

(by C Granger)

(CWJ Granger, Investigating Causal Relations by Econometric Models and Cross-spectral Methods. *Econometrica,* 1969)

$$\mathrm{PredAcc}[Y_t | \text{🌐}_{<t}] = \mathrm{PredAcc}[Y_t | \text{🌐}_{<t} \setminus X_{<t}]$$

Granger causality fails to predict the effects of interventions!

(N Ay and D Polani, Information flows in causal networks. *Advances in Complex Systems,* 2008)

## Common causal frameworks

▸ Potential Outcomes Framework

   may work under certain (untestable) assumptions

▸ Granger Causality

   problems with confounding

   may fail to predict effects of interventions

▸ Dynamic Causal Modelling

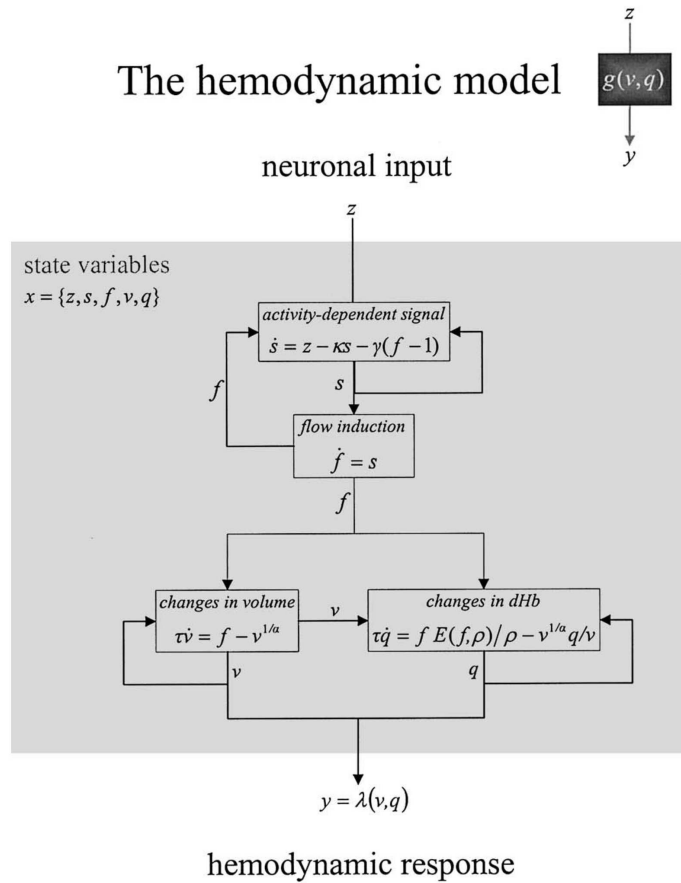▸ Causal Bayesian Networks and Structural Equation Models

# Dynamic Causal Modelling

Causality in DCM is used in a control theory sense and means that, under the model, activity in one brain area causes dynamics in another, and that these dynamics cause the observations.
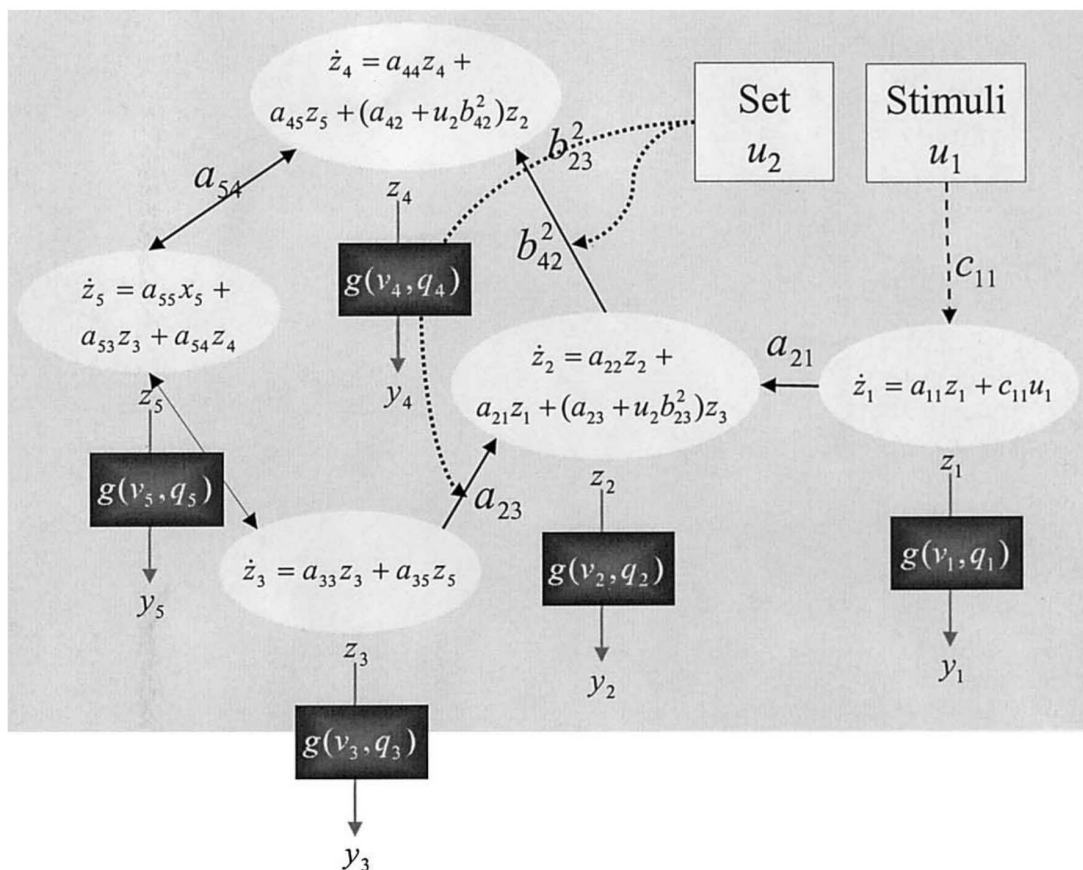
<div align="right">(Friston, <em>PLOS Biology,</em> 2009)</div>

Inference procedure:

- ▸ Observe

- ▸ Define models $\mathcal{M} = \{M_1, \ldots, M_N\}$

- ▸ Fit models to observed data

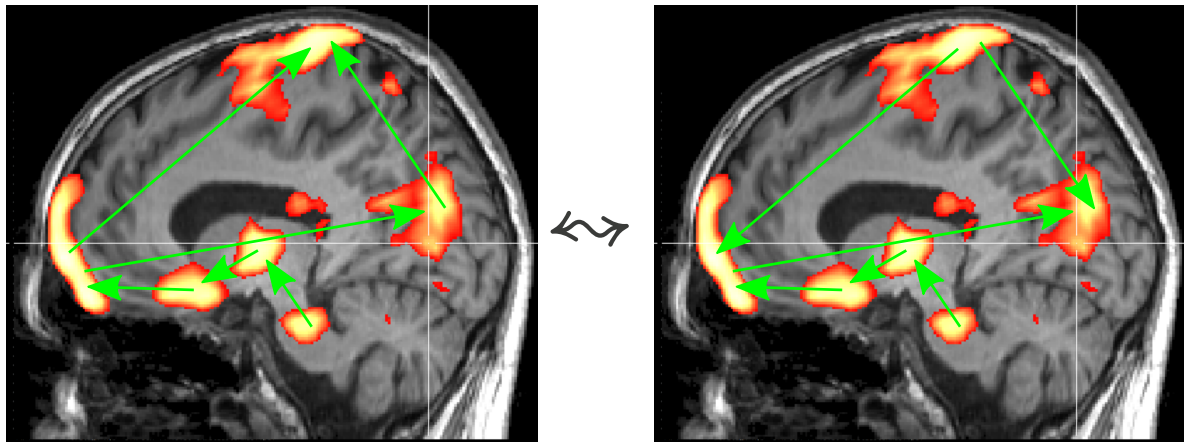- ▸ Best fitting model $\widehat{M}$ wins

## The hemodynamic model

$g(v,q)$

neuronal input

$z$

state variables
$x = \{z, s, f, v, q\}$

activity-dependent signal
$$\dot{s} = z - \kappa s - \gamma(f-1)$$

$s$

flow induction
$$\dot{f} = s$$

$f$

changes in volume
$$\tau\dot{v} = f - v^{1/\alpha}$$

$v$

changes in dHb
$$\tau\dot{q} = f\, E(f,\rho)/\rho - v^{1/\alpha} q/v$$

$q$

$$y = \lambda(v,q)$$

hemodynamic response

(KJ Friston et al., Dynamic Causal Modelling. *NeuroImage,* 2003)

16

$$\dot{z}_4 = a_{44}z_4 + a_{45}z_5 + (a_{42} + u_2 b_{42}^2)z_2$$

$a_{54}$

$z_4$

$g(v_4, q_4)$

$y_4$

$b_{23}^2$

$b_{42}^2$

Set $u_2$

Stimuli $u_1$

$c_{11}$

$$\dot{z}_5 = a_{55}x_5 + a_{53}z_3 + a_{54}z_4$$

$z_5$

$g(v_5, q_5)$

$y_5$

$$\dot{z}_2 = a_{22}z_2 + a_{21}z_1 + (a_{23} + u_2 b_{23}^2)z_3$$

$a_{21}$

$$\dot{z}_1 = a_{11}z_1 + c_{11}u_1$$

$z_1$

$g(v_1, q_1)$

$y_1$

$a_{23}$

$$\dot{z}_3 = a_{33}z_3 + a_{35}z_5$$

$z_2$

$g(v_2, q_2)$

$y_2$

$z_3$

$g(v_3, q_3)$

$y_3$

(KJ Friston et al., Dynamic Causal Modelling. *NeuroImage,* 2003)
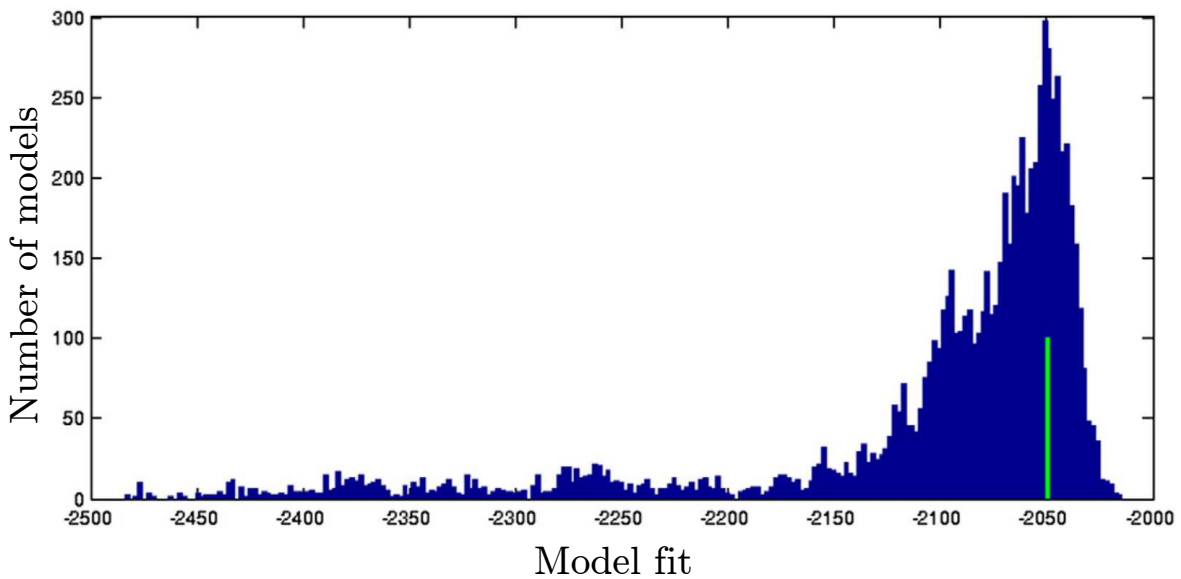
17

---

Dynamic Causal Modelling

Causality in DCM is used in a control theory sense and means that, under the model, activity in one brain area causes dynamics in another, and that these dynamics cause the observations.

(Friston, *PLOS Biology,* 2009)

Inference procedure:

- ▸ Observe

- ▸ Define models  $\mathcal{M} = \{M_1, \ldots, M_N\}$

- ▸ Fit models to observed data

- ▸ Best fitting model $\widehat{M}$ wins

Is $\widehat{M}$ guaranteed to reflect the true connectivities?



$\Longrightarrow$ Similar model fit does not translate into similar connectivities!

(Lohmann et al., Critical comments on dynamic causal modelling. *NeuroImage,* 2012)                                                    20

---

Common causal frameworks

---

‣ Potential Outcomes Framework

   may work under certain (untestable) assumptions

‣ Granger Causality

   problems with confounding

   may fail to predict effects of interventions

‣ Dynamic Causal Modelling

   unclear how it predicts interventional setting

   inference procedure provably correct?

‣ Causal Bayesian Networks and Structural Equation Models

# Causal Bayesian Networks and Structural Equation Models

A Structural Equation Model (SEM) $\mathcal{M}_X = (\mathcal{S}_X, \mathcal{I}_X, \mathbb{P}_{E_X})$ with

- ▸ structural equations $\mathcal{S}_X$;
- ▸ a set of interventions $\mathcal{I}_X$;
- ▸ exogenous variables distributed according to $\mathbb{P}_{E_X}$

induces distributions $\mathbb{P}_X$ over the $X$ variables for each $i \in \mathcal{I}_X$.

(J Pearl, *Causality: Models, reasoning, and inference,* 2000; P Spirtes et al., *Causation, Prediction, and Search,* 2001)

$$\mathcal{M}_X = (\mathcal{S}_X, \mathcal{I}_X, \mathbb{P}_{E_X})$$

‣ $\mathcal{S}_X = \begin{cases} X_1 = E_1 \\ \\ X_2 = X_1 + E_2 \end{cases}$

‣ $\mathcal{I}_X = \{\varnothing, \ \mathrm{do}(X_1 = 5), \ \mathrm{do}(X_2 = 3)\}$

‣ $E \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

| observational | intervention on $X_1$ | intervention on $X_2$ |
|---|---|---|
| $\mathbb{P}_{X_1}^{\varnothing} \sim \mathbb{N}(0, 1)$ | $\mathbb{P}_{X_1}^{\mathrm{do}(X_1=5)} \equiv 5$ | $\mathbb{P}_{X_1}^{\mathrm{do}(X_2=3)} \sim \mathbb{N}(0, 1)$ |
| $\mathbb{P}_{X_2}^{\varnothing} \sim \mathbb{N}(0, 2)$ | $\mathbb{P}_{X_2}^{\mathrm{do}(X_1=5)} \sim \mathbb{N}(5, 1)$ | $\mathbb{P}_{X_2}^{\mathrm{do}(X_2=3)} \equiv 3$ |

(J Pearl, *Causality: Models, reasoning, and inference,* 2000; P Spirtes et al., *Causation, Prediction, and Search,* 2001)

23

## Causal Bayesian Networks

Definition of Cause and Effect
$$X \to Y \iff \mathbb{P}_Y^{\mathrm{do}(X=x)} \neq \mathbb{P}_Y^{\varnothing} \text{ for some } x$$
Causal Markov Condition
$$\text{d-separation} \rightsquigarrow \text{independence}$$
Faithfulness
$$\text{d-separation} \leftsquigarrow \text{independence}$$

| chain | fork | collider |
|---|---|---|
| $X \to Y \to Z$ | $X \leftarrow Y \to Z$ | $X \to Y \leftarrow Z$ |
| $X \not\perp\!\!\!\perp Z$ | $X \not\perp\!\!\!\perp Z$ | $X \perp\!\!\!\perp Z$ |
| $X \perp\!\!\!\perp Z \mid Y$ | $X \perp\!\!\!\perp Z \mid Y$ | $X \not\perp\!\!\!\perp Z \mid Y$ |

(J Pearl, *Causality: Models, reasoning, and inference,* 2000; P Spirtes et al., *Causation, Prediction, and Search,* 2001)

24

- ▸ Randomised stimulus $S$

- ▸ Observe neural activity $X$ and $Y$

- ⤳ Estimate $\mathbb{P}^{\varnothing}_{S,X,Y}$

- ▸ Assume we find

  - ○ $S \not\perp X \implies$ existence of path between $S$ and $X$ w/o collider
  - ○ $S \not\perp Y \implies$ existence of path between $S$ and $Y$ w/o collider
  - ○ $S \perp Y|X \implies$ all paths between $S$ and $Y$ blocked by $X$

- ▸ Can rule out cases such as $S \to X \leftarrow h \to Y$

- ▸ Can formally prove that $X$ indeed is a cause of $Y$

- $\implies$ Robust against hidden confounding

## Application: Neural Dynamics of Probabilistic Reward Prediction

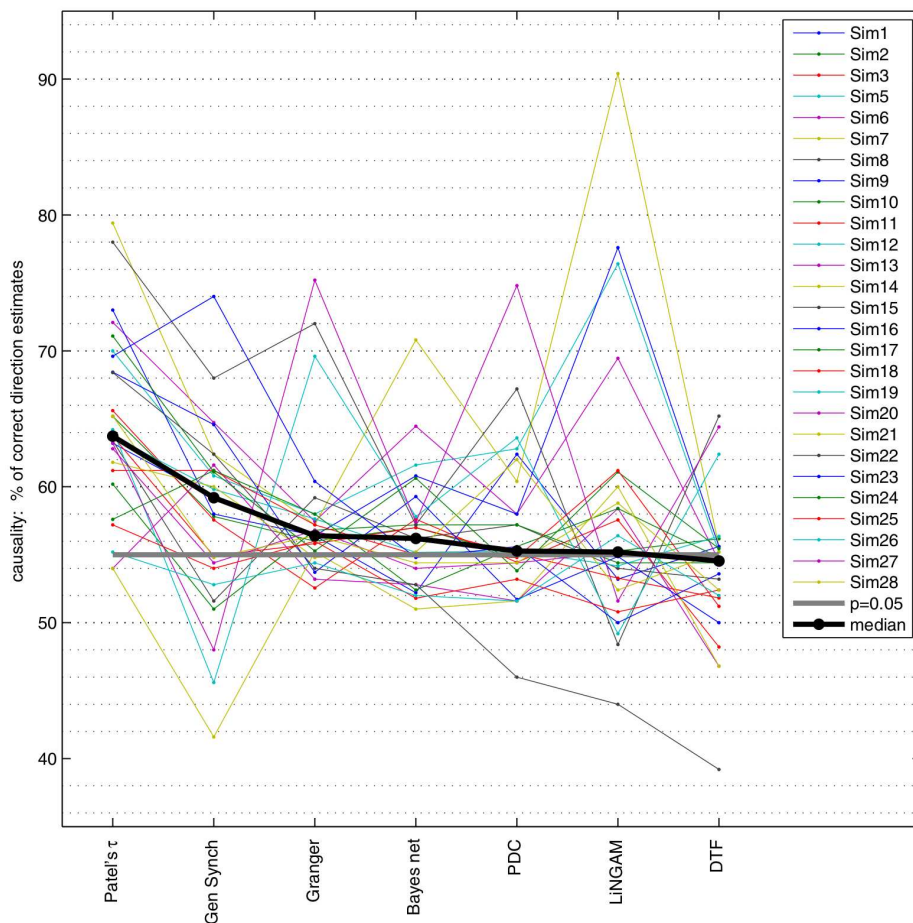Bach et al. ● Probabilistic Reward Prediction

(Bach, Symmonds, Barnes, and Dolan, Whole-brain neural dynamics of probabilistic reward prediction. *Journal of Neuroscience,* 2017)     27

---

## Common causal frameworks

---

‣ Potential Outcomes Framework

  may work under certain (untestable) assumptions

‣ Granger Causality

  problems with confounding

  may fail to predict effects of interventions

‣ Dynamic Causal Modelling

  unclear how it predicts interventional setting

  inference procedure provably correct?

‣ Causal Bayesian Networks and Structural Equation Models

  may work under certain (untestable) assumptions
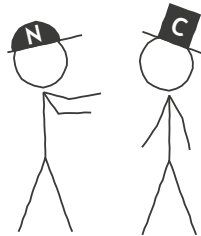
  not finding dependence is not evidence for independence

# Wrap-Up

- ▸ (Causal) Inference rests on *untestable* assumptions.

- ▸ Causal inference algorithms appear to perform above chance-level.

- ▸ Causal inference may be useful to guide the design of interventional studies.
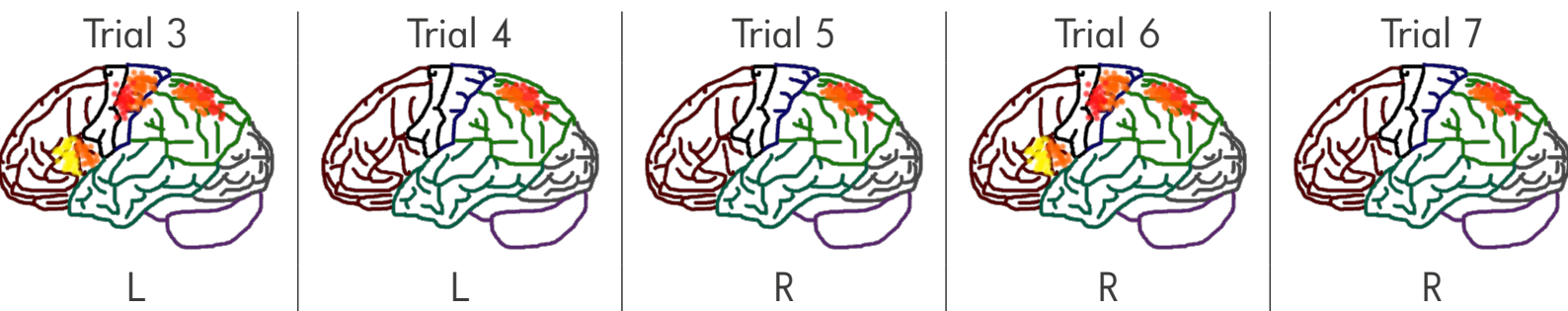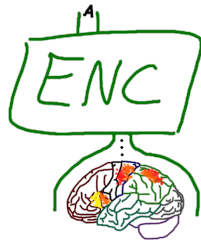
sweichwald.de/prni2017

neural.engineering

# ADDENDA

# Causal interpretation of encoding and decoding models

## Relevance in encoding and decoding models



| Trial 3 | Trial 4 | Trial 5 | Trial 6 | Trial 7 |
|---------|---------|---------|---------|---------|
| L | L | R | R | R |

"Significant variation explained by experimental condition?"

$$X_i \not\perp C$$

$$X_i \not\perp C | \vec{X} \setminus X_i$$

"Does removal impair decoding performance?"
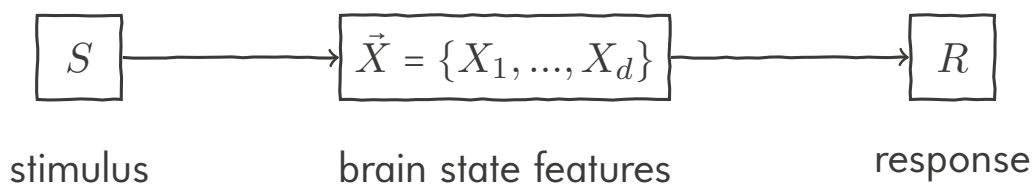
$$\text{relevant feature} \overset{?}{\leftrightsquigarrow} \text{cognitive process}$$

(S Weichwald et al., Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage,* 2015)　　32
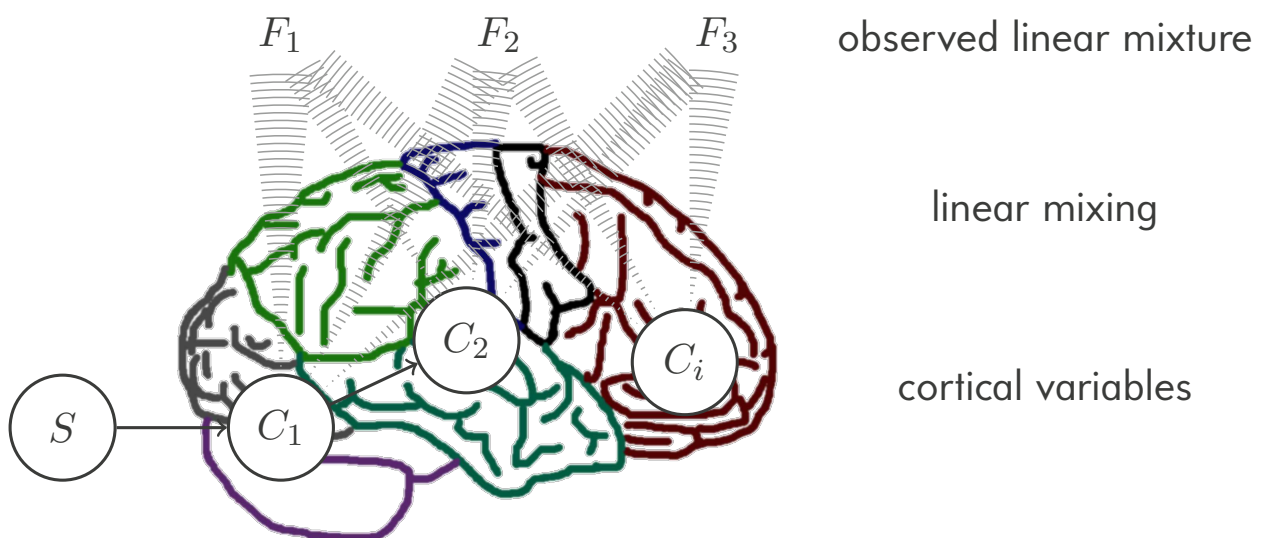
---

A new distinction: stimulus- vs response-based



$$S \longrightarrow \vec{X} = \{X_1, ..., X_d\} \longrightarrow R$$

stimulus　　　brain state features　　　response

| stimulus-based | | response-based |
|---|---|---|
| causal | encoding | *anti*-causal |
| *anti*-causal | decoding | causal |

(S Weichwald et al., Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage,* 2015)　　33

| | Feature $X_i$ relevant? | | Causal interpretation |
|---|---|---|---|
| | Encoding | Decoding | |
| Stimulus-based | × | | no effect of $S$ |
| | √ | | effect of $S$ |
| | | × | inconclusive |
| | | √ | inconclusive |
| Response-based | × | | no cause of $R$ |
| | √ | | inconclusive |
| | | × | inconclusive |
| | | √ | inconclusive |

| | Feature $X_i$ relevant? | | Causal interpretation |
|---|---|---|---|
| | Encoding | Decoding | |
| Stimulus-based | √ | √ | effect of $S$ |
| | √ | × | indirect effect of $S$ |
| | × | √ | provides context |
| | × | × | no effect of $S$ |
| Response-based | √ | √ | inconclusive |
| | √ | × | no direct cause of $R$ |
| | × | √ | provides context |
| | × | × | no cause of $R$ |

# MERL★iN

$F_1$ $F_2$ $F_3$     observed linear mixture

linear mixing

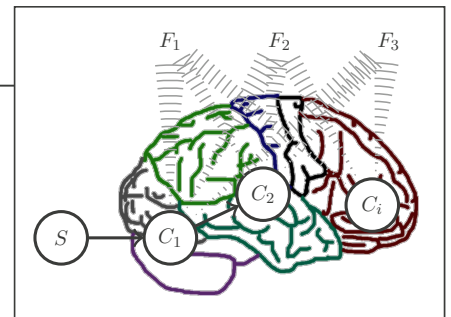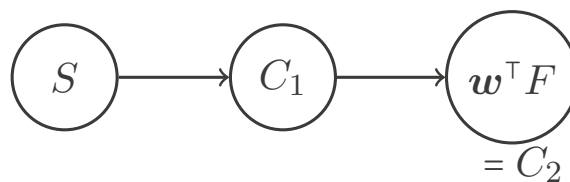$S$ $C_1$ $C_2$ $C_i$     cortical variables

*Given*

samples of $S, C_1$ and $F$

$$F = \begin{bmatrix} F_1 \\ \vdots \\ F_d \end{bmatrix} = \boldsymbol{A} \begin{bmatrix} C_1 \\ \vdots \\ C_d \end{bmatrix} = \boldsymbol{A}C$$

*Goal*

find linear combination $\boldsymbol{w}$ such that

*Idea*

Optimise $\boldsymbol{w}$ such that

(a)  $\mathrm{dep}\,(C_1, \boldsymbol{w}^\top F)$     is high
(b)  $\mathrm{dep}\,(S,\ \boldsymbol{w}^\top F\,|\,C_1)$ is low

*Implementation*

Optimise $\boldsymbol{w}$ and $\sigma, \theta$ such that

$$\mathrm{HSIC}\,(C_1, \boldsymbol{w}^\top F) \qquad\qquad \text{is high}$$
$$-\ \mathrm{HSIC}\big(\ \boldsymbol{w}^\top F - \mathrm{krr}_{\sigma,\theta}(C_1)\ ,\ (S, C_1)\ \big) \text{ is low}$$

is being maximised.