

Why metacognition matters for (computational) psychiatry

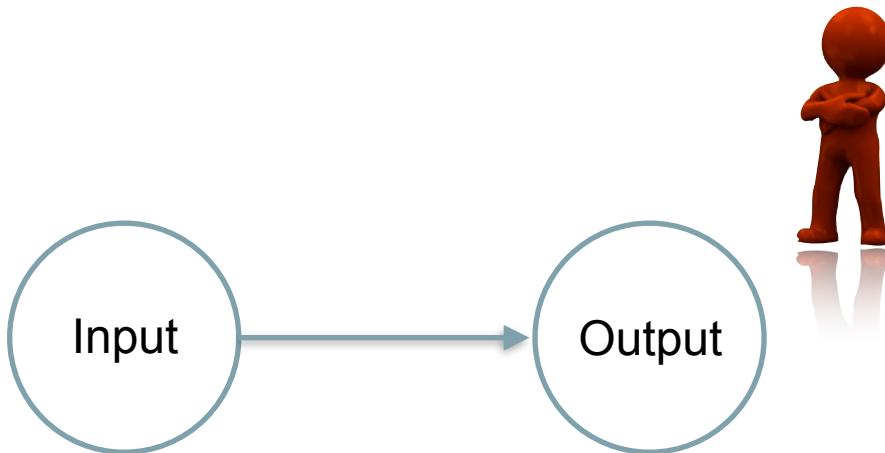
Steve Fleming

Wellcome Centre for Human Neuroimaging, UCL

stephen.fleming@ucl.ac.uk

metacoglab.org

Defining metacognition



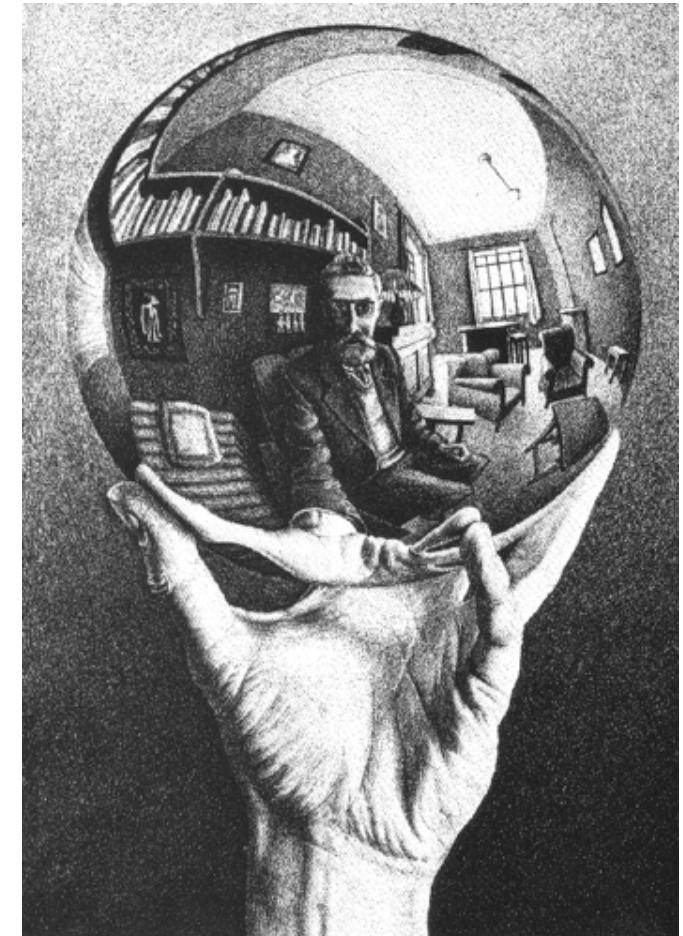
Defining metacognition

Self-reflection

Recursive thought

Introspection

etc...



Defining metacognition

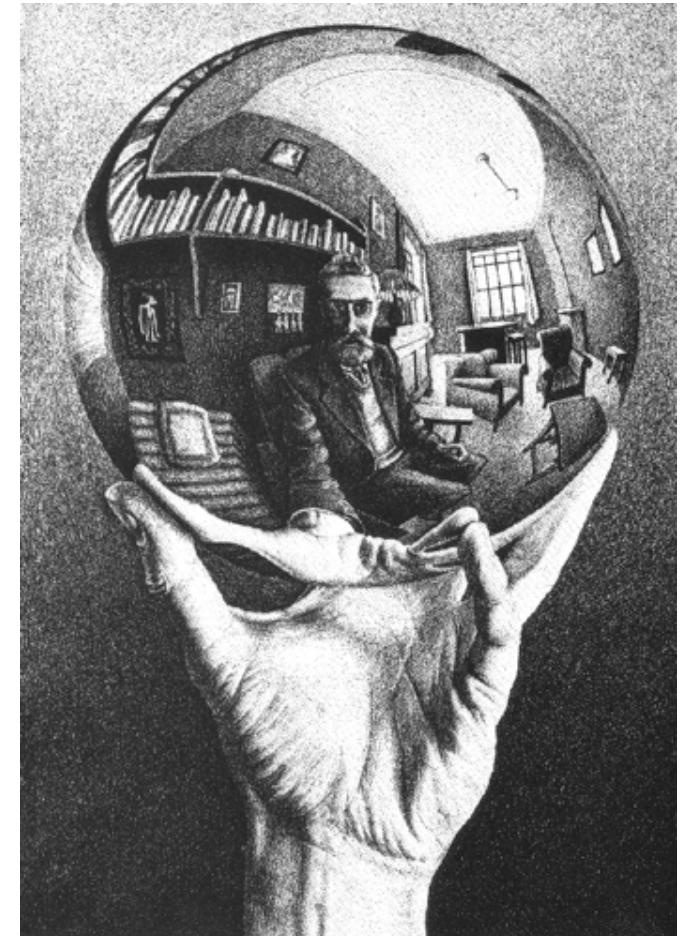
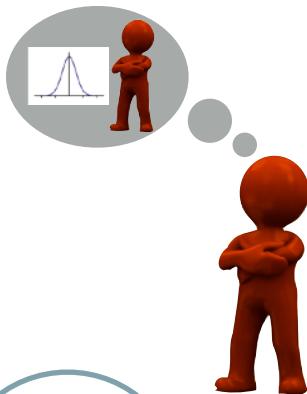
- “cognition about cognitive phenomena...” (**Flavell, 1979**)

Self-reflection

Recursive thought

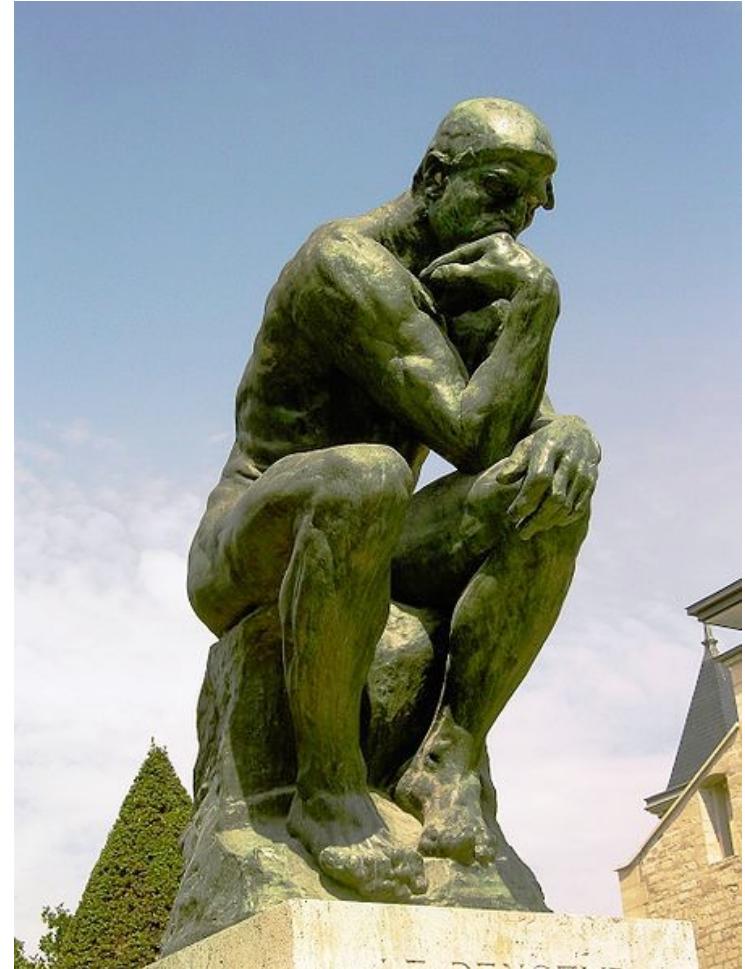
Introspection

etc...

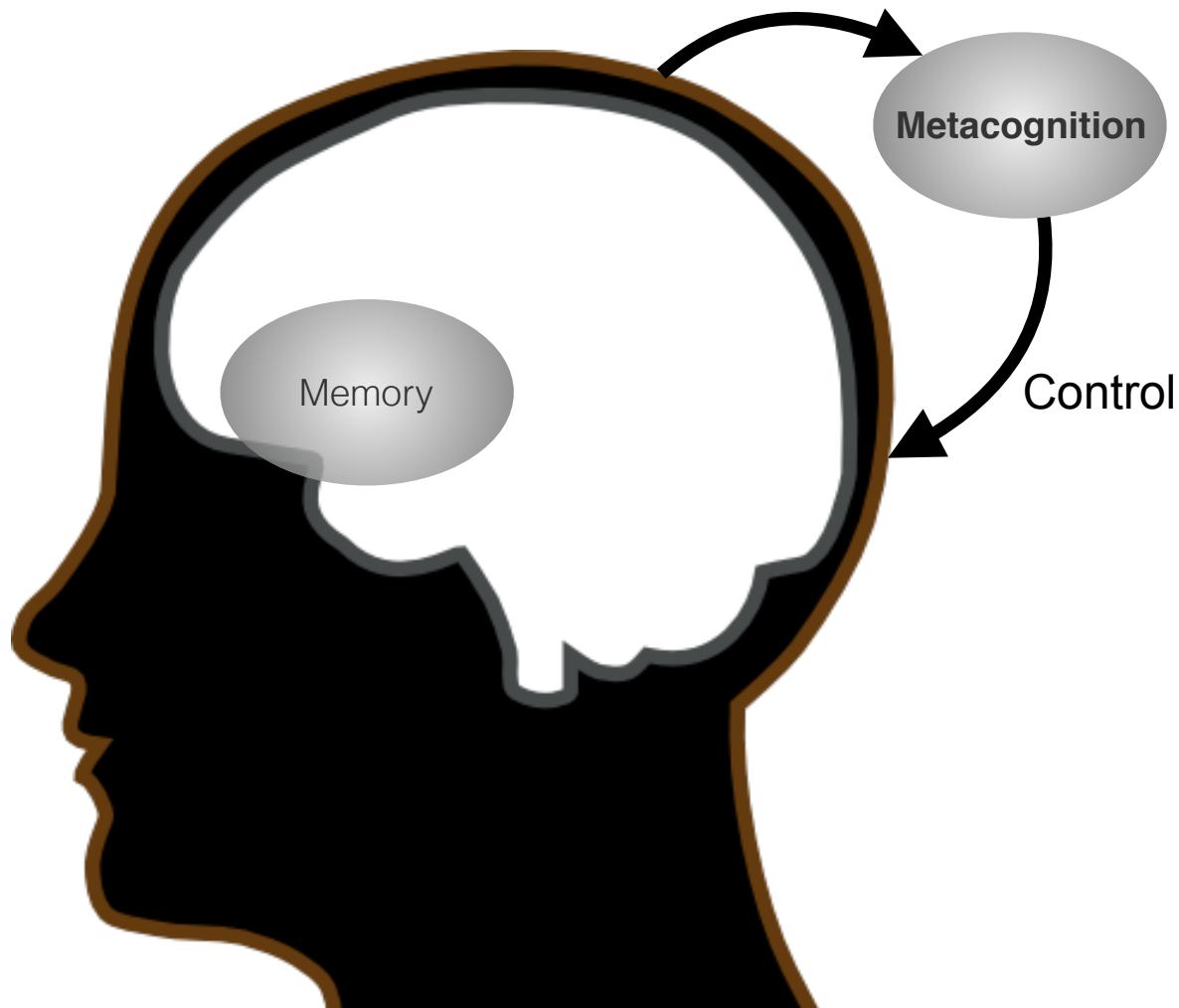


Defining metacognition

- “Metacognition research concerns the processes by which people self-reflect on their own cognitive and memory processes (monitoring) and how they put their metaknowledge to use in regulating their information processing and behaviour (control)” (**Koriat, 2007**)



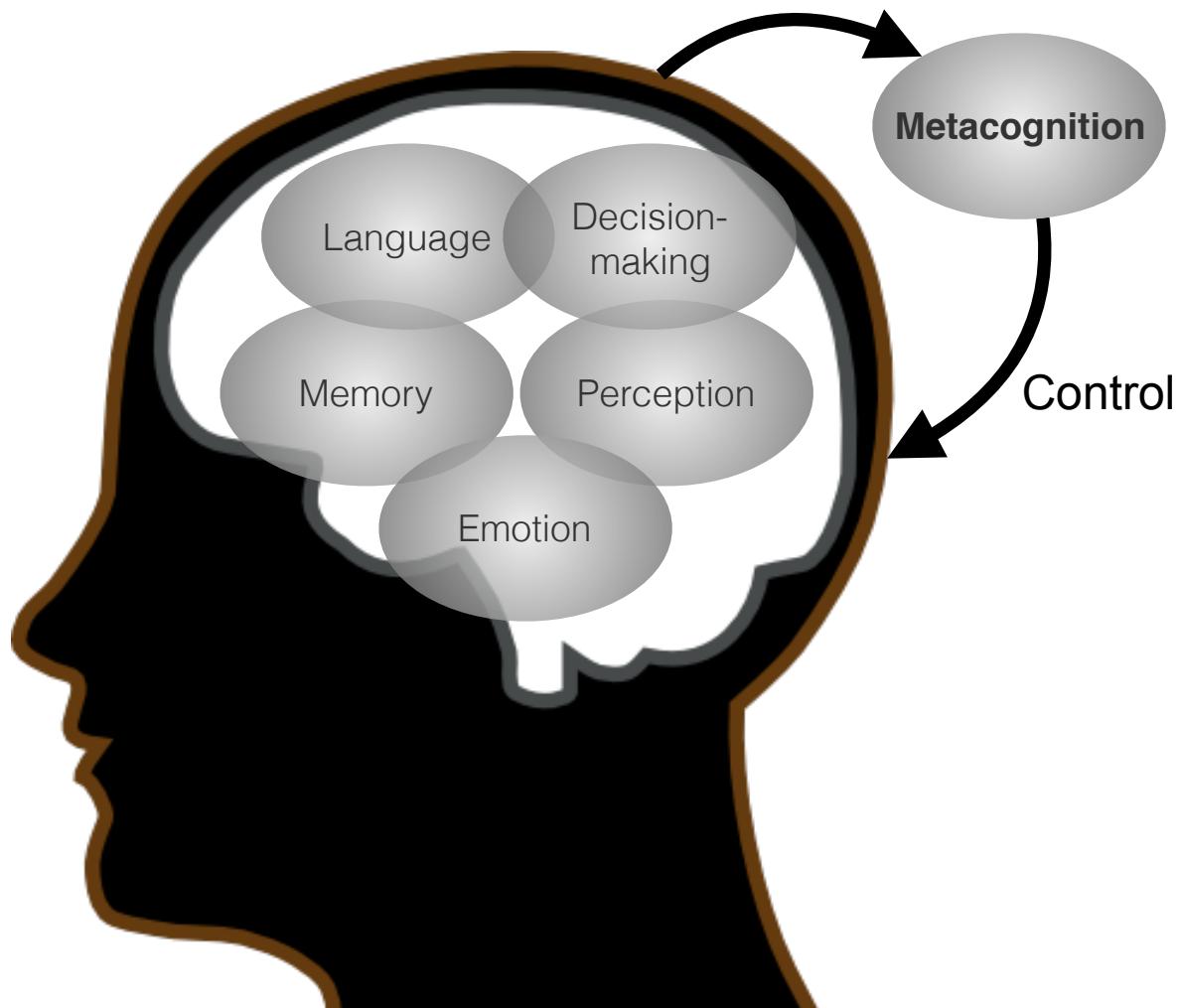
Monitoring / introspection



Metacognition

= reflecting on and
controlling other
cognitive processes

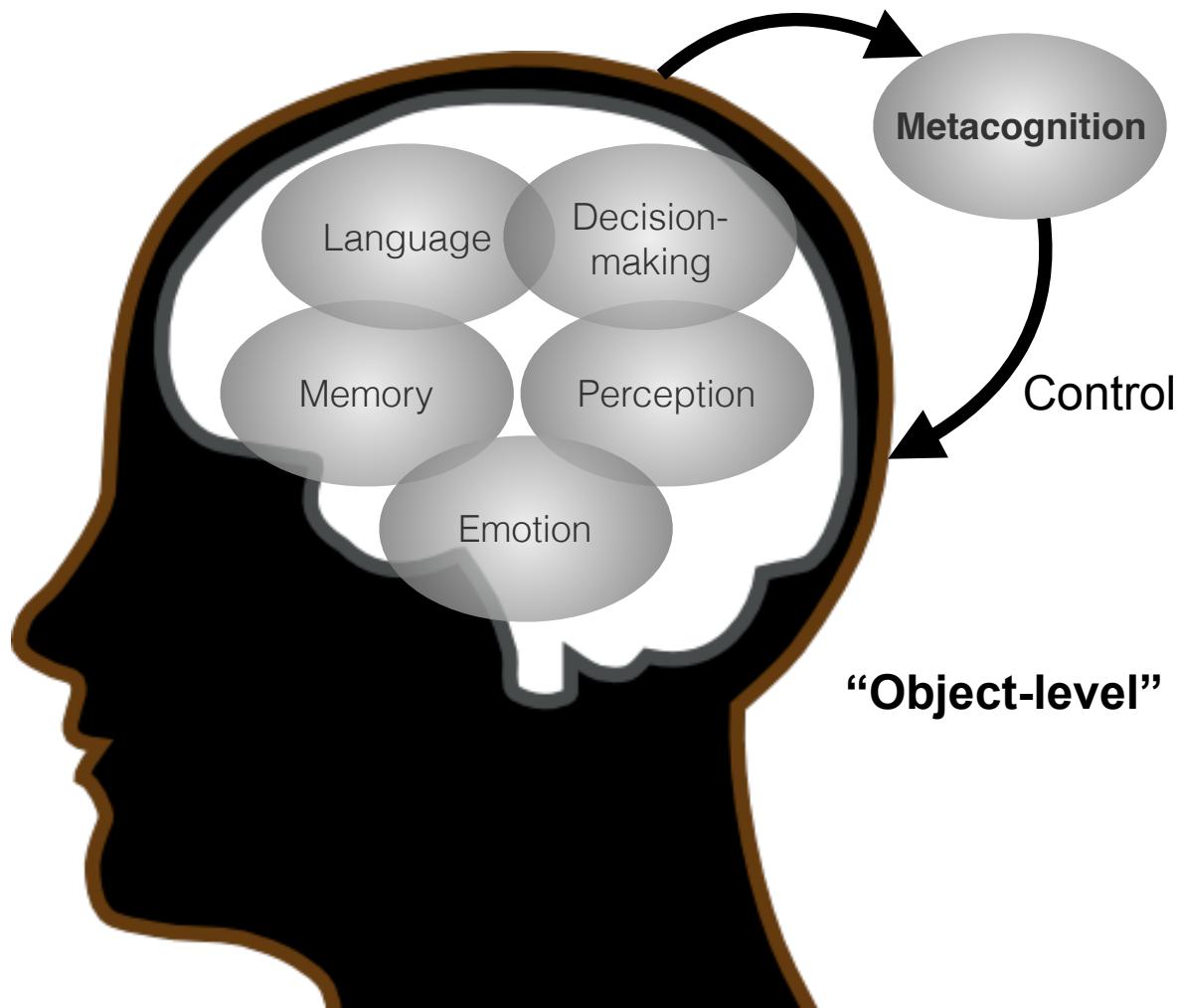
Monitoring / introspection



Metacognition

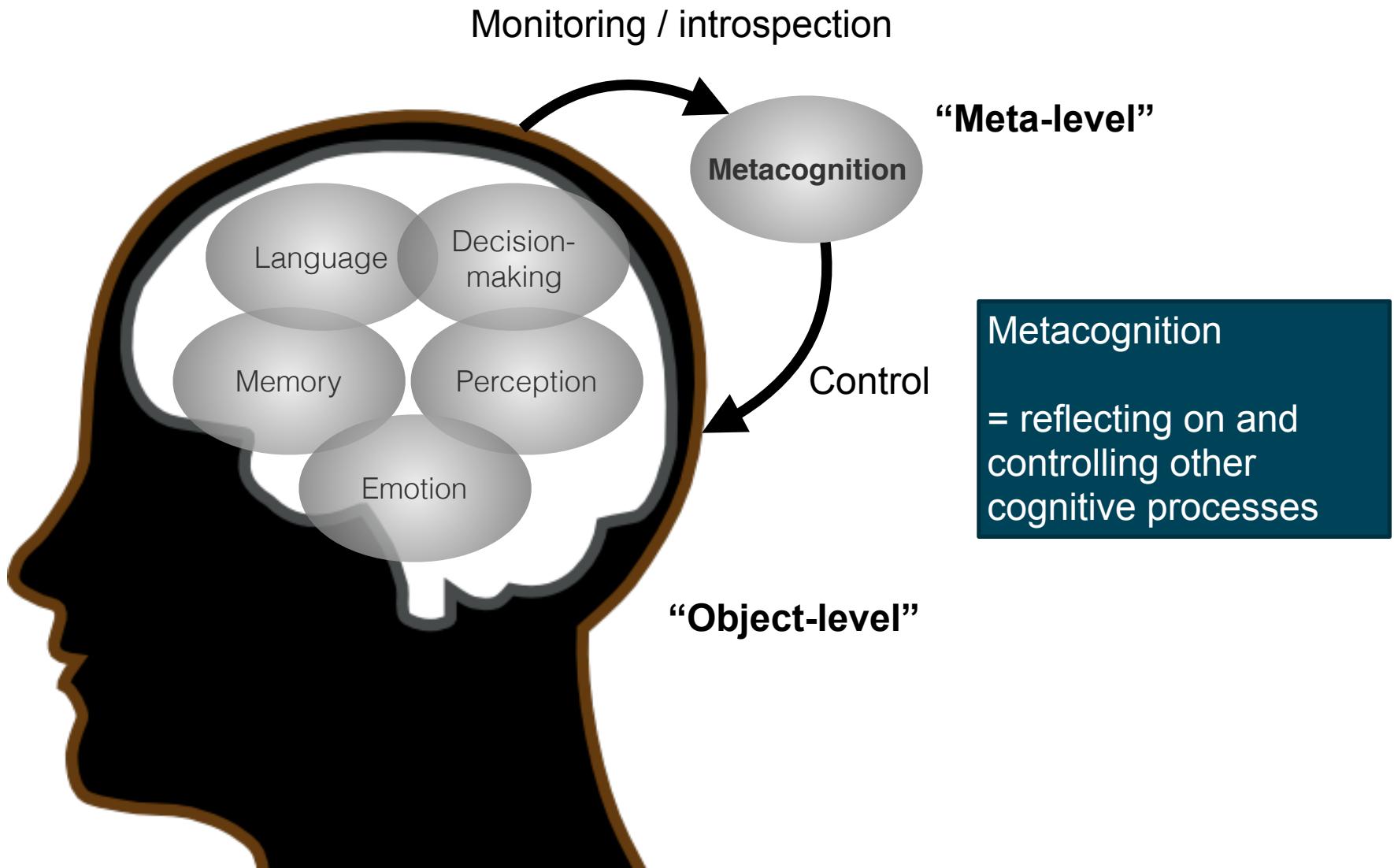
= reflecting on and controlling other cognitive processes

Monitoring / introspection



Metacognition

= reflecting on and controlling other cognitive processes



Metacognition
= reflecting on and
controlling other
cognitive processes

Questions for today



Questions for today

- How can we begin to study / quantify metacognition?



Questions for today

- How can we begin to study / quantify metacognition?
- What neurocognitive architecture supports metacognition?



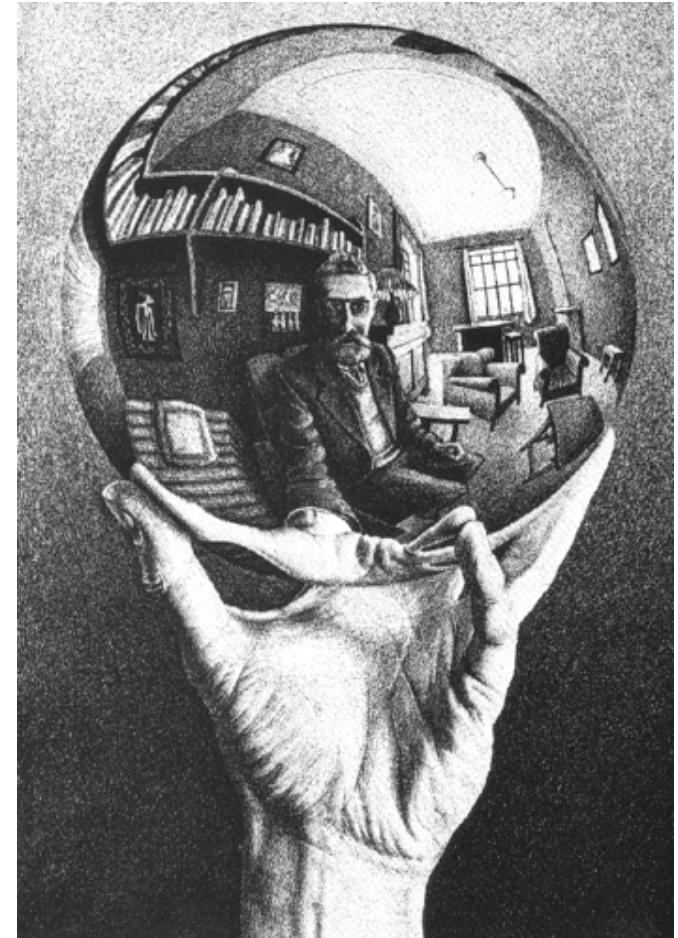
Questions for today

- How can we begin to study / quantify metacognition?
- What neurocognitive architecture supports metacognition?
- Is metacognition related to psychopathology?

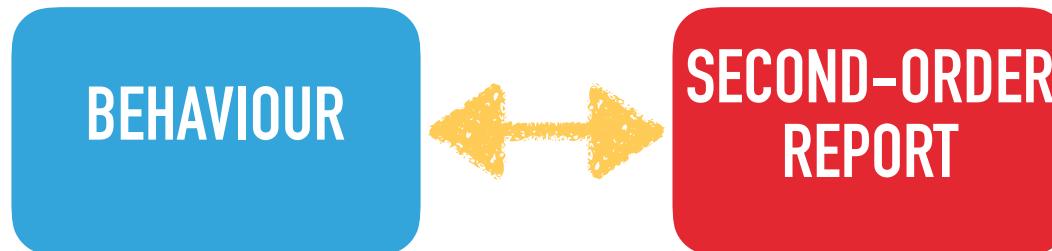


Questions for today

- How can we begin to study / quantify metacognition?
- What neurocognitive architecture supports metacognition?
- Is metacognition related to psychopathology?
- Can we design interventions to modulate (domain-general) metacognition?



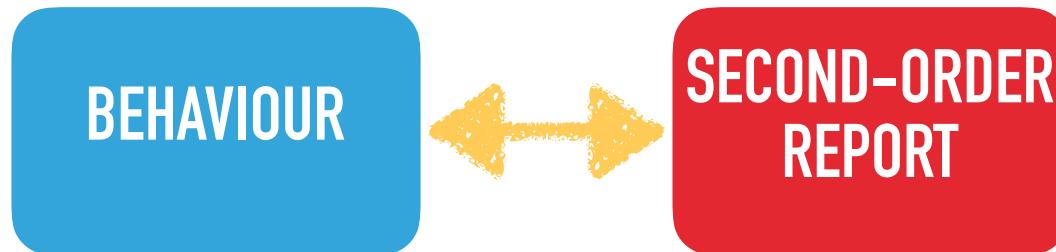
How can we measure metacognition?



E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

How can we measure metacognition?

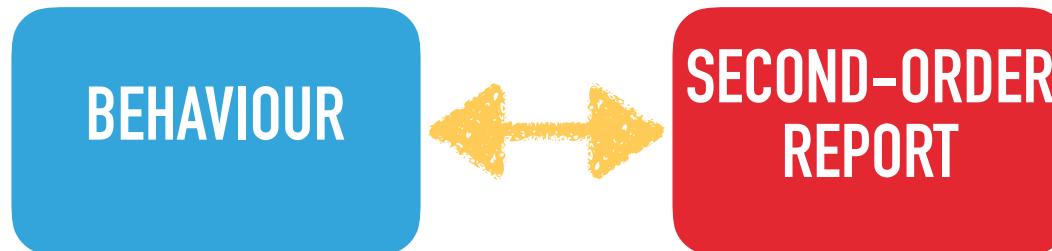


E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

- **Second-order reports**

How can we measure metacognition?

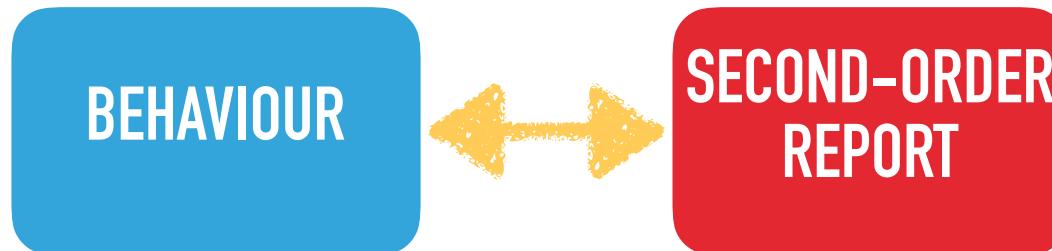


E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

- **Second-order reports**

How can we measure metacognition?

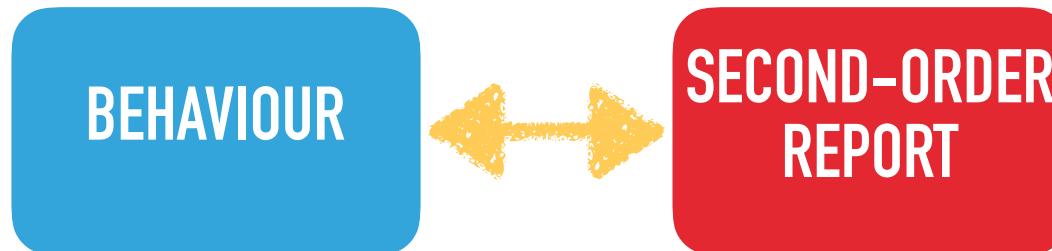


E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

- **Second-order reports**
 - Judgment or report about self-performance

How can we measure metacognition?



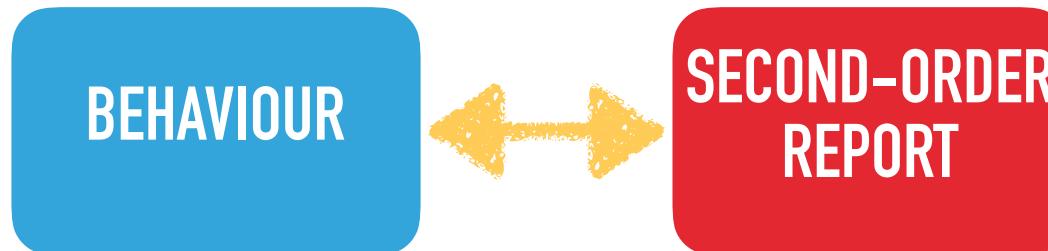
E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

- **Second-order reports**

- Judgment or report about self-performance
- Many different variants, prospective (before the task) or retrospective (after the task)

How can we measure metacognition?



E.g. answer to
exam question

E.g. **confidence** in
getting the answer
right

- **Second-order reports**

- Judgment or report about self-performance
- Many different variants, prospective (before the task) or retrospective (after the task)
- Metacognitive accuracy = relationship between second-order reports and behaviour

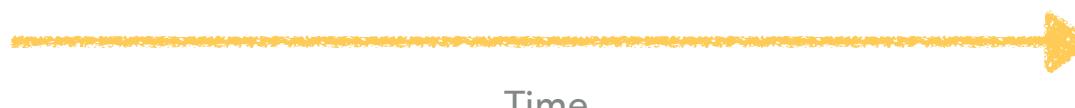
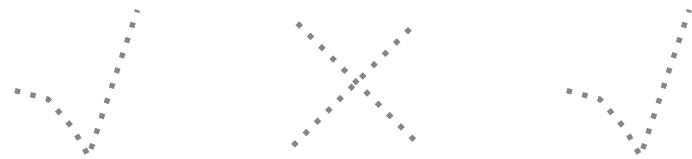
How do we study metacognition?

Quantify **statistical association** between behaviour and metacognitive judgments:

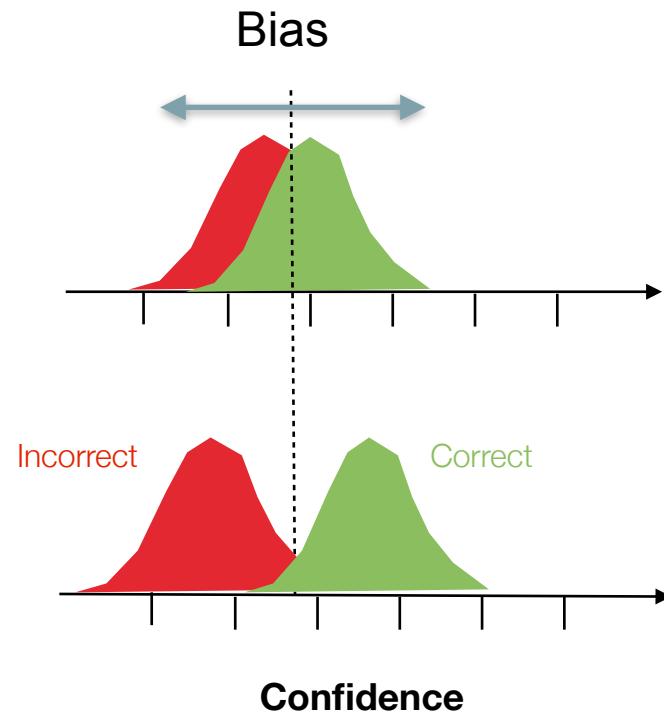
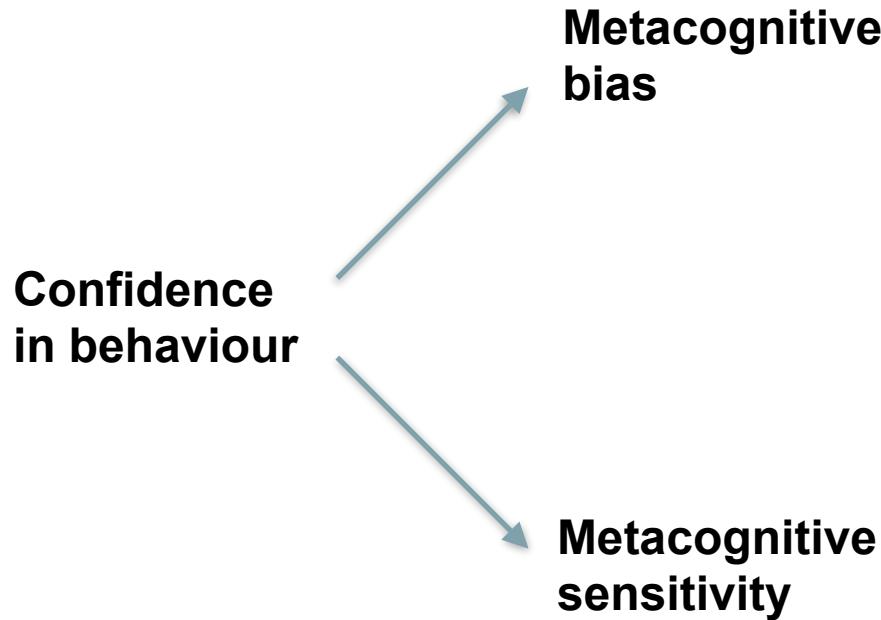
Subjective
confidence



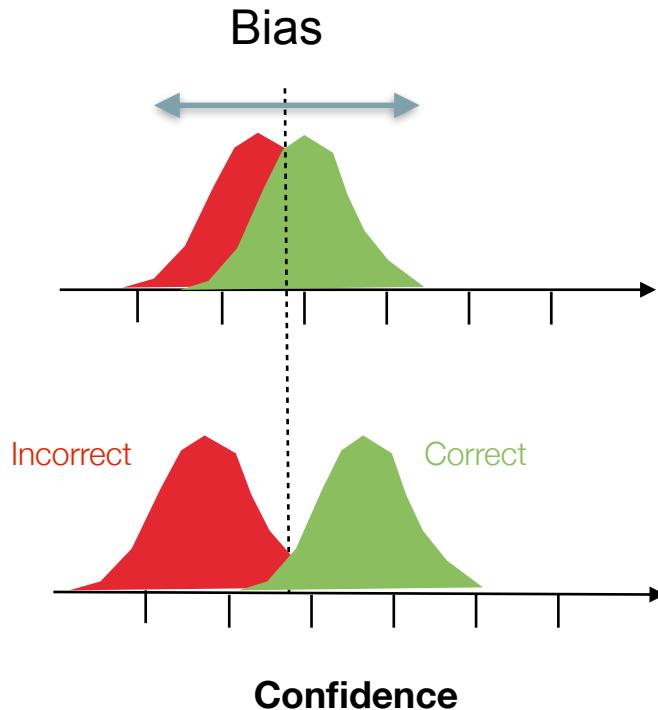
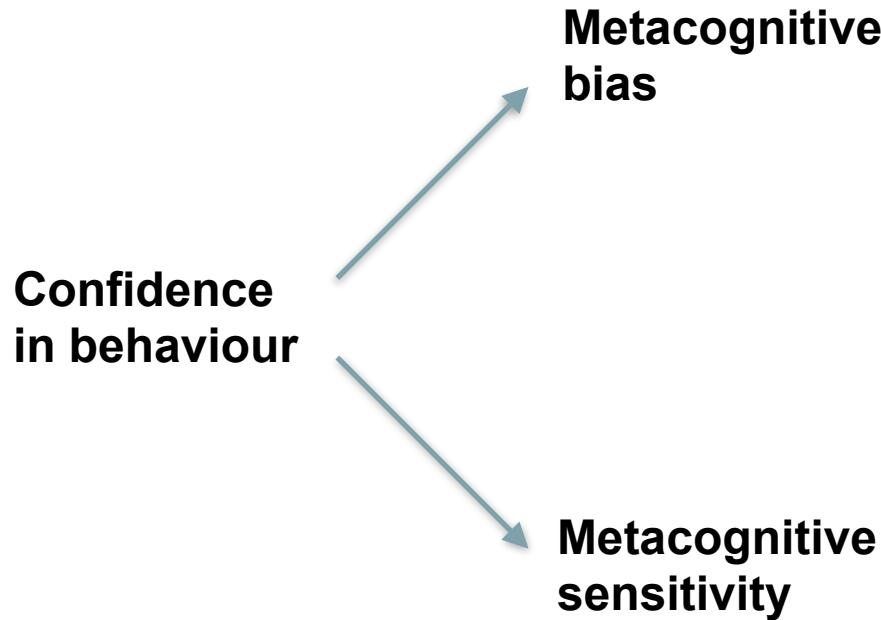
Objective
performance



Metacognitive bias and sensitivity



Metacognitive bias and sensitivity

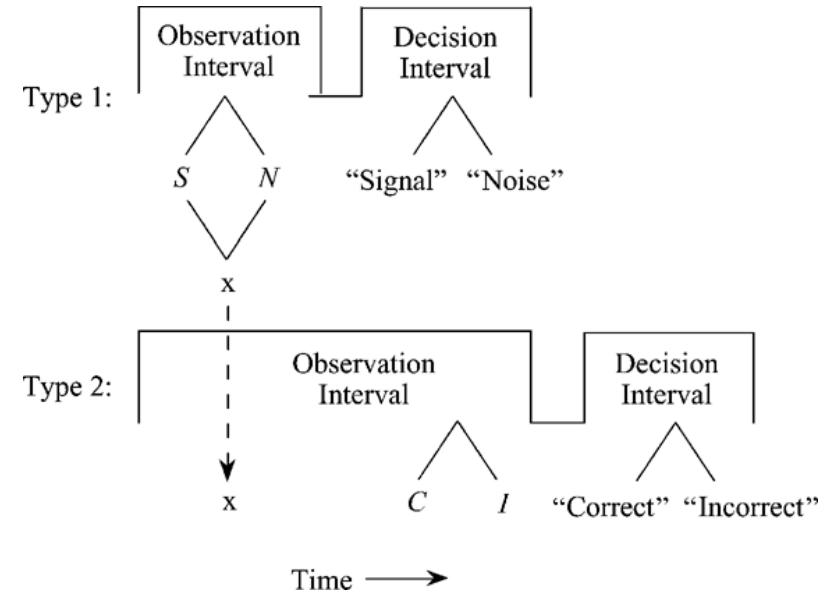
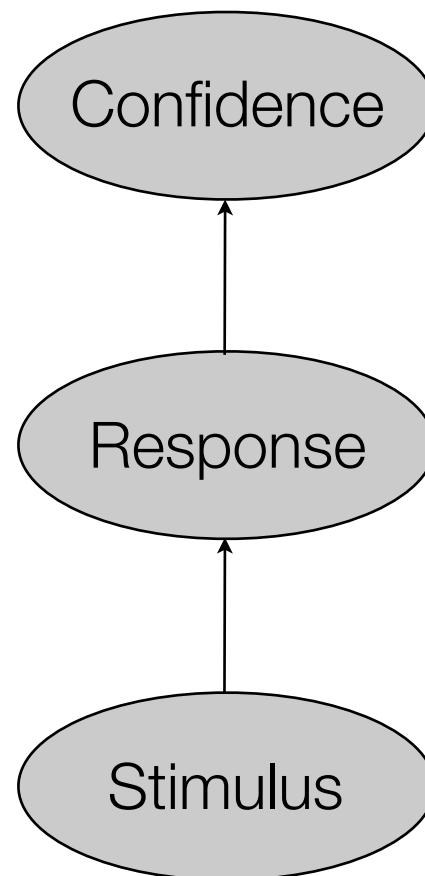


Other terminology in the literature:

Bias: calibration, confidence level, self-perceived ability, self-belief...

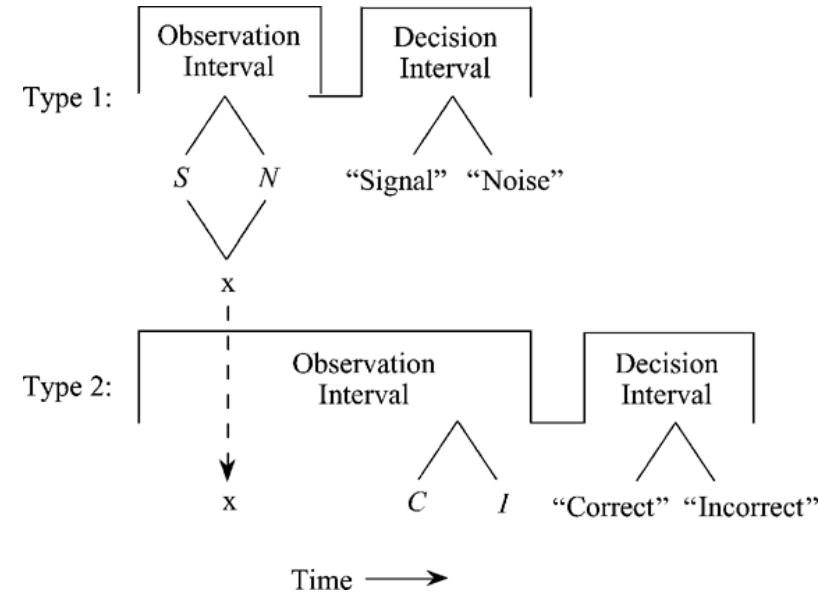
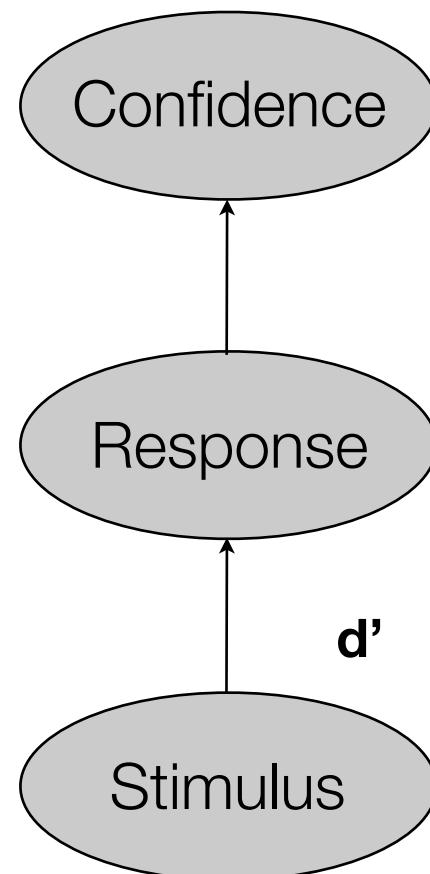
Sensitivity: discrimination, resolution, metacognitive awareness, insight...

Quantifying metacognitive sensitivity



Quantifying metacognitive sensitivity

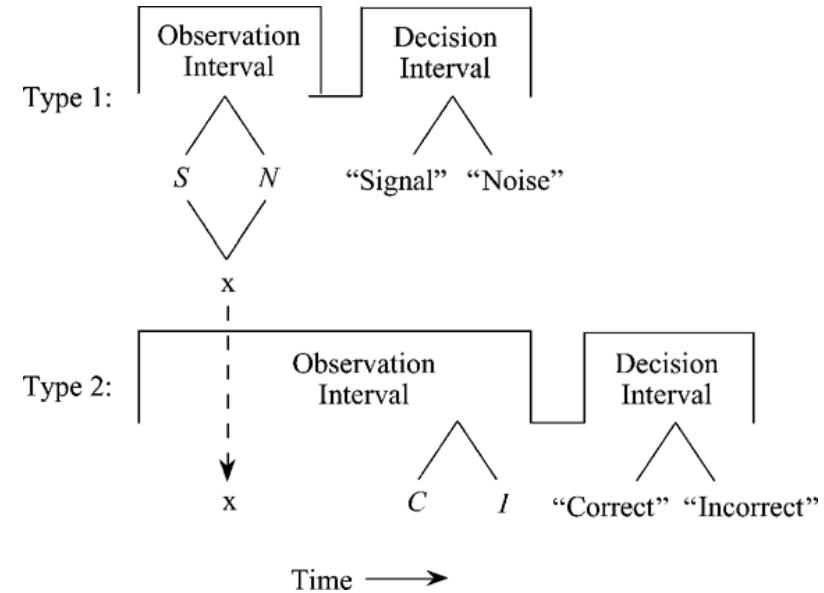
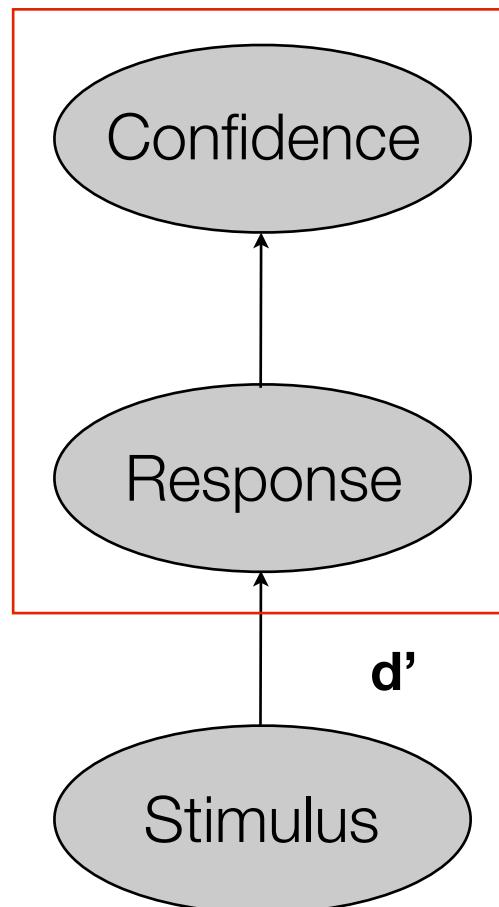
First-order
sensitivity



Quantifying metacognitive sensitivity

Second-order
sensitivity

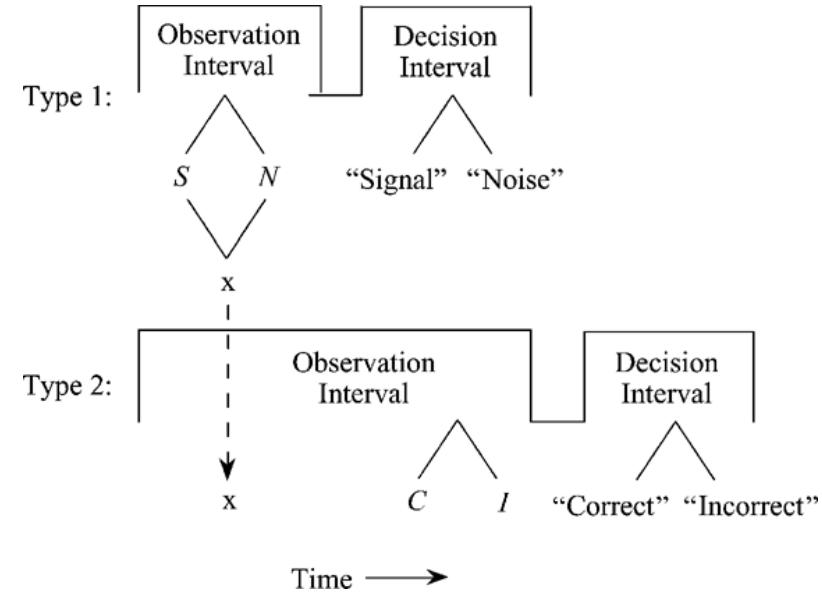
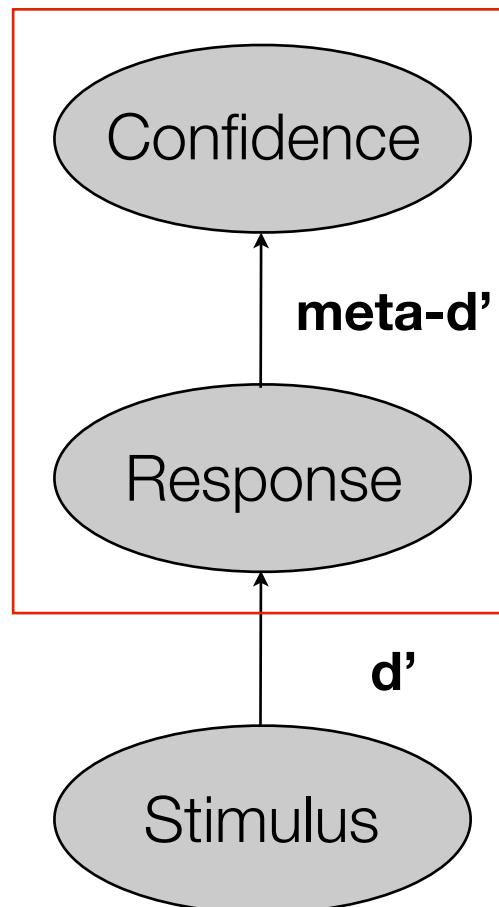
First-order
sensitivity



Quantifying metacognitive sensitivity

Second-order
sensitivity

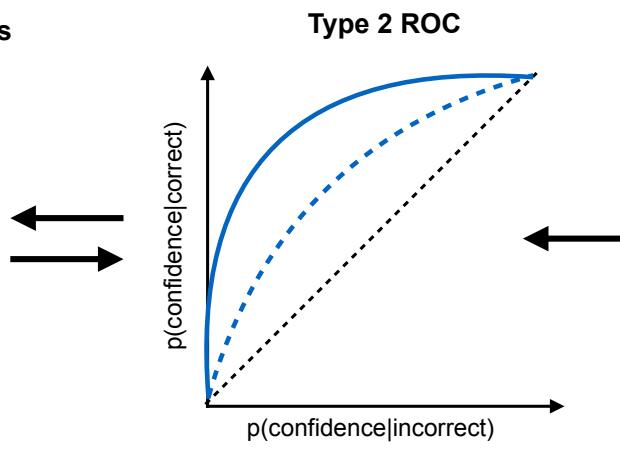
First-order
sensitivity



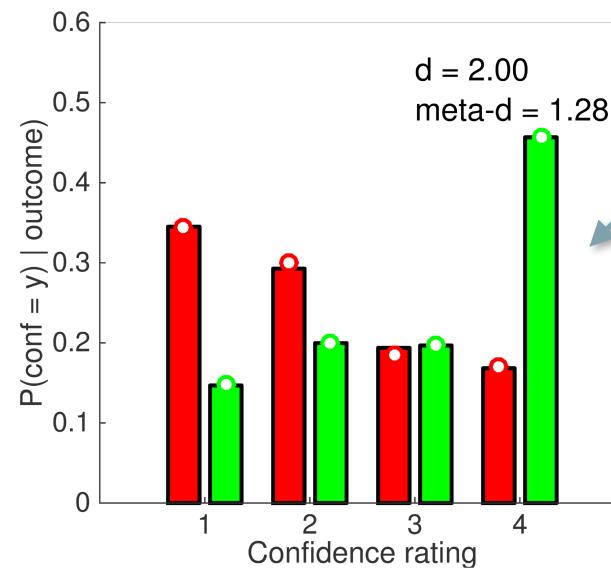
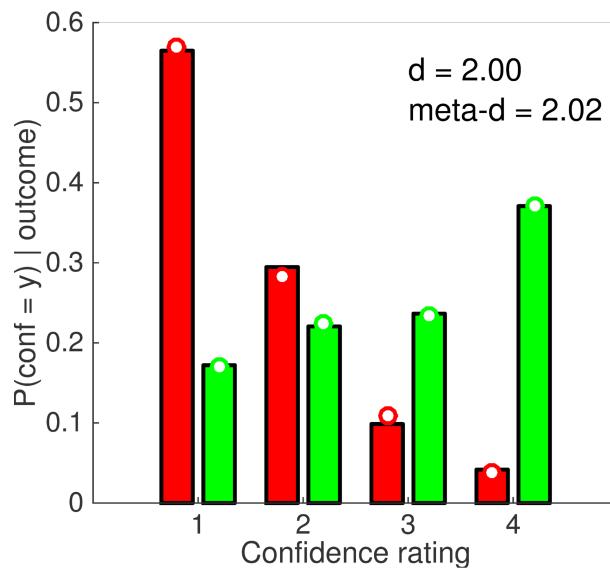
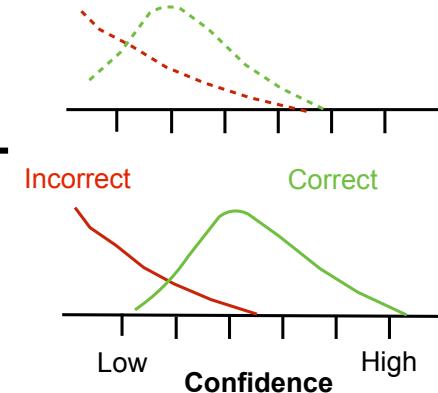
ROC analysis / meta- d'

Type 1 SDT parameters

meta- d' (fitted to type 2 ROC) compared to observed d'



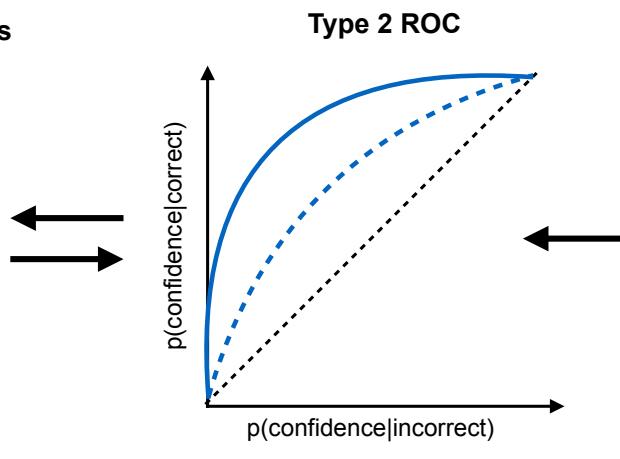
Observed confidence distributions



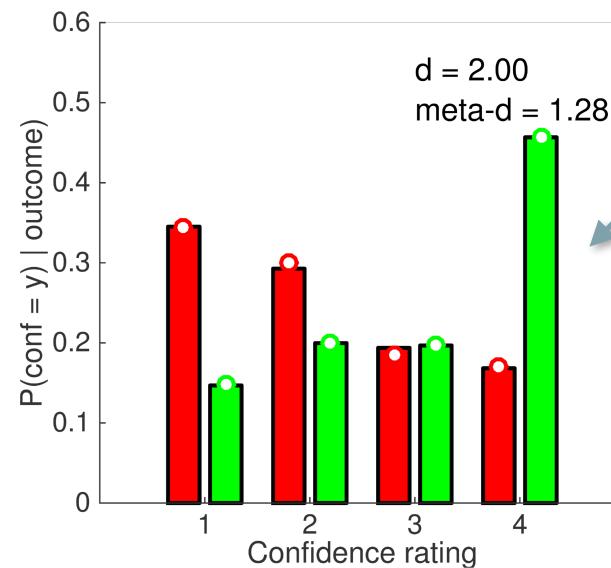
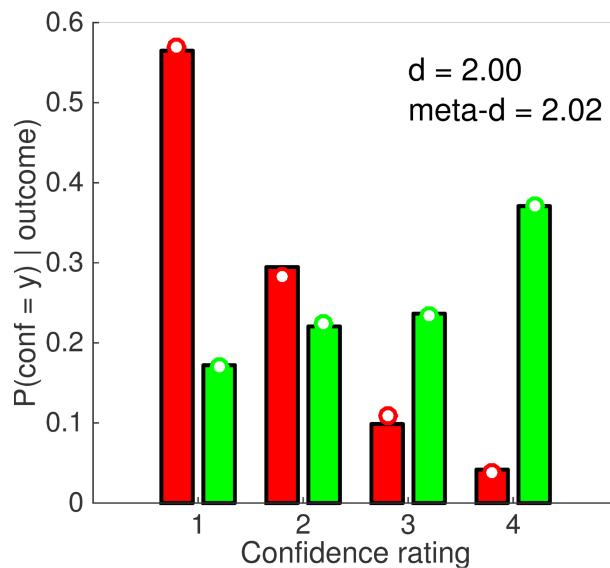
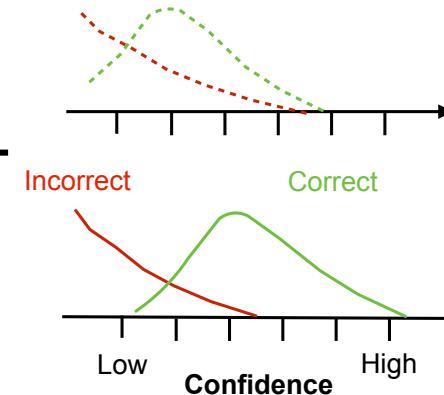
ROC analysis / meta- d'

Type 1 SDT parameters

meta- d' (fitted to type 2 ROC) compared to observed d'



Observed confidence distributions



Gaussian noise added to confidence ratings

meta- $d'/d' =$ metacognitive efficiency

Bayesian hierarchical model for meta-d' (HMeta-d)

$$\mu_{c2} \sim \mathcal{N}(0, 10)$$

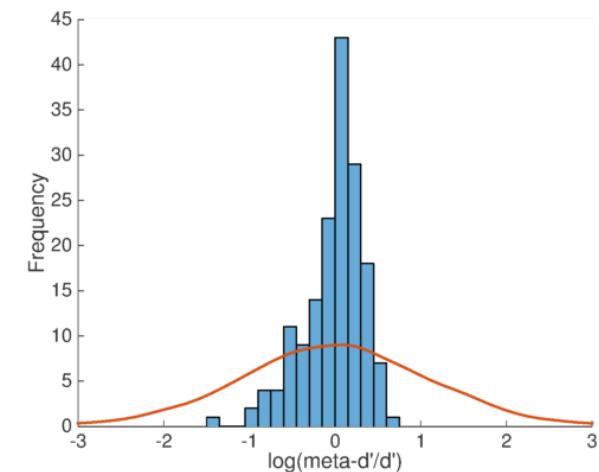
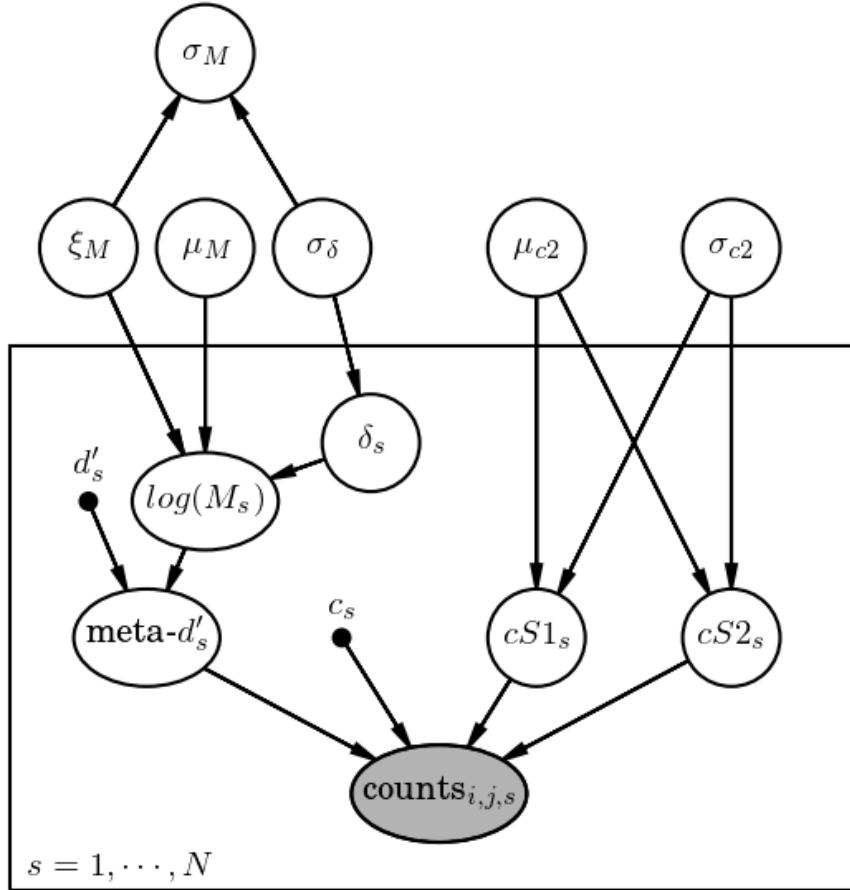
$$\sigma_{c2} \sim \mathcal{HN}(10)$$

$$\mu_M \sim N(0, 1)$$

$$\sigma_M = |\xi_M| \times \delta_s$$

$$\xi_M \sim Beta(1, 1)$$

$$\sigma_\delta \sim \mathcal{HN}(1)$$



Bayesian hierarchical model for meta-d' (HMeta-d)

$$\mu_{c2} \sim \mathcal{N}(0, 10)$$

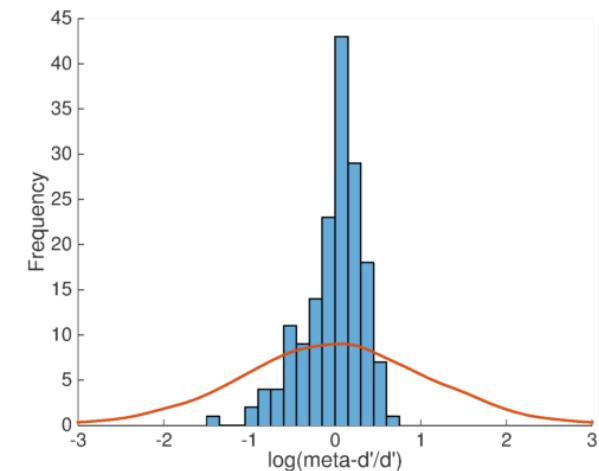
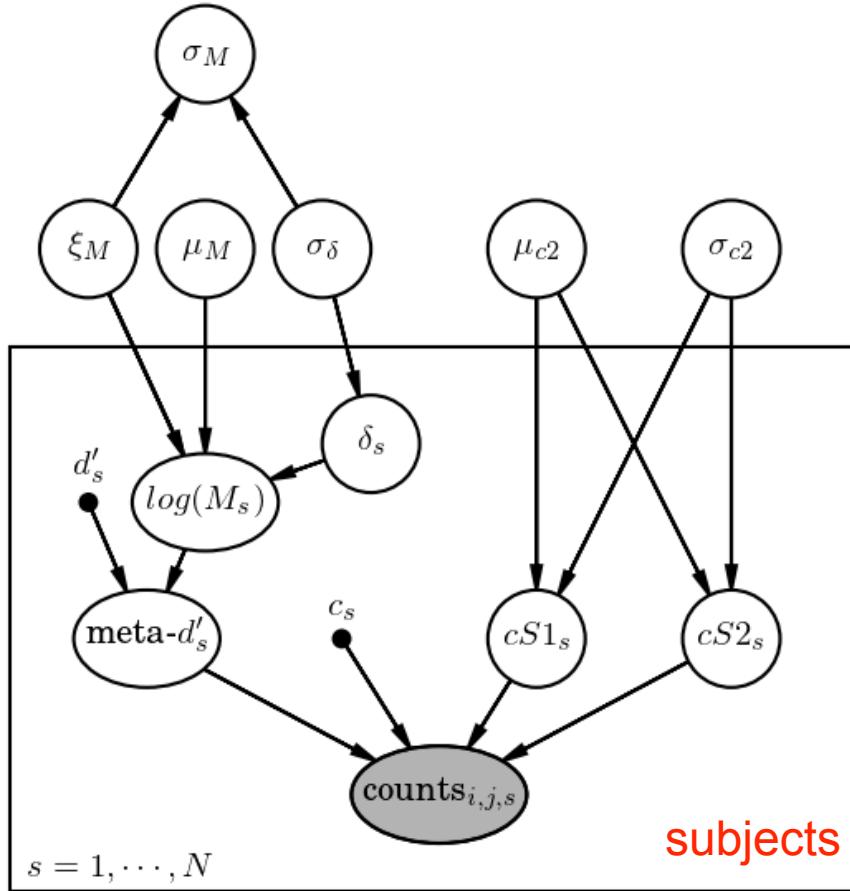
$$\sigma_{c2} \sim \mathcal{HN}(10)$$

$$\mu_M \sim N(0, 1)$$

$$\sigma_M = |\xi_M| \times \delta_s$$

$$\xi_M \sim Beta(1, 1)$$

$$\sigma_\delta \sim \mathcal{HN}(1)$$



Bayesian hierarchical model for meta-d' (HMeta-d)

$$\mu_{c2} \sim \mathcal{N}(0, 10)$$

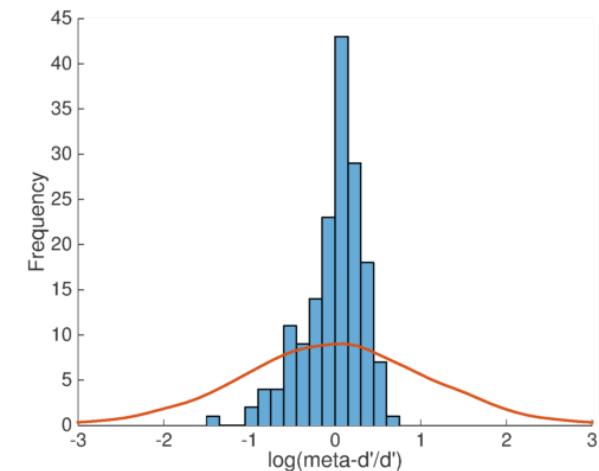
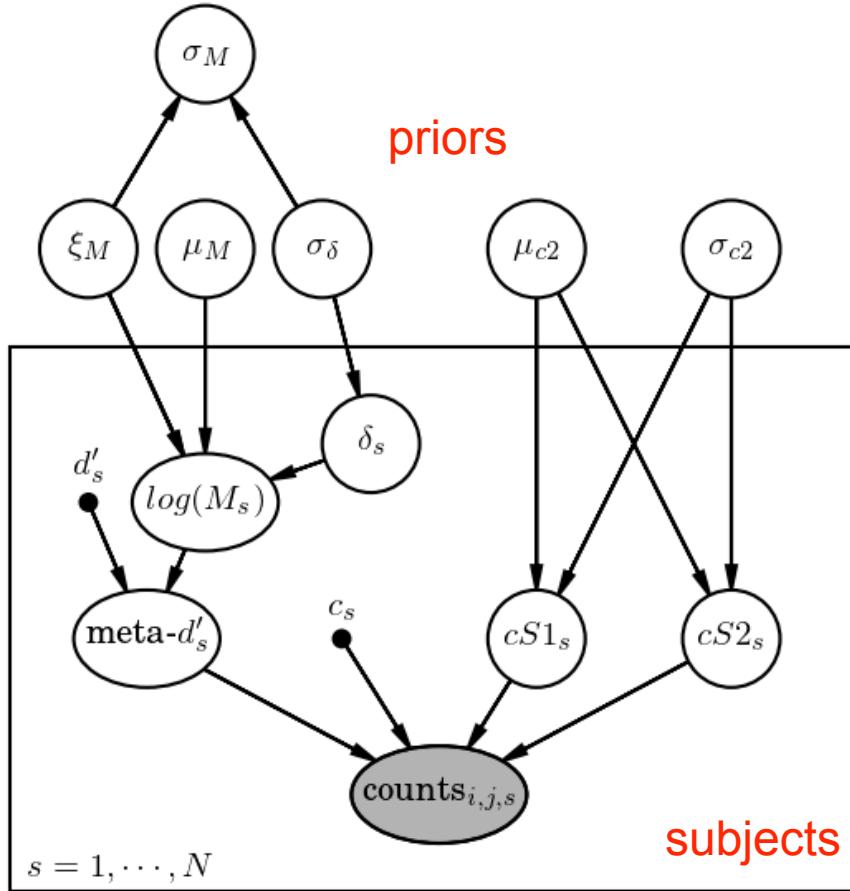
$$\sigma_{c2} \sim \mathcal{HN}(10)$$

$$\mu_M \sim N(0, 1)$$

$$\sigma_M = |\xi_M| \times \delta_s$$

$$\xi_M \sim Beta(1, 1)$$

$$\sigma_\delta \sim \mathcal{HN}(1)$$



Bayesian hierarchical model for meta-d' (HMeta-d)

$$\mu_{c2} \sim \mathcal{N}(0, 10)$$

$$\sigma_{c2} \sim \mathcal{HN}(10)$$

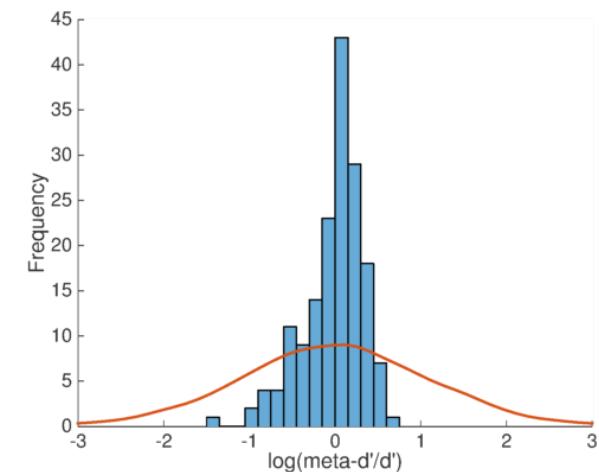
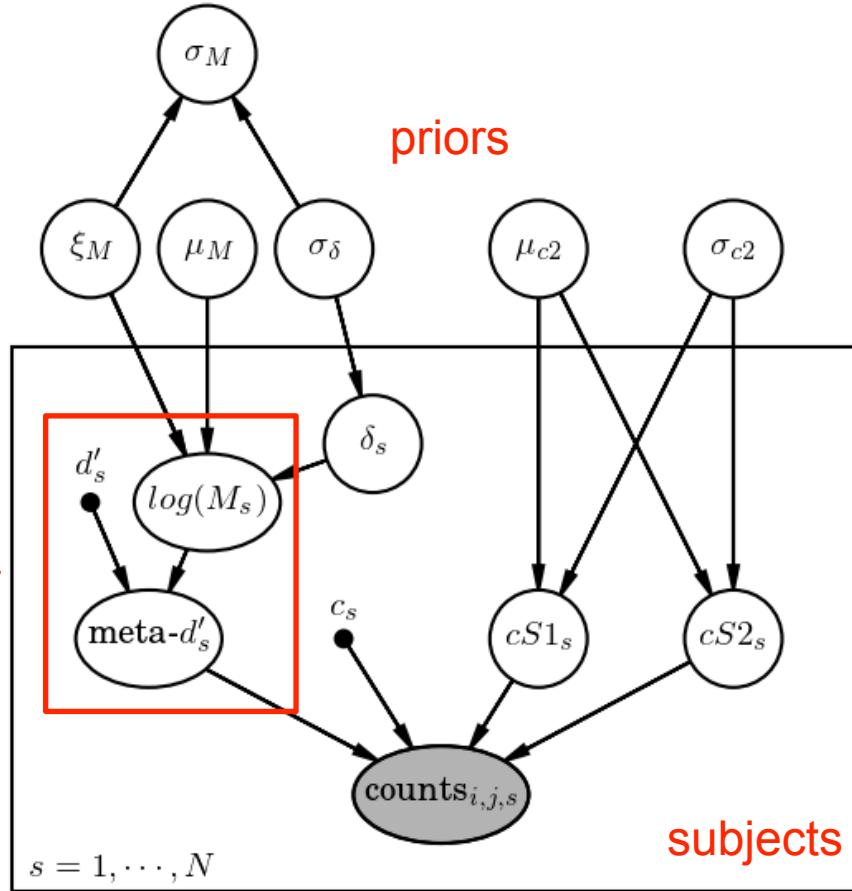
$$\mu_M \sim N(0, 1)$$

$$\sigma_M = |\xi_M| \times \delta_s$$

$$\xi_M \sim Beta(1, 1)$$

$$\sigma_\delta \sim \mathcal{HN}(1)$$

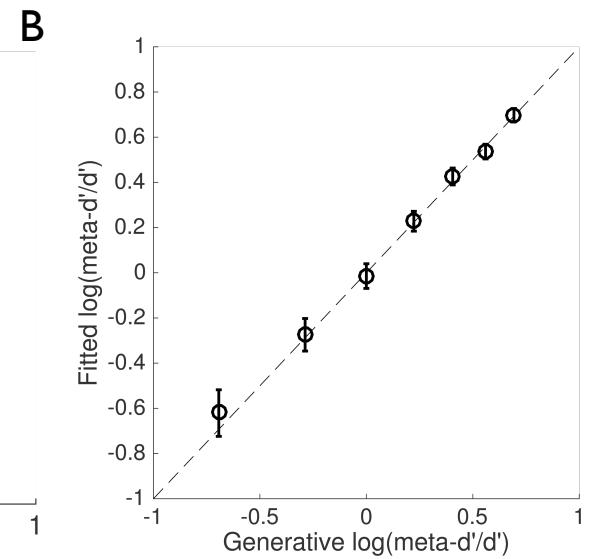
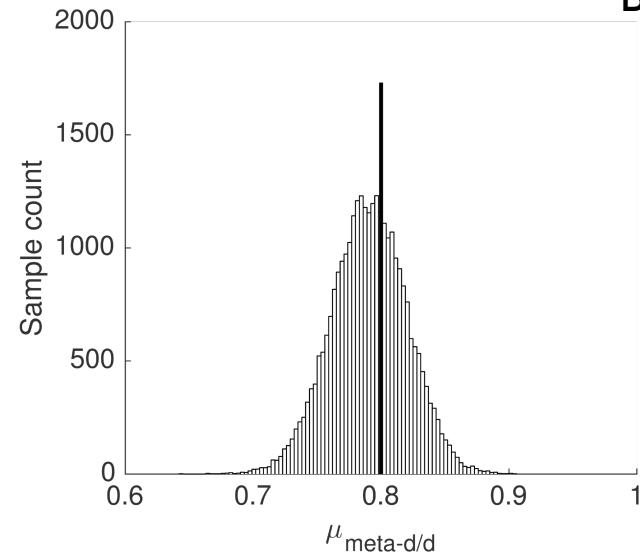
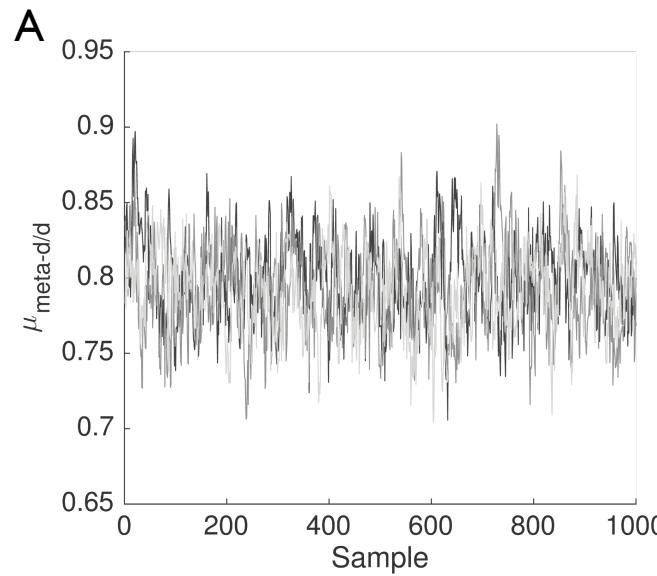
direct
estimation of
meta-d'/d'



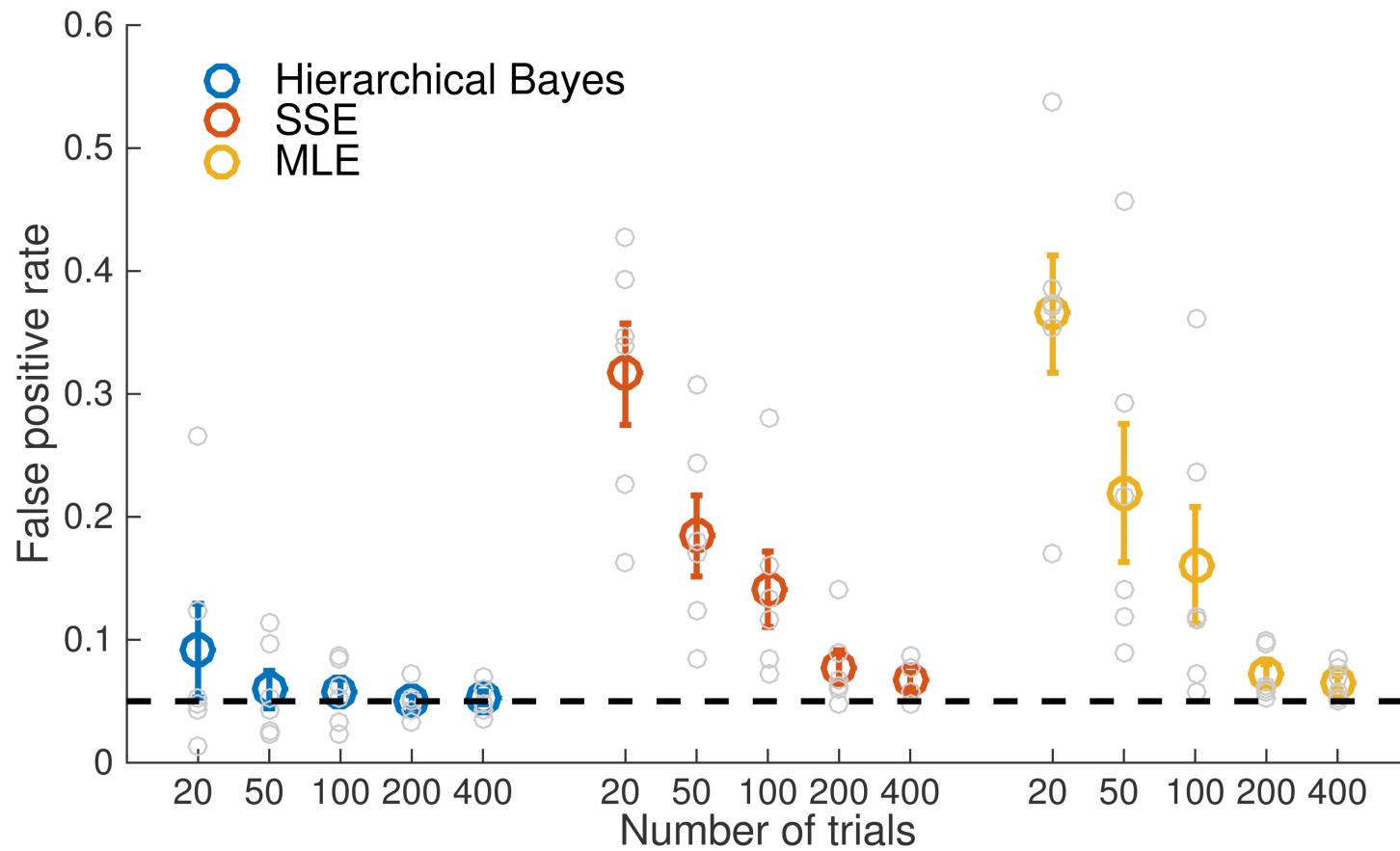
Hierarchical model for meta- d' (HMeta-d)

Code and tutorial available at <https://github.com/smfleming/HMeta-d>

MCMC samples of group-level metacognitive efficiency:

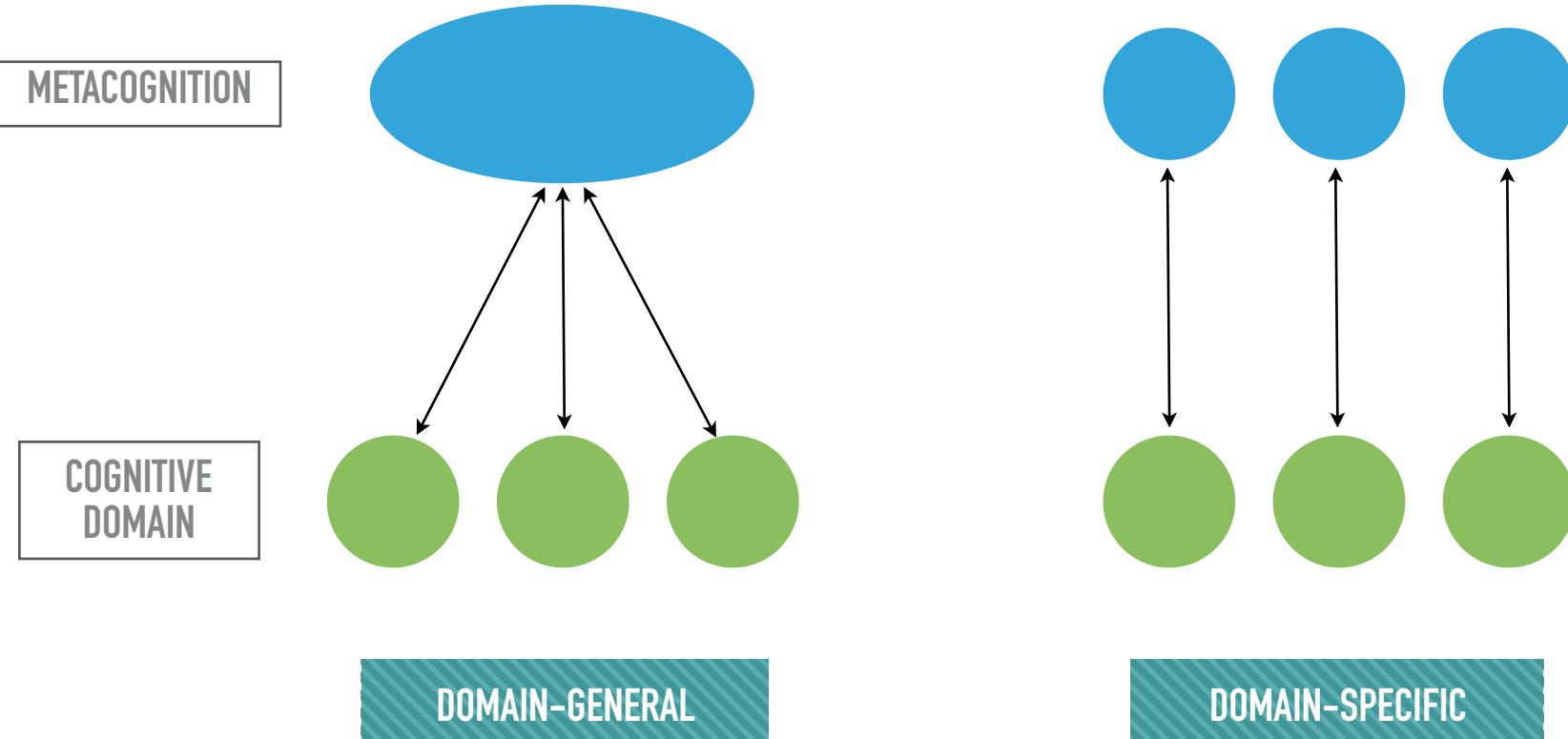


Hierarchical model for meta- d' (HMeta-d)



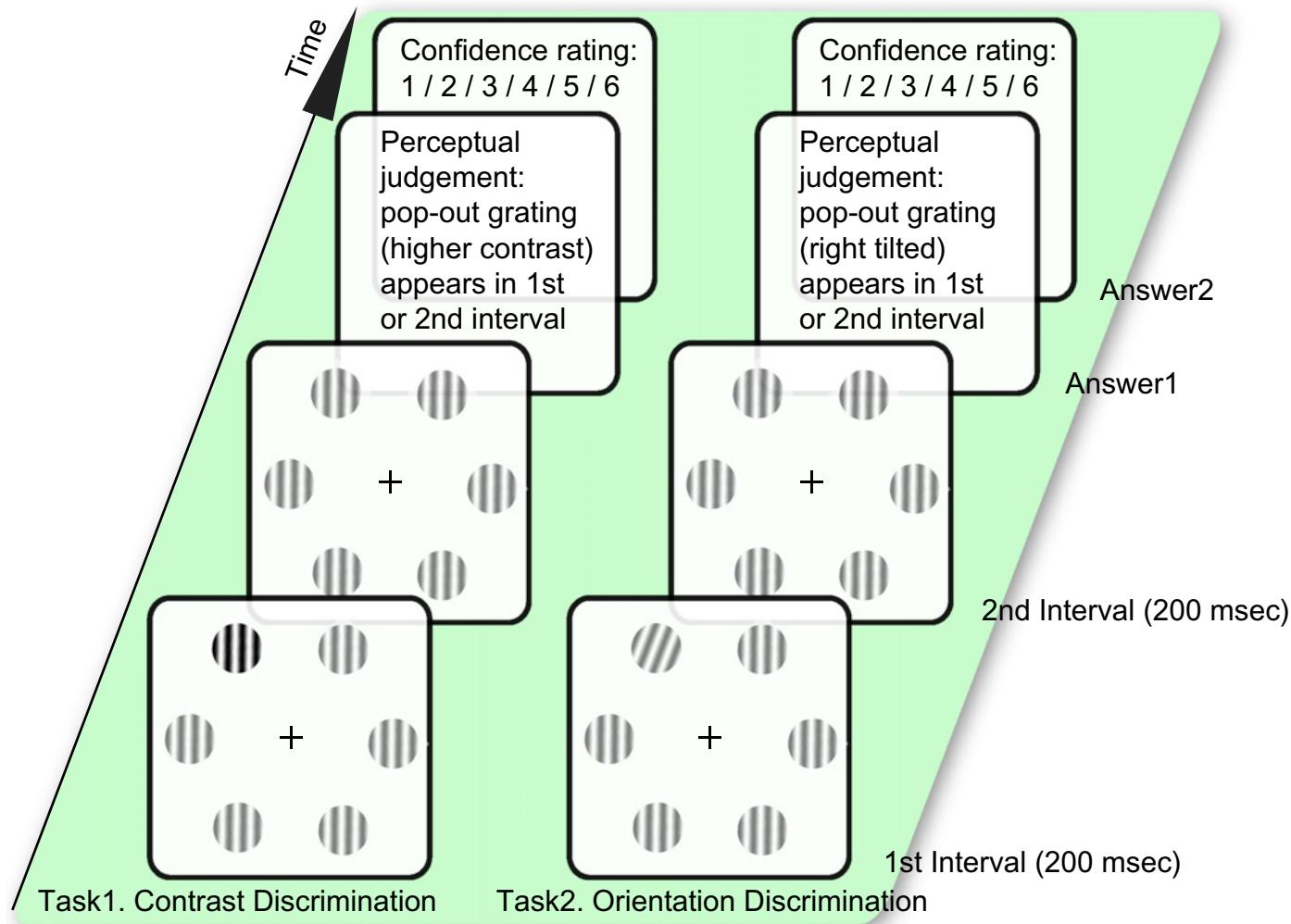
Accurate estimates of **group-level metacognitive efficiency** possible
with relatively low trial numbers when using HMeta-d

Domain-general or domain-specific?

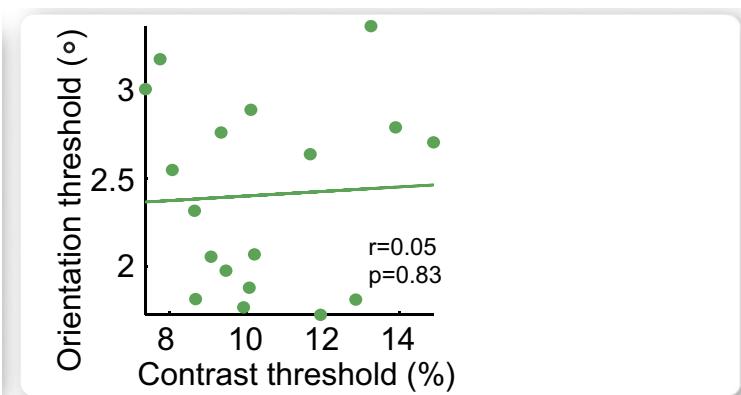


Two broad approaches:

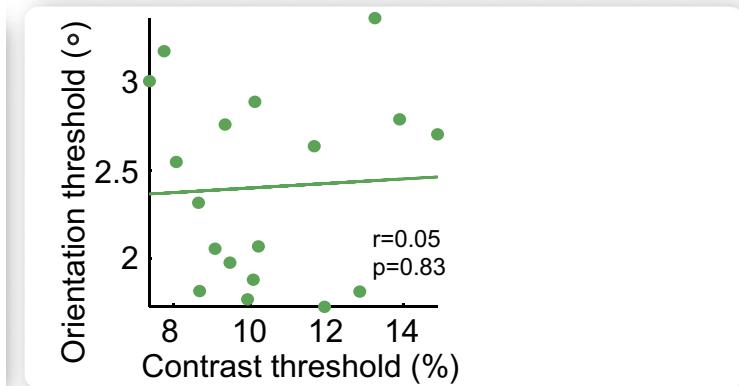
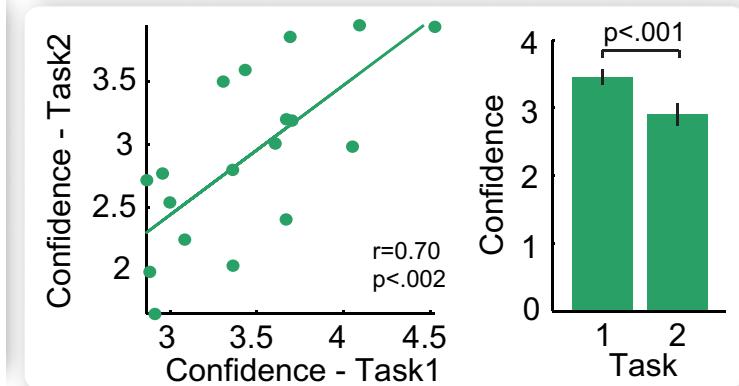
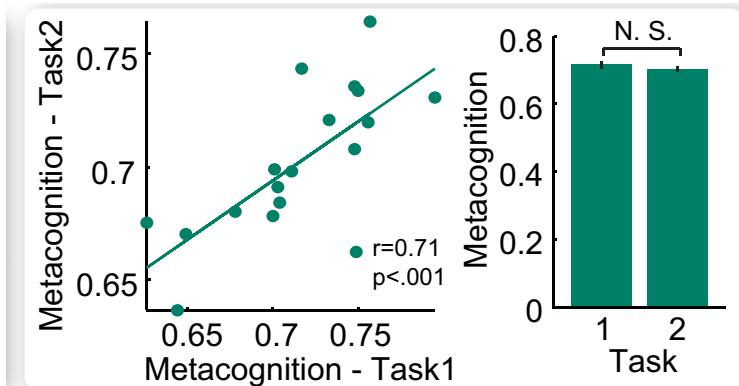
- Analysis of shared variance in individual differences
- Identification of shared/distinct neural resources supporting metacognitive ability



Does having good metacognition on task 1 predict good metacognition on task 2?



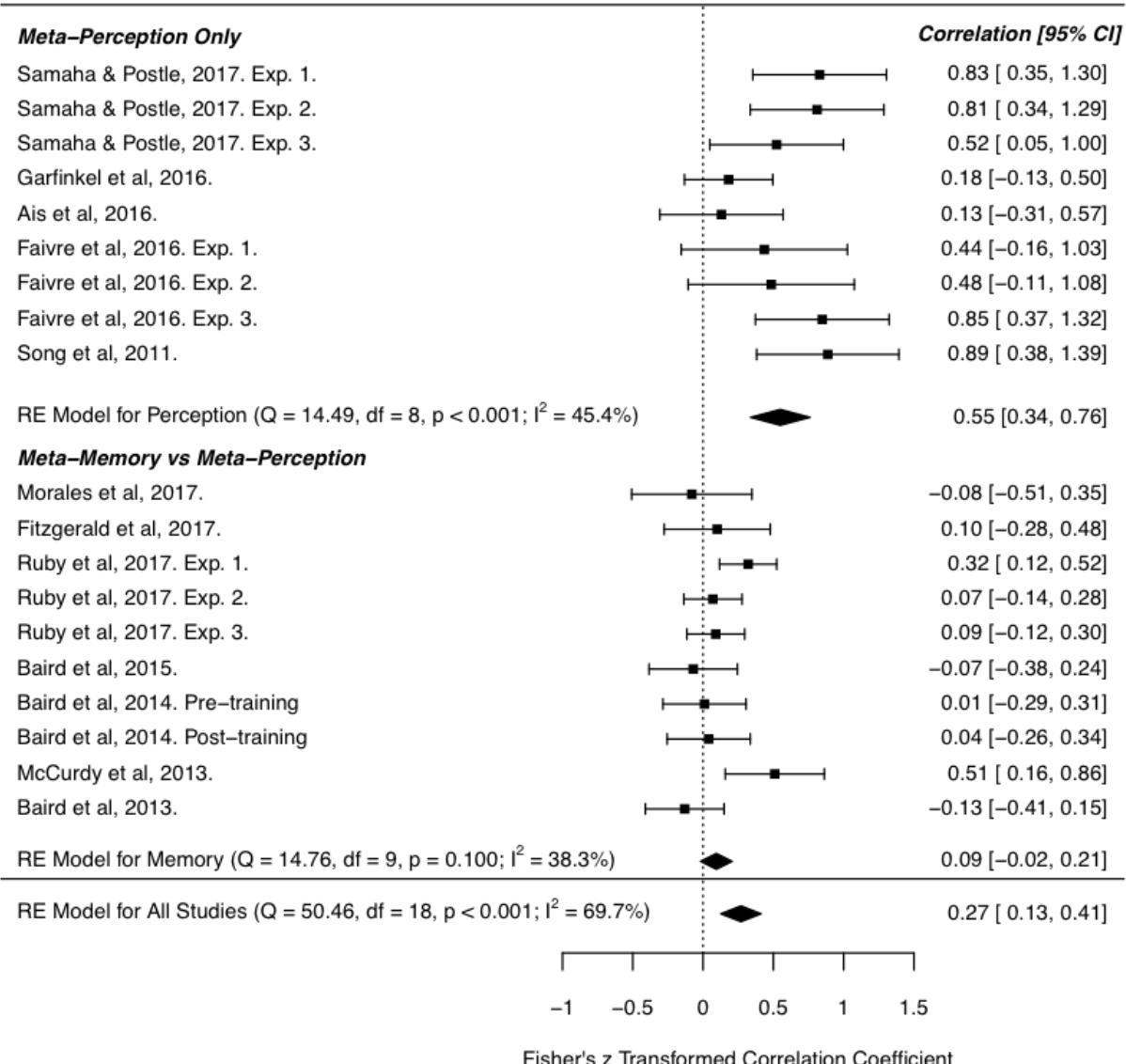
Performance
thresholds uncorrelated



Significant correlations
in metacognition
(sensitivity and bias)
across tasks

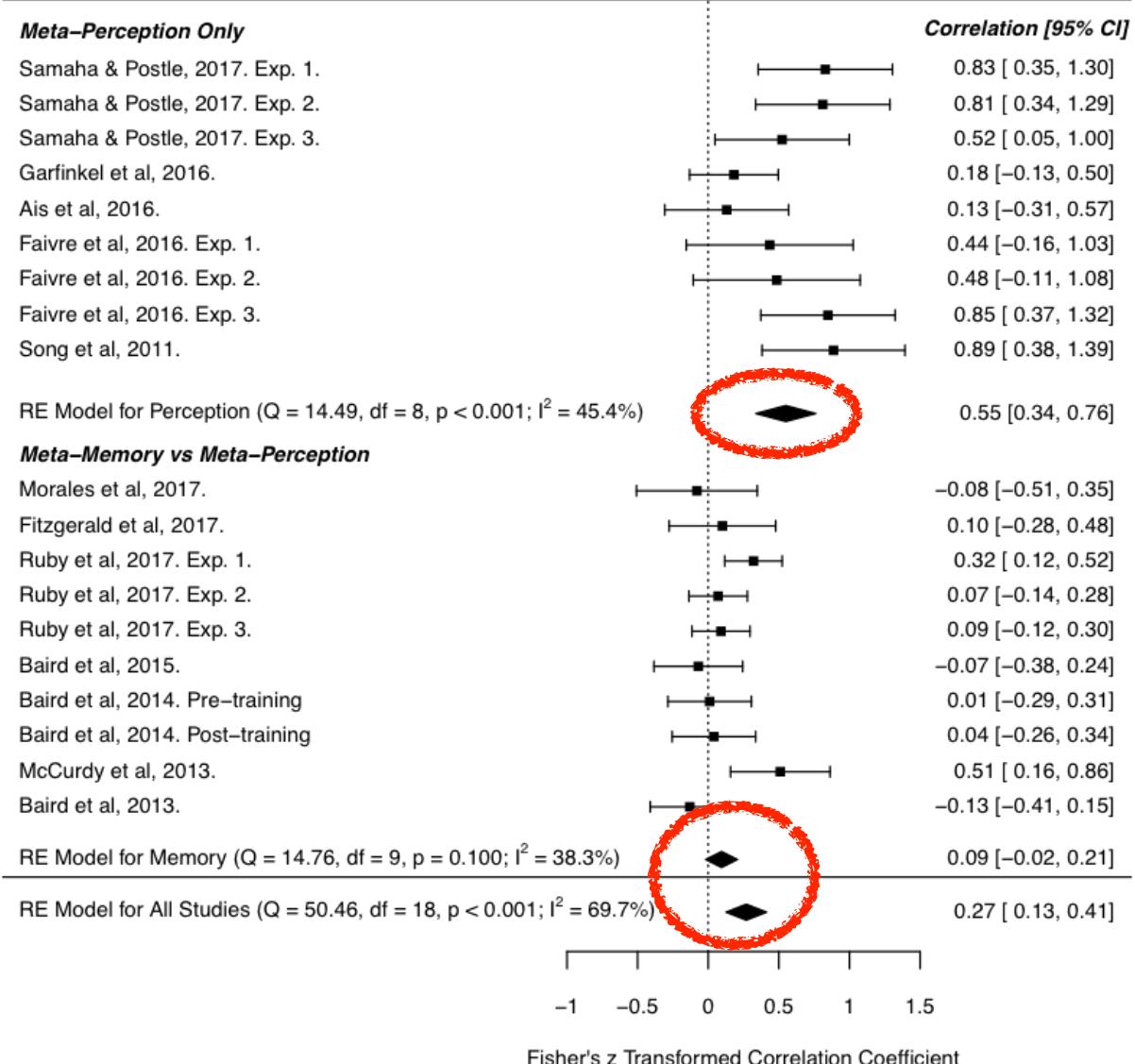
Performance
thresholds uncorrelated

Meta-analysis of efficiency correlations



Meta-analysis of efficiency correlations

Some evidence for domain-generality, but lack of consistency across task designs / low power



N=181,
hierarchical
modelling of
metacognitive
efficiency across
4 2AFC tasks

N=181,
hierarchical
modelling of
metacognitive
efficiency across
4 2AFC tasks

$$[\log(M1_s) \log(M2_s) \log(M3_s) \log(M4_s)] \sim N \left(\begin{bmatrix} \mu_{M1} \\ \mu_{M2} \\ \mu_{M3} \\ \mu_{M4} \end{bmatrix}, \begin{bmatrix} \sigma^2_{M1} & \rho_{M1M2}\sigma_{M1}\sigma_{M2} & \rho_{M1M3}\sigma_{M1}\sigma_{M3} & \rho_{M1M4}\sigma_{M1}\sigma_{M4} \\ \rho_{M1M2}\sigma_{M1}\sigma_{M2} & \sigma^2_{M2} & \rho_{M2M3}\sigma_{M2}\sigma_{M3} & \rho_{M2M4}\sigma_{M2}\sigma_{M4} \\ \rho_{M1M3}\sigma_{M1}\sigma_{M3} & \rho_{M2M3}\sigma_{M2}\sigma_{M3} & \sigma^2_{M3} & \rho_{M3M4}\sigma_{M3}\sigma_{M4} \\ \rho_{M1M4}\sigma_{M1}\sigma_{M4} & \rho_{M2M4}\sigma_{M2}\sigma_{M4} & \rho_{M3M4}\sigma_{M3}\sigma_{M4} & \sigma^2_{M4} \end{bmatrix} \right)$$

Priors were specified as follows:

$$\mu_{M1}, \mu_{M2}, \mu_{M3}, \mu_{M4} \sim N(0, 1)$$

$$\sigma_{M1}, \sigma_{M2}, \sigma_{M3}, \sigma_{M4} \sim \text{InvSqrtGamma}(0.001, 0.001)$$

$$\rho_{M1M2}, \rho_{M1M3}, \rho_{M1M4}, \rho_{M2M3}, \rho_{M2M4}, \rho_{M3M4} \sim \text{Uniform}(-1, 1)$$

N=181,
hierarchical
modelling of
metacognitive
efficiency across
4 2AFC tasks

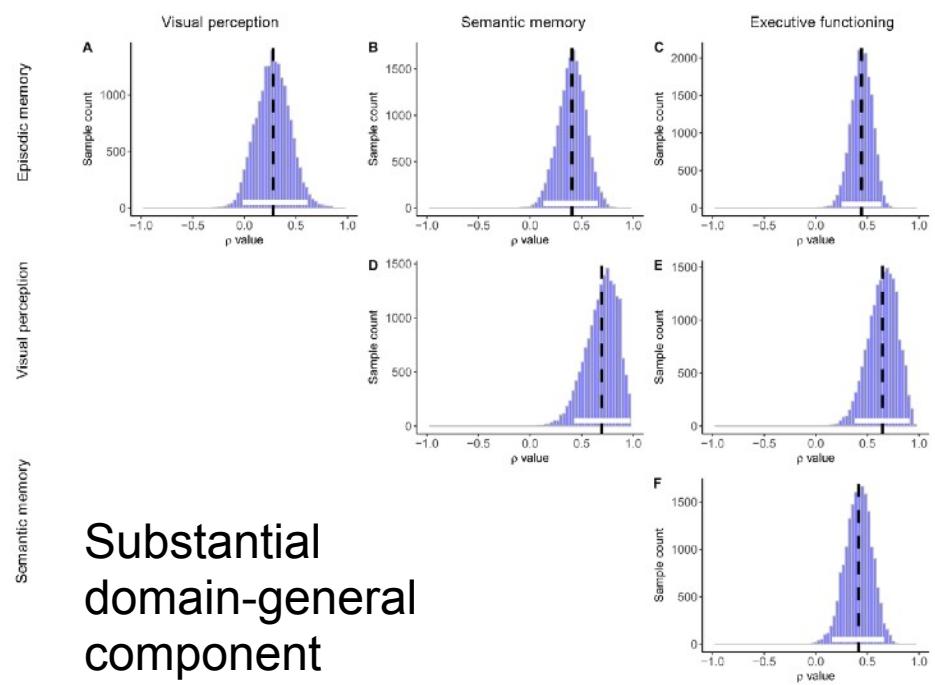
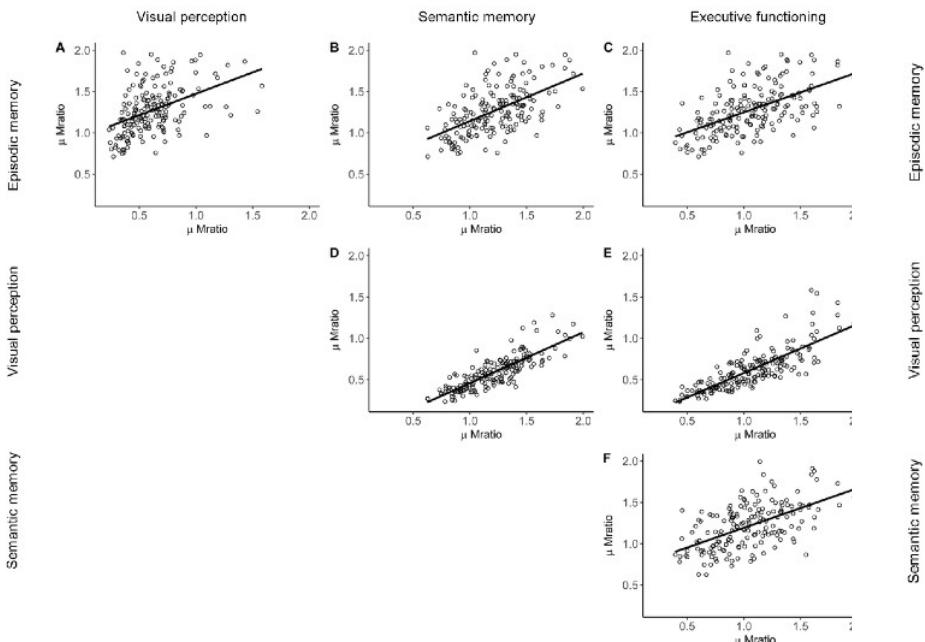
$$[\log(M1_s) \log(M2_s) \log(M3_s) \log(M4_s)] \sim N \left(\begin{bmatrix} \mu_{M1} \\ \mu_{M2} \\ \mu_{M3} \\ \mu_{M4} \end{bmatrix}, \begin{bmatrix} \sigma^2_{M1} & \rho_{M1M2}\sigma_{M1}\sigma_{M2} & \rho_{M1M3}\sigma_{M1}\sigma_{M3} & \rho_{M1M4}\sigma_{M1}\sigma_{M4} \\ \rho_{M1M2}\sigma_{M1}\sigma_{M2} & \sigma^2_{M2} & \rho_{M2M3}\sigma_{M2}\sigma_{M3} & \rho_{M2M4}\sigma_{M2}\sigma_{M4} \\ \rho_{M1M3}\sigma_{M1}\sigma_{M3} & \rho_{M2M3}\sigma_{M2}\sigma_{M3} & \sigma^2_{M3} & \rho_{M3M4}\sigma_{M3}\sigma_{M4} \\ \rho_{M1M4}\sigma_{M1}\sigma_{M4} & \rho_{M2M4}\sigma_{M2}\sigma_{M4} & \rho_{M3M4}\sigma_{M3}\sigma_{M4} & \sigma^2_{M4} \end{bmatrix} \right)$$

Priors were specified as follows:

$$\mu_{M1}, \mu_{M2}, \mu_{M3}, \mu_{M4} \sim N(0, 1)$$

$$\sigma_{M1}, \sigma_{M2}, \sigma_{M3}, \sigma_{M4} \sim \text{InvSqrtGamma}(0.001, 0.001)$$

$$\rho_{M1M2}, \rho_{M1M3}, \rho_{M1M4}, \rho_{M2M3}, \rho_{M2M4}, \rho_{M3M4} \sim \text{Uniform}(-1, 1)$$

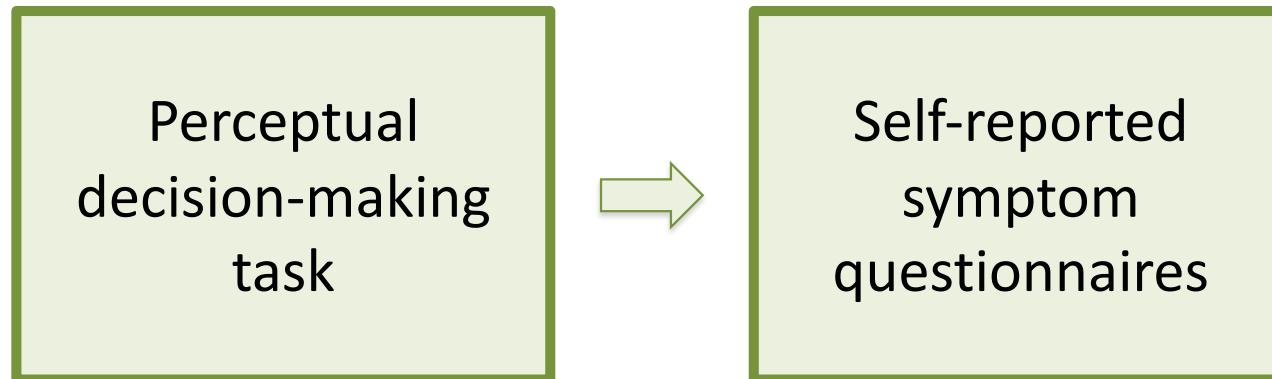


Substantial
domain-general
component

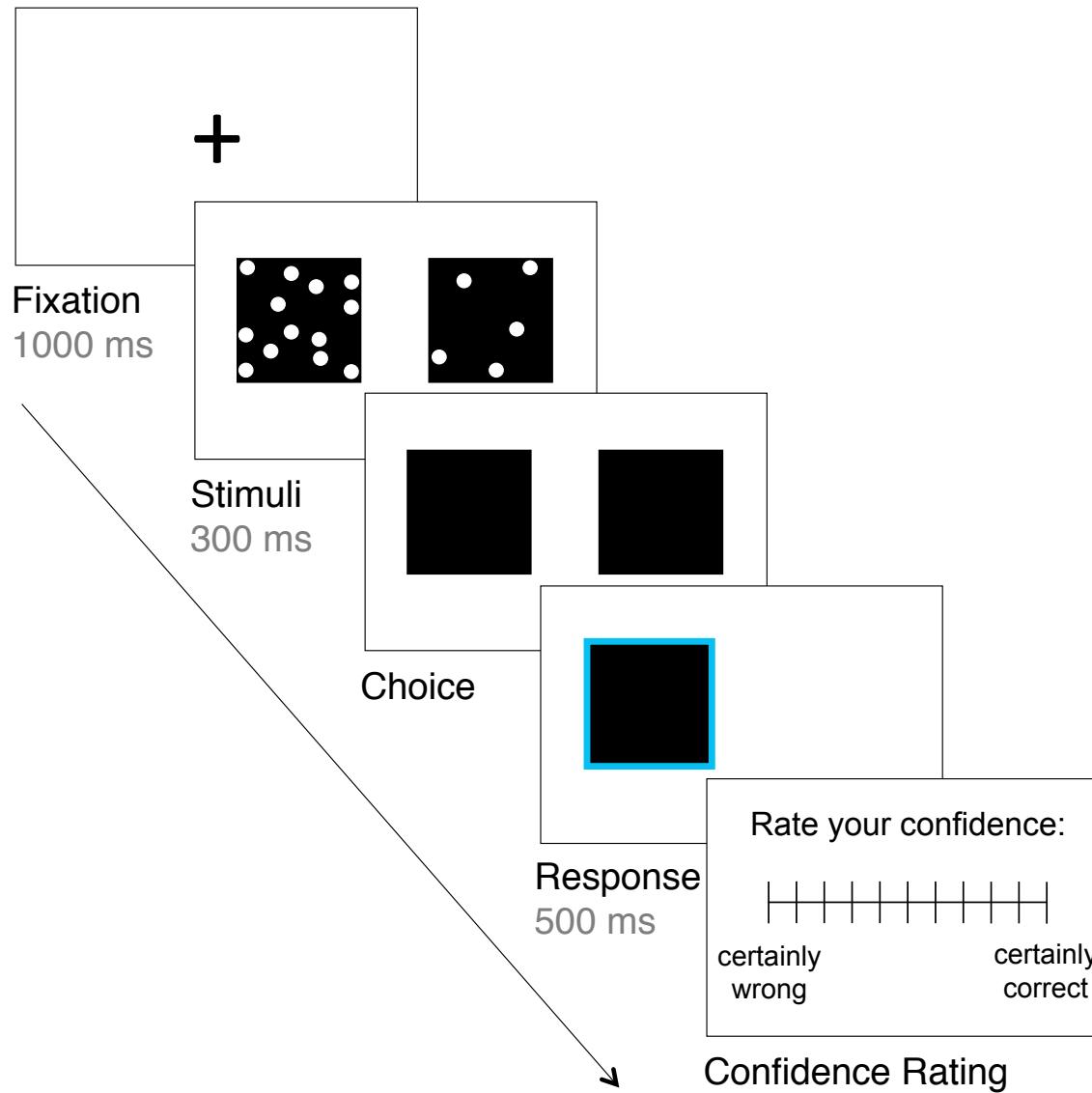
Is metacognition related to psychopathology?



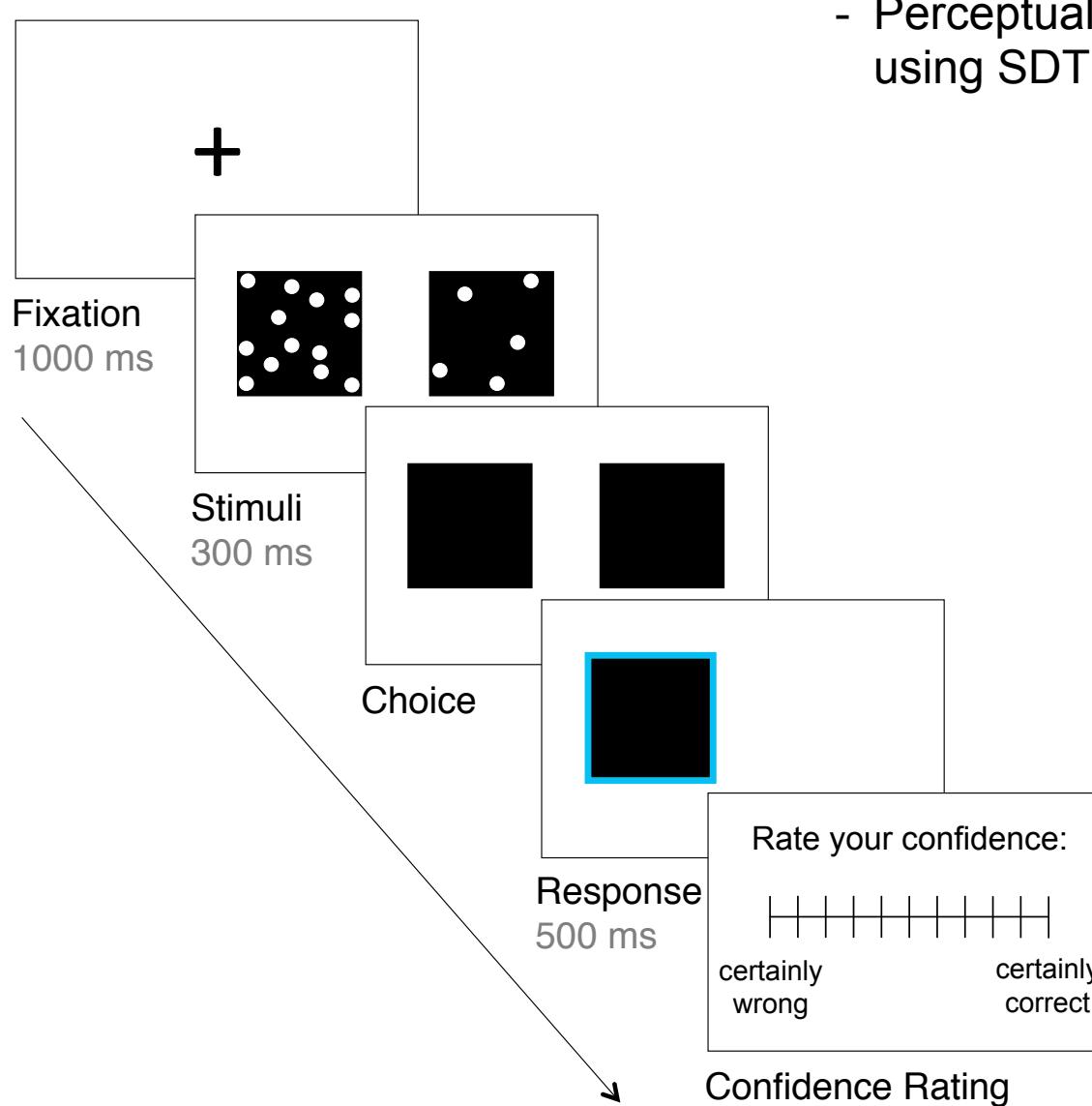
Experiment 1 - N=498 participants
Experiment 2 - N=497 participants



Decision-making + confidence task

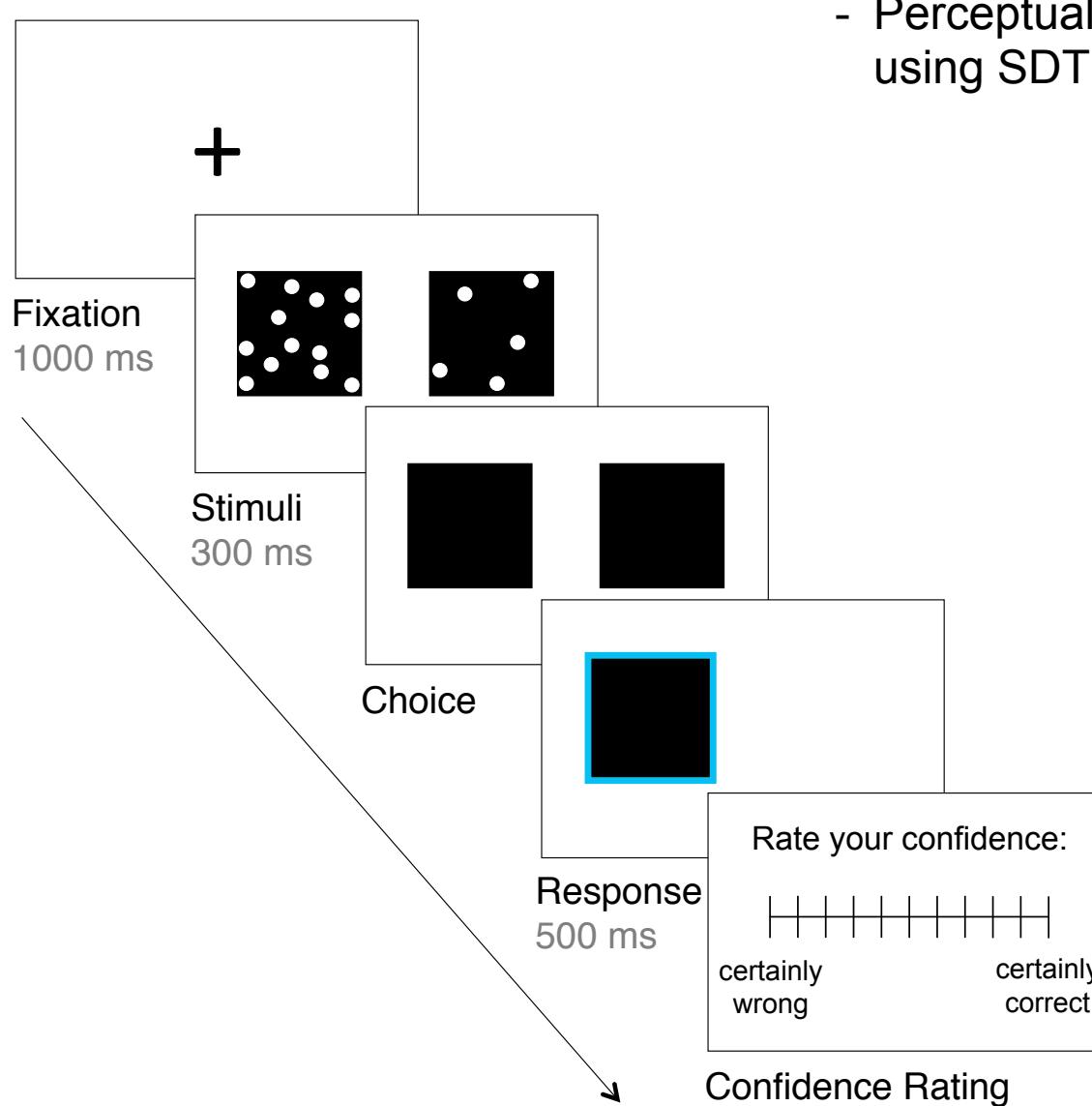


Decision-making + confidence task



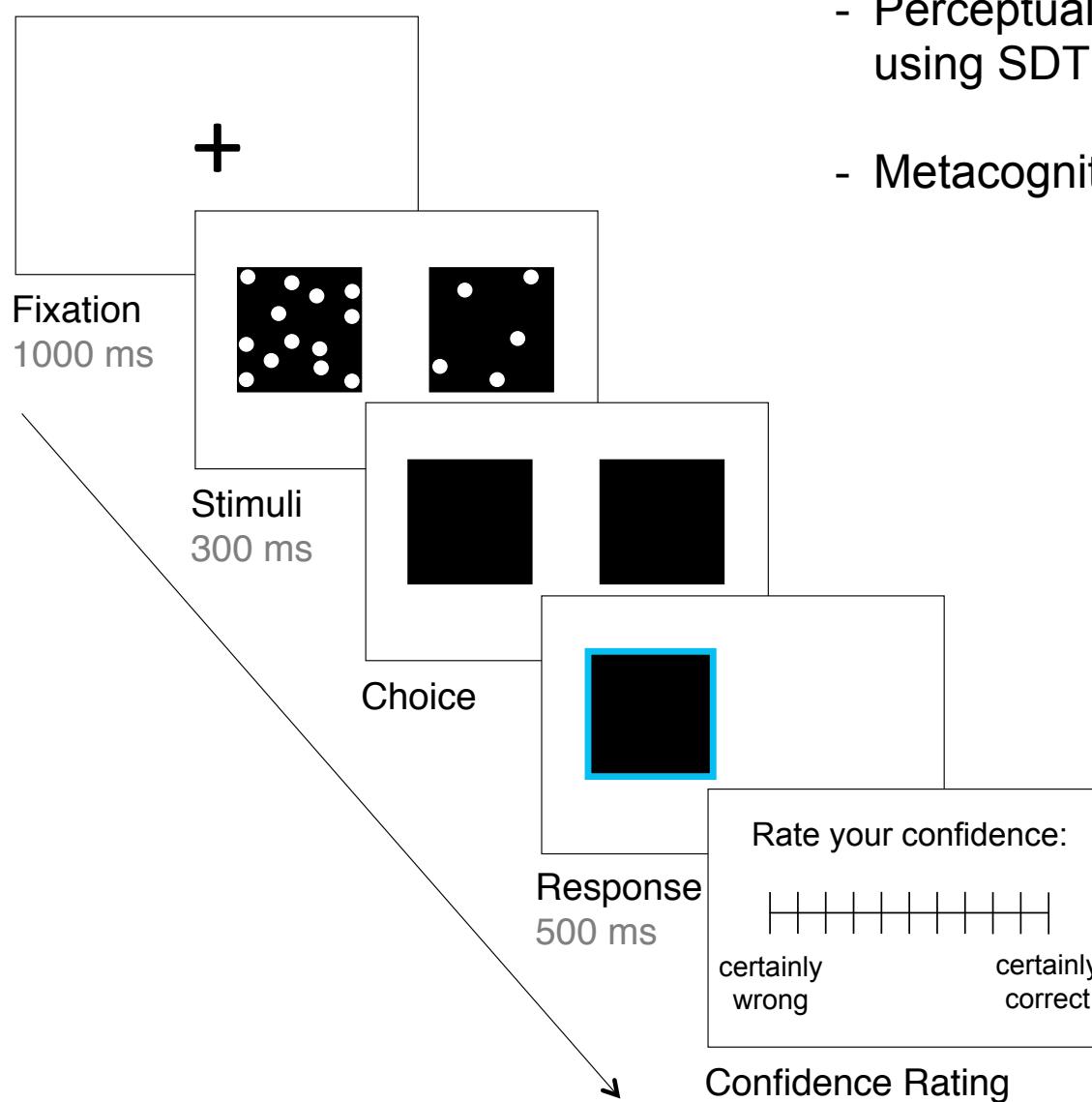
- Perceptual decision-making quantified using SDT and drift-diffusion modelling

Decision-making + confidence task



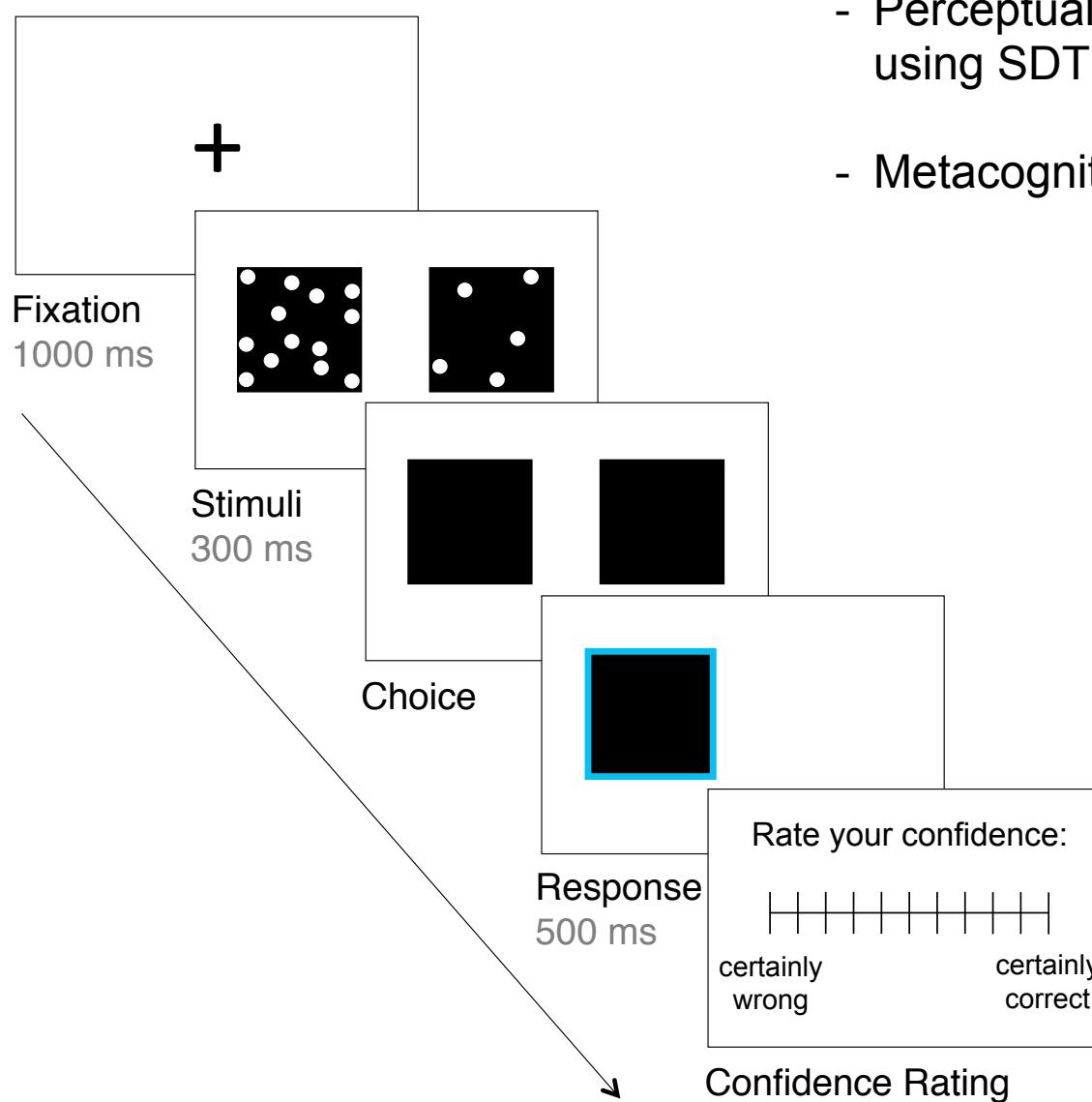
- Perceptual decision-making quantified using SDT and drift-diffusion modelling

Decision-making + confidence task



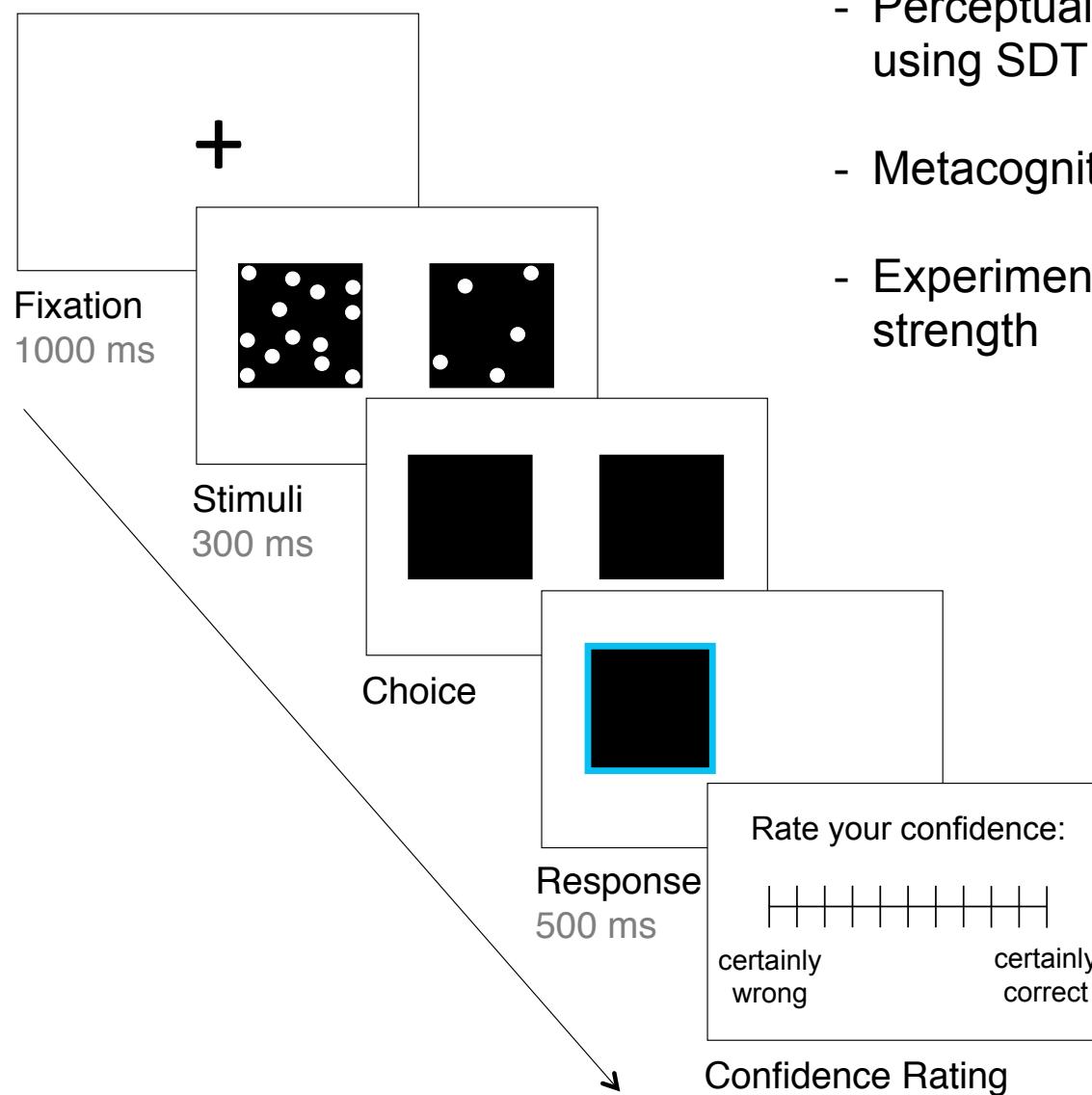
- Perceptual decision-making quantified using SDT and drift-diffusion modelling
- Metacognition quantified using meta-d'

Decision-making + confidence task



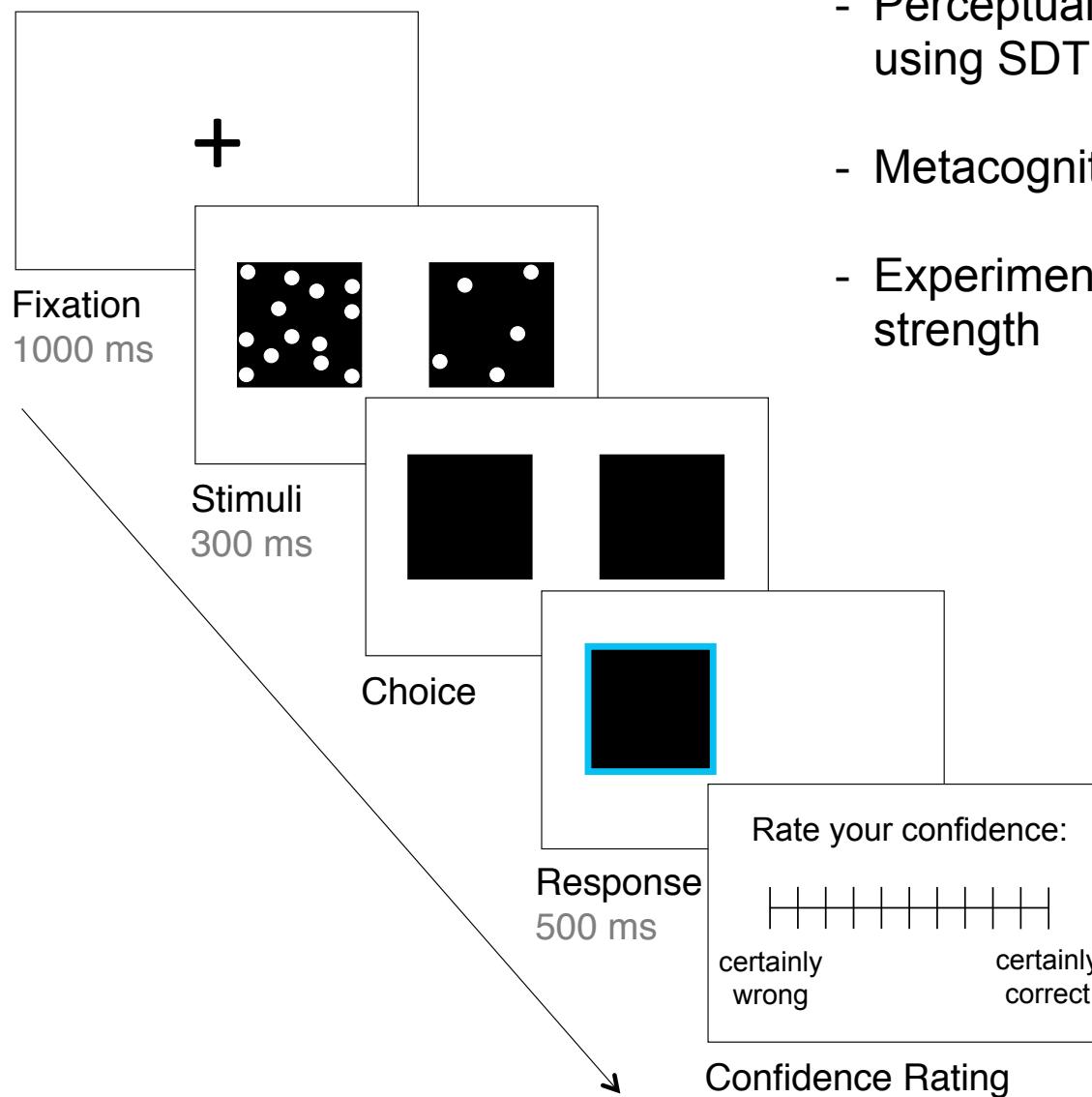
- Perceptual decision-making quantified using SDT and drift-diffusion modelling
- Metacognition quantified using meta-d'

Decision-making + confidence task



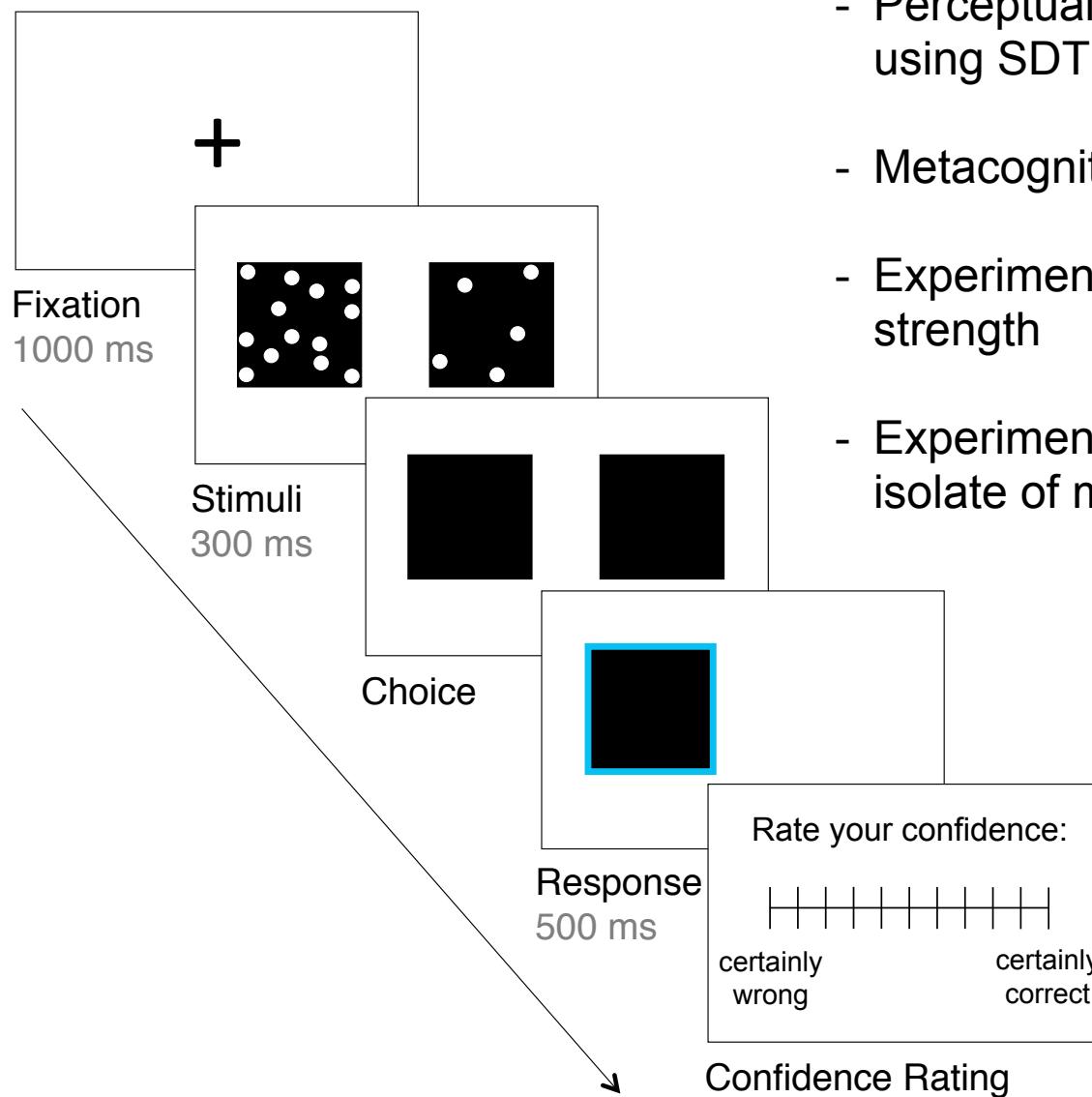
- Perceptual decision-making quantified using SDT and drift-diffusion modelling
- Metacognition quantified using meta-d'
- Experiment 1 = variable stimulus strength

Decision-making + confidence task



- Perceptual decision-making quantified using SDT and drift-diffusion modelling
- Metacognition quantified using meta-d'
- Experiment 1 = variable stimulus strength

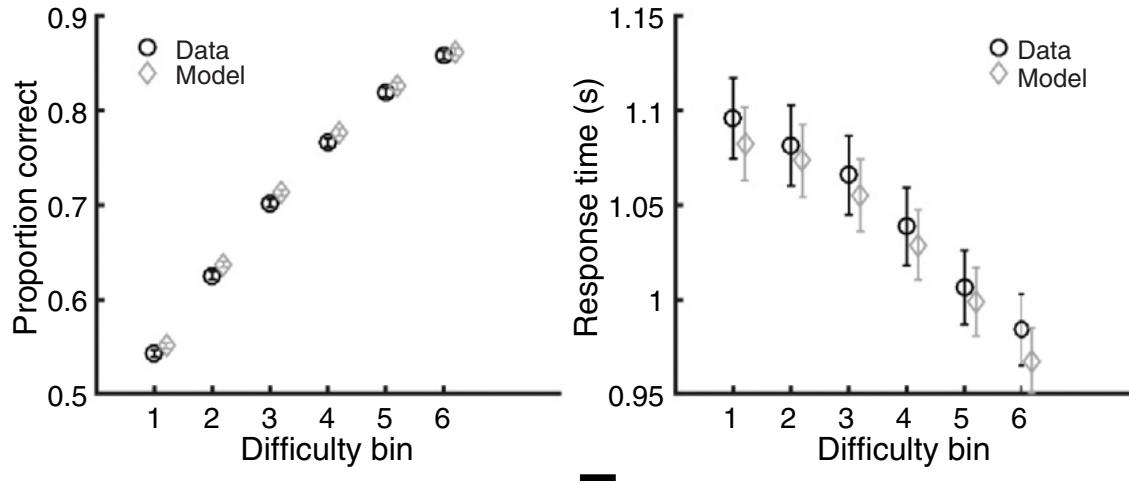
Decision-making + confidence task



- Perceptual decision-making quantified using SDT and drift-diffusion modelling
- Metacognition quantified using meta-d'
- Experiment 1 = variable stimulus strength
- Experiment 2 = staircase used to isolate of metacognitive variability

Decision-making + confidence task

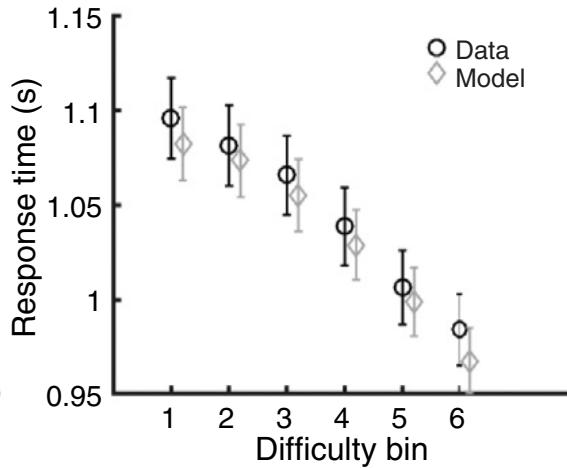
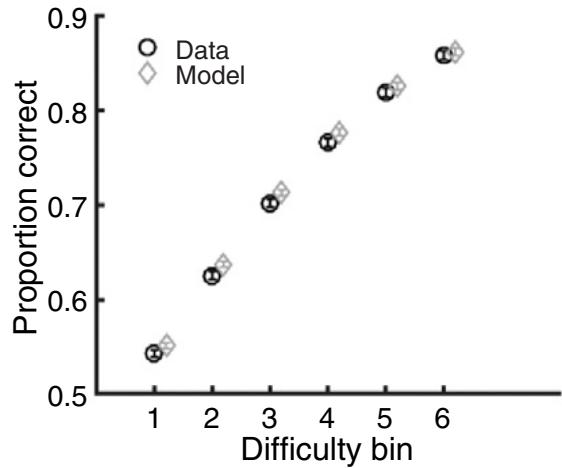
Perceptual decision performance



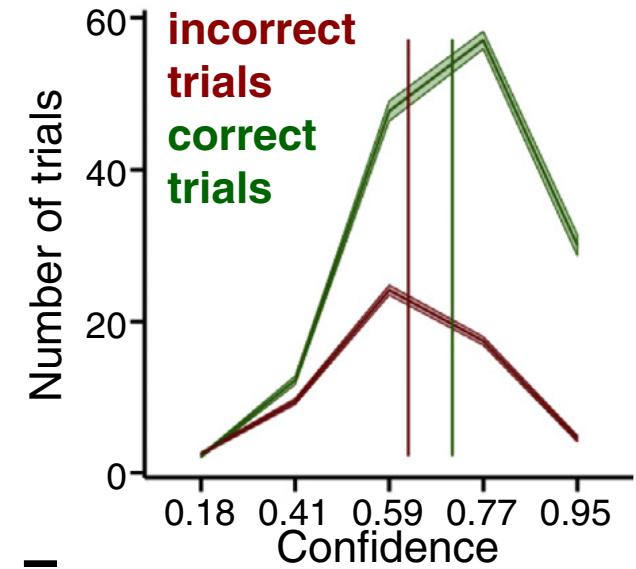
Experiment 1, N=498

Decision-making + confidence task

Perceptual decision performance

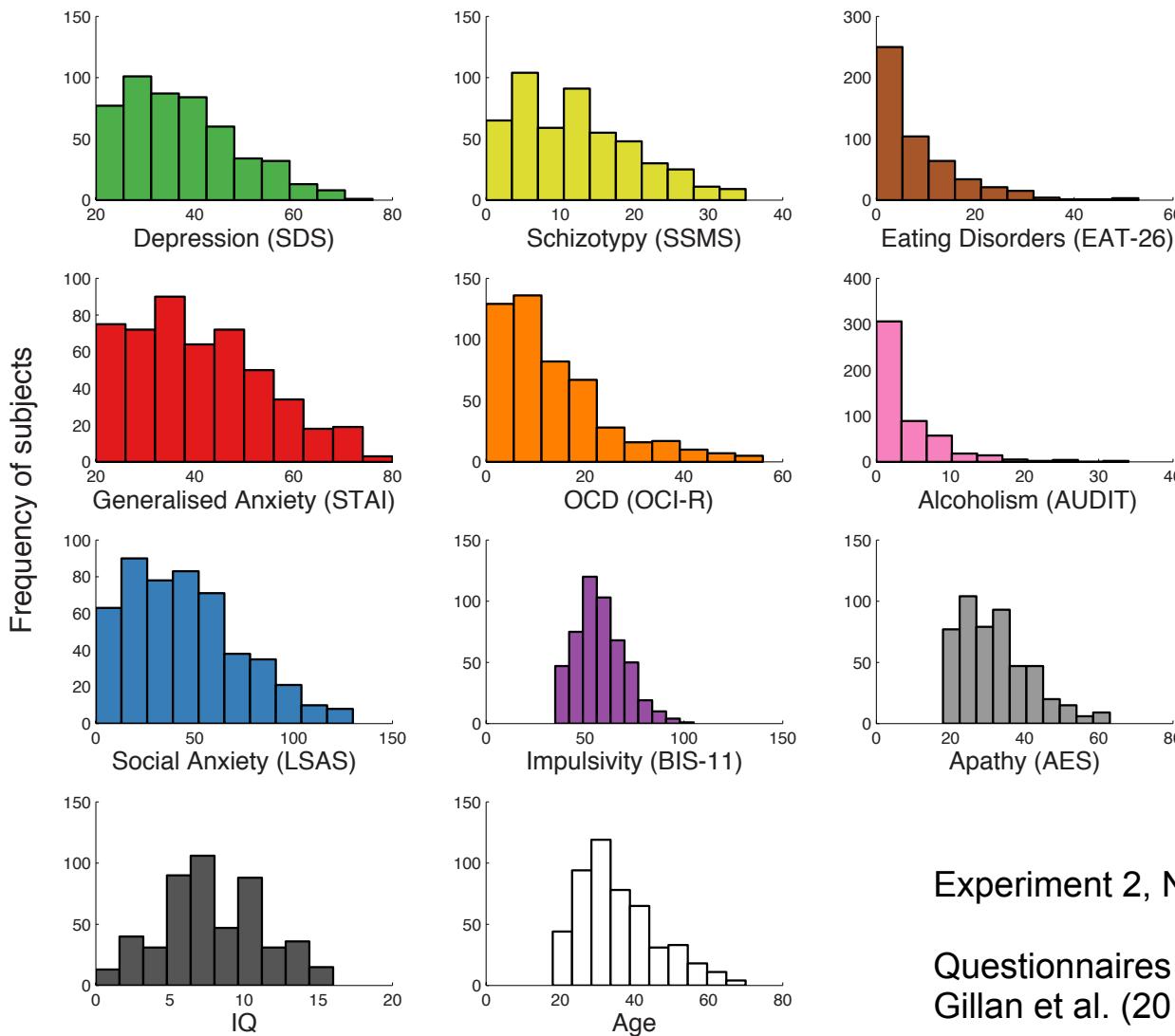


Confidence ratings



Experiment 1, N=498

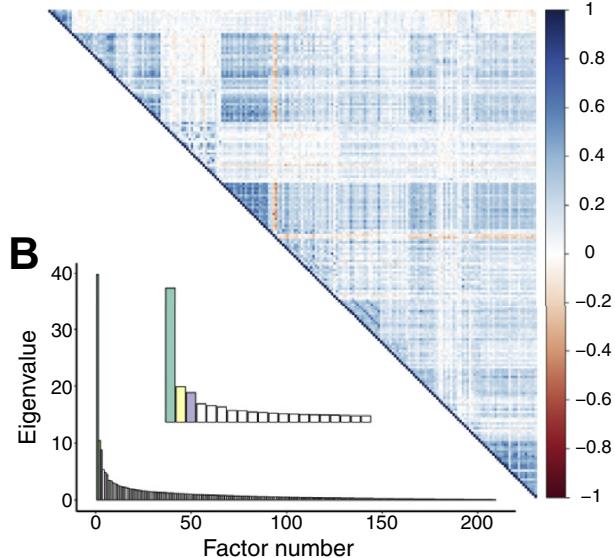
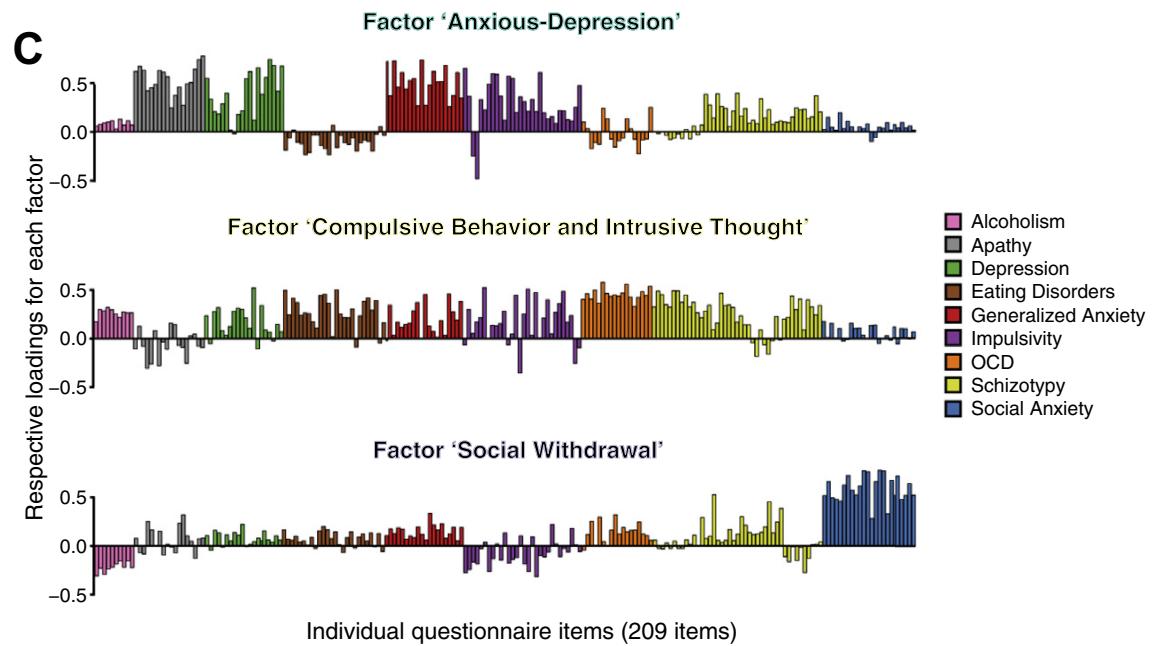
Self-reported symptom questionnaires



Experiment 2, N=497

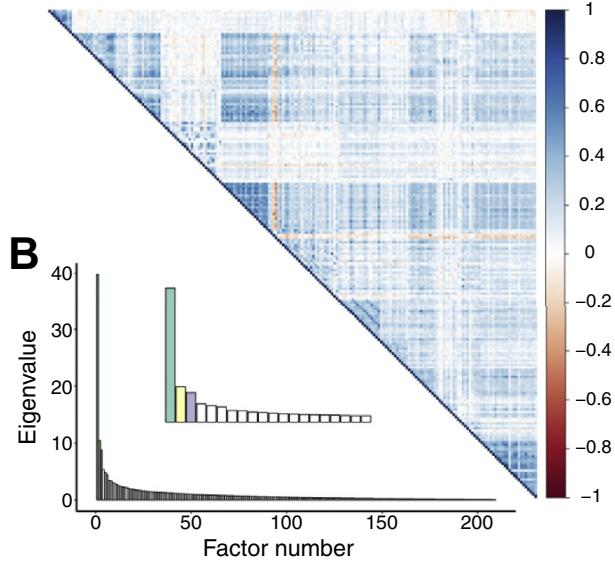
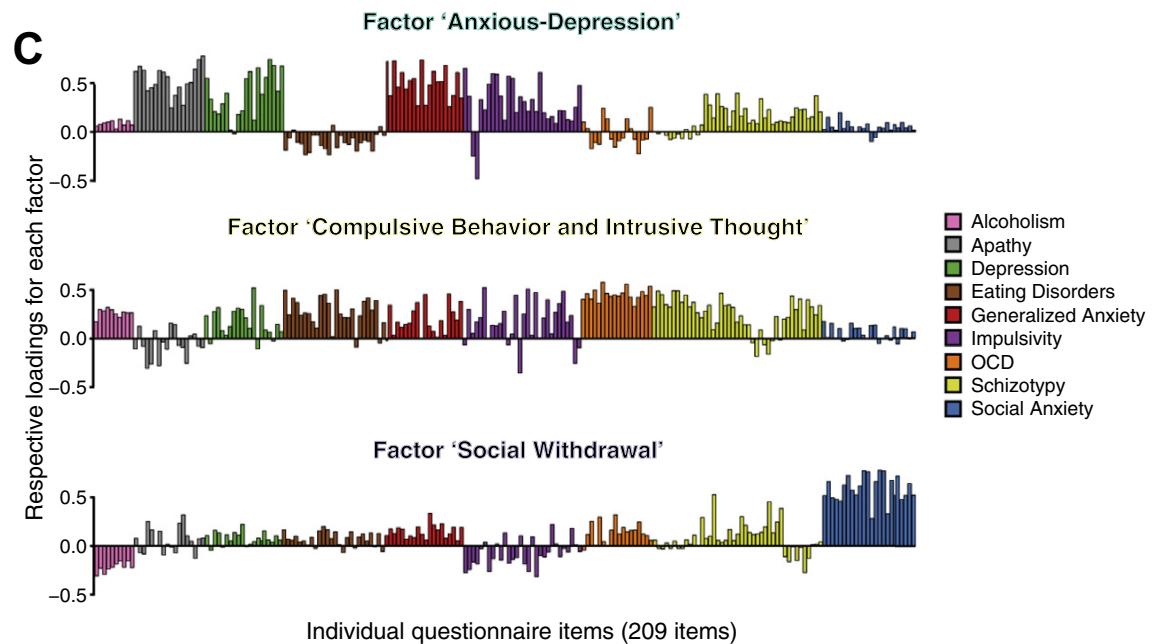
Questionnaires based on
Gillan et al. (2016) eLife

Identifying latent transdiagnostic dimensions

A**C**

See also Gillan et al. (2016) eLife

Identifying latent transdiagnostic dimensions

A**C**

- Alcoholism
- Apathy
- Depression
- Eating Disorders
- Generalized Anxiety
- Impulsivity
- OCD
- Schizotypy
- Social Anxiety

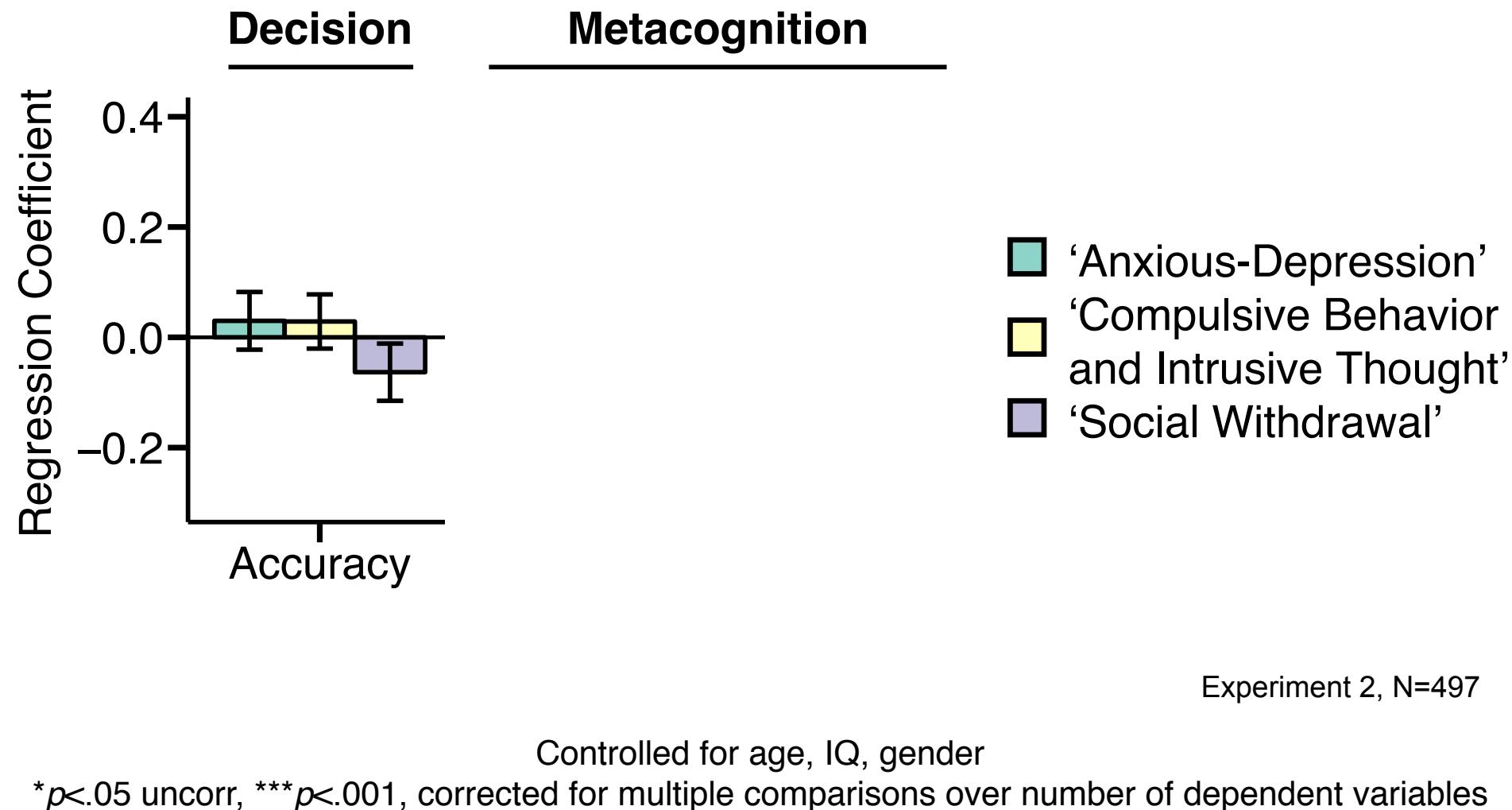
Anxious/
Depression

Compulsive
Behavior and
Intrusive
Thought

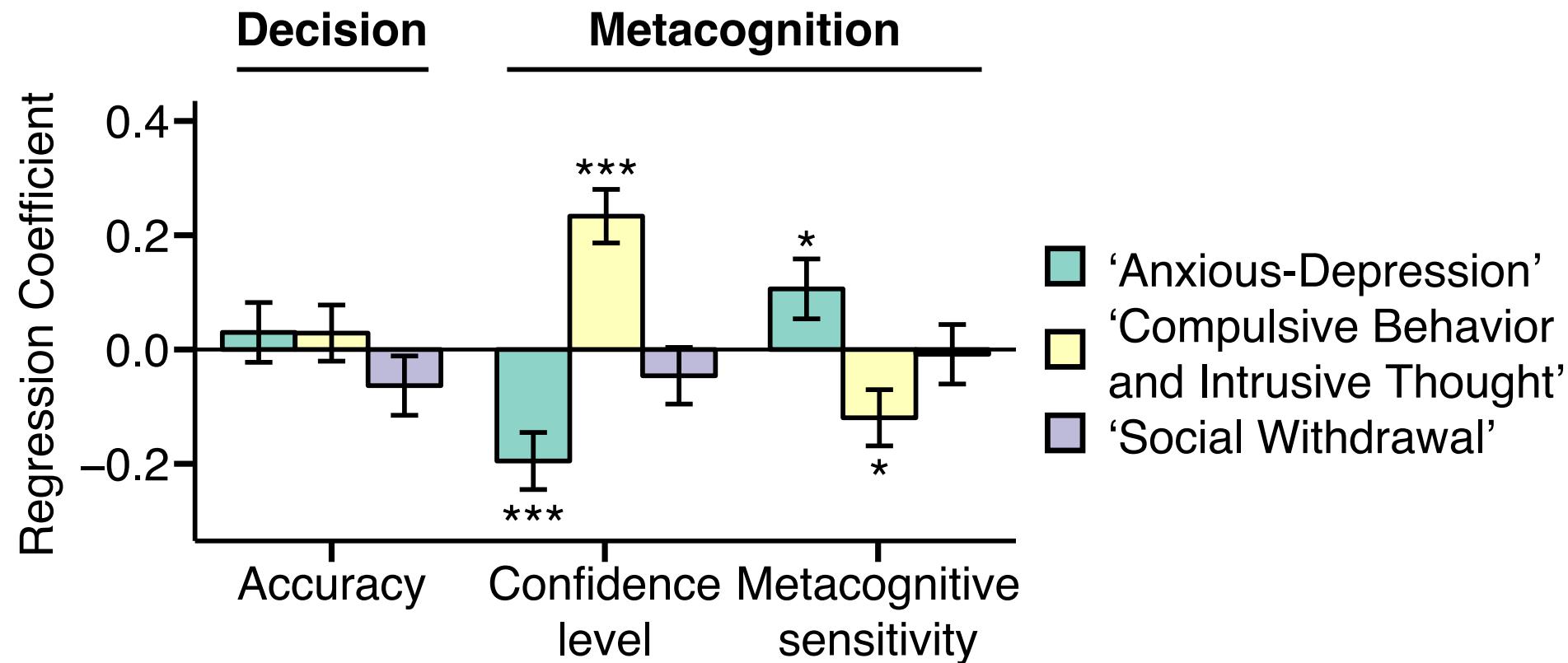
Social
Withdrawal

See also Gillan et al. (2016) eLife

Metacognition (but not decision performance) is associated with latent symptom dimensions



Metacognition (but not decision performance) is associated with latent symptom dimensions

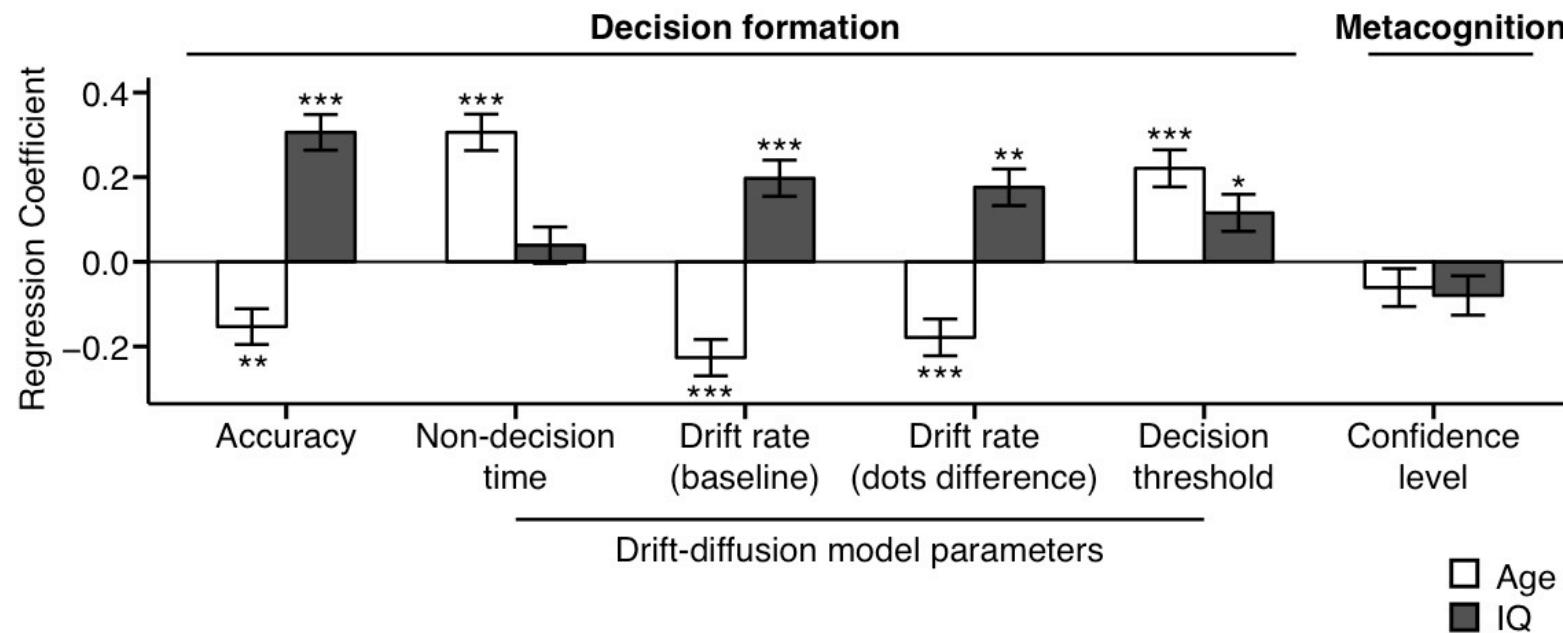


Experiment 2, N=497

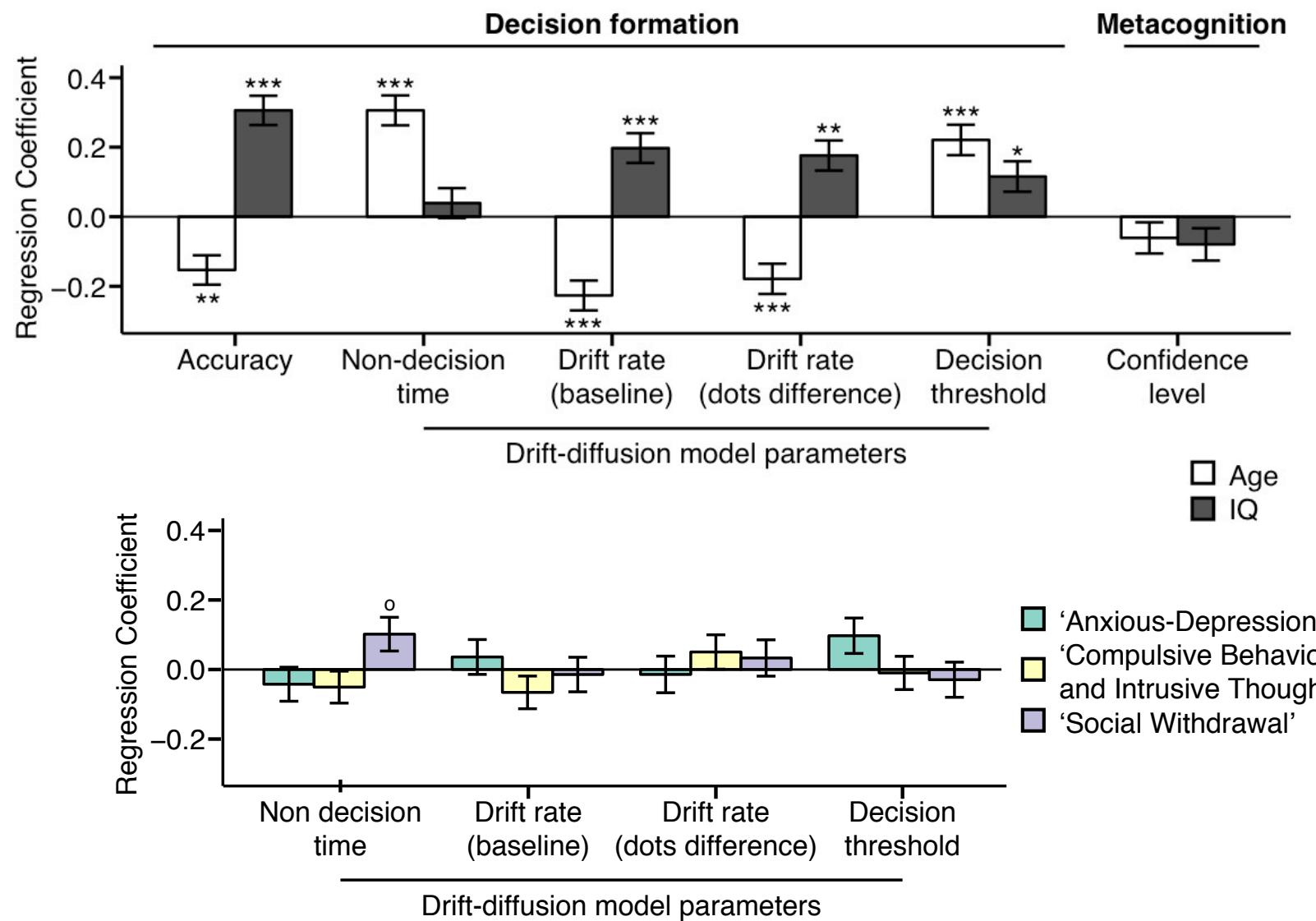
Controlled for age, IQ, gender

* $p < .05$ uncorr, *** $p < .001$, corrected for multiple comparisons over number of dependent variables

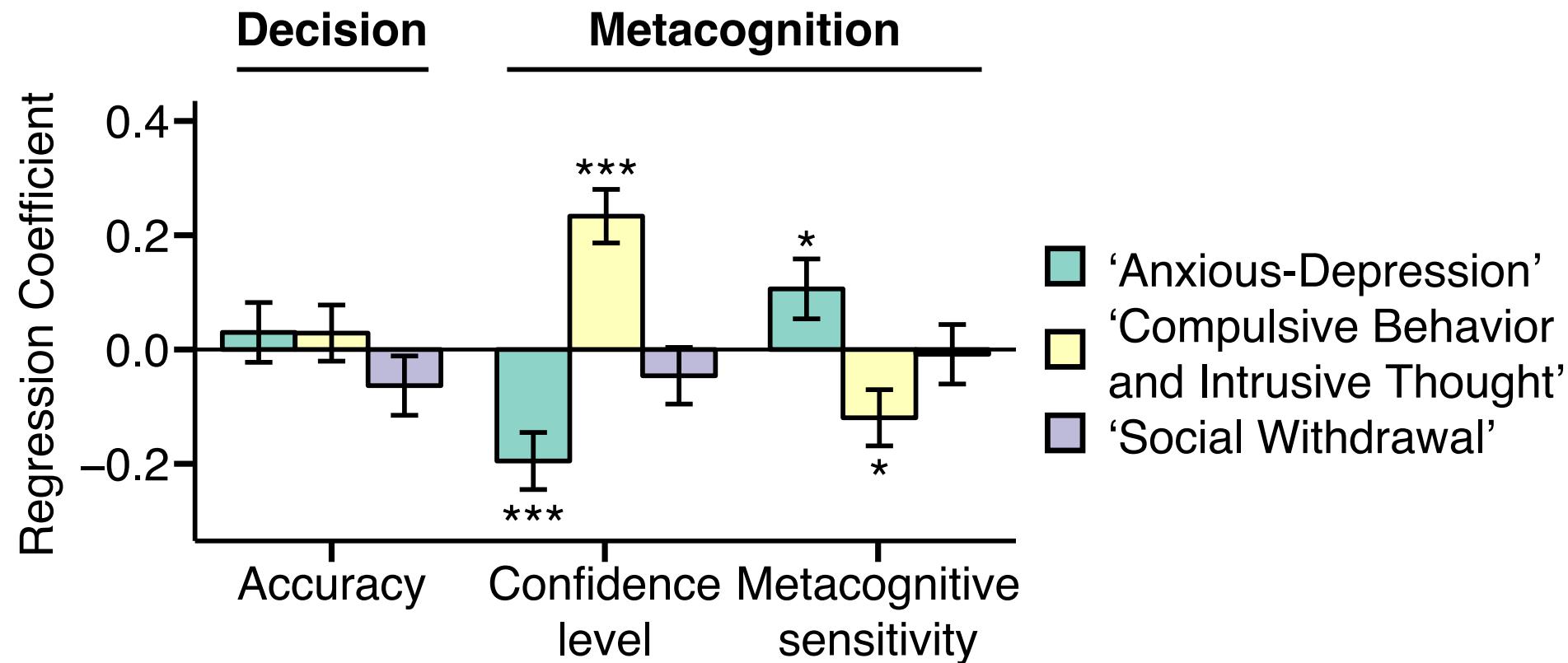
Dissociating metacognition and decision performance



Dissociating metacognition and decision performance



Metacognition (but not decision performance) is associated with latent symptom dimensions



Experiment 2, N=497

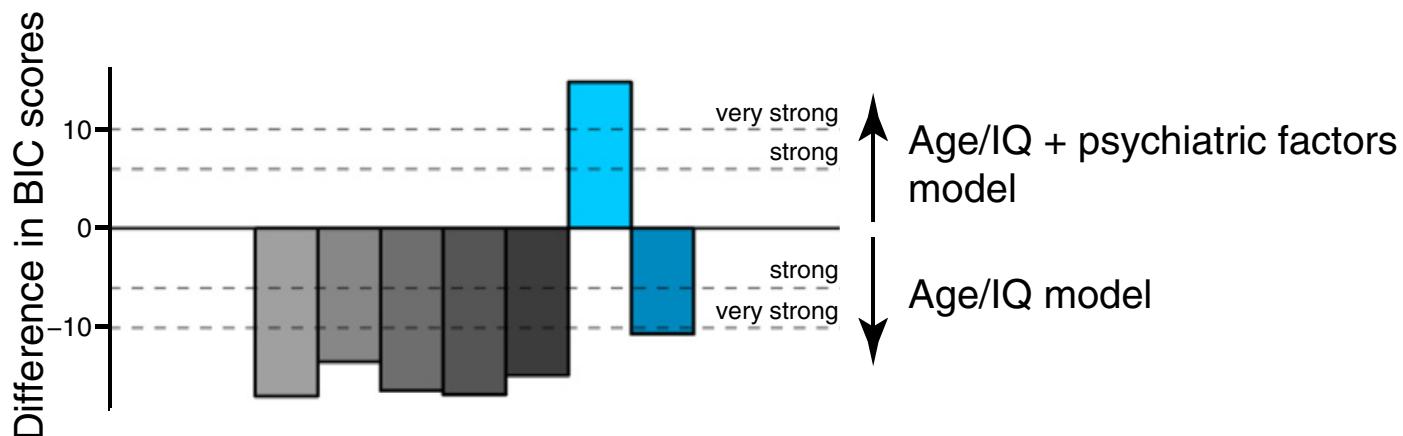
Controlled for age, IQ, gender

* $p < .05$ uncorr, *** $p < .001$, corrected for multiple comparisons over number of dependent variables

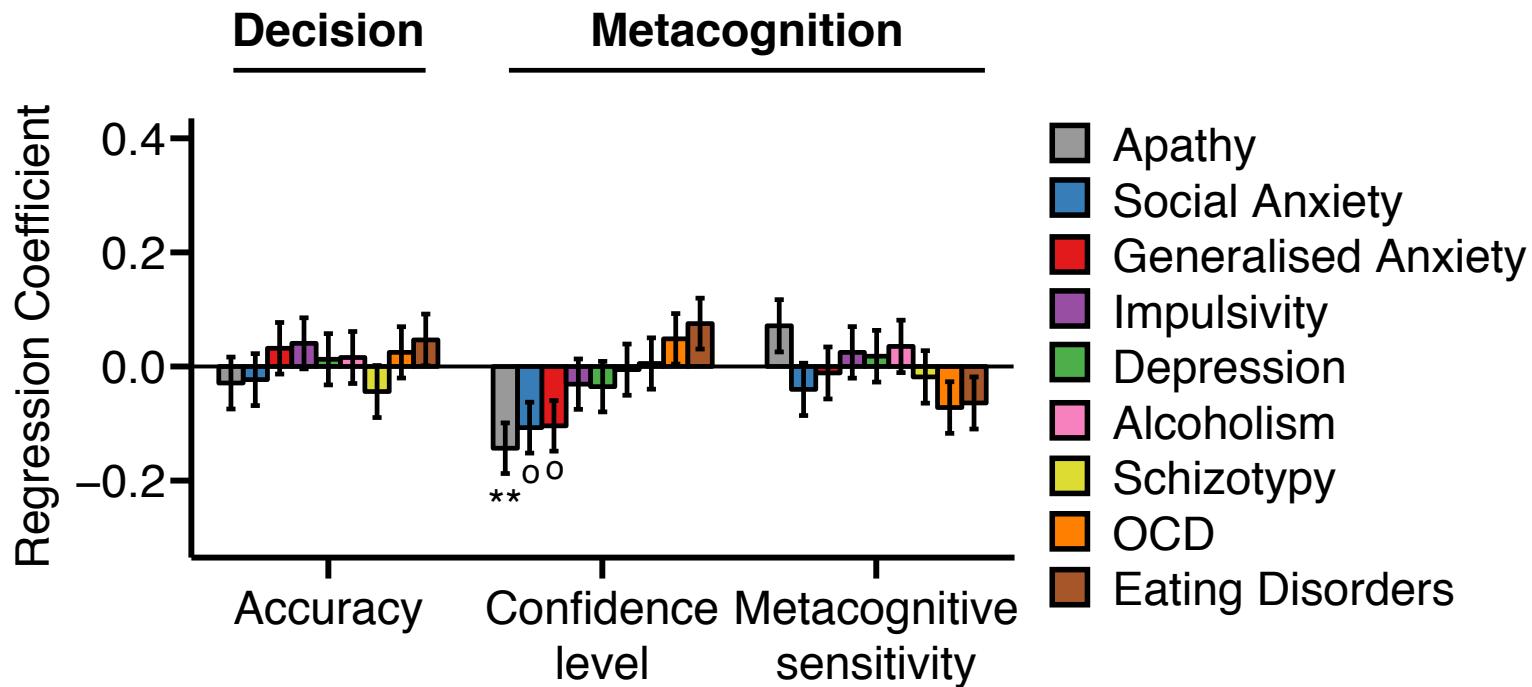
Metacognition (but not decision performance) is associated with latent symptom dimensions

Variables list:

Accuracy
Non-decision time
Drift rate (baseline)
Drift rate (dots difference)
Decision threshold
Confidence level
Metacognitive efficiency



Latent symptom dimensions are stronger predictors than individual scores



Controlled for age, IQ, gender

$^{\circ}p < .05$ uncorr, $^{**}p < .01$, corrected for multiple comparisons over number of dependent variables

Interim summary

Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks

Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks

Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks
- Psychiatric symptom dimensions are associated with distinct aspects of metacognitive ability, despite generally unaffected decision-making (at least in this task)

Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks
- Psychiatric symptom dimensions are associated with distinct aspects of metacognitive ability, despite generally unaffected decision-making (at least in this task)

Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks
- Psychiatric symptom dimensions are associated with distinct aspects of metacognitive ability, despite generally unaffected decision-making (at least in this task)
- What **brain systems** might support domain-general metacognition?

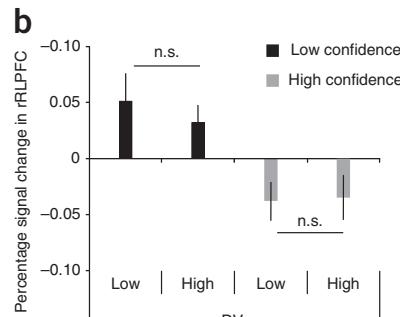
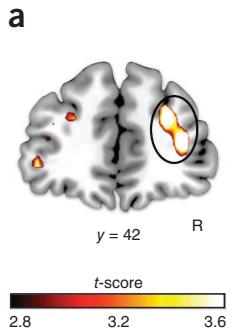
Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks
- Psychiatric symptom dimensions are associated with distinct aspects of metacognitive ability, despite generally unaffected decision-making (at least in this task)
- What **brain systems** might support domain-general metacognition?

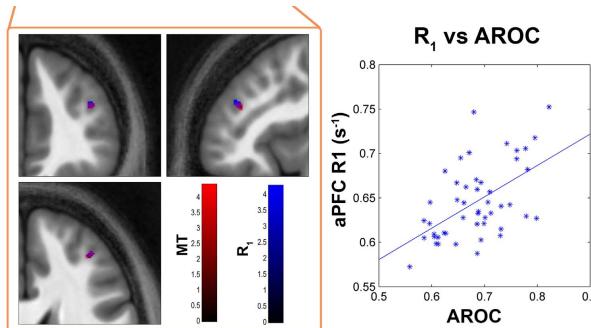
Interim summary

- Behavioural evidence indicates a domain-general resource supporting metacognition in simple tasks
- Psychiatric symptom dimensions are associated with distinct aspects of metacognitive ability, despite generally unaffected decision-making (at least in this task)
- What **brain systems** might support domain-general metacognition?
- Can we develop **neural/behavioural interventions** for modulating metacognition?

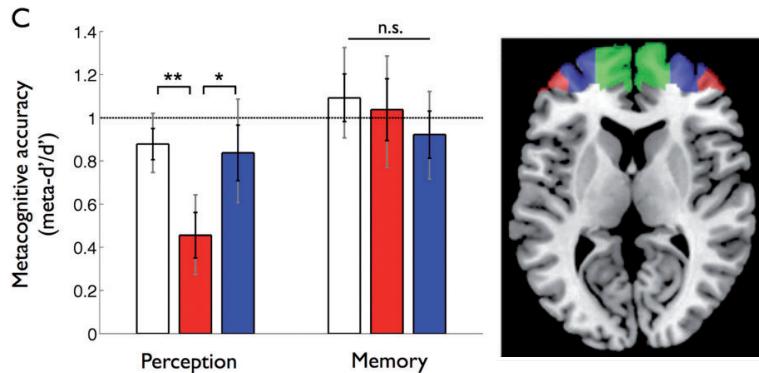
Metacognitive evaluation and aPFC



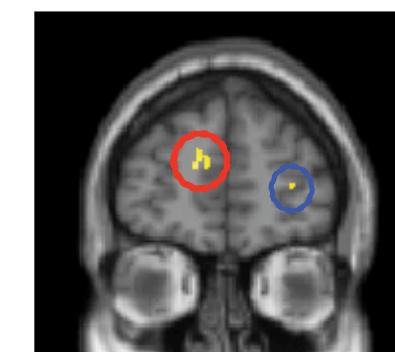
De Martino, Fleming et al. 2013 *Nat Neuro*



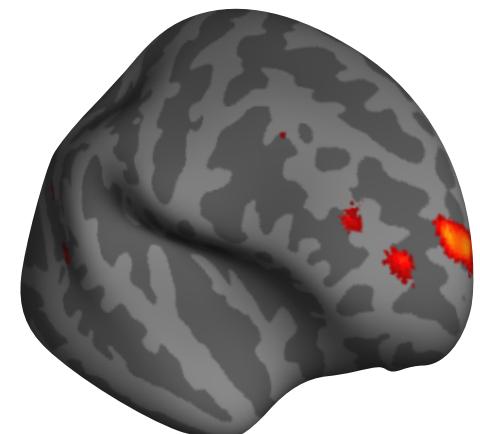
Allen et al. 2017 *Neuroimage*



Fleming et al. 2014 *Brain*



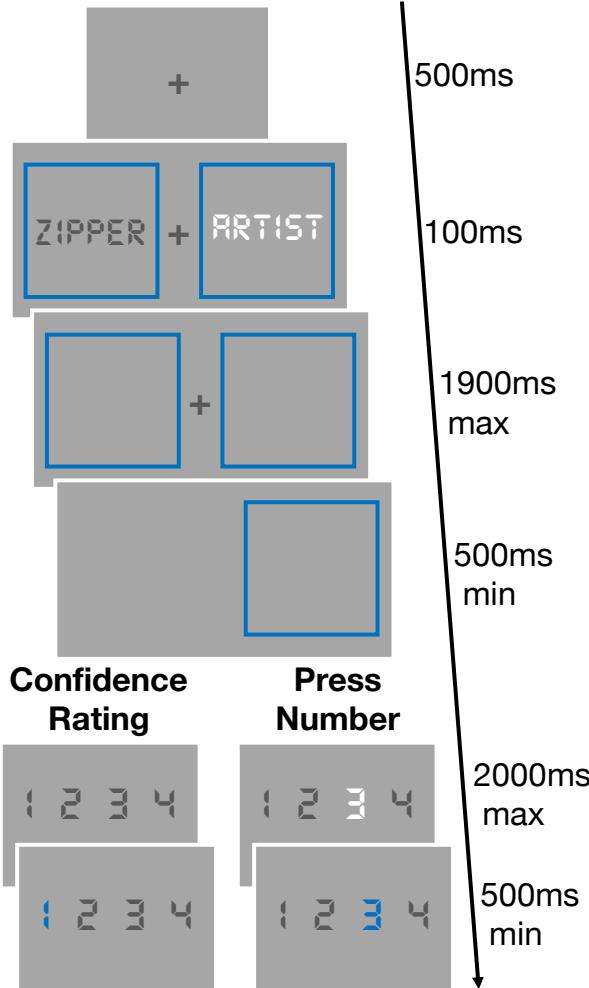
McCurdy et al. 2013 *J Neuro*;
Fleming et al. 2010 *Science*



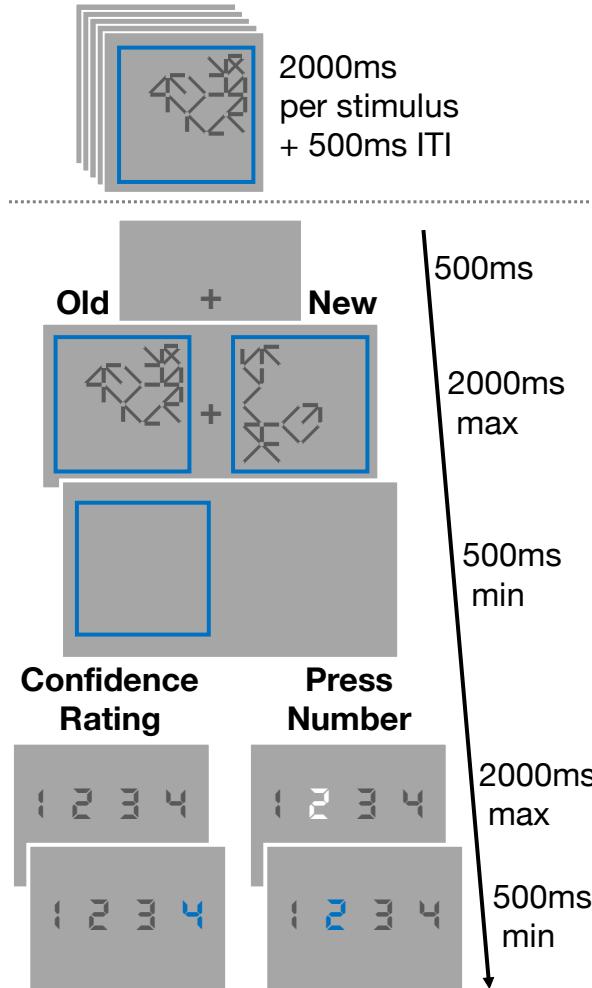
anterior prefrontal cortex (aPFC)

Shared signals for confidence across tasks?

B Perceptual trial



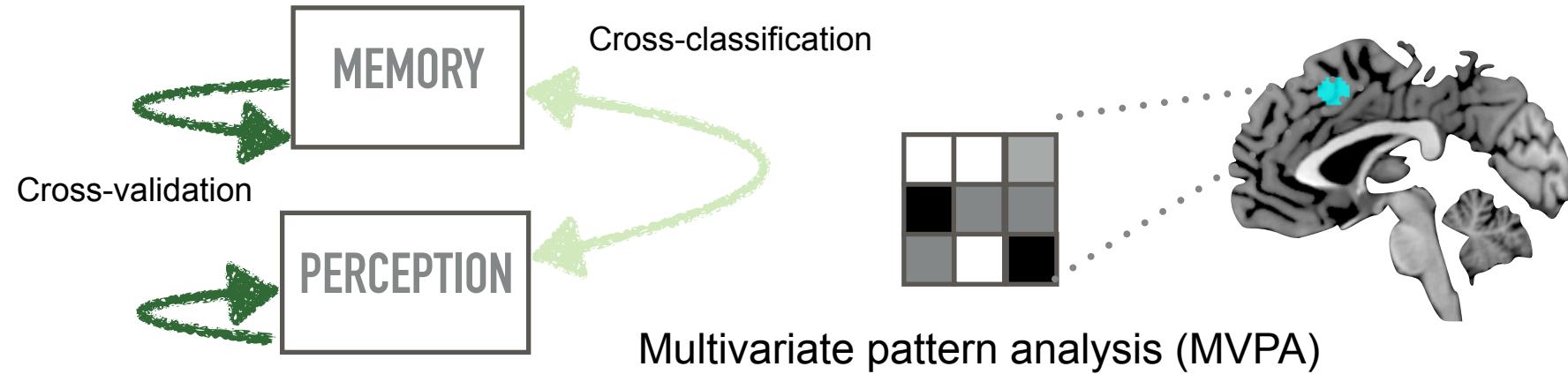
C Memory trial



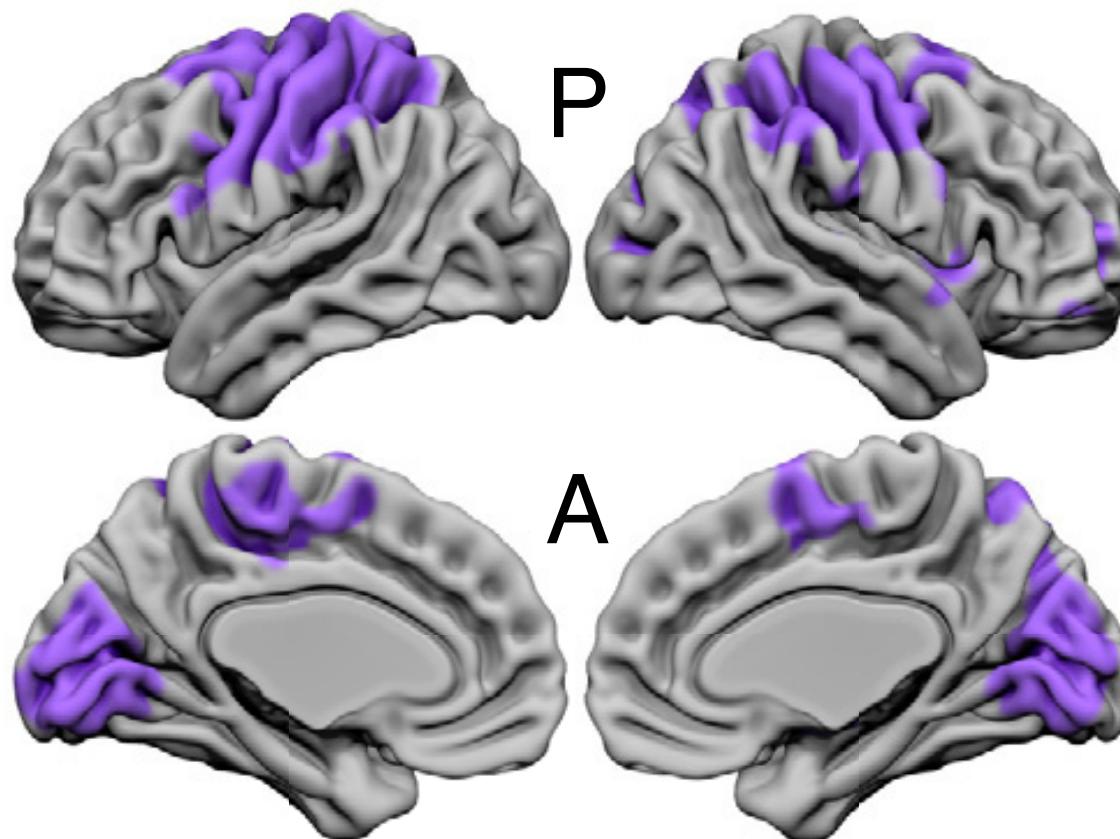
2 (stimulus) x
2 (domain)
block design



Shared signals for confidence across tasks?



Excluding button-press related activity patterns



Cross-classification

Confidence Rating Press Number

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

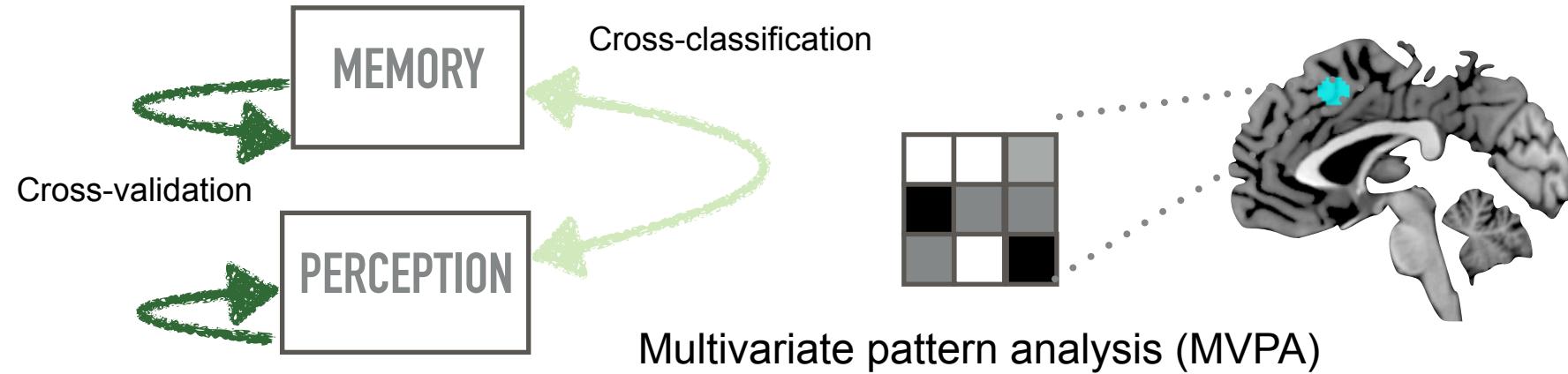
1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

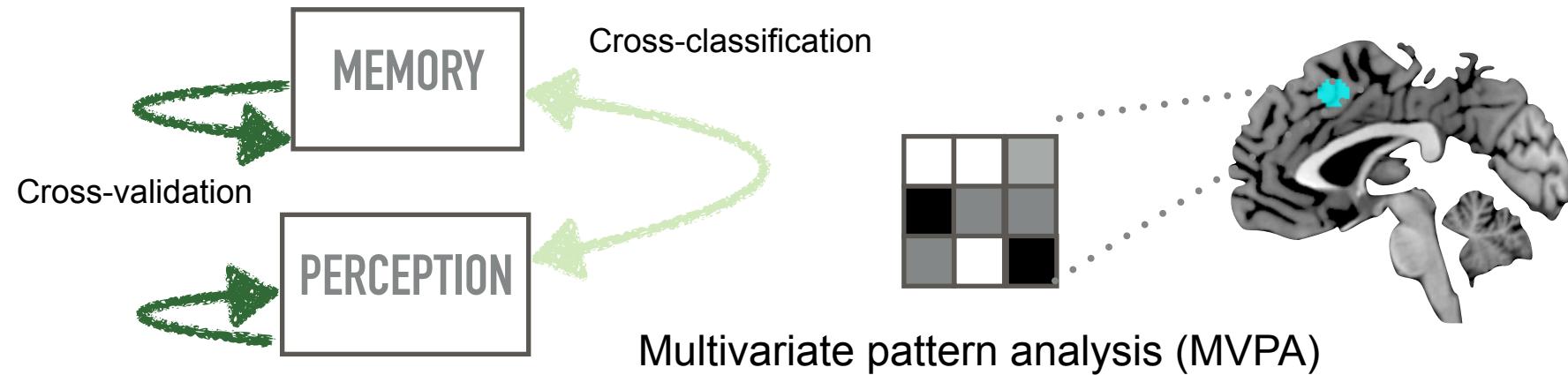
1	2	3	4
---	---	---	---

1	2	3	4
---	---	---	---

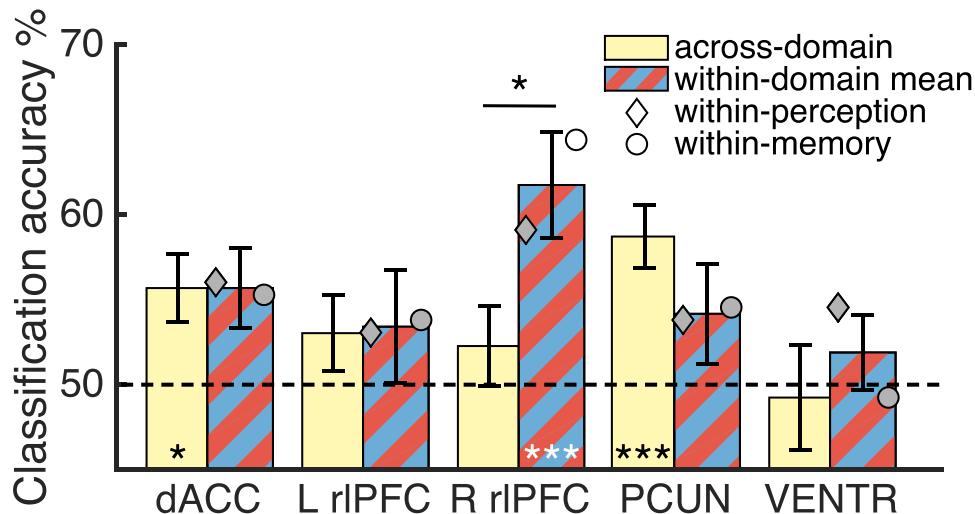
Shared signals for confidence across tasks?



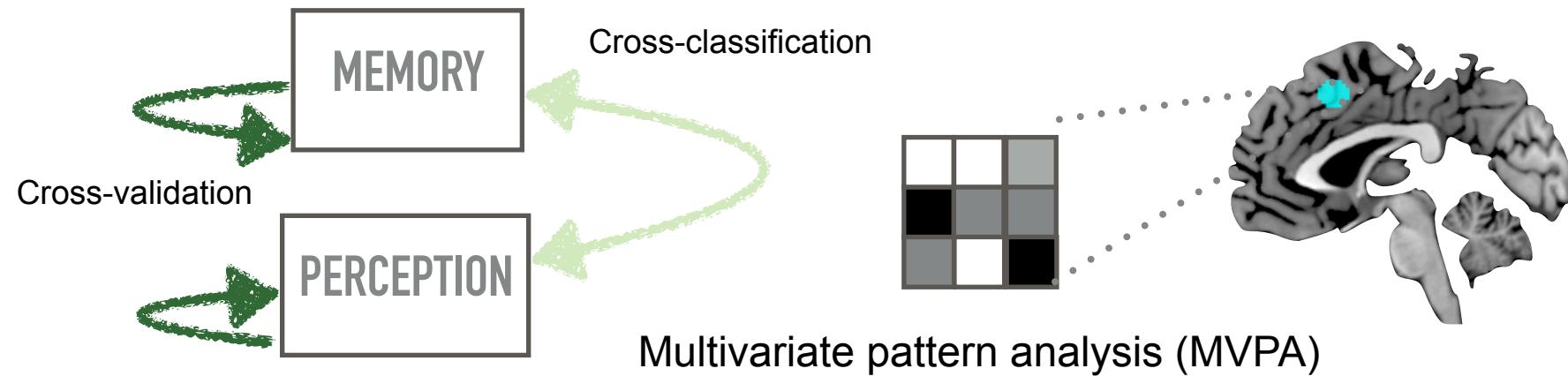
Shared signals for confidence across tasks?



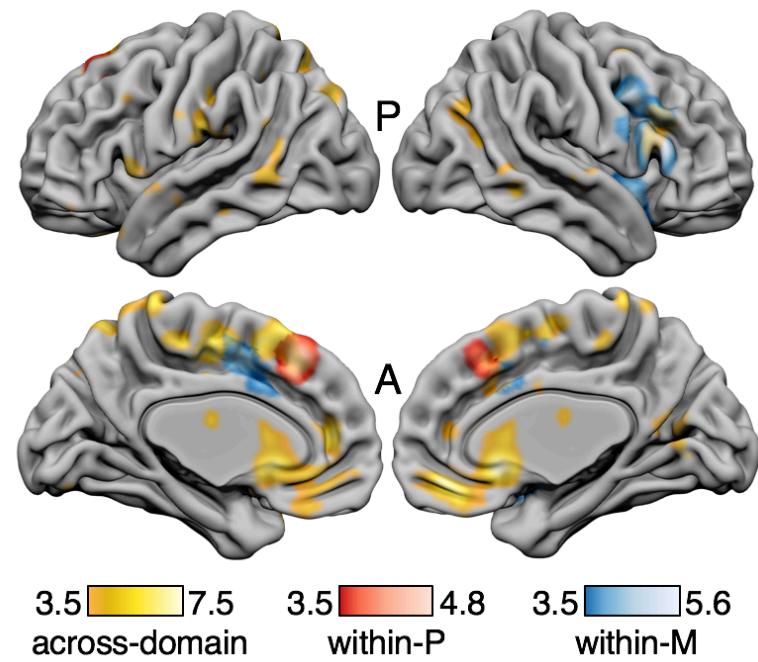
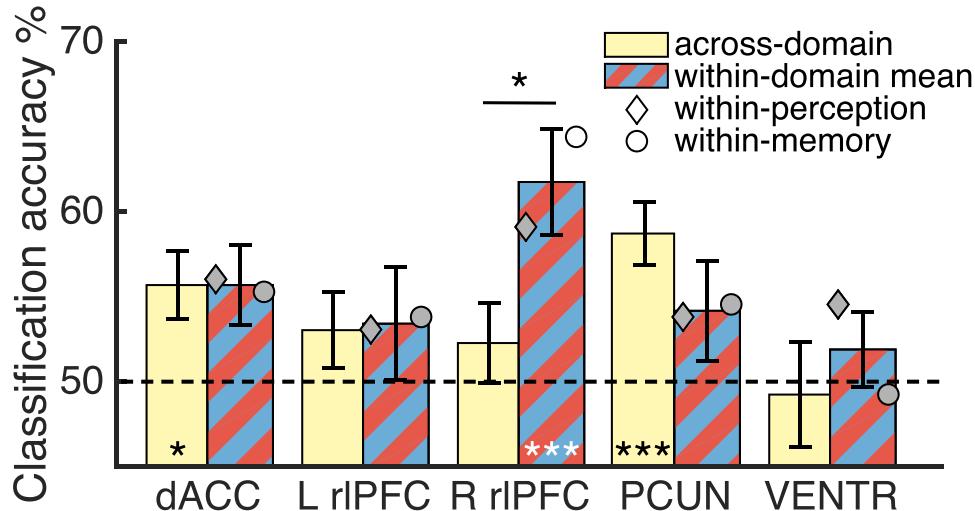
Confidence-related activity patterns



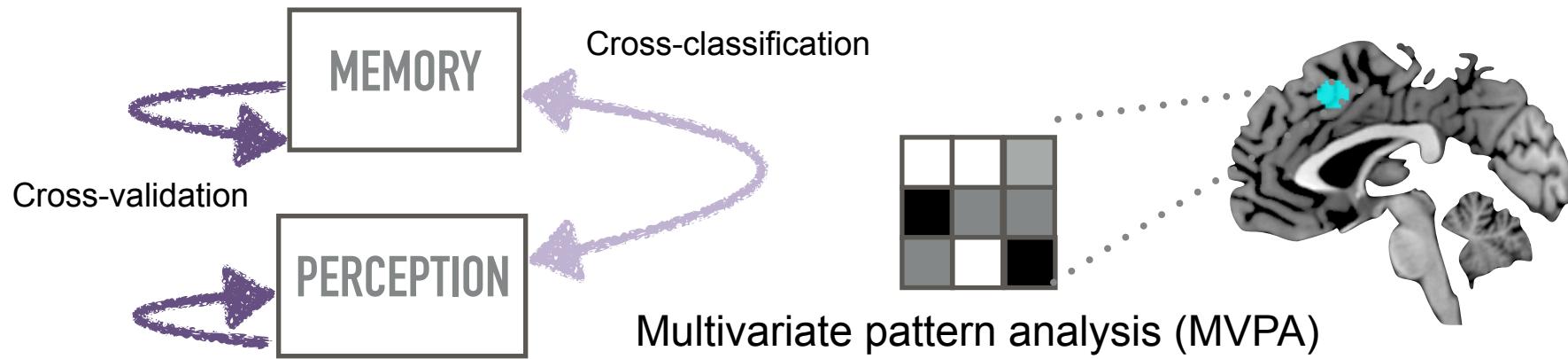
Shared signals for confidence across tasks?



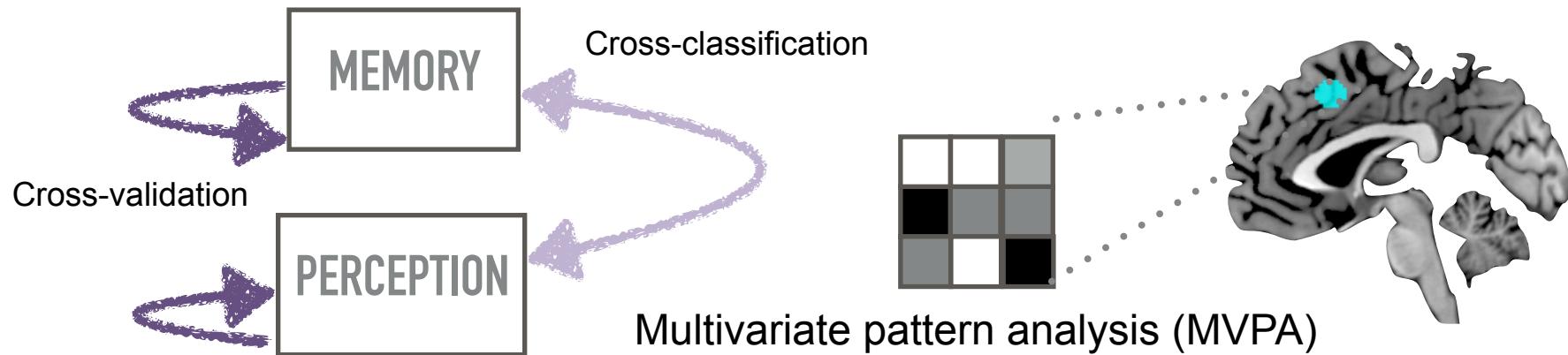
Confidence-related activity patterns



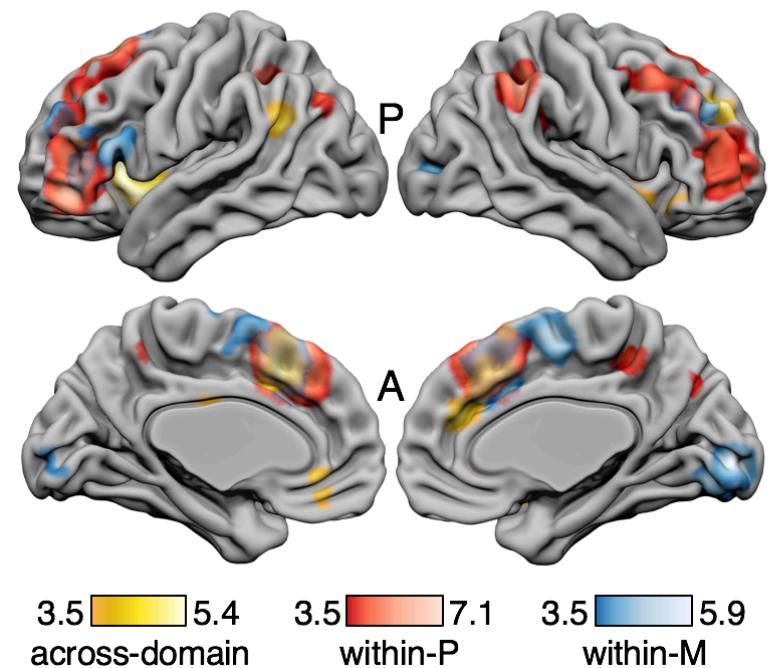
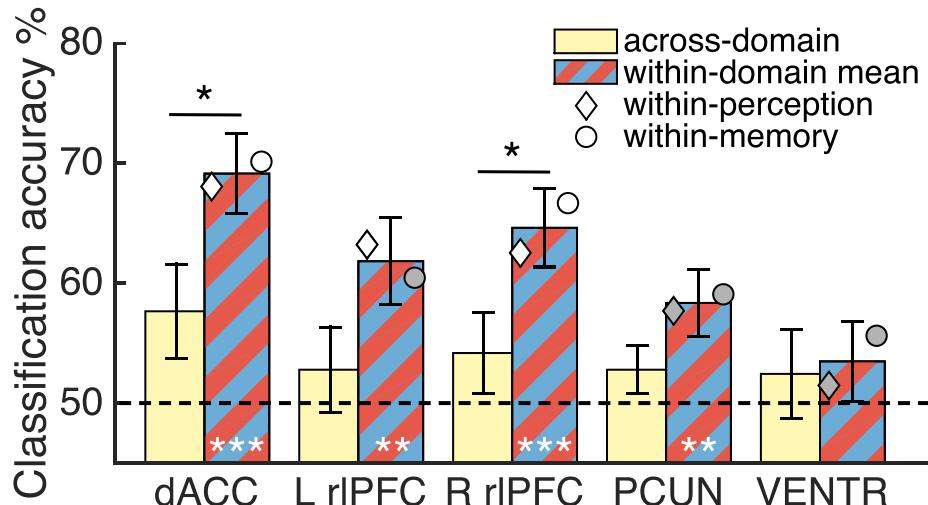
Judgment-related patterns



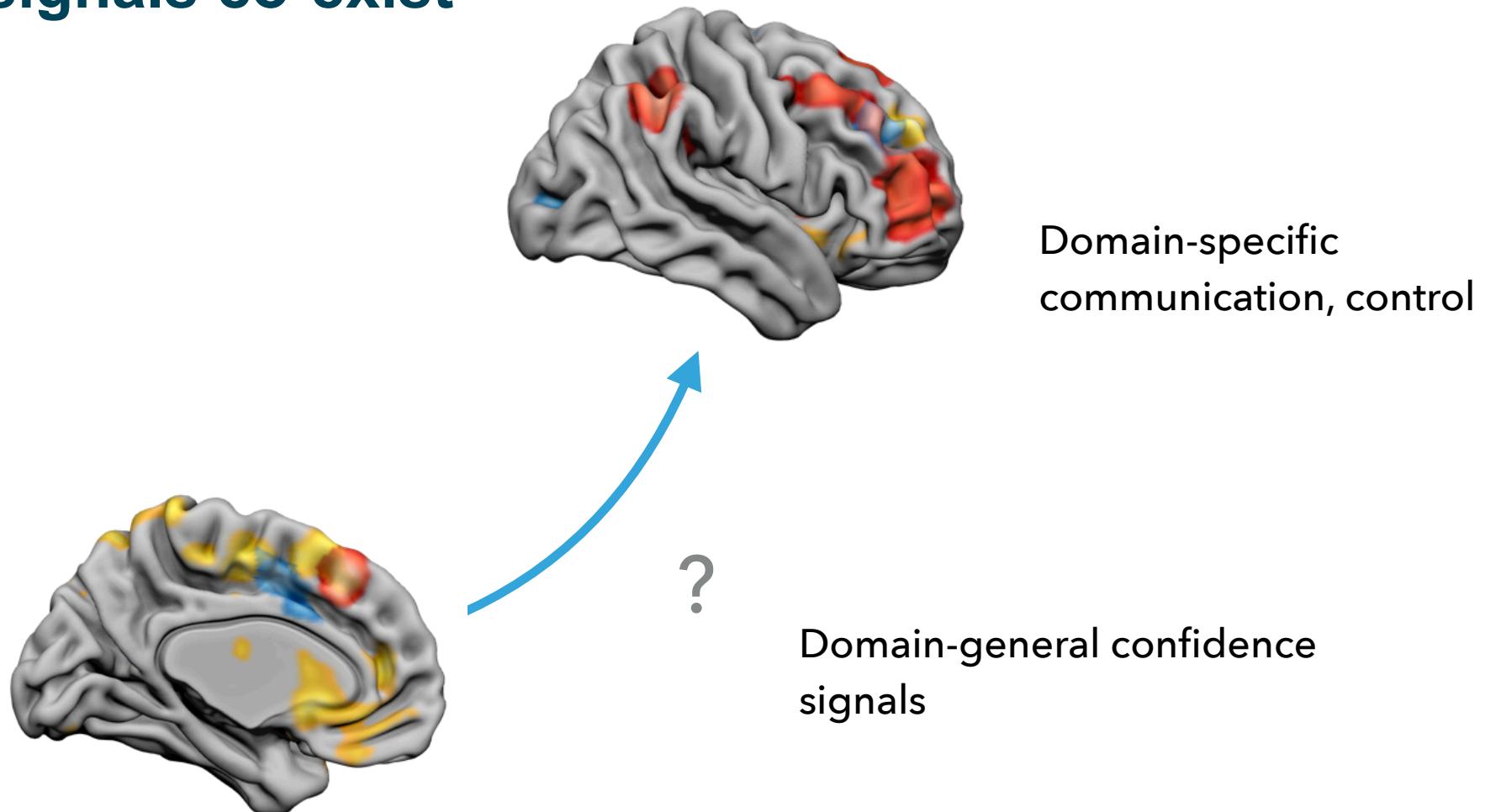
Judgment-related patterns



Judgment-related activity patterns



Domain-general and domain-specific confidence signals co-exist



Summary

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)
- Adopting a signal detection theory framework allows simultaneous estimation of both type 1 (d') and type 2 (meta- d') sensitivity

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)
- Adopting a signal detection theory framework allows simultaneous estimation of both type 1 (d') and type 2 (meta- d') sensitivity

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)
- Adopting a signal detection theory framework allows simultaneous estimation of both type 1 (d') and type 2 (meta- d') sensitivity
- Psychiatric symptom dimensions are associated with changes in metacognition over and above differences in behavioural performance

Summary

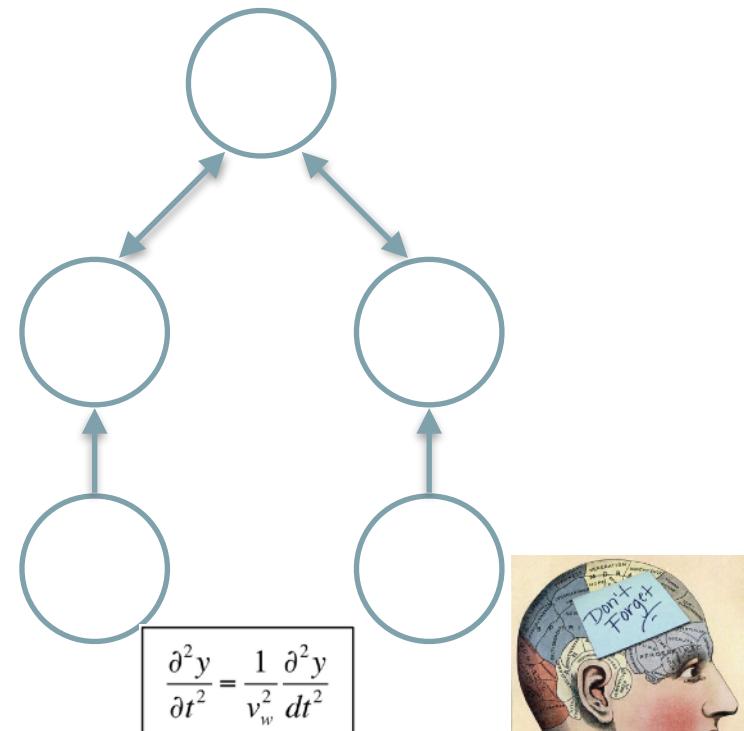
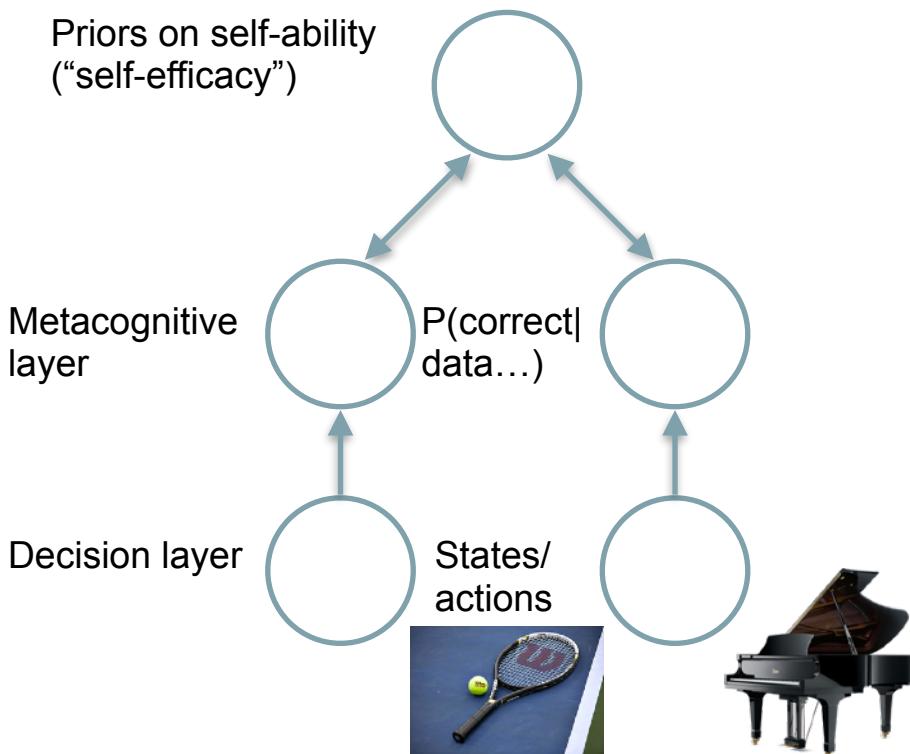
- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)
- Adopting a signal detection theory framework allows simultaneous estimation of both type 1 (d') and type 2 (meta- d') sensitivity
- Psychiatric symptom dimensions are associated with changes in metacognition over and above differences in behavioural performance

Summary

- We can measure metacognition across different tasks as the statistical association between behaviour and self-evaluation (confidence)
- Adopting a signal detection theory framework allows simultaneous estimation of both type 1 (d') and type 2 (meta- d') sensitivity
- Psychiatric symptom dimensions are associated with changes in metacognition over and above differences in behavioural performance
- Domain-general aspects of metacognition may prove relevant targets for computational psychiatry; potential to generalise to new tasks/experiences

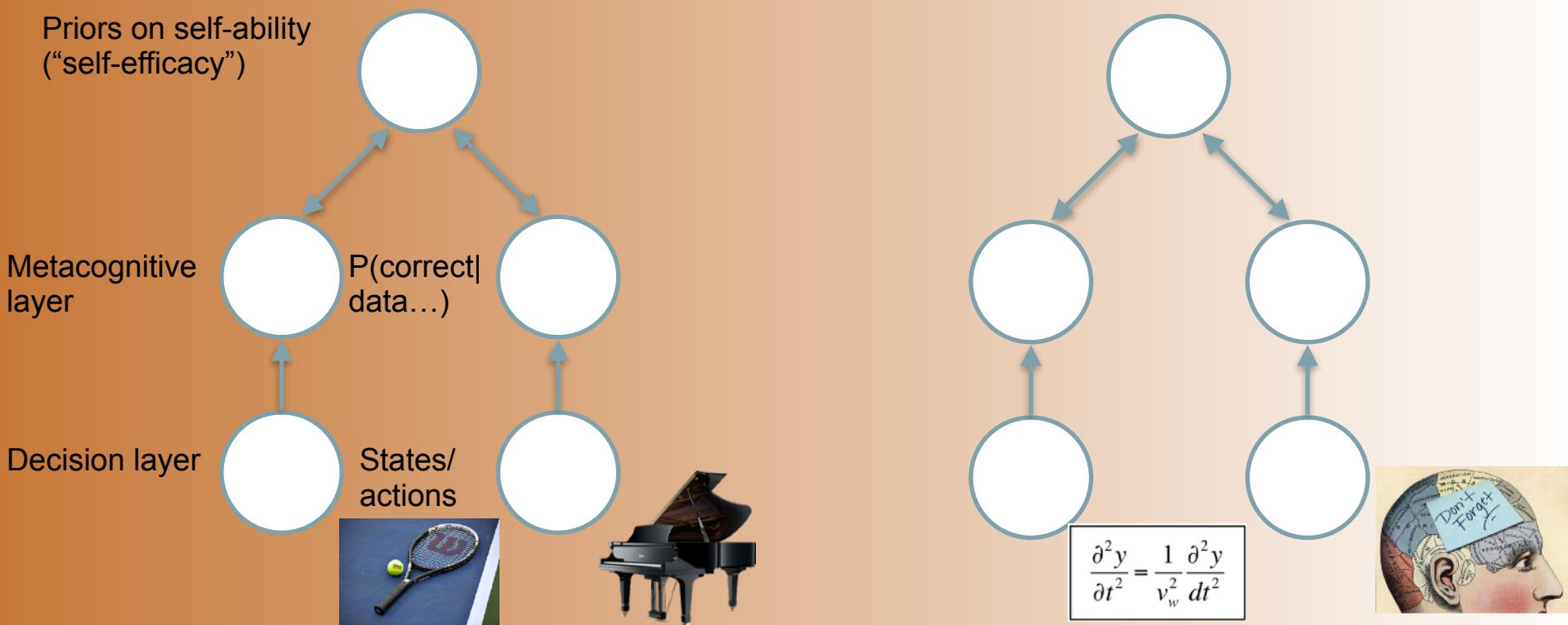
Metacognitive generalisation

The extent to which we **generalise** across domains may itself be an important/interesting individual difference



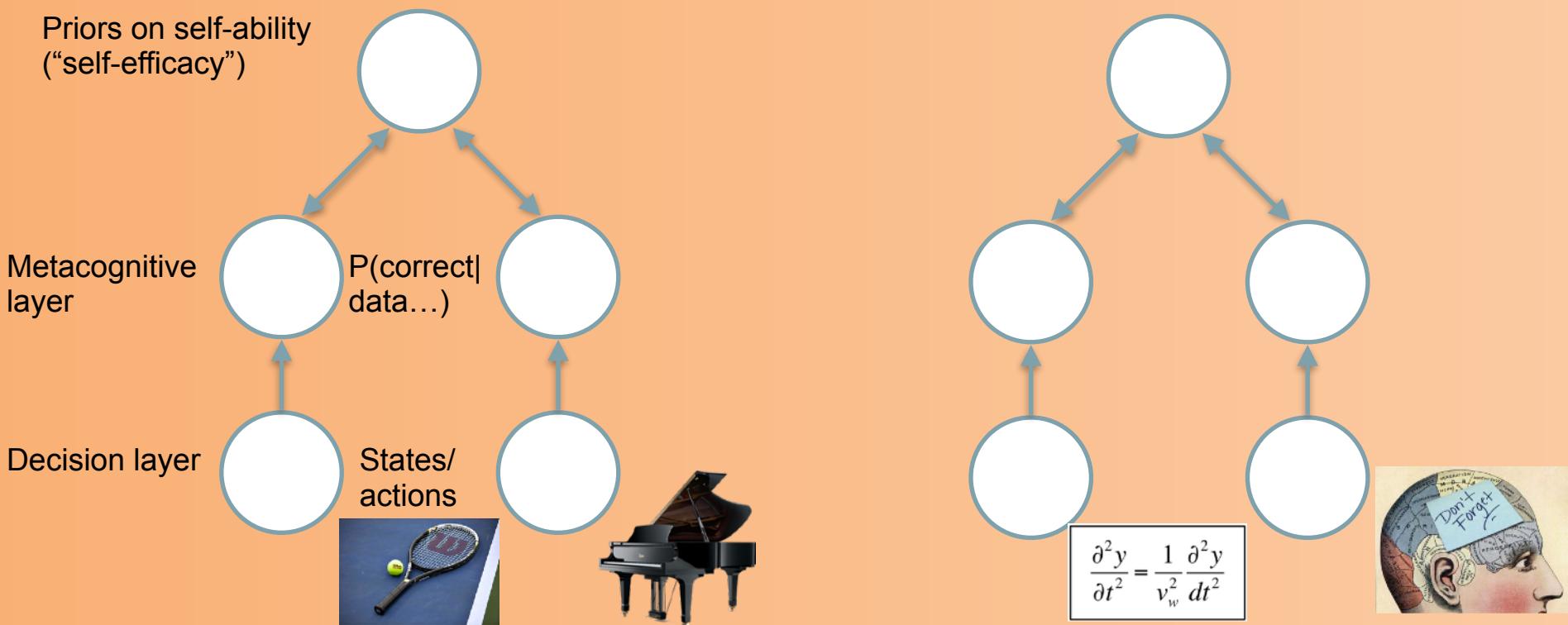
Metacognitive generalisation

The extent to which we **generalise** across domains may itself be an important/interesting individual difference

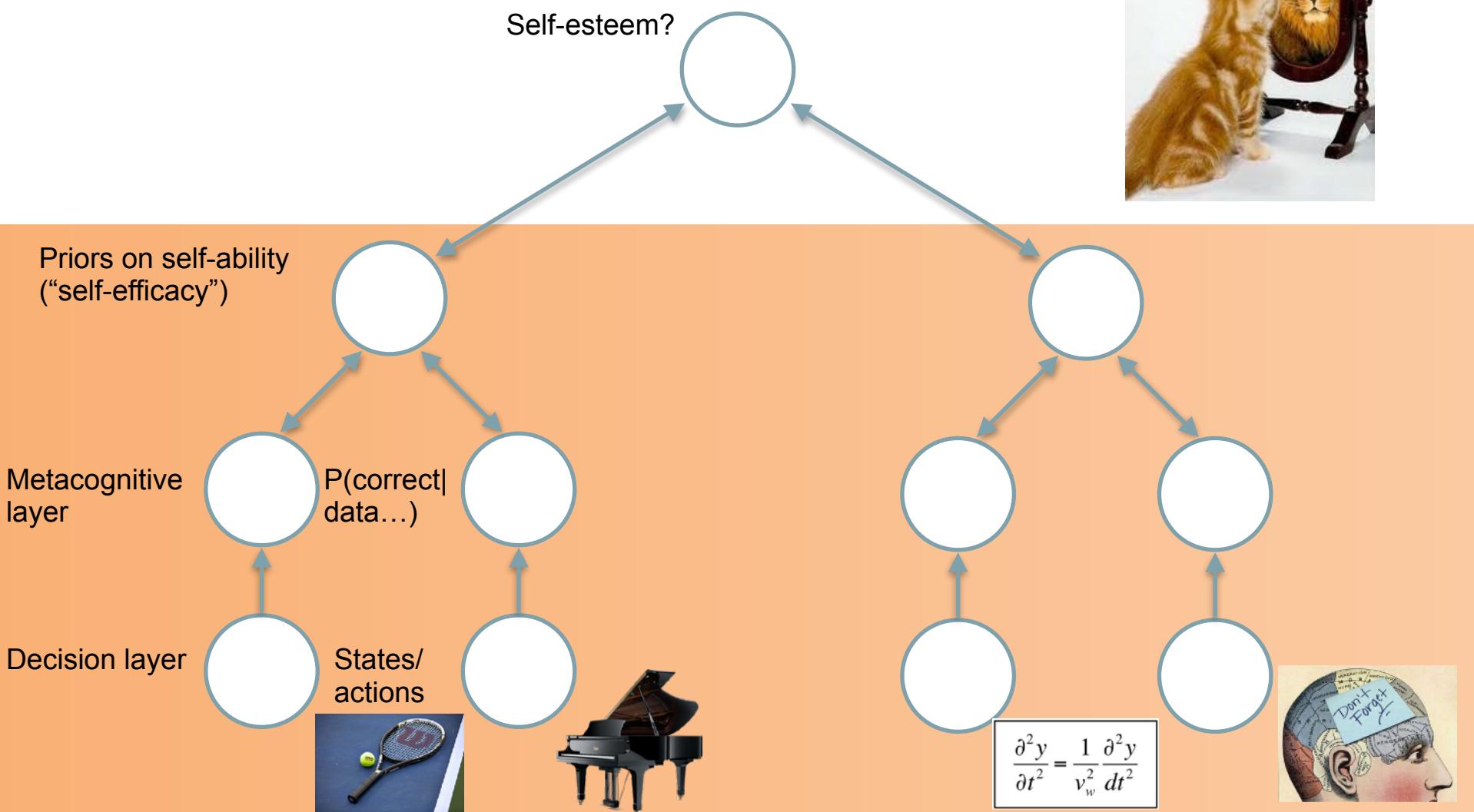


Metacognitive generalisation

The extent to which we **generalise** across domains may itself be an important/interesting individual difference

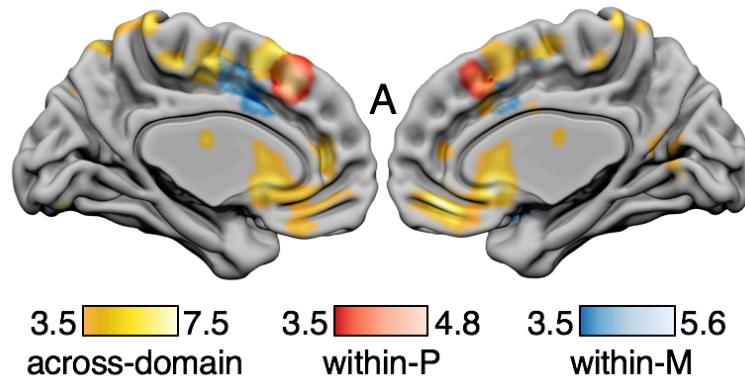


Metacognitive generalisation

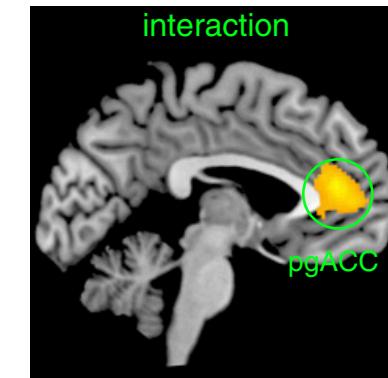


Metacognitive generalisation and pgACC

Confidence formation:



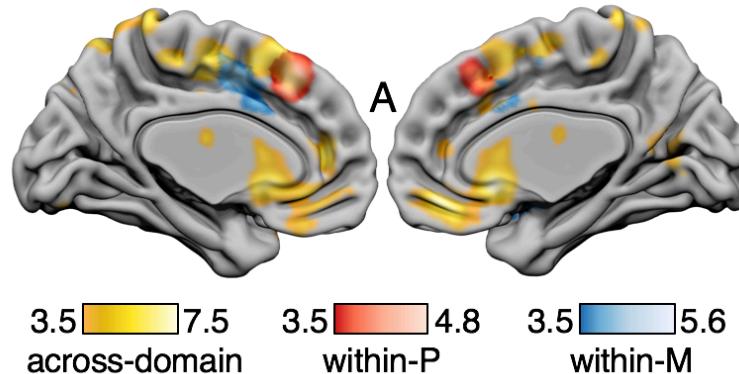
Morales et al. (2018) *J Neuro*



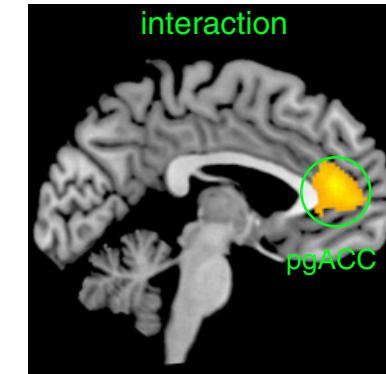
Bang & Fleming (2018) *PNAS*

Metacognitive generalisation and pgACC

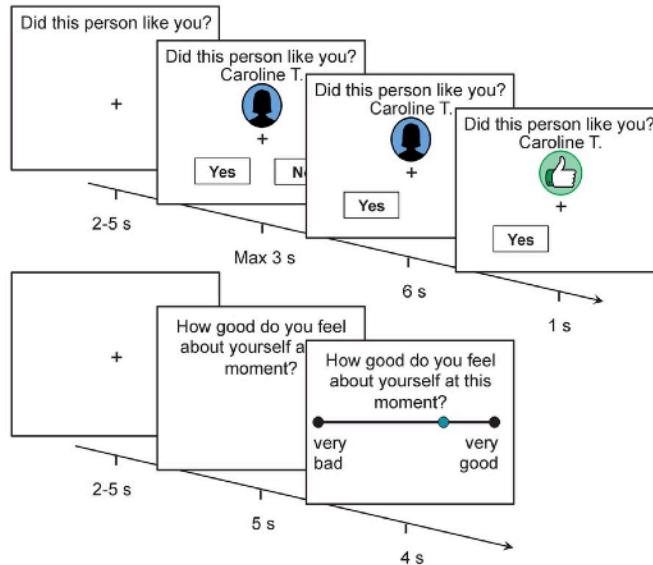
Confidence formation:



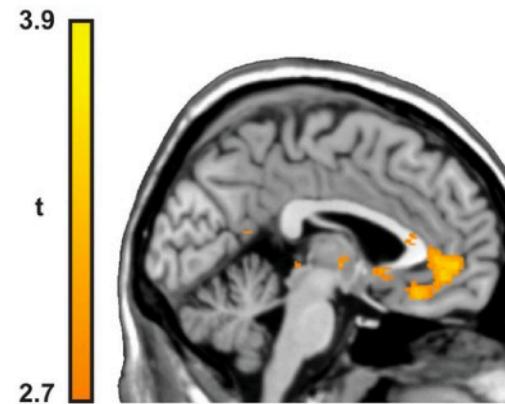
Morales et al. (2018) *J Neuro*



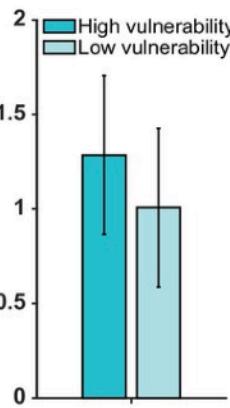
Bang & Fleming (2018) *PNAS*



Self-esteem updates:



Effect self-esteem update
vmPFC (a.u.)



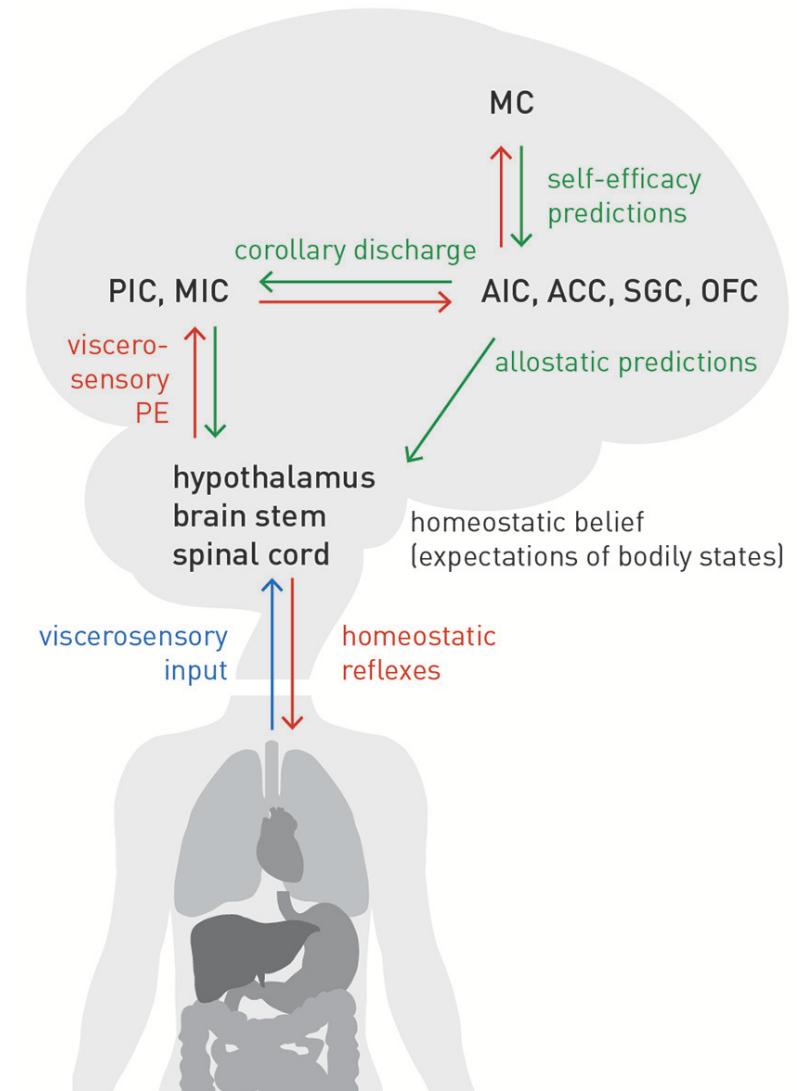
Will et al. (2017) *eLife*

Metacognitive generalisation

Generalisation may operate not only across exteroceptive domains, but also interoception

Allostatic Self-efficacy: A Metacognitive Theory of Dyshomeostasis-Induced Fatigue and Depression

Klaas E. Stephan^{1,2,3*}, Zina M. Manjaly^{1,4}, Christoph D. Mathys², Lilian A. E. Weber¹, Saeed Paliwal¹, Tim Gard^{1,5}, Marc Tittgemeyer³, Stephen M. Fleming², Helene Haker¹, Anil K. Seth⁶ and Frederike H. Petzschner¹



In this Issue

Metacognition: computation, neurobiology and function

Papers of a Theme Issue organized and edited by Stephen M. Fleming, Raymond J. Dolan,
Christopher D. Frith



The world's first science journal

ISSN 0962-8436

volume 367

number 1594

pages 1279–1438

Stephen M. Fleming
Christopher D. Frith *Editors*

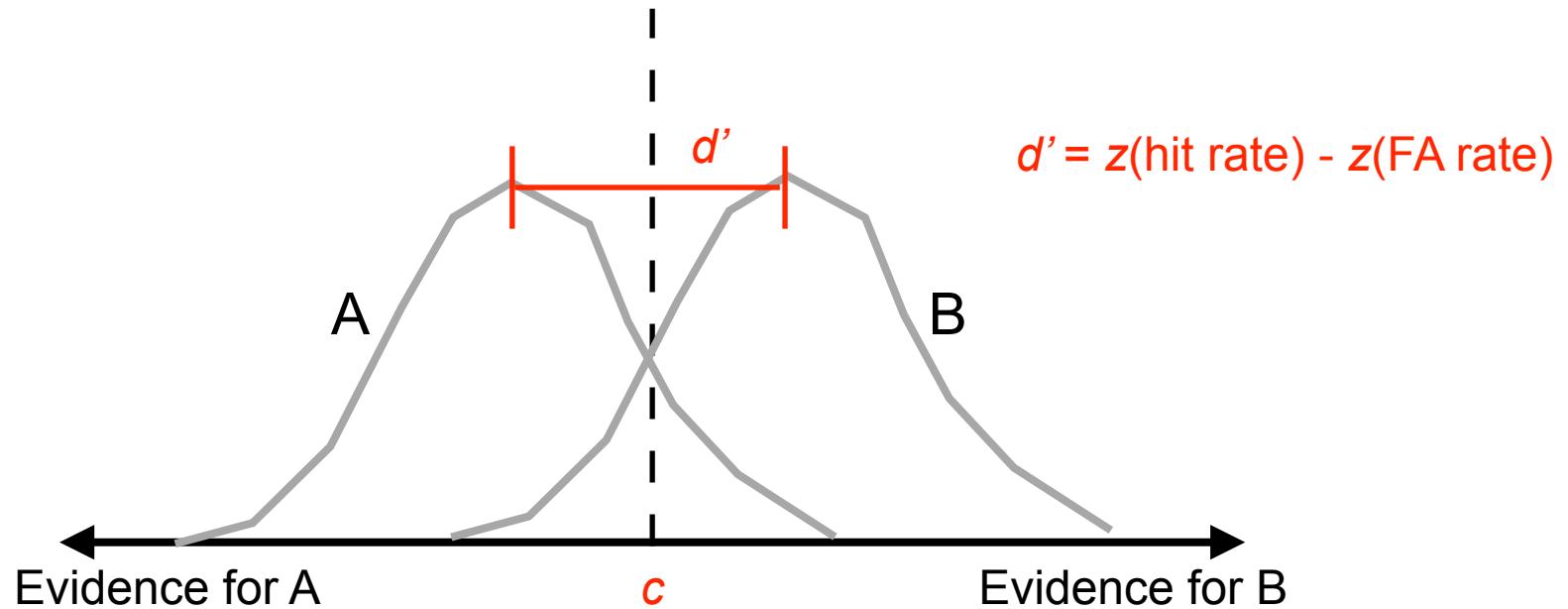
The Cognitive Neuroscience of Metacognition

Thank you

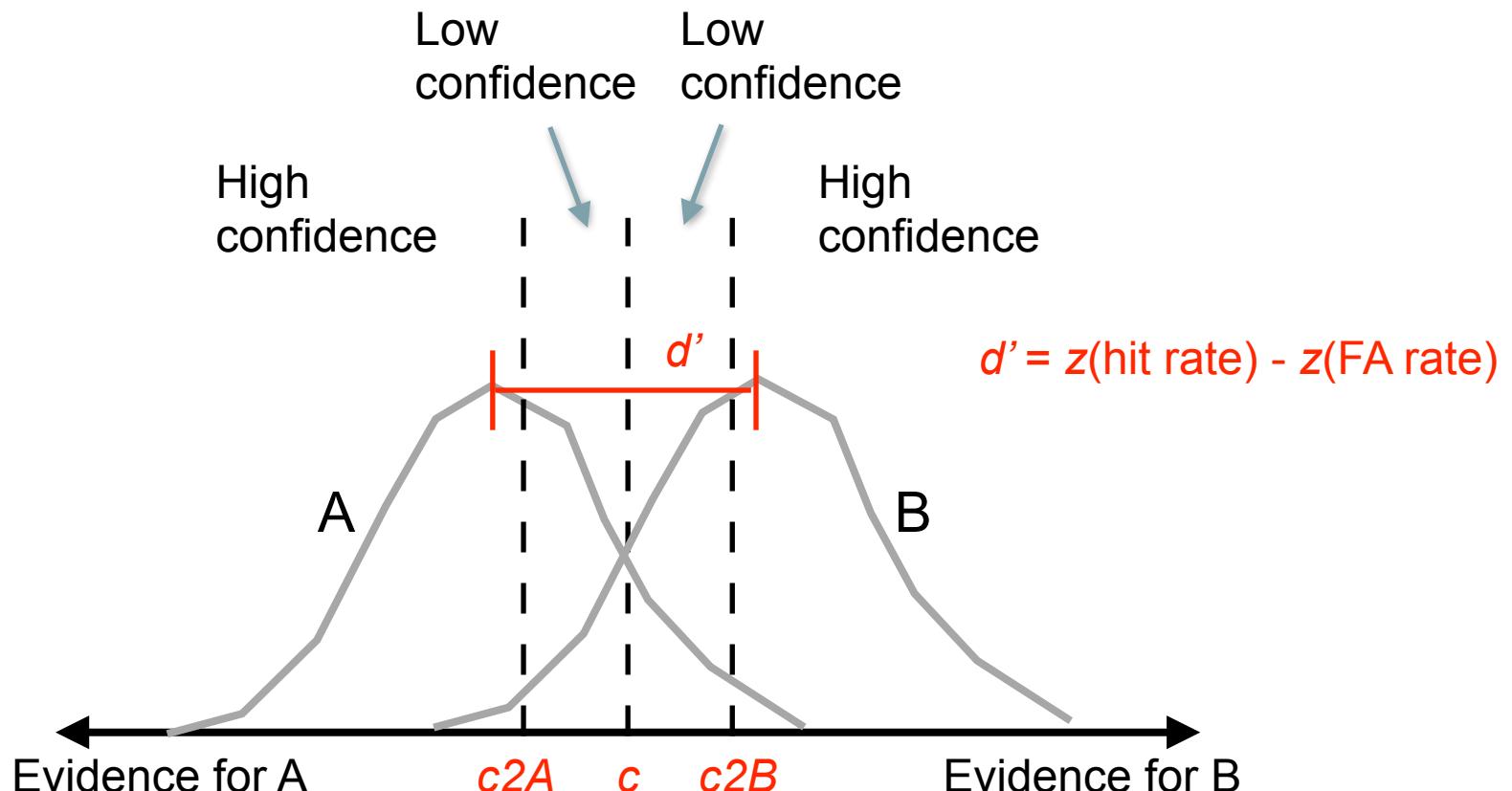


metacoglab.org

Signal detection theory (SDT)

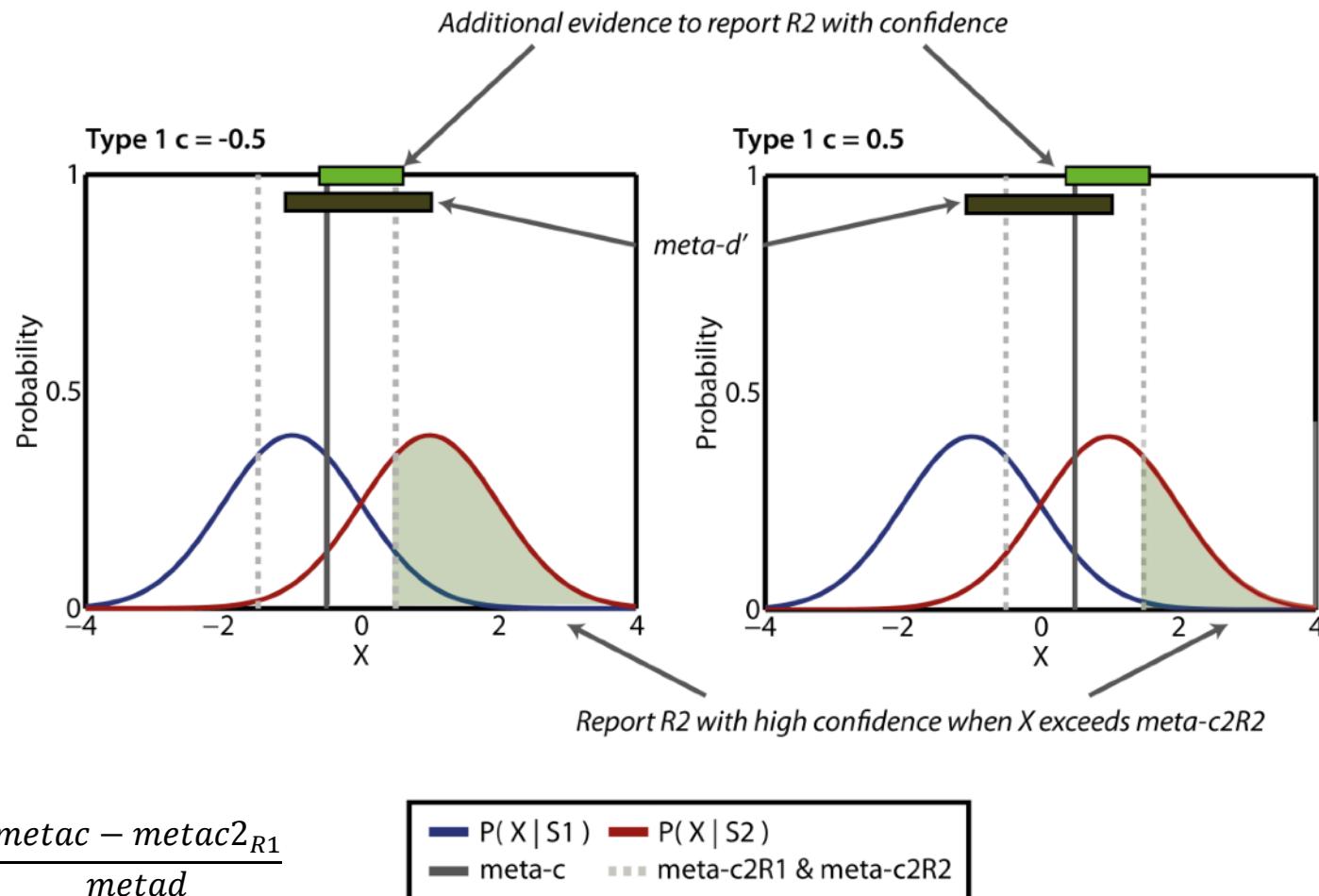


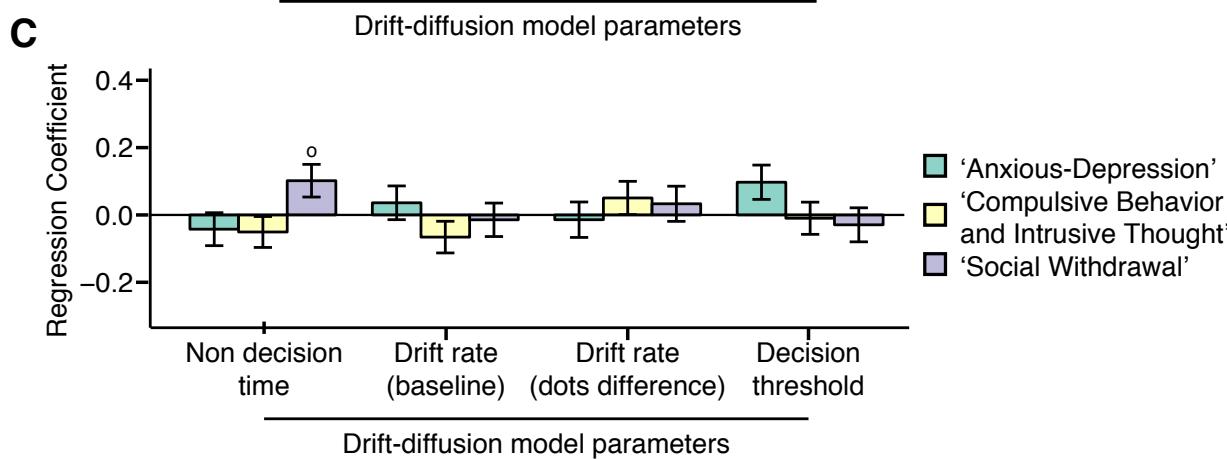
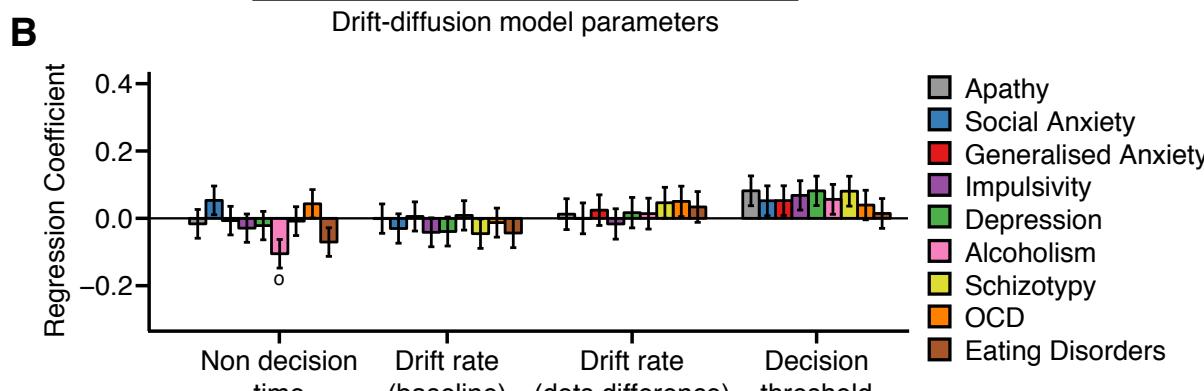
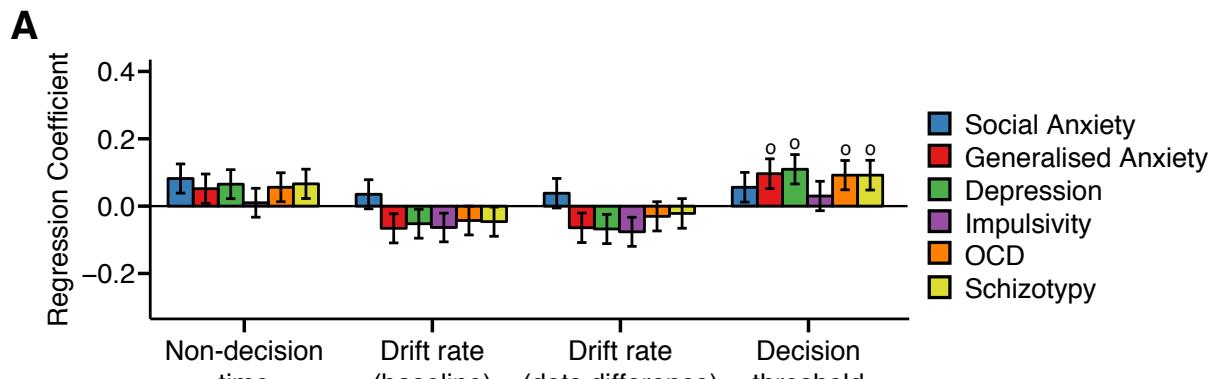
Signal detection theory (SDT)

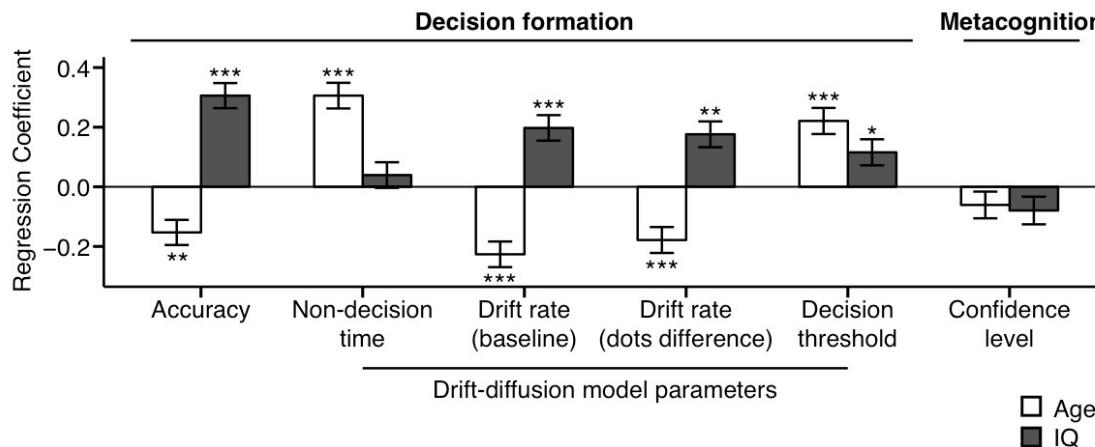
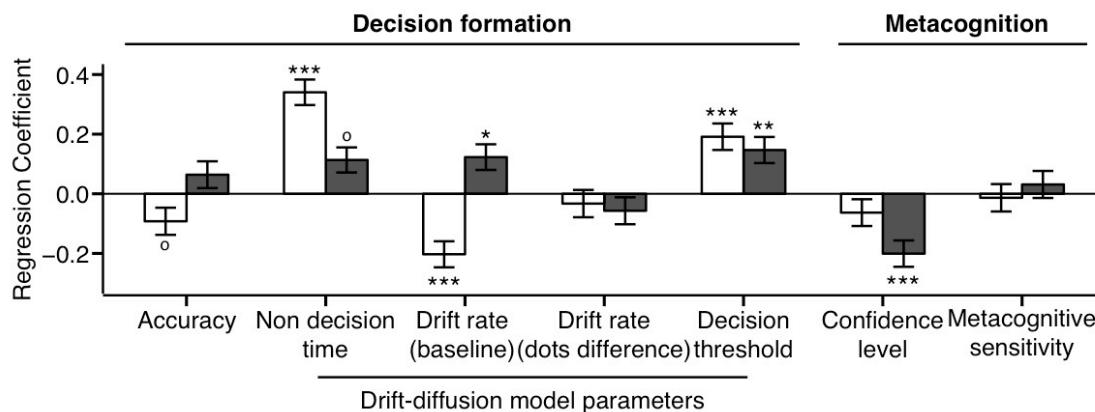


Confidence related to absolute
distance to decision criterion

Quantifying metacognitive bias





a
EXPERIMENT 1

b
EXPERIMENT 2

c

Variables list:

- Accuracy
- Non-decision time
- Drift rate (baseline)
- Drift rate (dots difference)
- Decision threshold
- Confidence level
- Metacognitive sensitivity

EXPERIMENT 2
