# Using neuroimaging to infer mental states: A guided tour through the minefield

Russell Poldrack

Department of Psychology
Stanford University

# Can neuroimaging tell us anything about the mind?

Max Coltheart

"No amount of knowledge about the hardware of a computer will tell you anything serious about the nature of the software that the computer runs. In the same way, no facts about the activity of the brain could be used to confirm or refute some information-processing model of cognition." (Coltheart, 2004, p. 22)
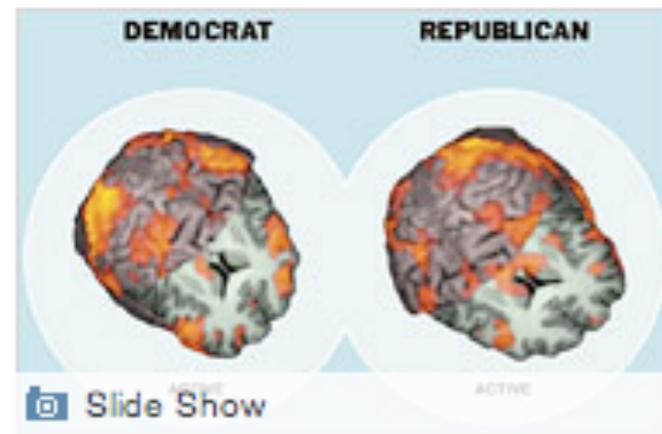
OP-ED CONTRIBUTORS
# This Is Your Brain on Politics

*This article was written by Marco Iacoboni, Joshua Freedman and Jonas Kaplan of the University of California, Los Angeles, Semel Institute for Neuroscience; Kathleen Hall Jamieson of the Annenberg Public Policy Center at the University of Pennsylvania; and Tom Freedman, Bill Knapp and Kathryn Fitzgerald of FKF Applied Research.*

**Multimedia**



DEMOCRAT    REPUBLICAN

Slide Show

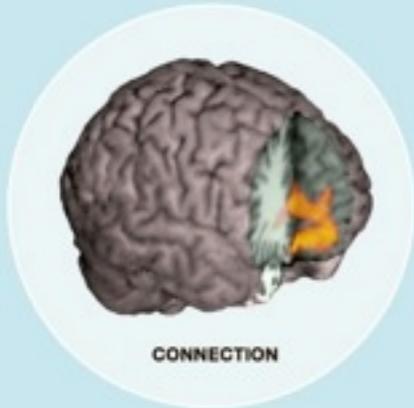This Is Your Brain on Politics

IN anticipation of the 2008 presidential election, we used functional magnetic resonance imaging to watch the brains of a group of swing voters as they responded to the leading presidential candidates. Our results reveal some voter impressions on which this election may well turn.

Our 20 subjects — registered voters who stated that they were open to choosing a candidate from either party next November — included 10 men and 10 women. In late summer, we asked them to answer a list of questions about their political preferences, then observed their brain activity for nearly an hour in the scanner at the Ahmanson Lovelace Brain Mapping Center at the University of California, Los Angeles. Afterward, each subject filled out a second questionnaire.
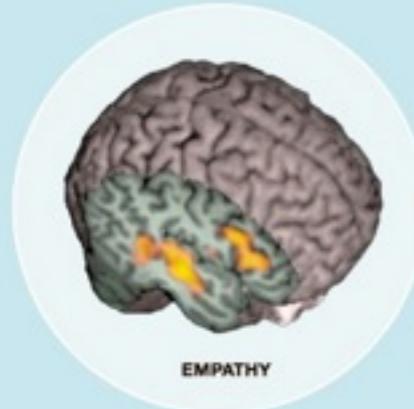
"In response to images of Democratic candidates, men exhibited activity in the medial orbital prefrontal cortex, indicating emotional connection and positive feelings."

"Images of Fred Thompson led to increased activity in the inferior frontal cortex, a brain structure associated with empathy."

"Subjects who had an unfavorable view of John Edwards responded to pictures of him with feelings of disgust, evidenced by increased activity in the insula, a brain area associated with negative emotions."

# Do you really love your iPhone?

buy·ology

Truth and Lies About
Why We Buy

MARTIN LINDSTROM

Foreword by Paco Underhill

- "Earlier this year, I carried out an fMRI experiment to find out whether iPhones were really, truly addictive, no less so than alcohol, cocaine, shopping or video games. In conjunction with the San Diego-based firm MindSign Neuromarketing, I enlisted eight men and eight women between the ages of 18 and 25. Our 16 subjects were exposed separately to audio and to video of a ringing and vibrating iPhone...most striking of all was the flurry of activation in the insular cortex of the brain, which is associated with feelings of love and compassion. The subjects' brains responded to the sound of their phones as they would respond to the presence or proximity of a girlfriend, boyfriend or family member.  In short, the subjects didn't demonstrate the classic brain-based signs of addiction. Instead, they loved their iPhones.

**To the Editor:**

"[You Love Your iPhone. Literally,](#)" by Martin Lindstrom (Op-Ed, Oct. 1), purports to show, using brain imaging, that our attachment to digital devices reflects not addiction but instead the same kind of emotion that we feel for human loved ones.

However, the evidence the writer presents does not show this.

The brain region that he points to as being "associated with feelings of love and compassion" (the insular cortex) is active in as many as one-third of all brain imaging studies.

Further, in studies of decision making the insular cortex is more often associated with negative than positive emotions.

The kind of reasoning that Mr. Lindstrom uses is well known to be flawed, because there is rarely a one-to-one mapping between any brain region and a single mental state; insular cortex activity could reflect one or more of several psychological processes.

We find it surprising that The Times would publish claims like this that lack scientific validity.

RUSSELL POLDRACK
Austin, Tex., Oct. 3, 2011

*The writer is a professor of psychology and neurobiology at the University of Texas at Austin. His letter was signed by 44 other neuroscientists.*

Insula
activity

# Does reverse inference work?



effort

craving

pain

Insula activity

p(process|act)

p(act|process)

$$p(process|act) = \frac{p(process) \star p(act|process)}{p(act)}$$

Some voxels active in more than 20% of studies

Yarkoni et al., 2011

- Informal reverse inference provides relatively weak evidence

*TICS*, 2006

# Formalizing reverse inference

- How can we more formally test the predictive ability of fMRI?

- Answer: statistical methods for prediction

  - Machine learning/statistical learning/pattern recognition

**fMRI dataset**

train to classify mental states

test accuracy of decoding on untrained data

*Cross-validation*:
- Train for each split of size k
- Compute average predictive accuracy on left-out data

# 96% correct classification

Haxby et al., 2001, *Science*

Free task selection (addition versus subtraction)

select

Task stimuli

Response mapping

56
33

65          89

23          18

variable delay

Train on 7 runs, test on 8th

r=3

spherical cluster

$v_i$

decoding

DELAY

0.5    accuracy    0.75

EXECUTION

Decoding accuracy

0.7

0.6

0.5 — chance

0.4

MPFCa    MPFCp    LLFPC    LIFS    RMFG    LFO

Haynes et al., 2007, *Current Biology*

**Mind-reading machine knows what you see**

15:26 25 April 2005
NewScientist.com news service

April 25, 2005

**Brain Scans Helps Scientists "Read" Minds**

# BBC NEWS

OPEN The News in 2 minutes

News Front Page

Africa
Americas
Asia-Pacific
Europe
Middle East
South Asia
UK
Business
Health
Medical notes
Science/Nature
Technology
Entertainment

Last Updated: Monday, 25 April, 2005, 00:05 GMT 01:05 UK

✉ E-mail this to a friend          🖶 Printable version

## Brain scan 'sees hidden thoughts'

Scientists say they can read a person's unconscious thoughts using a simple brain scan.

Functional MRI scans plot brain activity by looking at brain blood flow and are already used by researchers.

A team at University College London found with fMRI they could tell what a person was thinking deep down even when the individual was unaware themselves.

The scan picks up subliminal thought activity

# "60 Minutes", January 4, 2009

# "60 Minutes", January 4, 2009

"It's tough to make predictions, especially about the future." - Yogi Berra

- Existing work has primarily examined ability to predict mental states using a classifier trained on data from the same person

- For many applications of interest, such training data would not exist for the individual being tested

- Can we accurately generalize to new individuals?

# Predicting risky decisions

## Balloon Analog Risk Task (BART)

Balloon #1

Choice: Pump

Choice: Pump

Pre-cashout trial

Choice: Cash Out

This Balloon: 10
Total: 10

Reward

Balloon #2

Choice: Pump

Choice: Pump

Pre-pump trial

Choice: Pump

Choice: Pump

This Balloon: 0
Total: 10

Explode

Helfinstein et al, 2014, PNAS

# Crossvalidation across subjects

18 subjects

18 subjects

18 subjects

18 subjects

18 subjects

18 subjects

Train on 5 folds

Test on left-out fold

Randomly assign to folds 50 times and average results

Helfinstein et al, 2014, PNAS

# Classification accuracy for risk-taking

## Searchlight classification accuracy



Whole-brain
classification:
72%
p<0.002 under
null hypothesis
(by randomization)

Helfinstein et al, 2014, PNAS

# Classifying based on activity balance



Logistic regression:
67% accuracy

Blue: Pre-Pump > Pre-Cashout
Red: Pre-Cashout > Pre-Pump
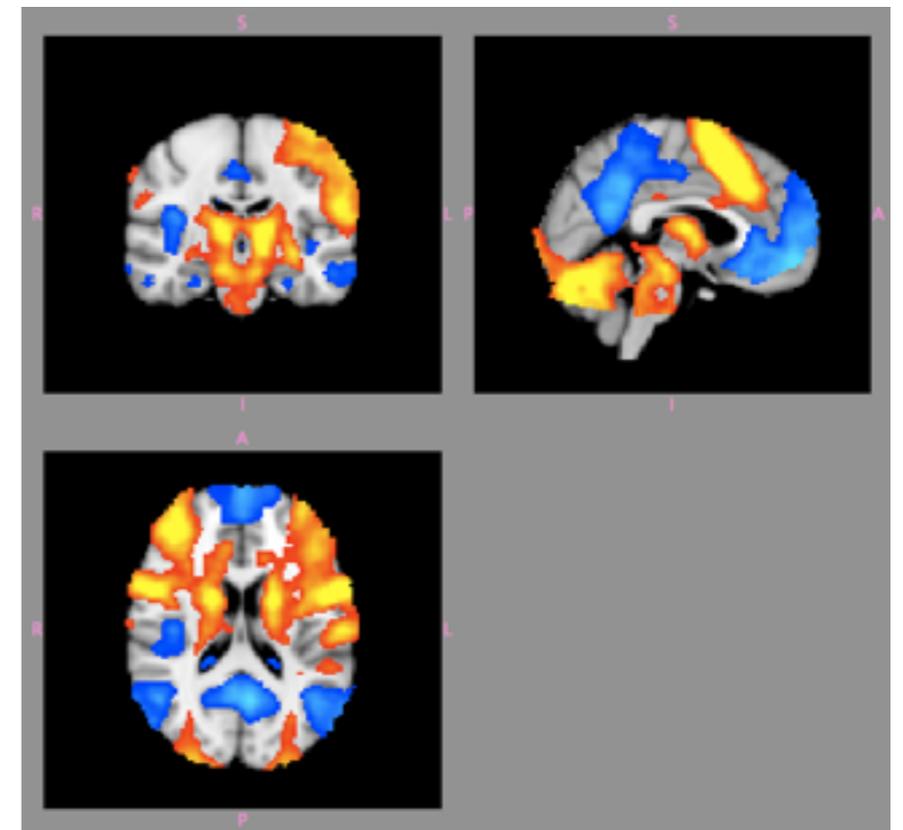
Helfinstein et al, 2014, PNAS

- Can we predict what task a subject was performing, using a classifier trained on other people?

  - 8 tasks, 130 subjects total

| Task # | Task description | # subjects | Design type |
|---|---|---|---|
| 1 | Risky decision making (Balloon analog risk task) (Stover et al., in preparation) | 16 | Event-related |
| 2 | Probabilistic classification (no feedback) (Aron et al., unpublished) | 20 | Event-related |
| 3 | Rhyme judgments on pseudowords (Xue et al., unpublished) | 13 | Event-related |
| 4 | Working memory (tone counting) (Foerde et al., 2006) | 17 | Event-related |
| 5 | 50/50 gain-loss gamble decisions (Tom et al., 2007) | 16 | Blocked |
| 6 | Living/nonliving decision on mirror-reversed words (Poldrack et al., unpublished) | 14 | Blocked |
| 7 | Reading pseudowords aloud (Xue et al., submitted) | 19 | Event-related |
| 8 | Response inhibition (successful stopping) (Aron & Poldrack, 2006) | 15 | Event-related |

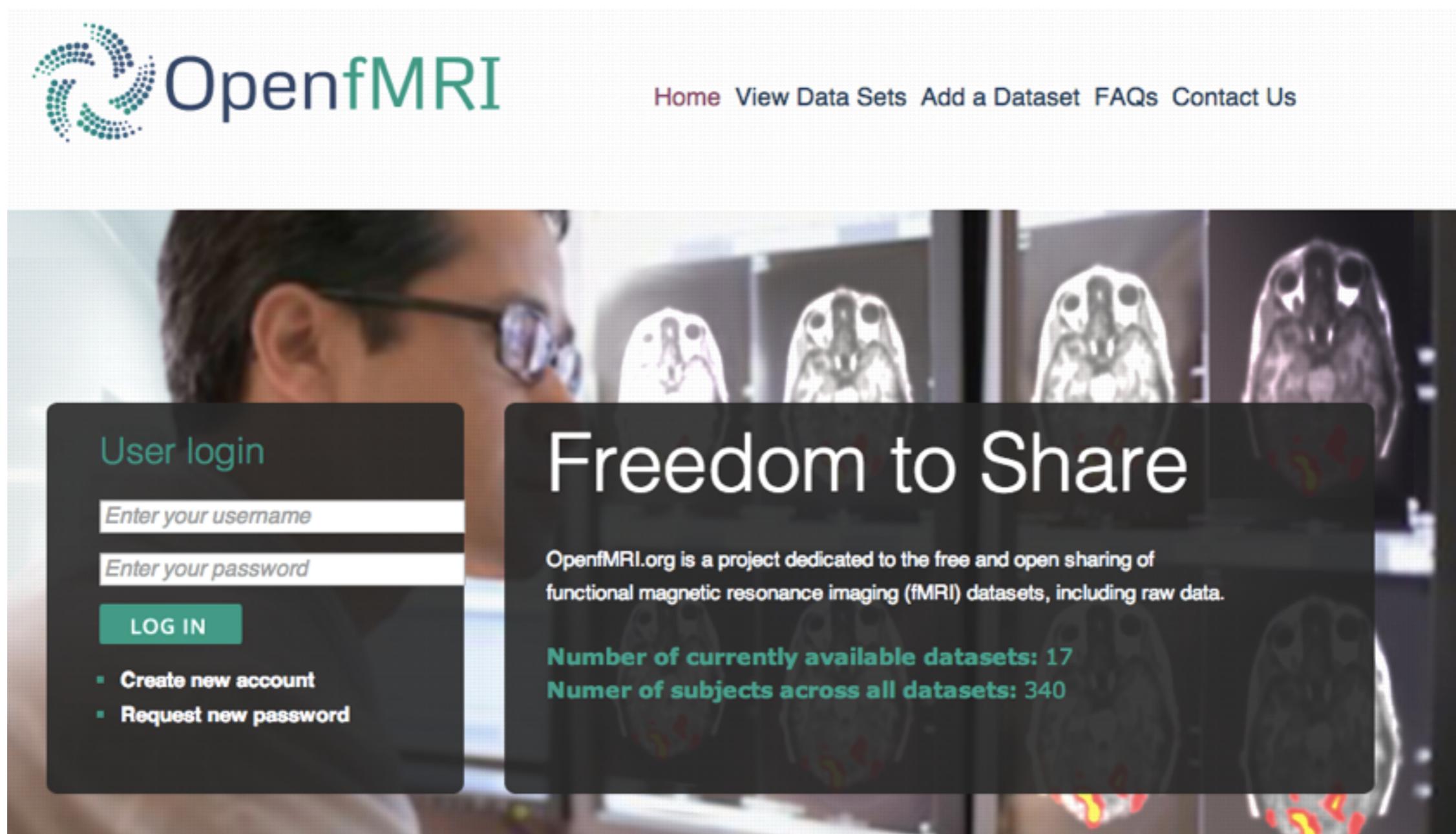| Analysis | Crossvalidated accuracy | # of voxels included |
|---|---|---|
| Union of all in-mask voxels across subjects (one-vs-one) | 74% | 417,231 |
| Intersection of in-mask voxels across subjects (one-vs-many) | 80.8% | 214,940 |
| Positively activated voxels only (across all 130 subjects, t > 3, p<.002) (one-vs-many) | 74.6% | 83,825 |
| Deactivated voxels only (t < -3, p<.002) (one-vs-many) | 50.8% | 23,736 |

Accuracy above 18.5% is significant at p<.05 by randomization



## Task chosen by classfier

| True task | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 |
|---|---|---|---|---|---|---|---|---|
| Task 1 | **87.5** | 6.0 | 0.0 | 0.0 | 6.0 | 0.0 | 0.0 | 0.0 |
| Task 2 | 0.0 | **90.0** | 0.0 | 0.0 | 0.0 | 0.0 | 5.0 | 5.0 |
| Task 3 | 8.0 | 23.0 | **61.5** | 0.0 | 0.0 | 8.0 | 0.0 | 0.0 |
| Task 4 | 0.0 | 0.0 | 0.0 | **82.4** | 0.0 | 0.0 | 0.0 | 18.0 |
| Task 5 | 0.0 | 38.0 | 0.0 | 0.0 | **43.8** | 18.2 | 0.0 | 0.0 |
| Task 6 | 0.0 | 28.0 | 0.0 | 0.0 | 0.0 | **71.4** | 0.0 | 0.0 |
| Task 7 | 0.0 | 11.0 | 0.0 | 0.0 | 0.0 | 0.0 | **84.0** | 5.0 |
| Task 8 | 0.0 | 0.0 | 7.0 | 0.0 | 0.0 | 0.0 | 27.0 | **63.0** |

26 tasks, 482 images from 338 subjects

# Classification results



Whole-brain:
47% accuracy

Poldrack et al., 2013, *Frontiers in Neuroinformatics*

- In neuroscience, correlations are often colloquially described as "prediction", but true prediction requires generalization to new samples
- The ability to predict quantitative variables for new individuals can be tested using crossvalidation
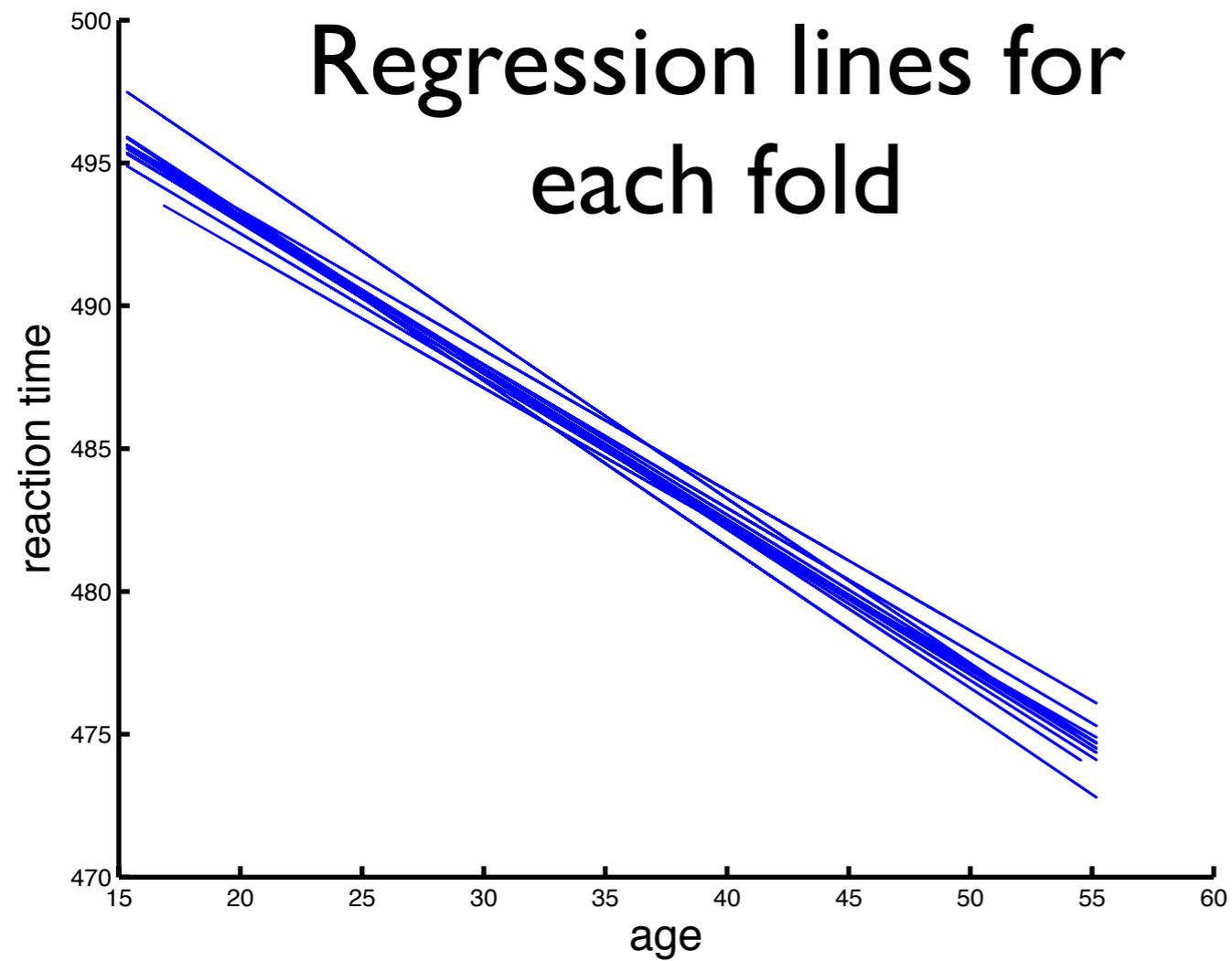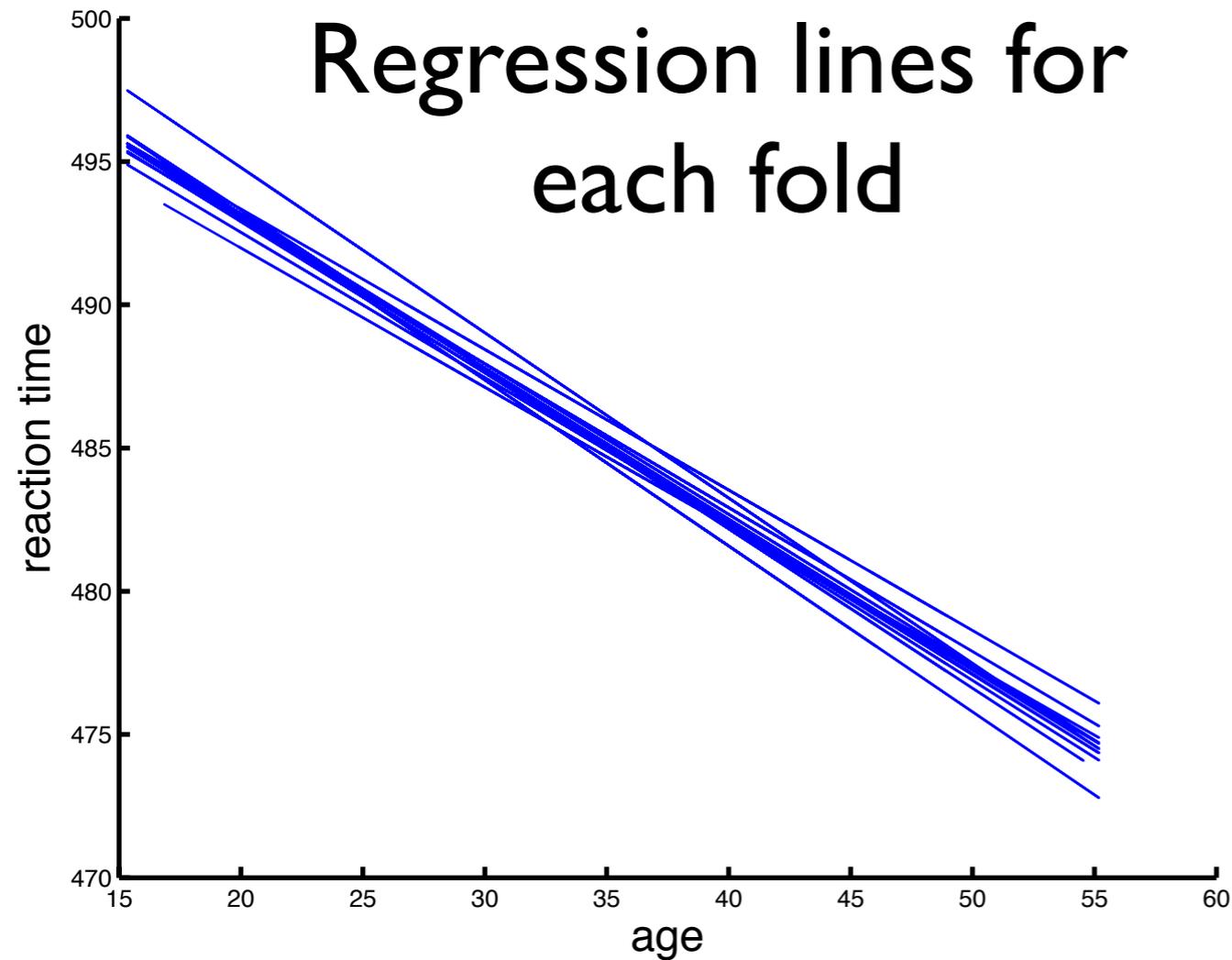
$$RT = 503.7 + age*{-}0.530$$

full-sample $R^2 = 0.694$

Regression lines for
each fold

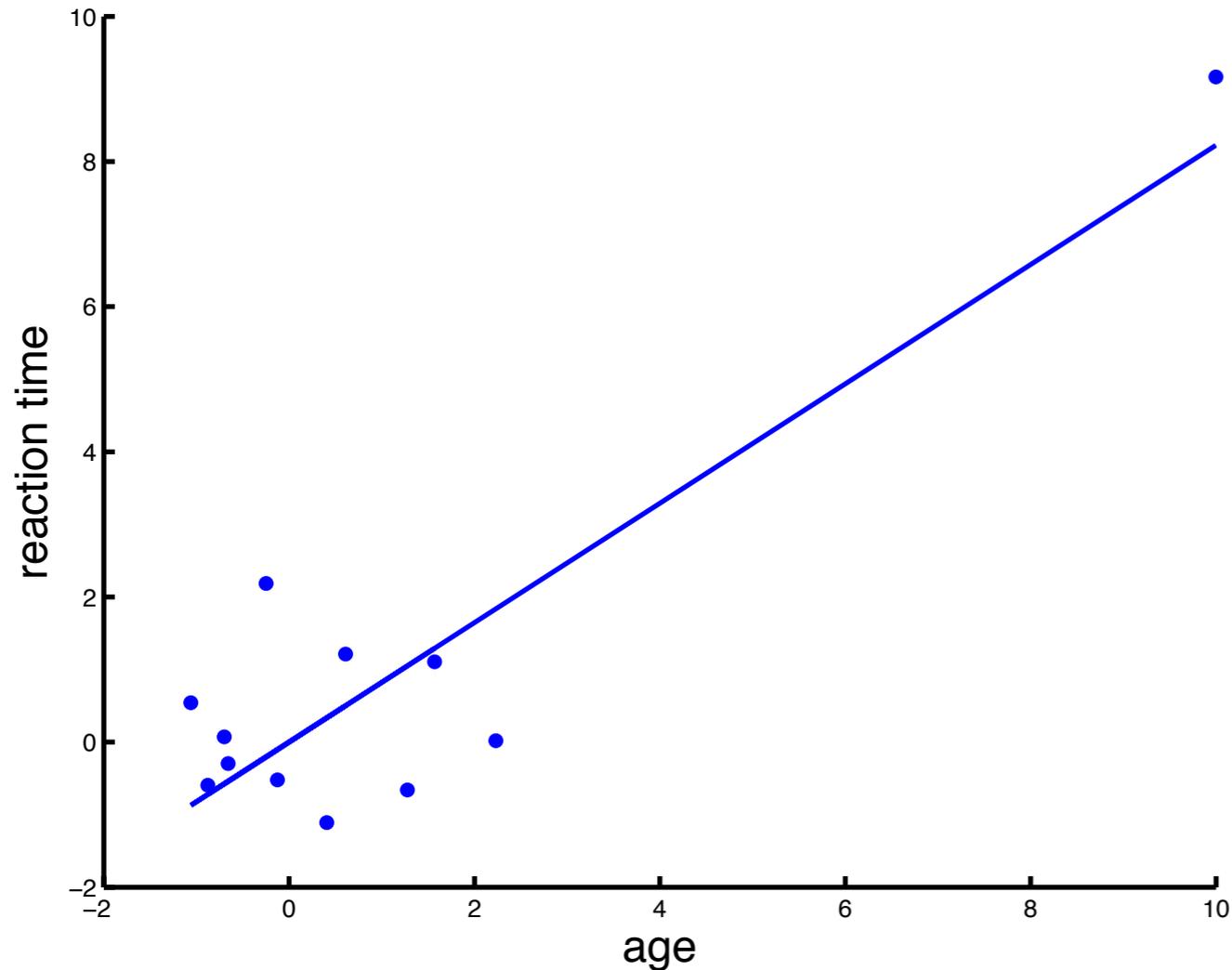full-sample $R^2 = 0.694$

Regression lines for each fold

full-sample $R^2 = 0.694$

LOO CV $R^2 = 0.586$
mean new sample $R^2 = 0.591$

# Correlation ≠ Prediction



- 11 datapoints sampled from normal distribution (no correlation)
- one outlier
- full-sample $R^2$ =0.785
- LOO CV $R^2$ = 0.025

- Highlights importance of examining the raw data!

⌂ > Early Edition > Eyal Aharoni

# Neuroprediction of future rearrest

Eyal Aharoni[a,b,1,2], Gina M. Vincent[c], Carla L. Harenski[a], Vince D. Calhoun[a,d], Walter Sinnott-Armstrong[e], Michael S. Gazzaniga[f], and Kent A. Kiehl[a,b,2]

Author Affiliations ⌃

"The present analysis shows that hemodynamic activity within the brain prospectively predicted rearrest in an offender sample."
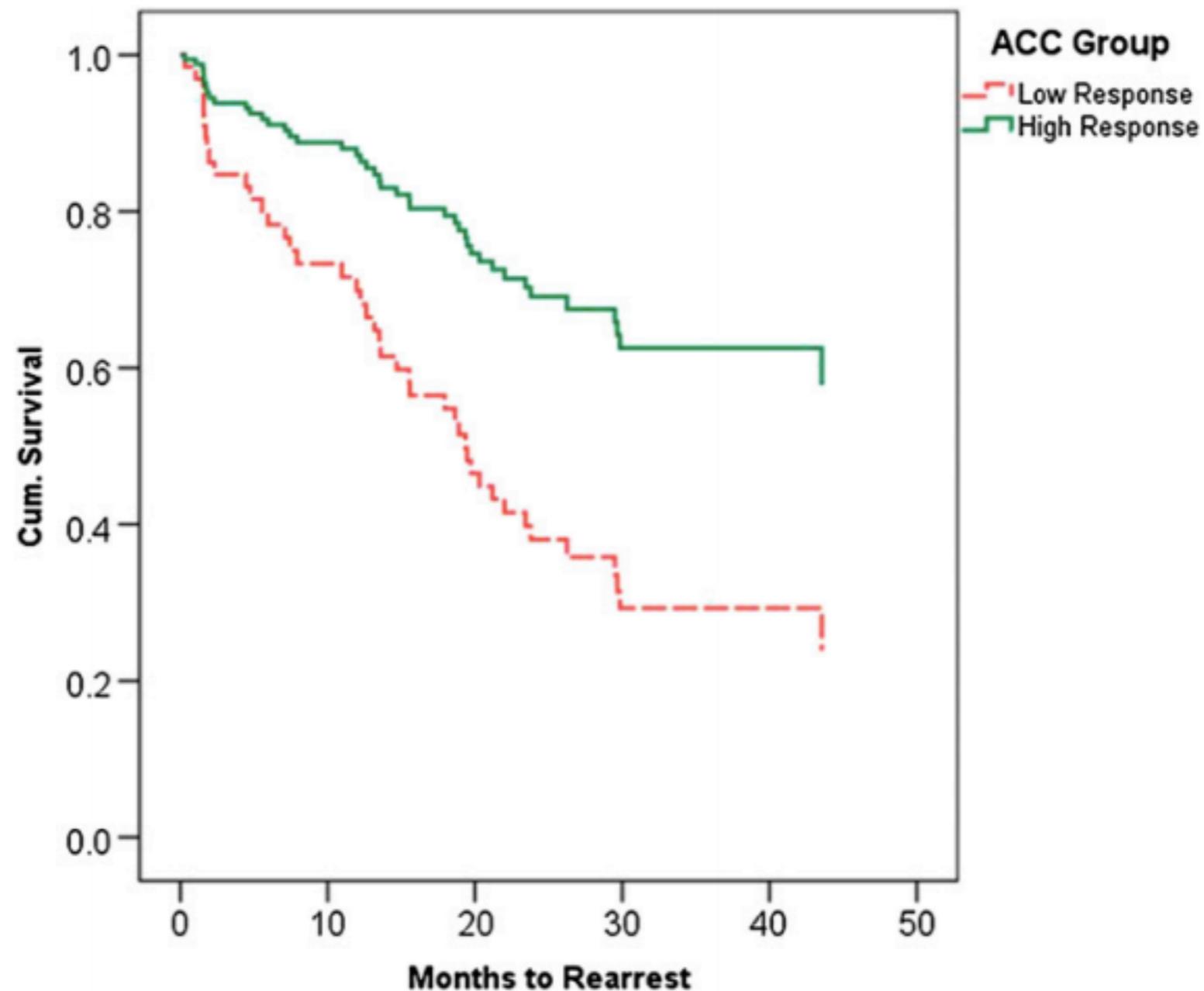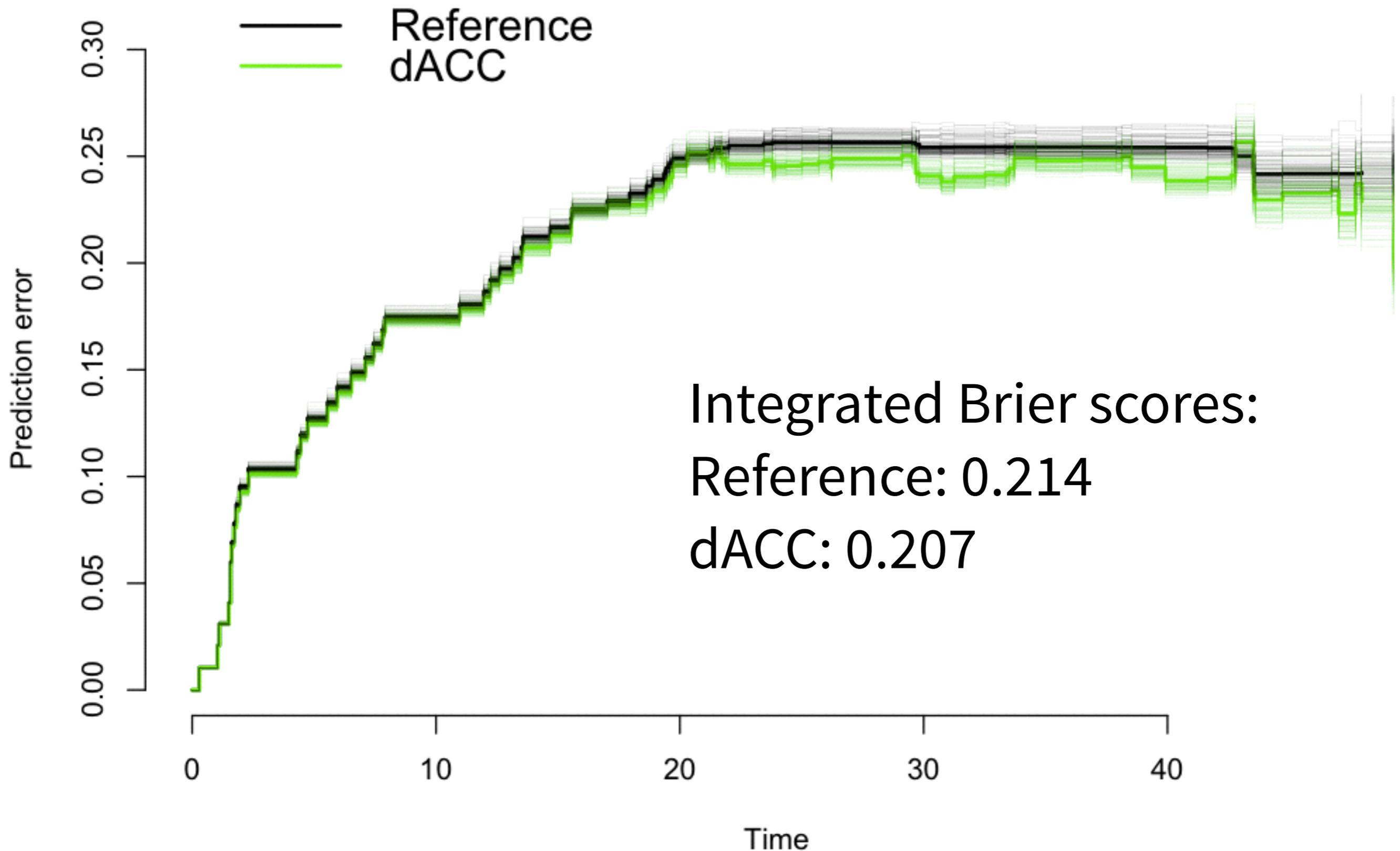
**Fig. 1.** Cox survival function showing proportional rearrest survival rates of high (solid green) vs. low (dashed red) ACC response groups for any crime over a 4-y period. Results of this median split analysis were equivalent to that of the parametric model: bootstrapped $B = 0.96$; SE $= 0.40$; $P < 0.01$; 95% CI, 0.29–1.84. The mean survival times to rearrest for the low and high ACC activity groups were 25.27 (2.80) mo and 32.42 (2.73) mo, respectively. The overall probabilities of rearrest were 60% for the low ACC group and 46% for the high ACC group.
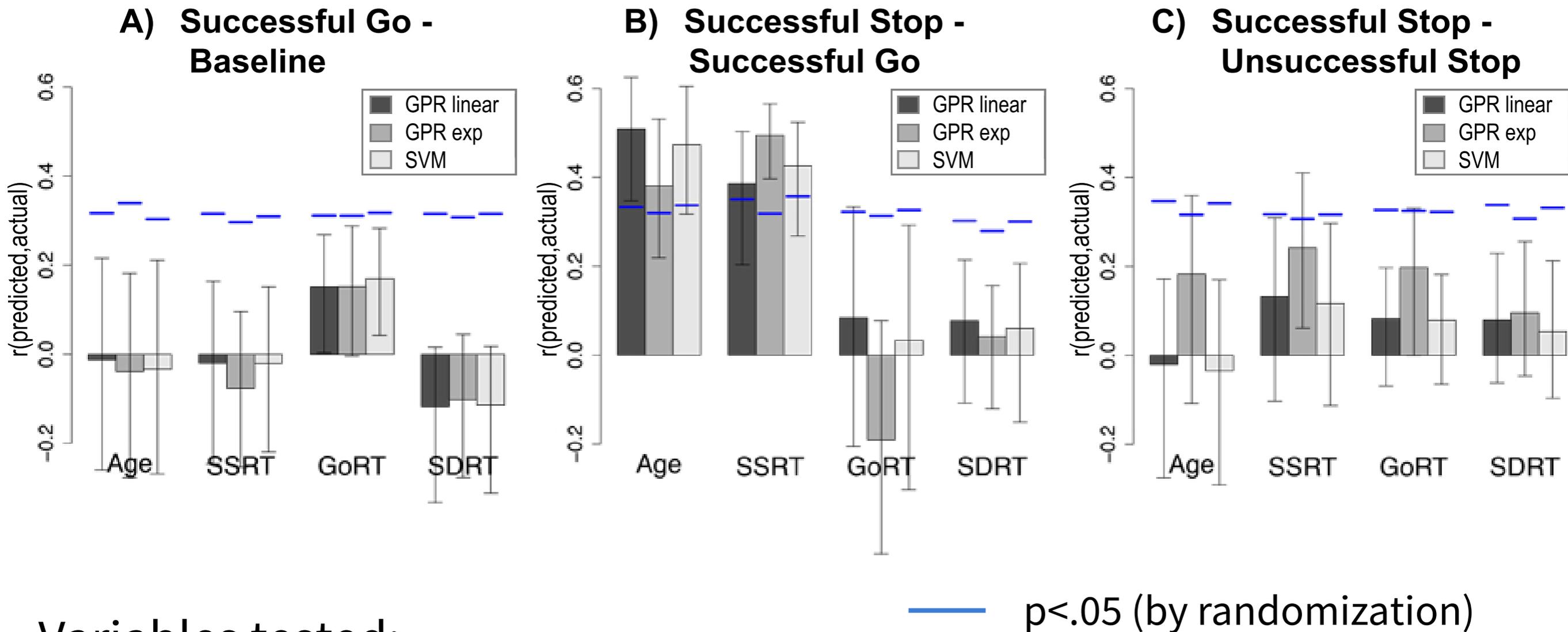
Aharoni et al., 2013

Prediction error using crossvalidation

Integrated Brier scores:
Reference: 0.214
dACC: 0.207

# Predicting individual differences from fMRI



Variables tested:

Age: subject's age

SSRT: stop signal reaction time

GoRT: go reaction time

SDRT: std. dev. of go reaction time

Cohen et al., 2010, *Frontiers in Human Neuroscience*

# Meta-analytic decoding

- All of the results to this point were based on fMRI data from individual subjects

- Can we push this even further?

  - Can we use meta-analytic data from papers?

# Activation locations

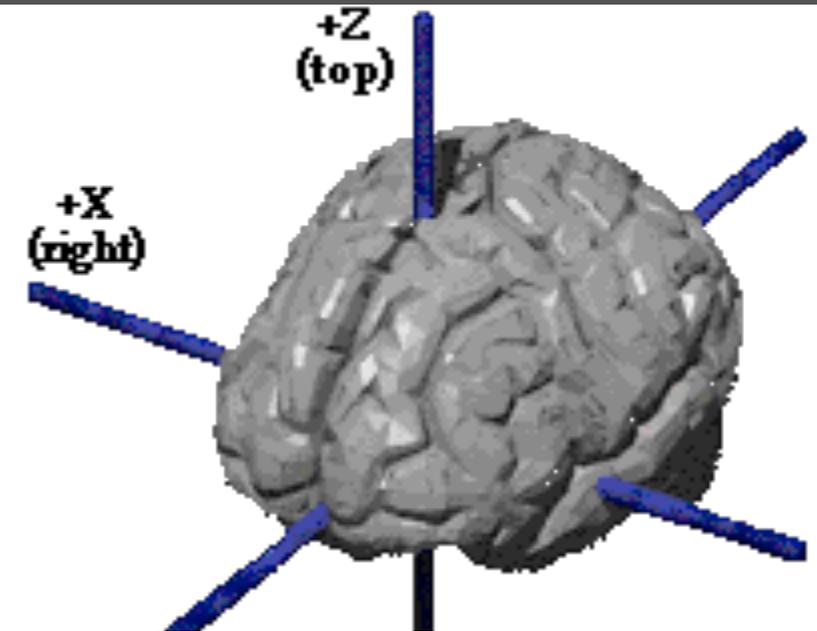▸ Brain activity is reported in (somewhat) standardized format



+Z (top)
+X (right)

Table 1
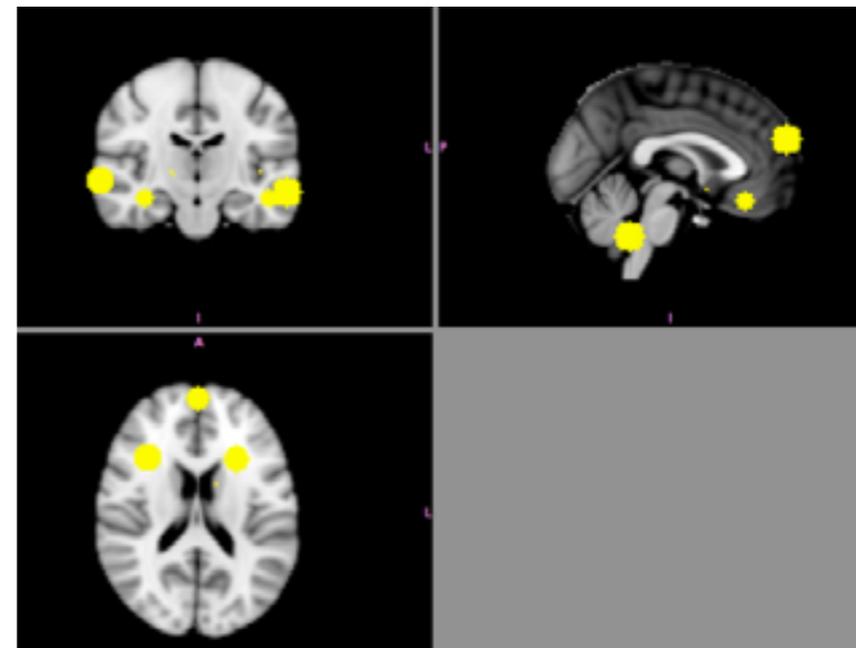Regions that showed a condition × time interaction in the ANOVA analysis

| No. | Region | Hemisphere | BA | x | y | z | mm³ |
|---|---|---|---|---|---|---|---|
| 1 | Middle/superior temporal gyrus | L | 21/22/37 | −52 | −54 | 9 | 13257 |
| 2 | Inferior frontal gyrus | L | 45/46/9 | −49 | 26 | 6 | 2781 |
| 3 | Posterior cerebellum | L | | −19 | −79 | −38 | 2214 |
| 4 | Dorsomedial PFC | L | 9/8 | −11 | 42 | 47 | 3051 |
| 5 | Left anterior PFC | L | 10 | −37 | 49 | 15 | 2025 |
| 6 | Inferior parietal cortex | L | 40/7 | −42 | −58 | 47 | 3132 |
| 7 | Dorsal premotor cortex | L | 6 | −43 | 0 | 50 | 1485 |
| 8 | Lingual gyrus | L | 17 | −10 | −95 | −2 | 378 |
| 9 | Middle /superior temporal gyrus | R | 21/22/37 | 52 | −40 | 5 | 16470 |
| 10 | Inferior frontal gyrus | R | 45/46 | 51 | 28 | 6 | 2241 |
| 11 | Posterior cerebellum | R | | 23 | −78 | −34 | 2808 |
| 12 | Dorsomedial PFC | R | 9 | 5 | 53 | 29 | 405 |
| 13 | Right anterior PFC | R | 10 | 38 | 42 | 21 | 5022 |
| 14 | Inferior parietal cortex | R | 40/7 | 42 | −53 | 48 | 9963 |
| 15 | Superior frontal gyrus | R | 6/8 | 10 | 28 | 60 | 297 |
| 16 | Anterior cingulate cortex | M | 32 | 0 | 26 | 35 | 5076 |
| 17 | Posterior cingulate cortex | M | 23/31/7 | 0 | −35 | 31 | 9612 |
| 18 | Precuneus | M | 7/19 | 1 | −76 | 36 | 10044 |

# Creating meta-analytic brain maps

- Automated Coordinate Extraction (Yarkoni et al, 2011, *Nature Methods*)

  - Automatically extracts activation tables from fMRI papers for 17 journals

  - Current database has 4,393 papers (with full text)

  - Good accuracy

    - 84% sensitivity, 97% specificity against SumsDB manual database

- Meta-analytic maps created for each paper

  - 10mm sphere placed at each focus

| X | Y | Z |
|---|---|---|
| 12 | 57 | -6 |
| 33 | 21 | 15 |
| 24 | 15 | 60 |
| 42 | 6 | 51 |
| 24 | -3 | 57 |

Automated coordinate extraction →

# Neurosynth.org

# Automated meta-analysis



Yarkoni et al., 2011, *Nature Methods*

# Automated meta-analysis



Previous meta-analyses

Automated meta-analysis

A — Forward Inference (P(Act|Term))

B — Reverse Inference (P(Term|Act))

C

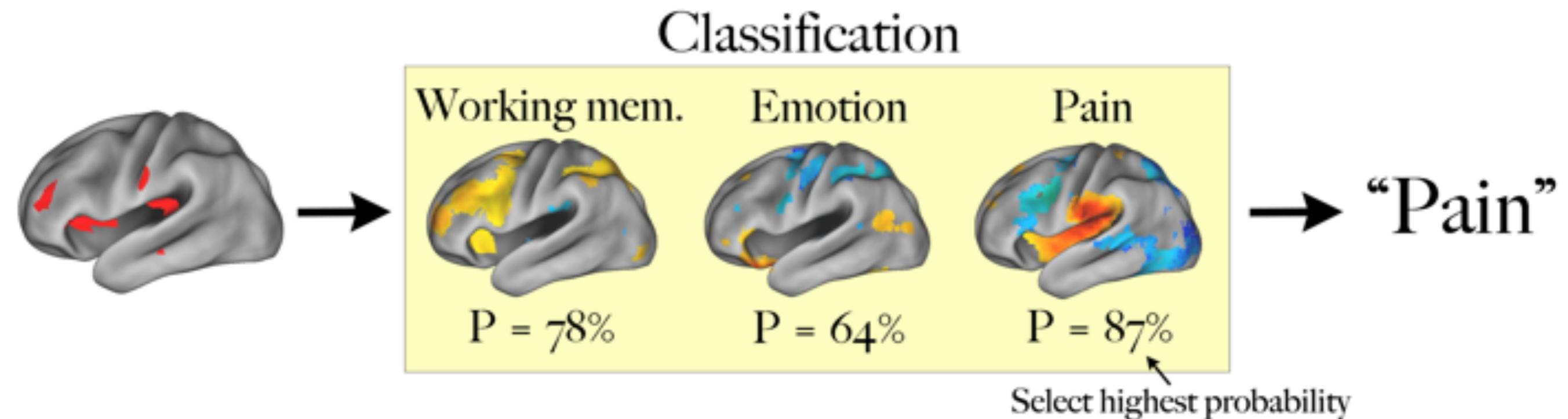Working Memory

Emotion

Pain

0   P(Act|Term)   0.4

0.1   P(Term|Act)   0.9

Yarkoni et al., 2011, *Nature Methods*

# Classification of cognitive states

- Given 2+ terms, can determine which is most likely given the data

- Naive Bayes classifier: assumes that all features (voxels) are independent; selects the most probable class

- Can apply this to any activation map—studies, individual subjects, etc.



Classification

Working mem.    Emotion    Pain

P = 78%    P = 64%    P = 87%

Select highest probability

→ "Pain"

Yarkoni et al, 2011, *Nature Methods*

- Cross-validated classification of all studies in database

- Select 25 high-frequency terms

- Pairwise classification: how well can we distinguish between each pair of terms?

Yarkoni et al, 2011, *Nature Methods*

Yarkoni et al, 2011, *Nature Methods*

# Using classification to understand mental structure

WM: working memory
TS: Task switching
RS: Response selection
RI: Response inhibition
CC: Cognitive control
BI: Bilingual language



A' (k=3)

Lenartowicz et al, 2010, *Topics in Cognitive Science*

- Can we identify cognitive states in individual (new) subjects?

- Difficult, because:

  - No opportunity for training

  - Data is of a fundamentally different type

- Tested in samples of subjects from working memory, emotion, and pain studies

  - Can we predict source study type?

Yarkoni et al, 2011, *Nature Methods*

# Classifying individual subjects



Yarkoni et al, 2011, *Nature Methods*

Table 2. Pearson correlations between searchlight classification map and NeuroSynth term-based reverse inference activation maps

| Term | Correlation (r) |
| --- | --- |
| Control | 0.1451 |
| Working | 0.1159 |
| Numerical | 0.1157 |
| Letter | 0.1081 |
| Attention | 0.1062 |
| Correct | 0.1060 |
| Cue | 0.0995 |
| Preparatory | 0.0970 |
| Load | 0.0959 |
| Hand | 0.0924 |

The 10 most highly correlated terms are listed. From Yarktoni et al. (26).

Helfinstein et al, 2014, PNAS

# Neurovault + Neurosynth = automated reverse inference



Gorgolewski et al., submitted

- We can decode mental states across individuals

- This can provide insights into the similarity space of mental processes

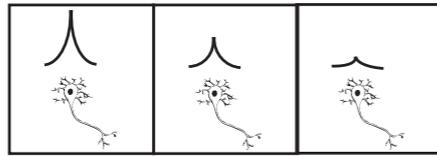  - And ultimately inform our ontology of mental processes

- Psychological theories rarely make clear predictions about activation
- But they often make predictions about similarity relations between stimuli

- We can test those against neuroimaging data
  - In principle we don't even have to care *where* the effects happen in the brain

# Representational similarity analysis
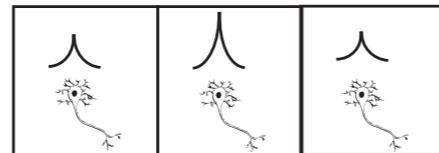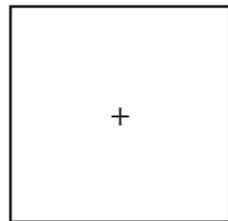


Stimulus Presentation

Neural Response

Activation Pattern

Time

## Similarity matrix

|  | 1 | med | low |
|---|---|---|---|
|  |  | 1 | med |
|  |  |  | 1 |

Davis & Poldrack, 2013

Kriegeskorte et al., 2008

- Some birds are "birdier" than others

# Typicality



Dimension 2: Predacity

Dimension 1: Size

After Smith, Shoben, & Rips (1974)    Photos via http://www.birdphotography.com/
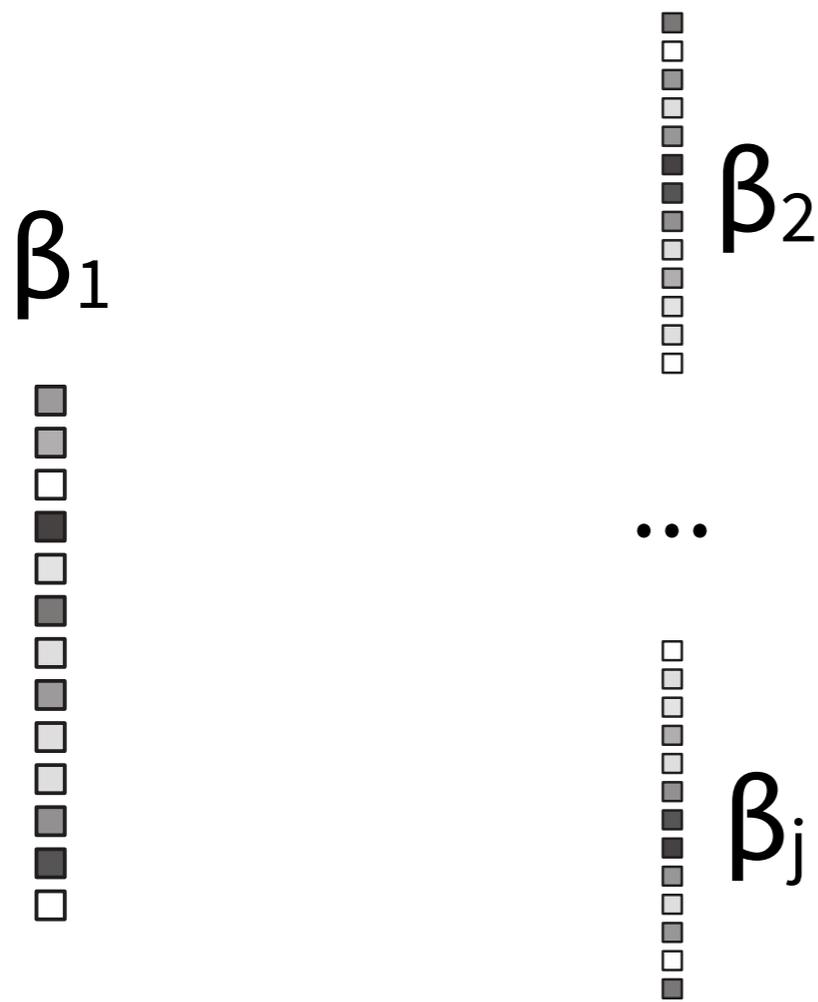
- May reflect:

  - Average similarity to other category members

  - Similarity to idealized members

    - Caricature effects

- Can we find a neural signature that is related to psychological typicality?

$\beta_2$

$\beta_1$

...

$\beta_j$

Davis & Poldrack, 2013, *Cerebral Cortex*

β₁

β₂

...

βⱼ

**DISTANCE IS DEFINED AS THE CORRELATION DISTANCE BETWEEN TWO BETA-SERIES ACTIVATION PATTERNS**

$$d_{ij} = \left[1 - corr(\beta_i, \beta_j)\right]/2$$
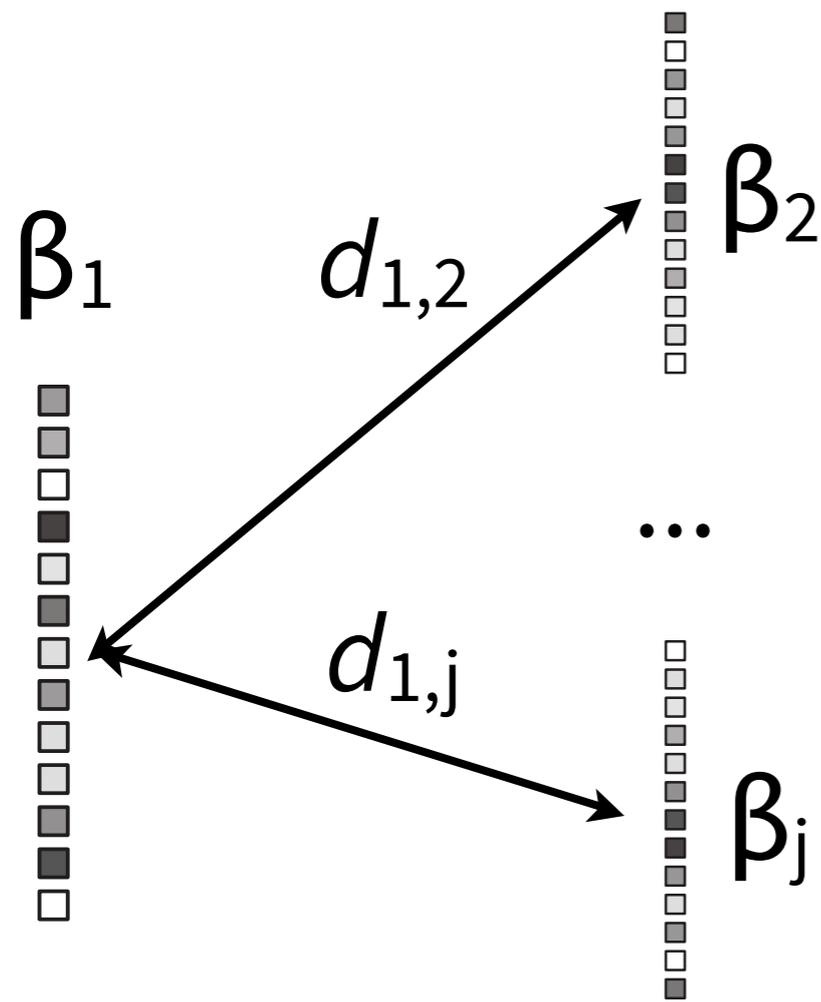
Davis & Poldrack, 2013, *Cerebral Cortex*

# Computing neural typicality



DISTANCE IS DEFINED AS THE CORRELATION DISTANCE BETWEEN TWO BETA-SERIES ACTIVATION PATTERNS
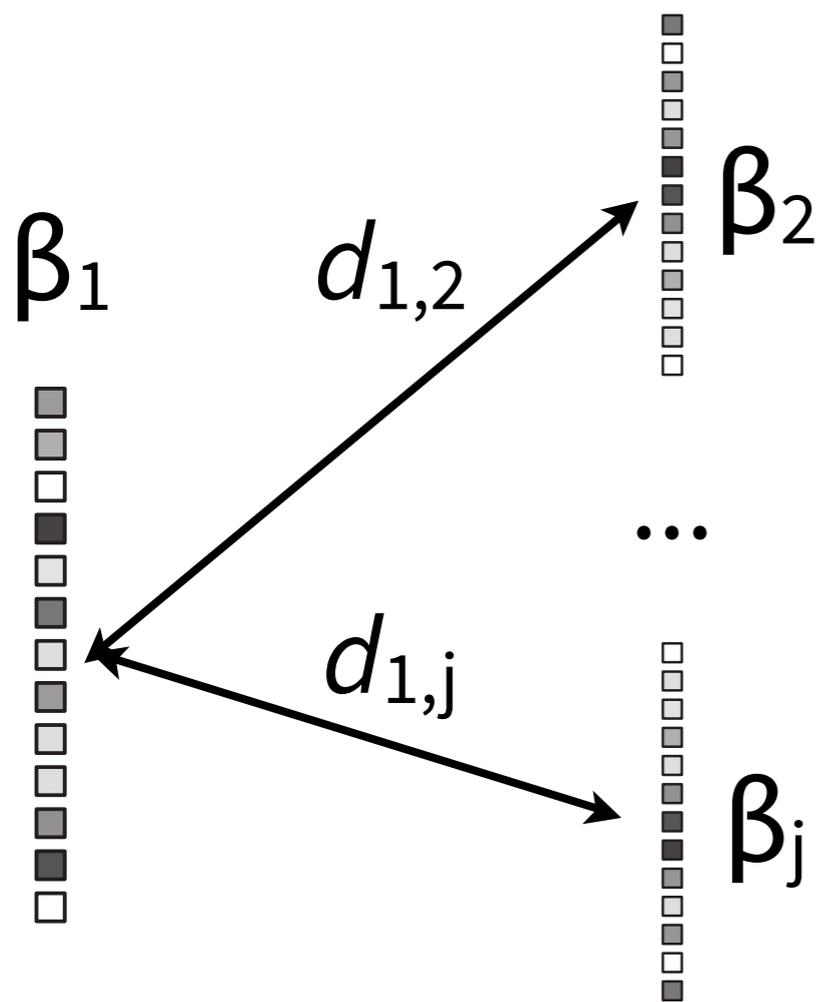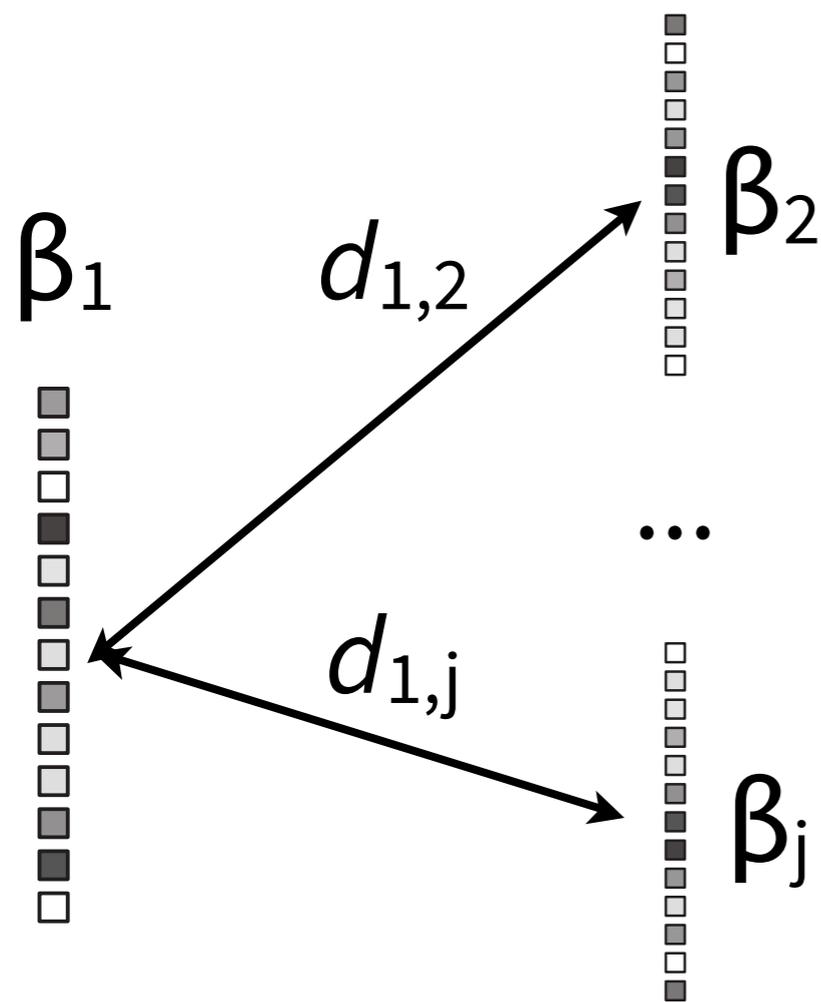
$$d_{ij} = \left[1 - corr(\beta_i, \beta_j)\right]/2$$

Davis & Poldrack, 2013, *Cerebral Cortex*

# Computing neural typicality



**DISTANCE IS DEFINED AS THE CORRELATION DISTANCE BETWEEN TWO BETA-SERIES ACTIVATION PATTERNS**

$$d_{ij} = \left[1 - corr(\beta_i, \beta_j)\right]/2$$

**SIMILARITY IS AN EXPONENTIAL FUNCTION OF THE DISTANCE BETWEEN TWO REPRESENTATIONS**
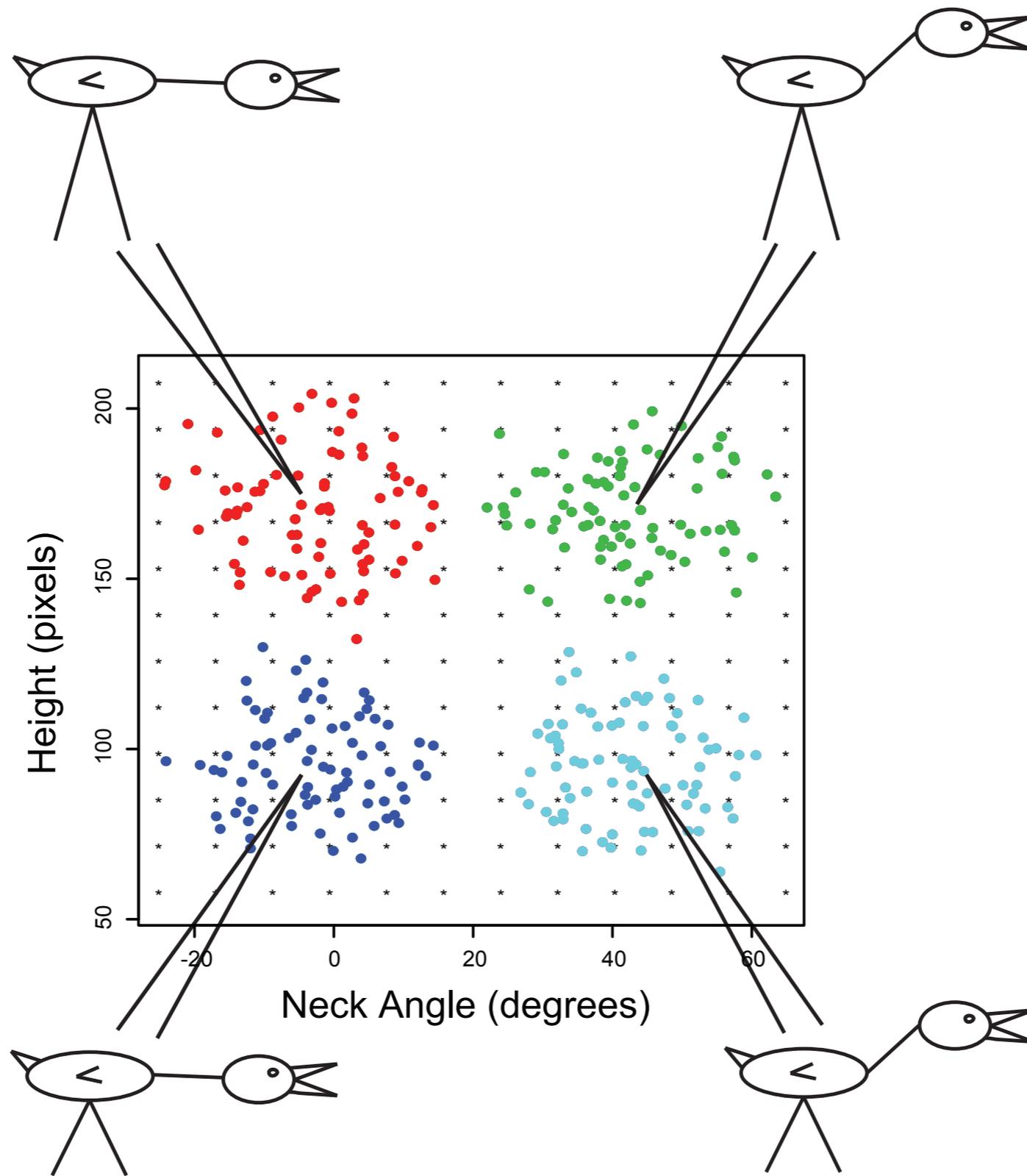
$$s_{ij} = \exp(-d_{ij})$$

Davis & Poldrack, 2013, *Cerebral Cortex*

# Computing neural typicality



β₁    $d_{1,2}$    β₂

...

$d_{1,j}$    βⱼ

$$typ(i \mid J) = \sum_{j \in J} s_{ij}$$

**DISTANCE IS DEFINED AS THE CORRELATION DISTANCE BETWEEN TWO BETA-SERIES ACTIVATION PATTERNS**

$$d_{ij} = \left[1 - corr(\beta_i, \beta_j)\right]/2$$

**SIMILARITY IS AN EXPONENTIAL FUNCTION OF THE DISTANCE BETWEEN TWO REPRESENTATIONS**

$$s_{ij} = \exp(-d_{ij})$$

**TYPICALITY IS BASED ON THE SUM OF SIMILARITIES BETWEEN REPRESENTATIONS OF AN OBJECT AND OTHER CATEGORY MEMBERS**

Davis & Poldrack, 2013, *Cerebral Cortex*

- Is neural typicality associated with subjective typicality ratings?

- Used a task in which subjective typicality and physical feature resemblance are dissociated
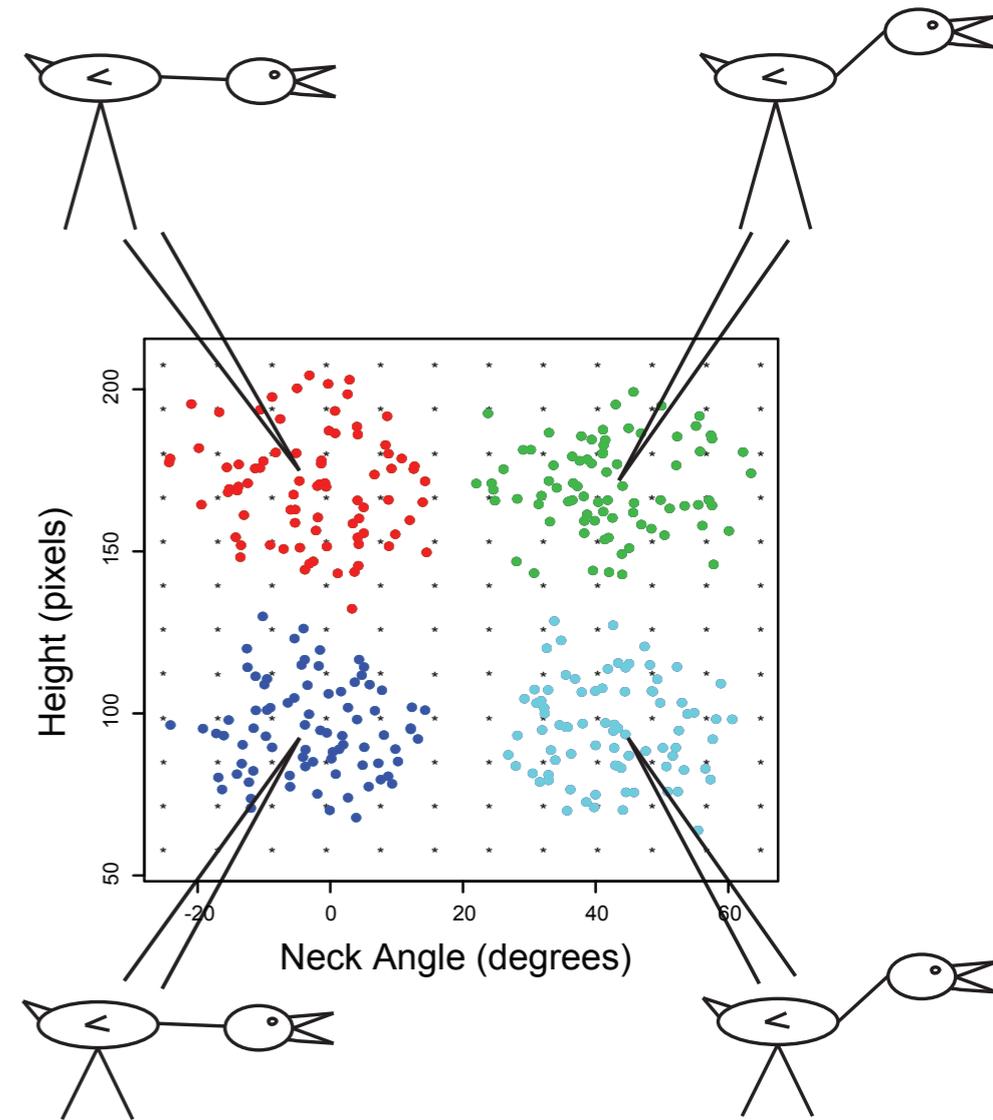
Davis & Poldrack, 2013, *Cerebral Cortex*

Does this bird live in nest
A or C ?

Incorrect,
This One Lives in Nest A

1.5 seconds

performed during
structural imaging

Height (pixels)

Neck Angle (degrees)

Davis & Poldrack, 2013, *Cerebral Cortex*

Does this bird live in nest
A, B, C, or D ?

performed while scanning

Davis & Poldrack, 2013, *Cerebral Cortex*

How Typical is this bird of Nest A (1-7)?

performed outside scanner
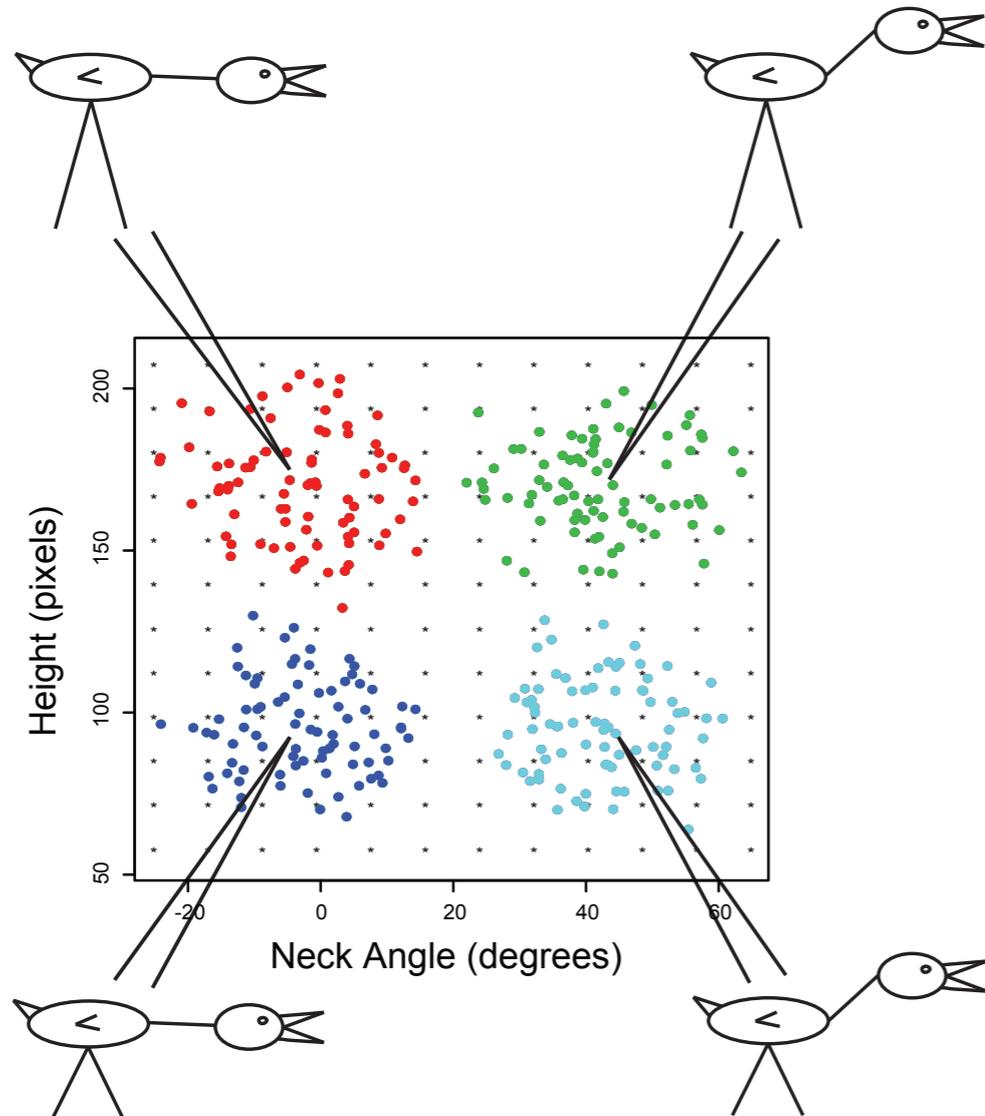
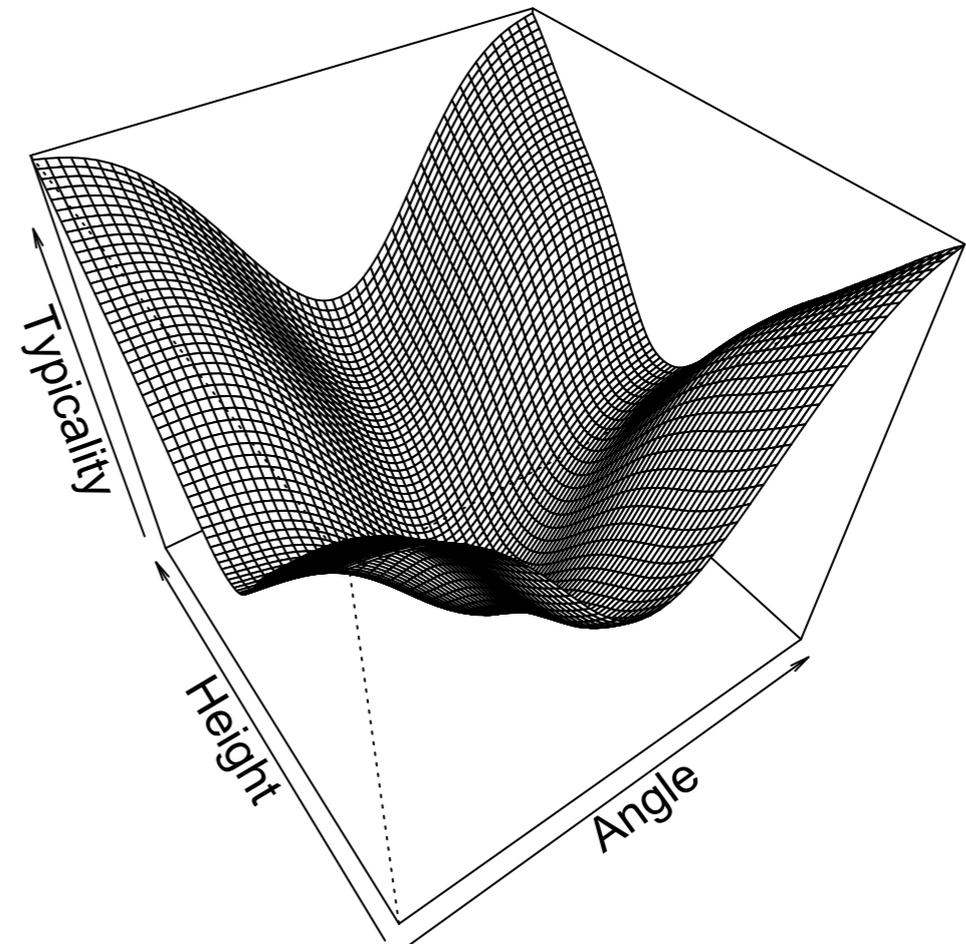Davis & Poldrack, 2013, *Cerebral Cortex*

beta-series

LS-single
(Mumford et al., 2012)

Design matrices for single-trial estimation

# Idealized stimuli are judged most typical



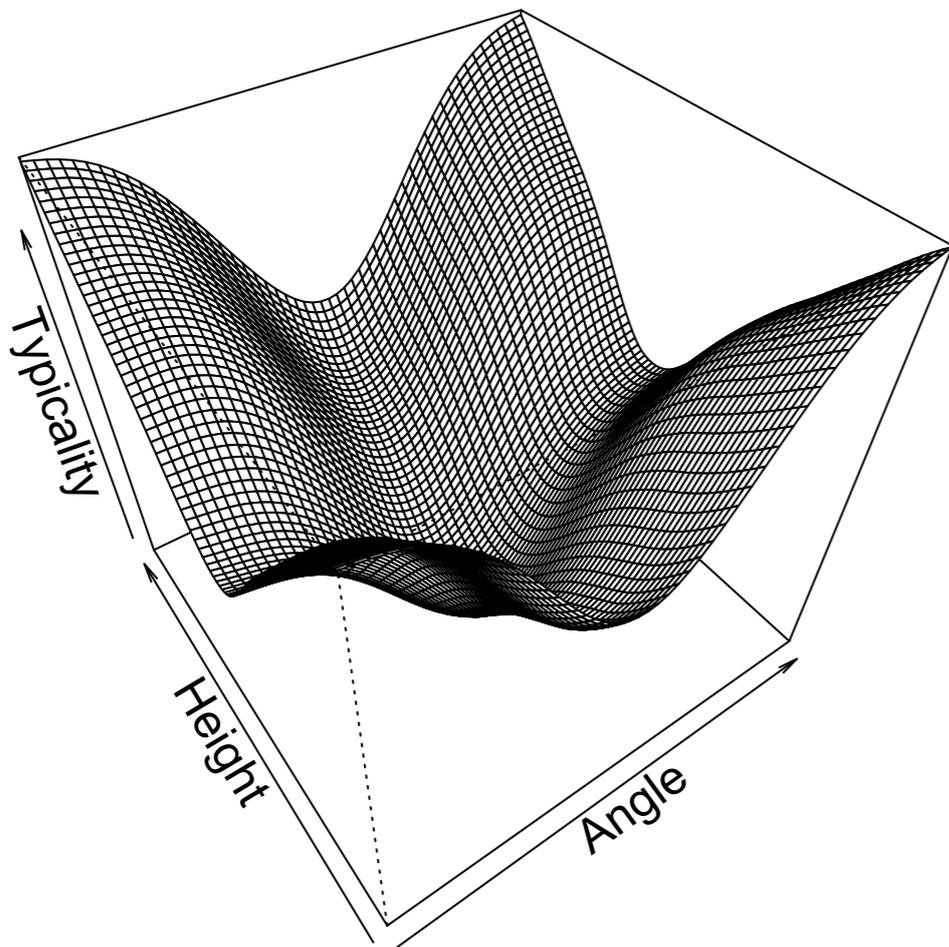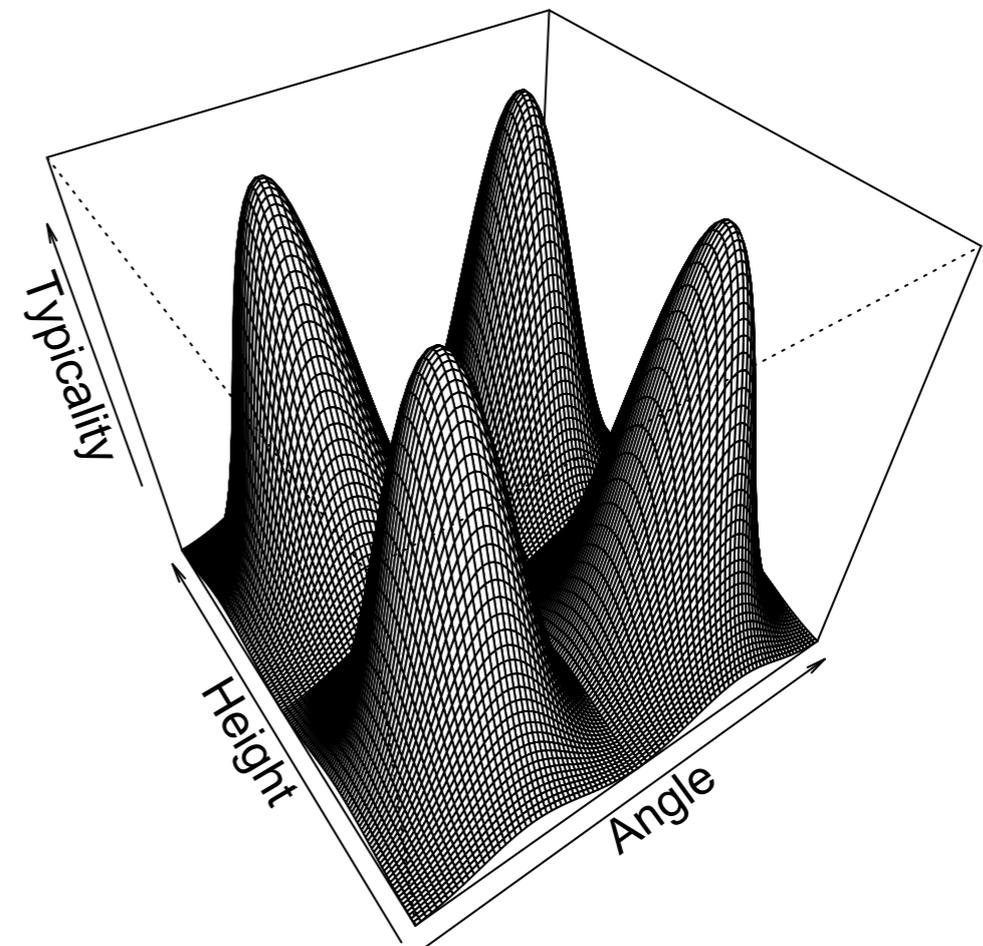GAM fit to typicality ratings

Davis & Poldrack, 2013, *Cerebral Cortex*

- Neural similarity space will reflect subjective typicality

- Neural similarity space will reflect physical typicality (likelihood given category)
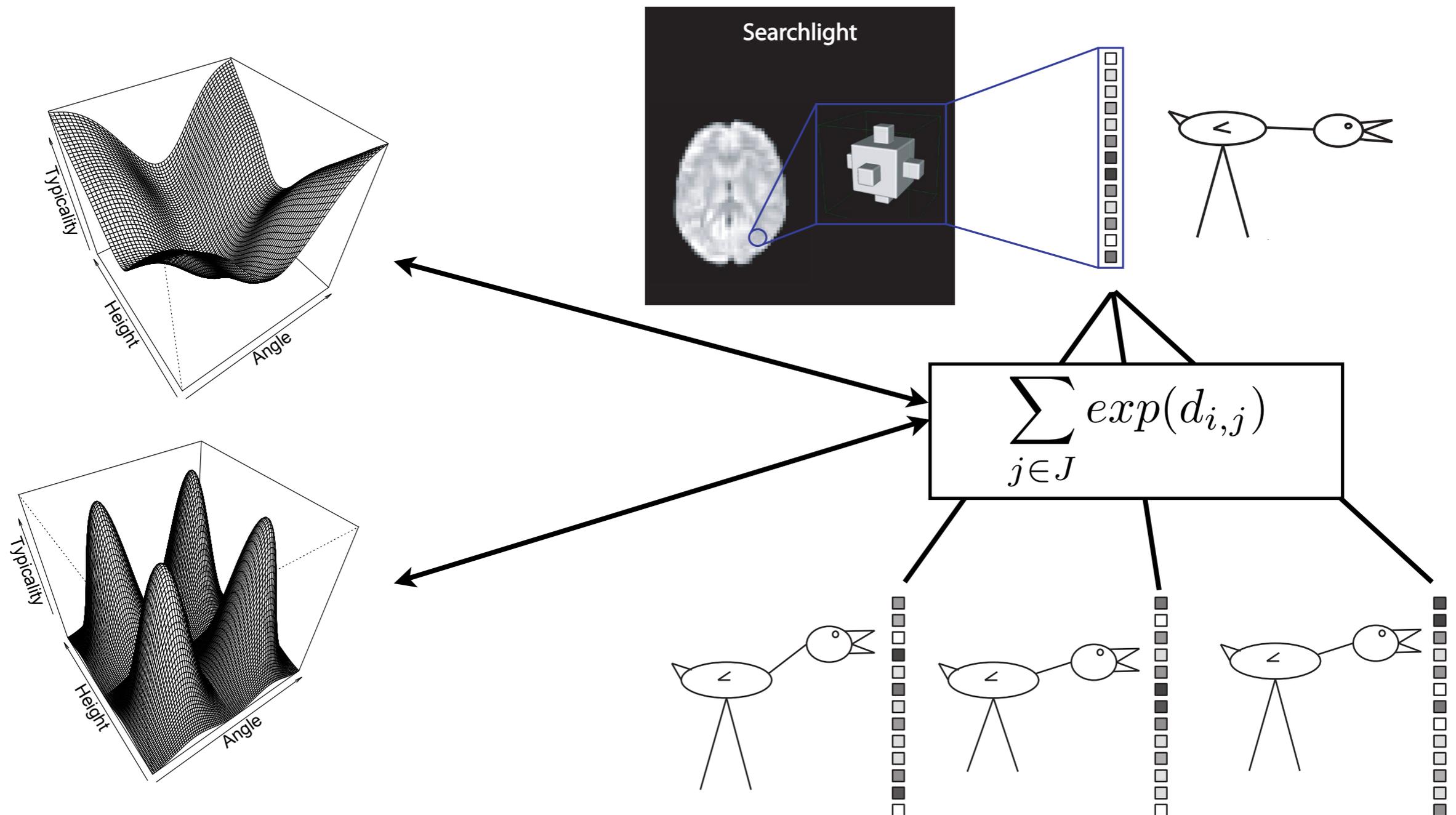


Obtained from behavior



Obtained using GCM

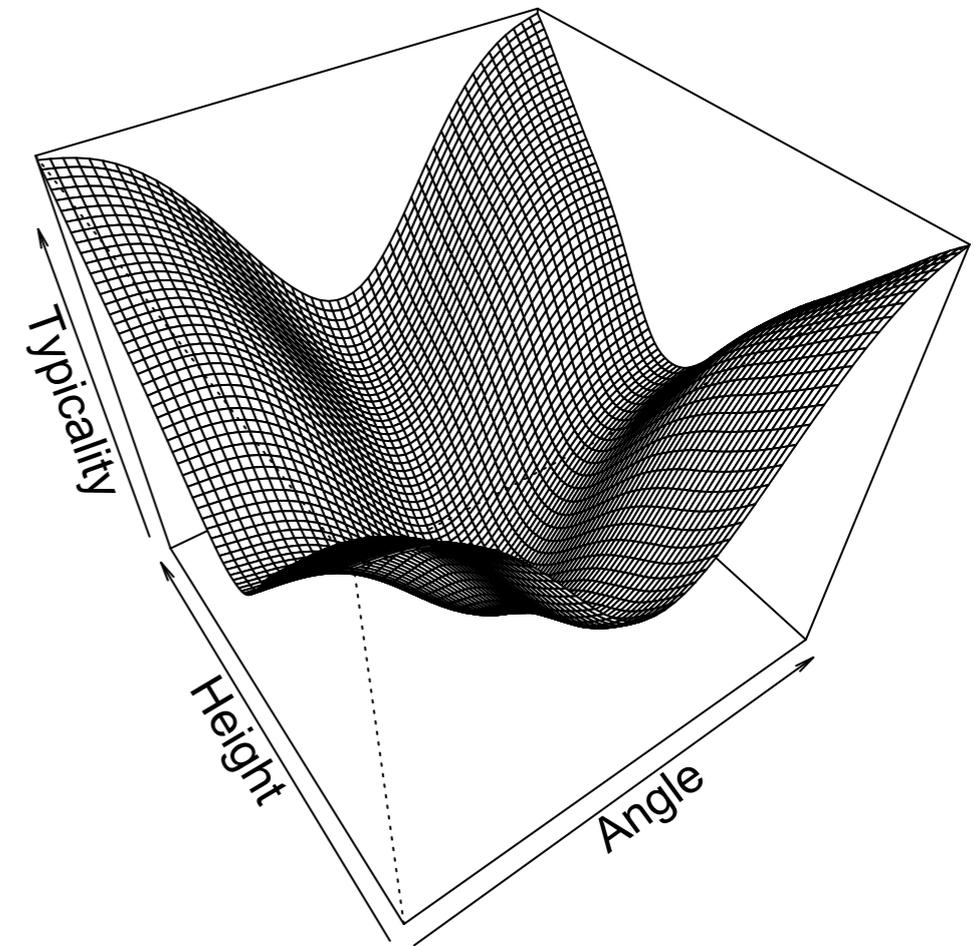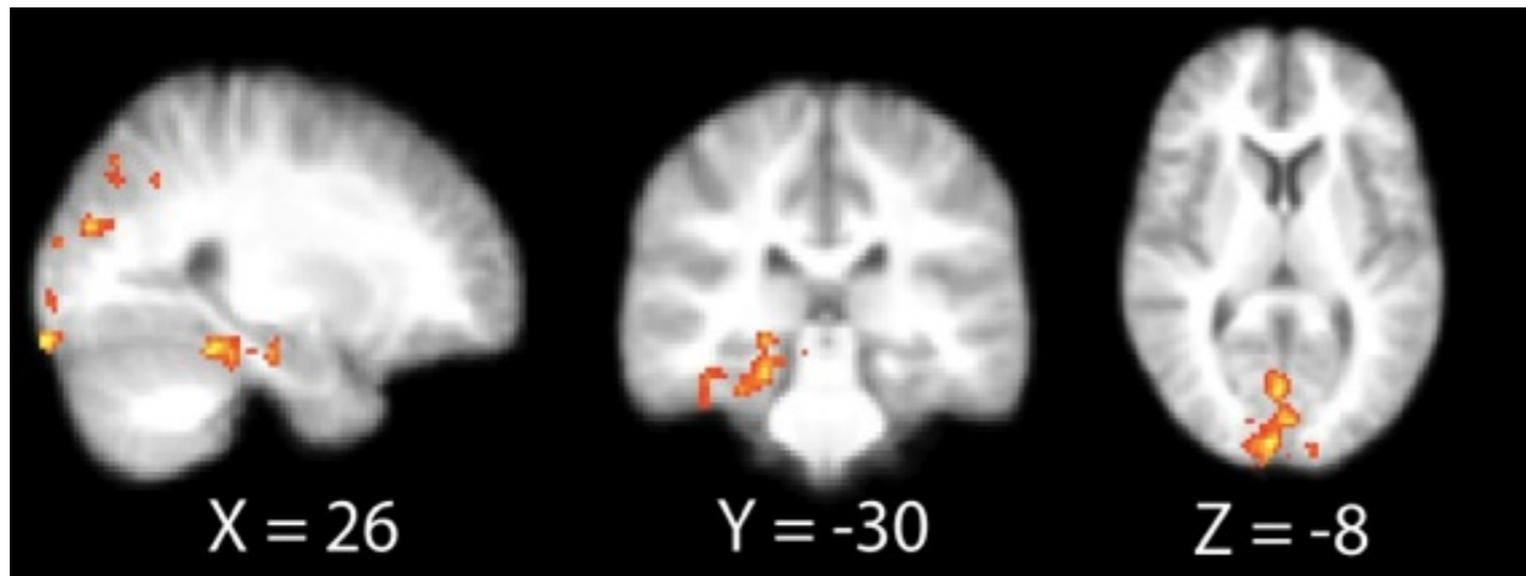Davis & Poldrack, 2013, *Cerebral Cortex*

- Correlate neural typicality with psychological and physical predictions



$$\sum_{j \in J} exp(d_{i,j})$$
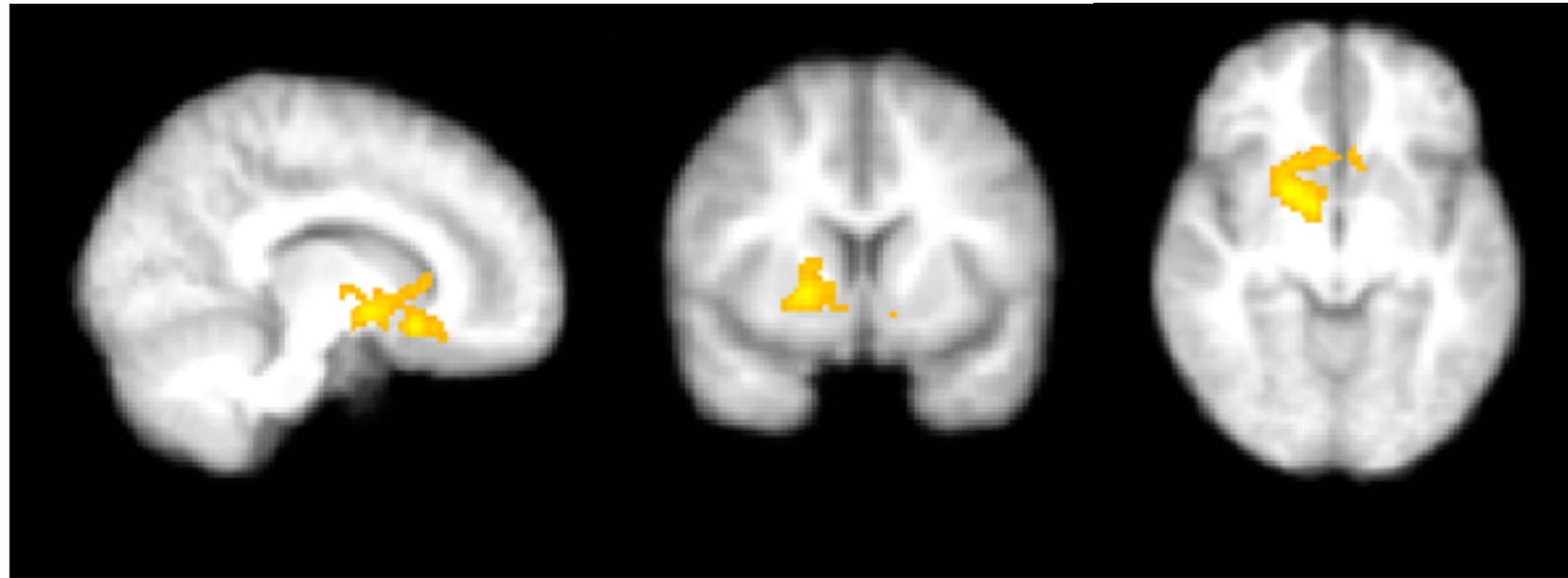
Regions in which neural typicality and psychological typicality are correlated



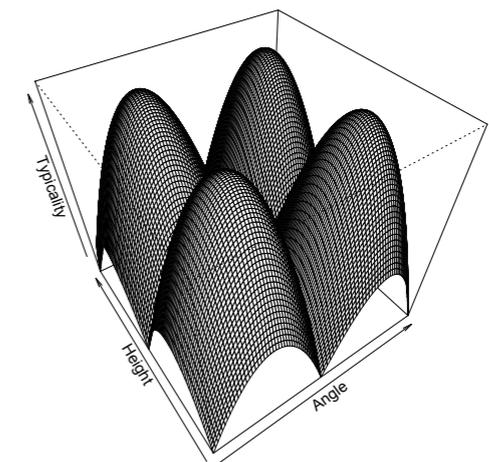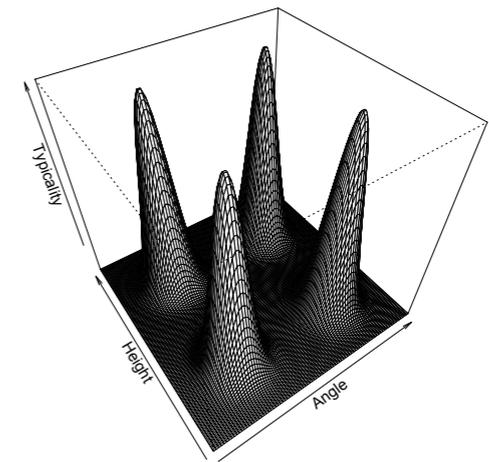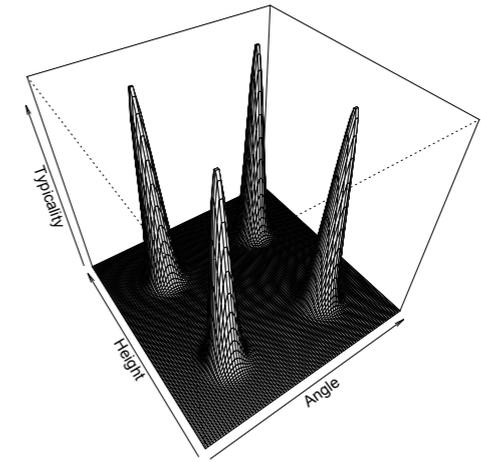Davis & Poldrack, 2013, *Cerebral Cortex*

Regions in which univariate activation
and psychological typicality are correlated

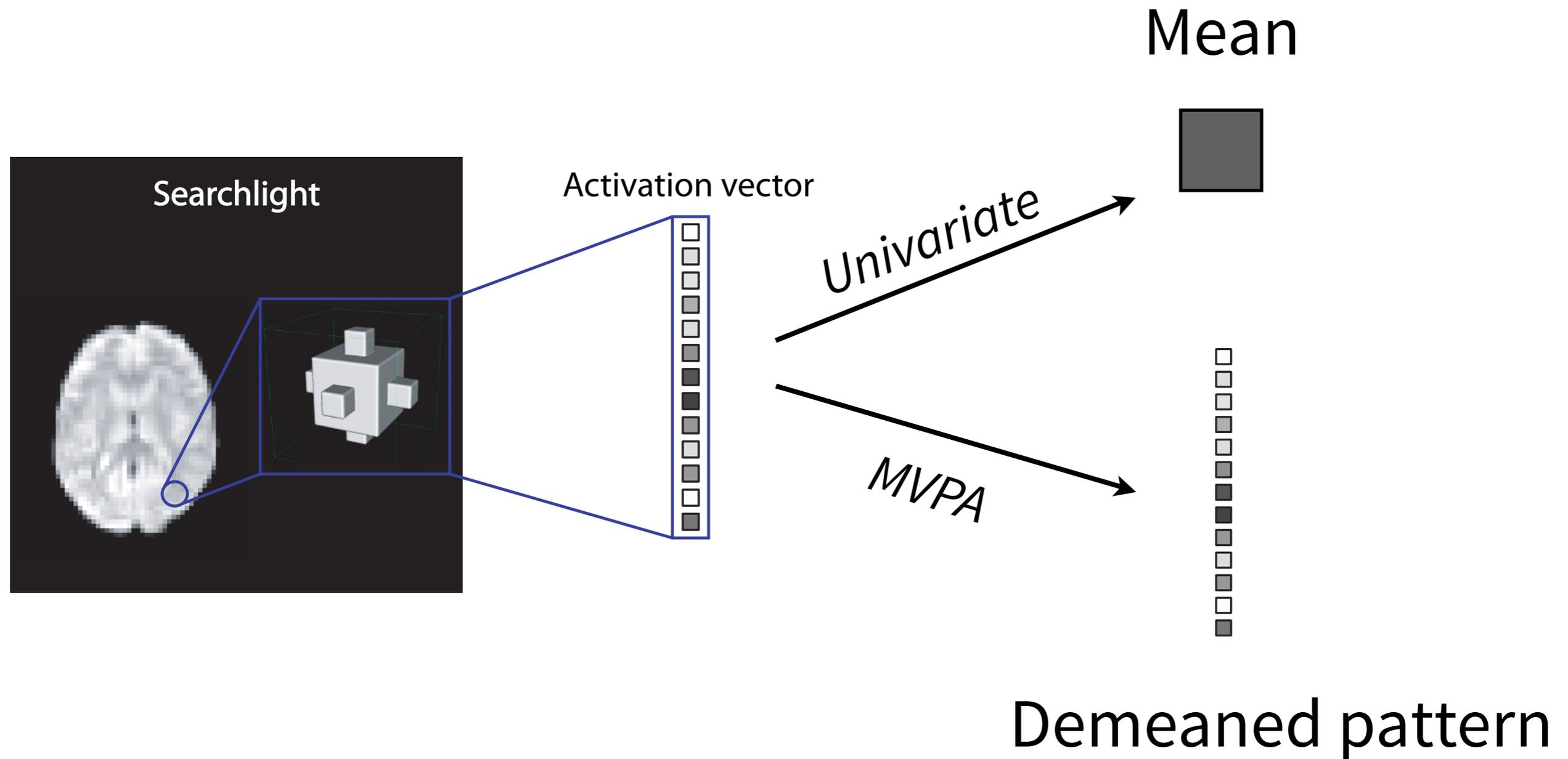Davis & Poldrack, 2013, *Cerebral Cortex*

- Physical similarity predictions obtained using GCM

  - Examined across multiple levels of variance

- There were no regions in which neural typicality or activation reflected physical typicality/
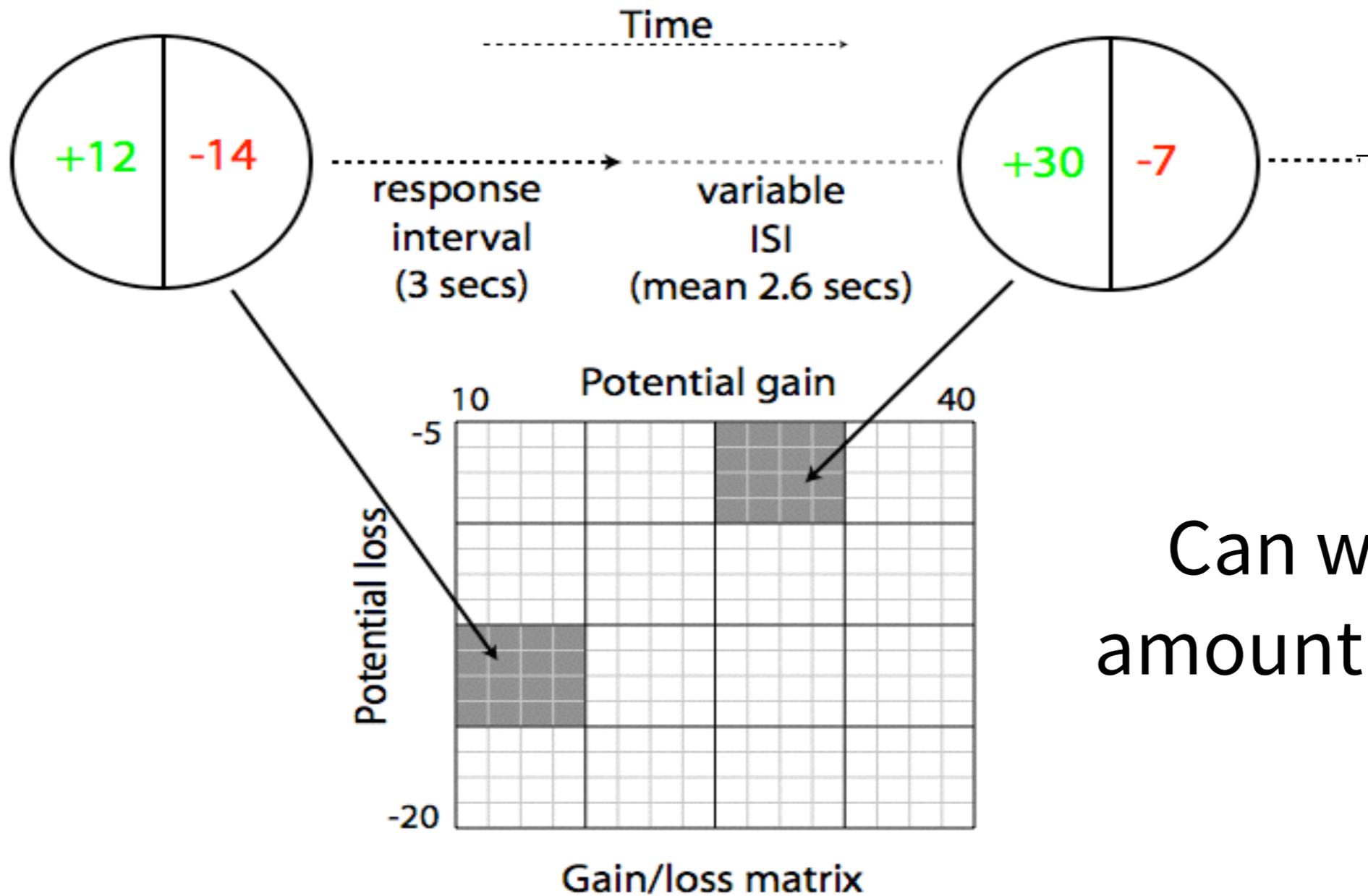


Davis & Poldrack, 2013, *Cerebral Cortex*

- Activation patterns are isomorphic to mental representations

- Subjective typicality reflected in neural typicality

  - Univariate analysis would have told a very different story

Mean

Searchlight

Activation vector

Univariate

MVPA

Demeaned pattern

Can we decode the amount of gain or loss?

Tom et al., 2007, *Science*

# The opportunistic nature of MVPA



Jimura & Poldrack, 2011, *Neuropsychologia*

# Differential sensitivity of MVPA



Gain

Loss

Lateral

Medial

a

b

M+

z = 2.0    z = 3.5

U+

z = 2.0    z = 3.5

Jimura & Poldrack, 2011, *Neuropsychologia*

Davis, Laroque et al., 2014, *Neuroimage*

## Similarity



## Classification



Davis, Laroque et al., 2014, *Neuroimage*

Univariate    Similarity    Classification

Davis, Laroque et al., 2014, *Neuroimage*

# Some things we have learned

- If classification results look too good, you have most likely done something wrong

- Always confirm results by randomizing data from the very beginning

  - run many times to get null distribution, make sure it's actually at chance

- Crossvalidation with regression is very tricky (don't use LOO)

- Differences between univariate and multivariate analyses can't be easily interpreted (Davis et al., 2014, *Neuroimage*)

- Trials orders must be separately randomized for each subject (Mumford et al., 2014, *Neuroimage*)

# Conclusions

- Neuroimaging data CAN provide evidence relevant to psychological questions

    - But informal reverse inference is not the way!

- Machine learning methods provide the means to decode and predict mental states from neuroimaging data

- Multivariate analyses can establish isomorphisms between neural and mental representations

# Acknowledgments

**Stanford**
Sanmi Koyejo
Chris Filo Gorgolewski
Karen LaRoque
Anthony Wagner

**UT Austin**
Jeanette Mumford
Sarah Helfinstein
Tom Schonberg
Tyler Davis
Tal Yarkoni

**UCLA**
Jessica Cohen
Robert Bilder
Eliza Congdon
Eydie London
Tyrone Cannon
Nelson Freimer

**Rutgers**
Stephen J. Hanson
Yaroslav Halchenko

**Princeton**
Ken Norman

NIMH
National Institute
of Mental Health

James S. McDonnell
Foundation

Data sets and code will be made available at www.openfmri.org