

# Bayesian model selection & averaging

Klaas Enno Stephan



Translational Neuromodeling Unit

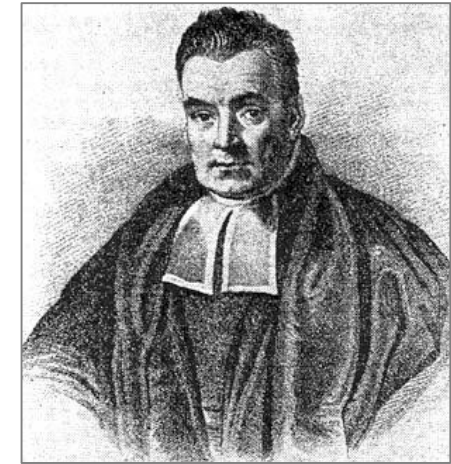
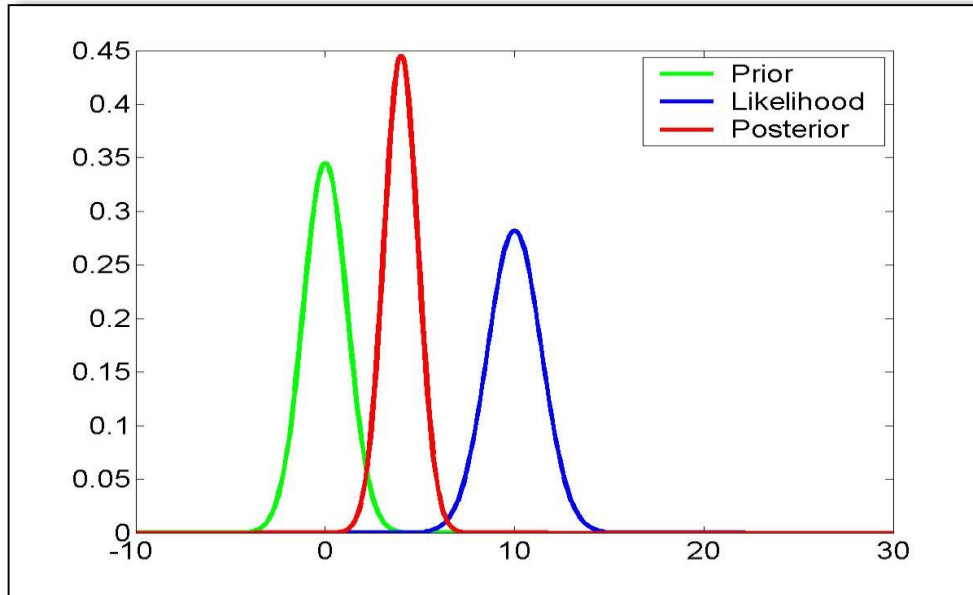


Universität  
Zürich<sup>UZH</sup>



Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

# Bayes' theorem

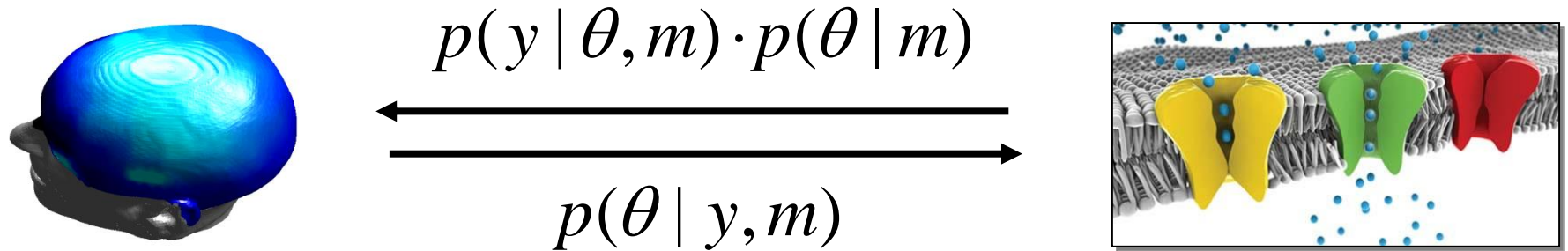


The Reverend Thomas Bayes  
(1702-1761)

$$p(\theta \mid y, m) = \frac{p(y \mid \theta, m) p(\theta \mid m)}{p(y \mid m)}$$

posterior = likelihood • prior / evidence

# Generative model

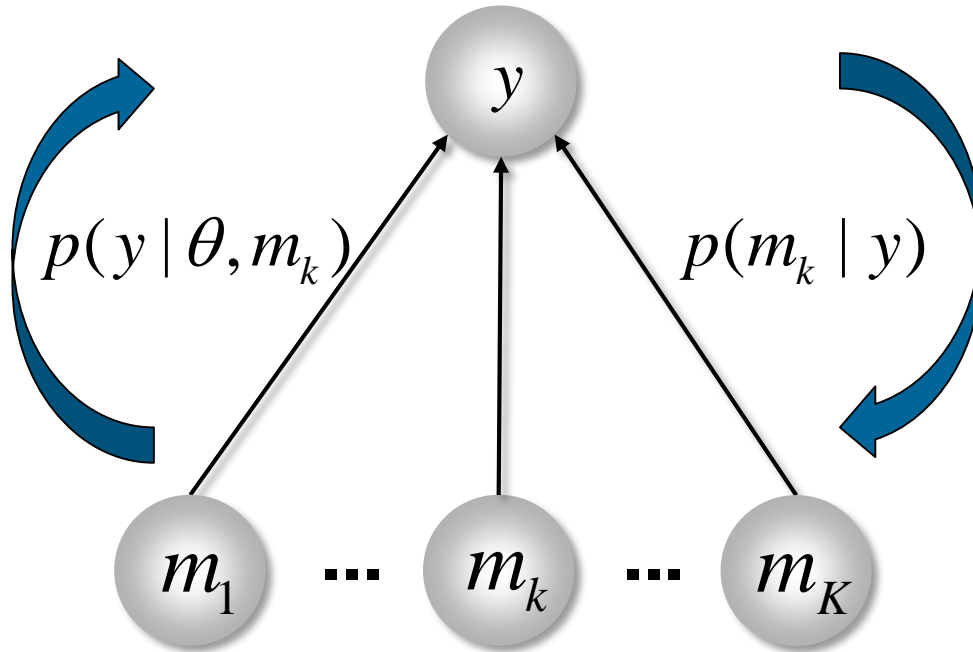


1. enforces mechanistic thinking: how could the data have been caused?
2. generate synthetic data (observations) by sampling from the prior – can model explain certain phenomena at all?
3. inference about parameters  $\rightarrow p(\theta|y)$
4. inference about model structure: formal approach to disambiguating mechanisms  $\rightarrow p(y|m)$  or  $p(m|y)$

# Long-term goal: Differential diagnosis based on generative models of disease symptoms

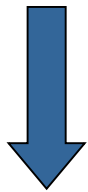
**SYMPTOM**  
(behaviour  
or physiology)

**HYPOTHETICAL  
MECHANISM**



# Model comparison and selection

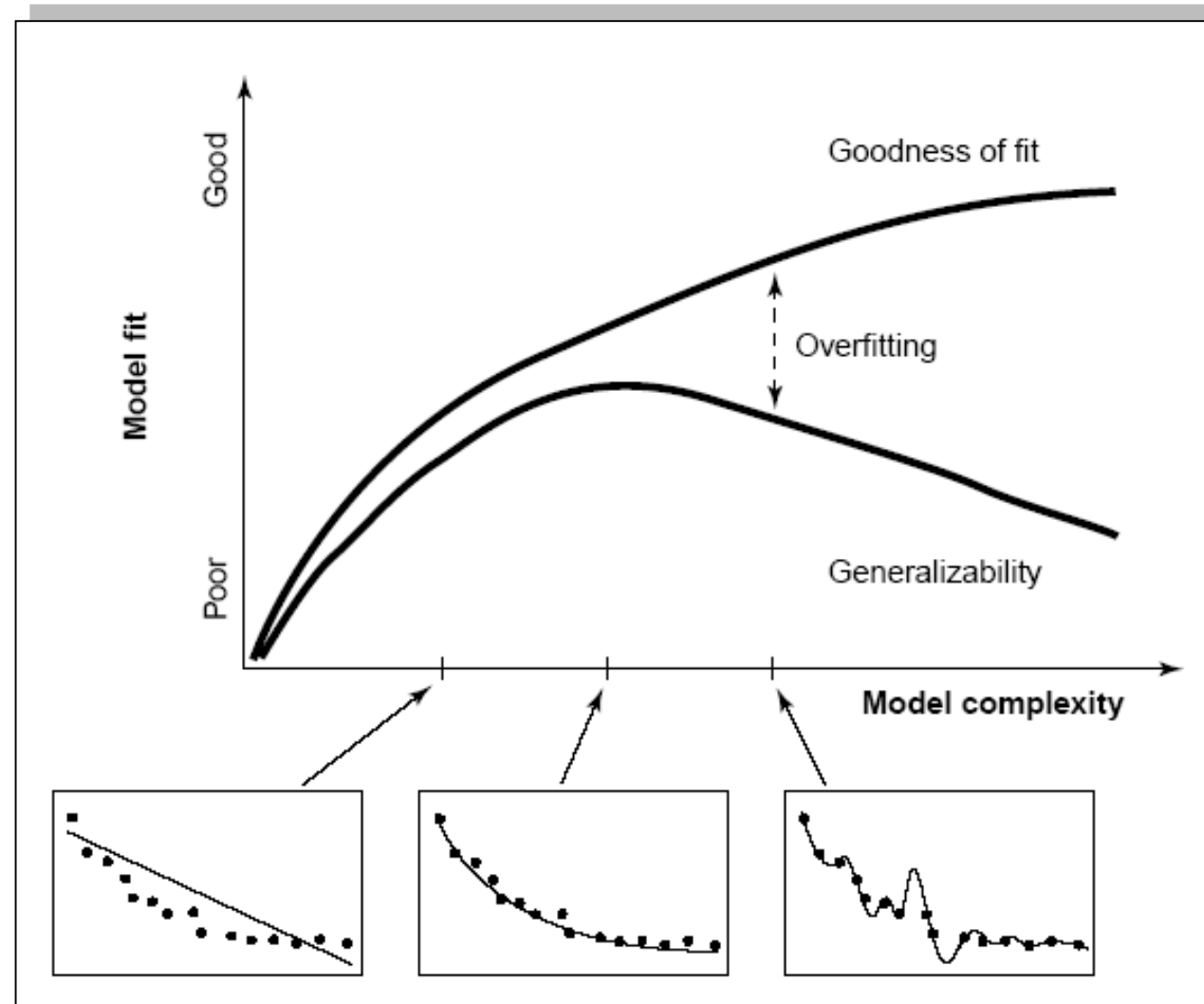
Given competing hypotheses on structure & functional mechanisms of a system, which model is the best?



Which model represents the best balance between model fit and model complexity?



For which model  $m$  does  $p(y|m)$  become maximal?



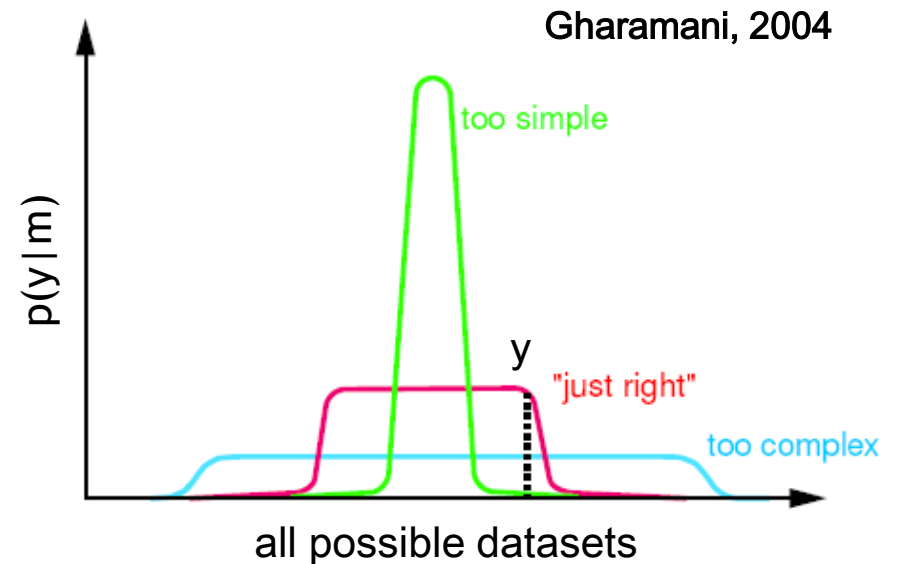
# Bayesian model selection (BMS)

Model evidence (marginal likelihood):

$$p(y | m) = \int p(y | \theta, m) p(\theta | m) d\theta$$

➡ accounts for both accuracy and complexity of the model

➡ “If I randomly sampled from my prior and plugged the resulting value into the likelihood function, how close would the predicted data be – on average – to my observed data?”



Various approximations, e.g.:

- negative free energy, AIC, BIC

McKay 1992, *Neural Comput.*  
Penny et al. 2004a, *NeuroImage*

# Model space (hypothesis set) $M$

Model space  $M$  is defined by prior on models.

Usual choice: flat prior over a small set of models.

$$p(m) = \begin{cases} 1/|M| & \text{if } m \in M \\ 0 & \text{if } m \notin M \end{cases}$$

In this case, the posterior probability of model  $i$  is:

$$p(m_i | y) = \frac{p(y | m_i) p(m_i)}{\sum_{j=1}^{|M|} p(y | m_j) p(m_j)} = \frac{p(y | m_i)}{\sum_{j=1}^{|M|} p(y | m_j)}$$

# Approximations to the model evidence in DCM

Logarithm is a  
monotonic function



Maximizing log model evidence  
= Maximizing model evidence

Log model evidence = balance between fit and complexity

$$\begin{aligned}\log p(y | m) &= \text{accuracy}(m) - \text{complexity}(m) \\ &= \log p(y | \theta, m) - \text{complexity}(m)\end{aligned}$$

Akaike Information Criterion:

$$AIC = \log p(y | \theta, m) - p$$

No. of  
parameters

No. of  
data points

Bayesian Information Criterion:

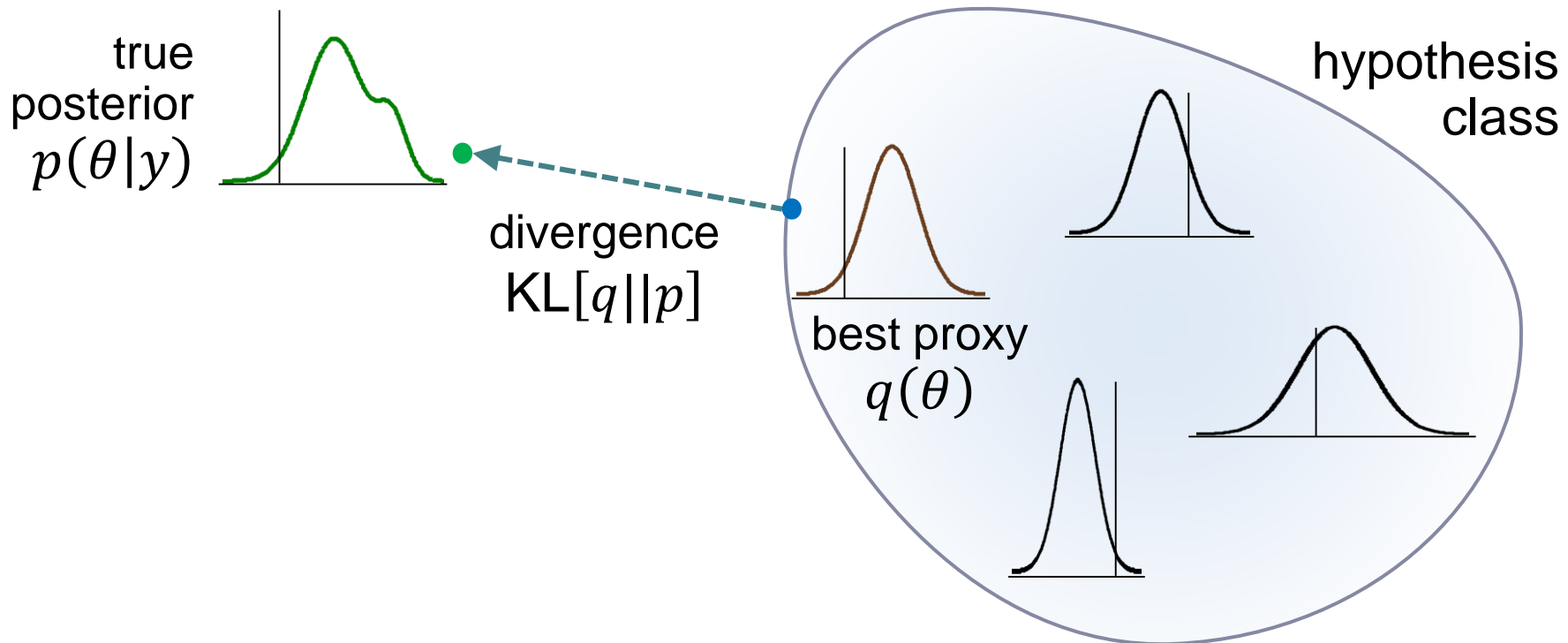
$$BIC = \log p(y | \theta, m) - \frac{p}{2} \log N$$



# Variational Bayes (VB)

Idea: find an approximate density  $q(\theta)$  that is maximally similar to the true posterior  $p(\theta|y)$ .

This is often done by assuming a particular form for  $q$  (fixed form VB) and then optimizing its sufficient statistics.



# The (negative) free energy approximation $F$

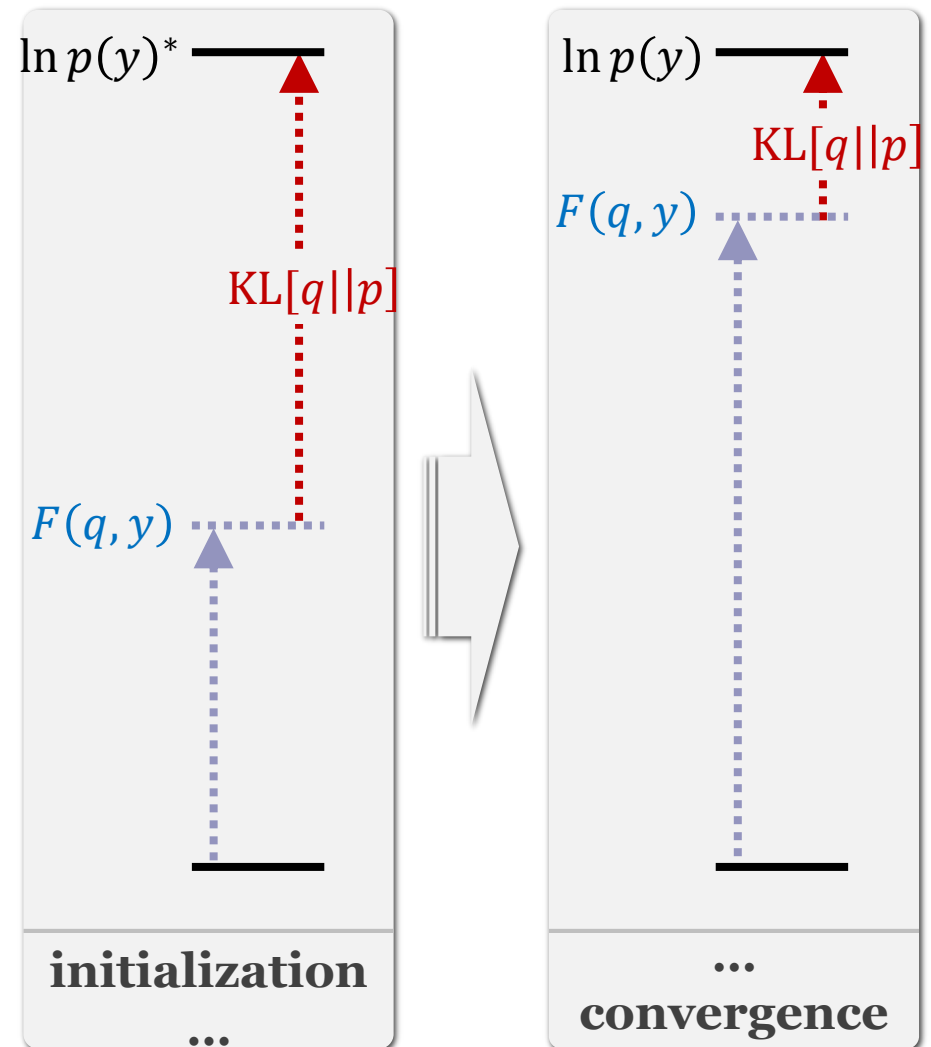
$$\ln p(y) = \underbrace{\text{KL}[q||p]}_{\substack{\text{divergence} \\ \geq 0 \\ \text{(unknown)}}} + \underbrace{F(q, y)}_{\substack{\text{neg. free} \\ \text{energy} \\ \text{(easy to evaluate} \\ \text{for a given } q)}}$$

$F(q, y)$  is a functional wrt. the approximate posterior  $q(\theta)$ .

Maximizing  $F(q, y)$  is equivalent to:

- minimizing  $\text{KL}[q||p]$
- tightening  $F(q, y)$  as a lower bound to the log model evidence

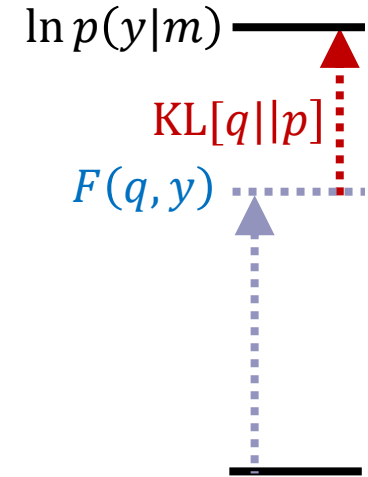
When  $F(q, y)$  is maximized,  $q(\theta)$  is our best estimate of the posterior.



# The (negative) free energy approximation $F$

$F$  is a lower bound on the log model evidence, where the bound is determined by the KL divergence between an approximate posterior  $q$  and the true posterior::

$$\log p(y | m) = F + KL[q(\theta), p(\theta | y, m)]$$



Like AIC/BIC,  $F$  is an accuracy/complexity tradeoff:

$$F = \underbrace{\langle \log p(y | \theta, m) \rangle}_{\text{accuracy}} - \underbrace{KL[q(\theta), p(\theta | m)]}_{\text{complexity}}$$

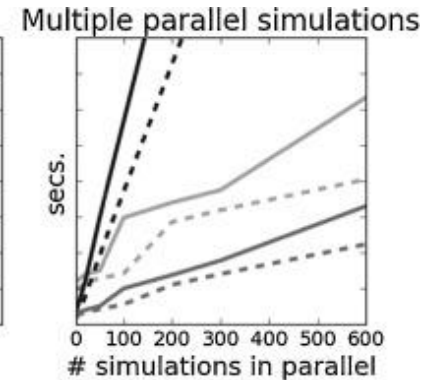
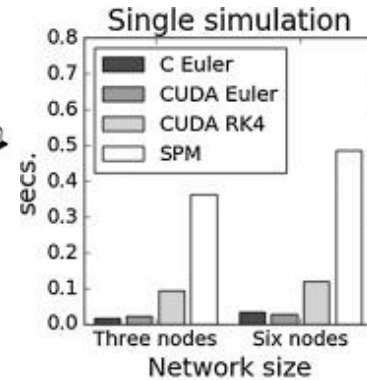
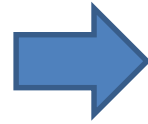
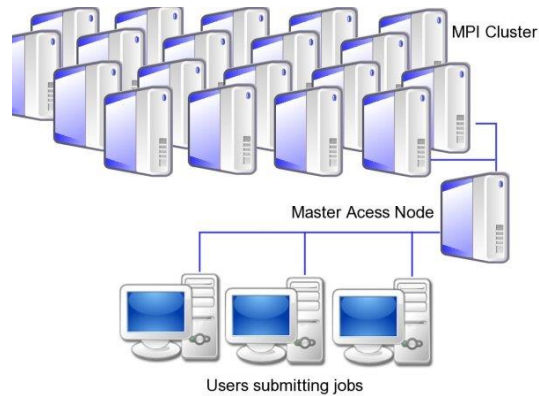
# The complexity term in $F$

- In contrast to AIC & BIC, the complexity term of the negative free energy  $F$  accounts for parameter interdependencies.

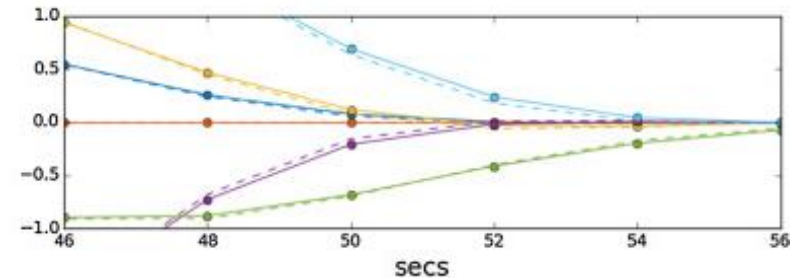
$$\begin{aligned} & KL[q(\theta), p(\theta | m)] \\ &= \frac{1}{2} \ln |C_\theta| - \frac{1}{2} \ln |C_{\theta|y}| + \frac{1}{2} (\mu_{\theta|y} - \mu_\theta)^T C_\theta^{-1} (\mu_{\theta|y} - \mu_\theta) \end{aligned}$$

- determinant = measure of “volume” (space spanned by the eigenvectors of the matrix)
- The complexity term of  $F$  is higher
  - the more independent the prior parameters ( $\uparrow$  effective DFs)
  - the more dependent the posterior parameters
  - the more the posterior mean deviates from the prior mean

# mpdcm: Computing the evidence by sampling



$$\left. \begin{array}{l} \dot{x} = f(x, u_1, \theta_1) \\ \dot{x} = f(x, u_2, \theta_2) \\ \vdots \\ \dot{x} = f(x, u_1, \theta_1) \end{array} \right\} \text{mpdcm\_integrate(dcms)} \left\{ \begin{array}{l} y_1 \\ y_2 \\ \vdots \\ y_3 \end{array} \right.$$



# Bayes factors

To compare two models, we could just compare their log evidences.

But: the log evidence is just some number – not very intuitive!

A more intuitive interpretation of model comparisons is made possible by Bayes factors:

positive value,  $[0; \infty[$

$$B_{12} = \frac{p(y | m_1)}{p(y | m_2)}$$

Kass & Raftery classification:

$B_{12}$	$p(m_1 y)$	Evidence
1 to 3	50-75%	weak
3 to 20	75-95%	positive
20 to 150	95-99%	strong
$\geq 150$	$\geq 99\%$	Very strong

## Fixed effects BMS at group level

**Group Bayes factor (GBF)** for  $1 \dots K$  subjects:

$$GBF_{ij} = \prod_k BF_{ij}^{(k)}$$

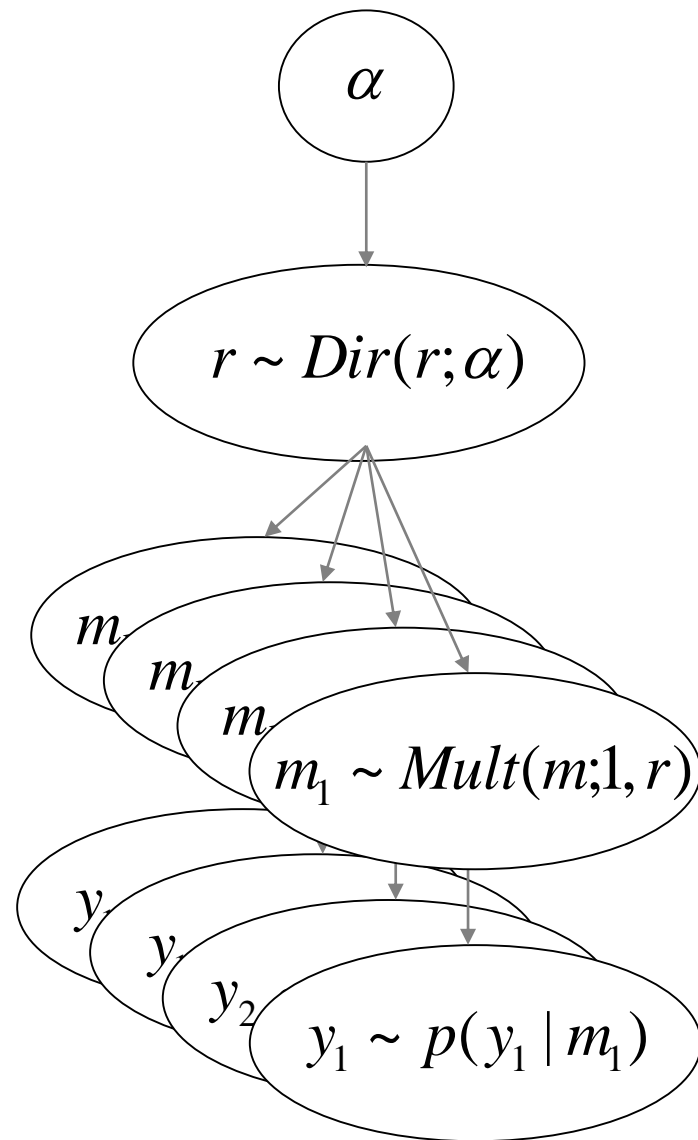
**Average Bayes factor (ABF):**

$$ABF_{ij} = \sqrt[K]{\prod_k BF_{ij}^{(k)}}$$

**Problems:**

- blind with regard to group heterogeneity
- sensitive to outliers

# Random effects BMS for heterogeneous groups



Dirichlet parameters  $\alpha$   
= “occurrences” of models in the population

Dirichlet distribution of model probabilities  $r$

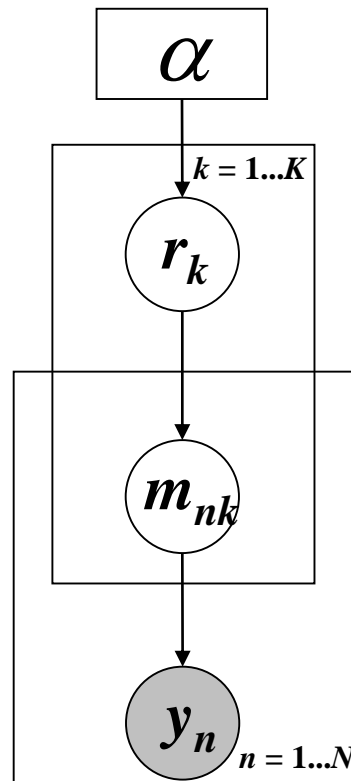
Multinomial distribution of model labels  $m$

Measured data  $y$

**Model inversion  
by Variational  
Bayes (VB) or  
MCMC**



# Random effects BMS for heterogeneous groups



Dirichlet parameters  $\alpha$   
= “occurrences” of models in the population

Dirichlet distribution of model probabilities  $r$

Multinomial distribution of model labels  $m$

Measured data  $y$

**Model inversion  
by Variational  
Bayes (VB) or  
MCMC**

# Reminder: VB in a nutshell (mean-field approximation)

- ❶ Neg. free-energy approx. to model evidence.

$$\ln p(y|m) = F + KL[q(\theta, \lambda), p(\theta, \lambda | y)]$$

$$F = \langle \ln p(y, \theta, \lambda) \rangle_q - KL[q(\theta, \lambda), p(\theta, \lambda | m)]$$

- ❷ Mean field approx.

$$p(\theta, \lambda | y) \approx q(\theta, \lambda) = q(\theta)q(\lambda)$$

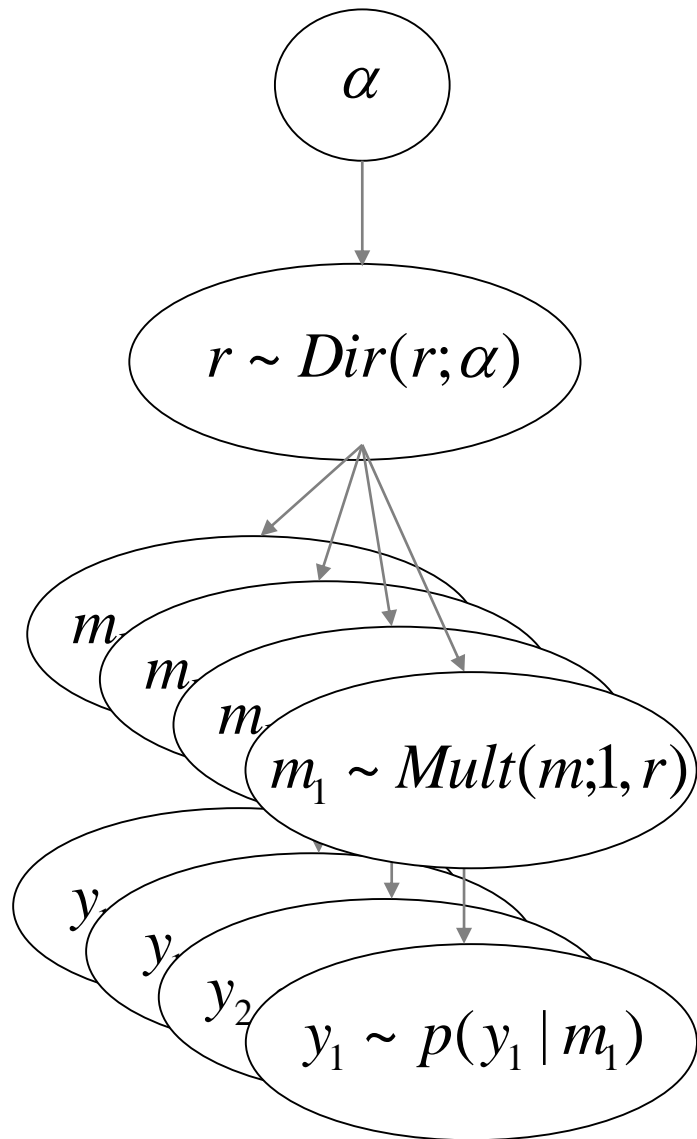
- ❸ Maximise neg. free energy wrt.  $q$  = minimise divergence, by maximising variational energies

$$q(\theta) \propto \exp(I_\theta) = \exp\left[\langle \ln p(y, \theta, \lambda) \rangle_{q(\lambda)}\right]$$

$$q(\lambda) \propto \exp(I_\lambda) = \exp\left[\langle \ln p(y, \theta, \lambda) \rangle_{q(\theta)}\right]$$

- ❹ Iterative updating of sufficient statistics of approx. posteriors (e.g., by gradient ascent).

# VB solution for random effects BMS



$$p(r | \alpha) = \text{Dir}(r, \alpha) = \frac{1}{Z(\alpha)} \prod_k r_k^{\alpha_k - 1}$$

$$Z(\alpha) = \prod_k \Gamma(\alpha_k) / \Gamma\left(\sum_k \alpha_k\right)$$

$$p(m_n | r) = \prod_k r_k^{m_{nk}}$$

$$p(y_n | m_{nk}) = \int p(y | \mathcal{G}) p(\mathcal{G} | m_{nk}) d\mathcal{G}$$

- ❶ Write down joint probability and take the log

$$\begin{aligned} p(y, r, m) &= p(y | m) p(m | r) p(r | \alpha_0) \\ &= p(r | \alpha_0) \left[ \prod_n p(y_n | m_n) p(m_n | r) \right] \\ &= \frac{1}{Z(\alpha_0)} \left[ \prod_k r_k^{\alpha_{0k}-1} \right] \left[ \prod_n p(y_n | m_n) \prod_k r_k^{m_{nk}} \right] \\ &= \frac{1}{Z(\alpha_0)} \prod_n \left[ \prod_k \left[ p(y_n | m_{nk}) r_k \right]^{m_{nk}} r_k^{\alpha_{0k}-1} \right] \end{aligned}$$

$$\ln p(y, r, m) = -\ln Z(\alpha_0) + \sum_n \sum_k \left( (\alpha_{0k} - 1) \ln r_k + m_{nk} (\log p(y_n | m_{nk}) + \ln r_k) \right)$$

② Mean field approx.

$$q(r, m) = q(r)q(m)$$

③ Maximise neg. free energy wrt.  $q$  =  
minimise divergence,  
by maximising  
variational energies

$$q(r) \propto \exp(I(r))$$

$$q(m) \propto \exp(I(m))$$

$$I(r) = \langle \log p(y, r, m) \rangle_{q(m)}$$

$$I(m) = \langle \log p(y, r, m) \rangle_{q(r)}$$

#### ④ Iterative updating of sufficient statistics of approx. posteriors

$$\alpha = \alpha_0$$

$$\alpha_0 = [1, \dots, 1]$$

**Until convergence**

$$u_{nk} = \exp \left( \ln p(y_n | m_{nk}) + \Psi(\alpha_k) - \Psi \left( \sum_k \alpha_k \right) \right)$$

$$g_{nk} = \frac{u_{nk}}{\sum_k u_{nk}}$$

$$\beta_k = \sum_n g_{nk}$$

$$\alpha = \alpha_0 + \beta$$

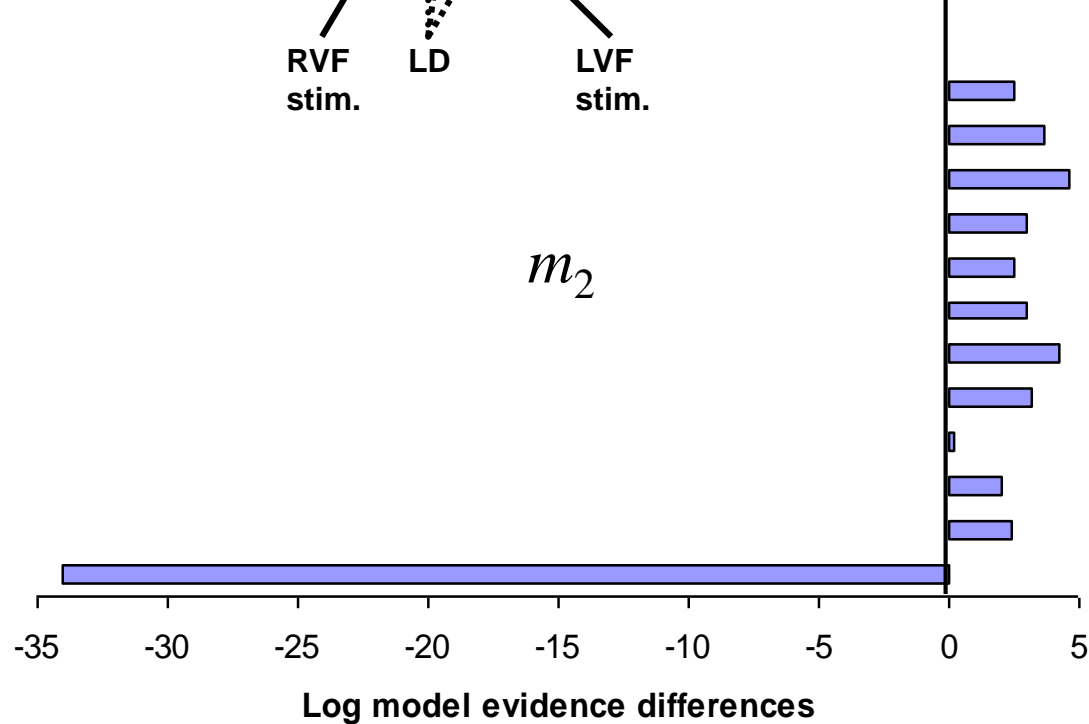
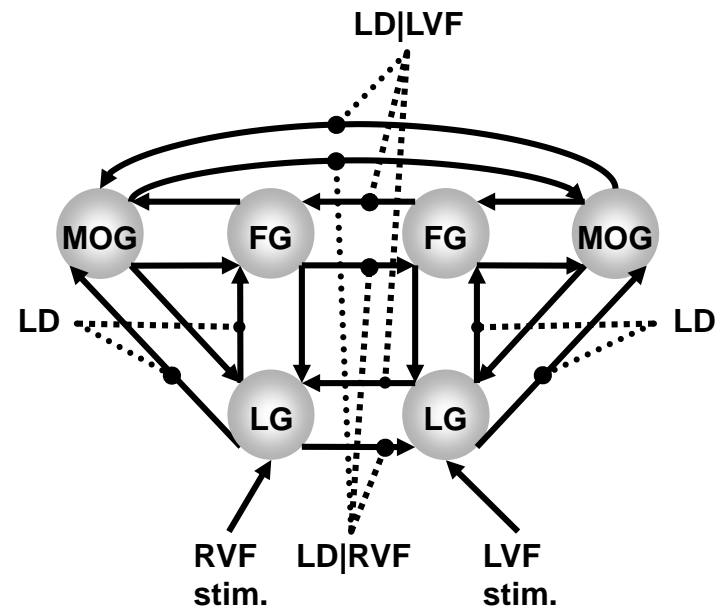
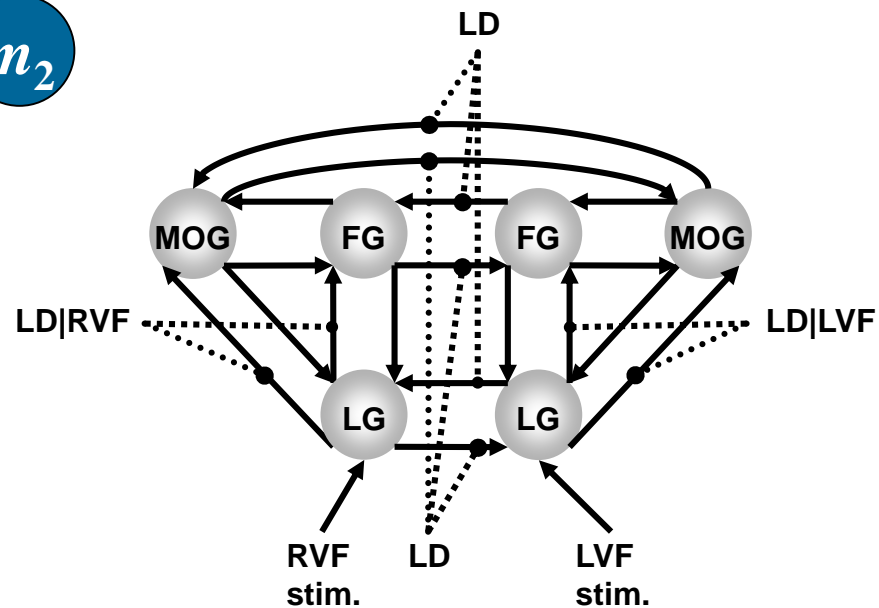
**end**

$$g_{nk} = q(m_{nk} = 1)$$

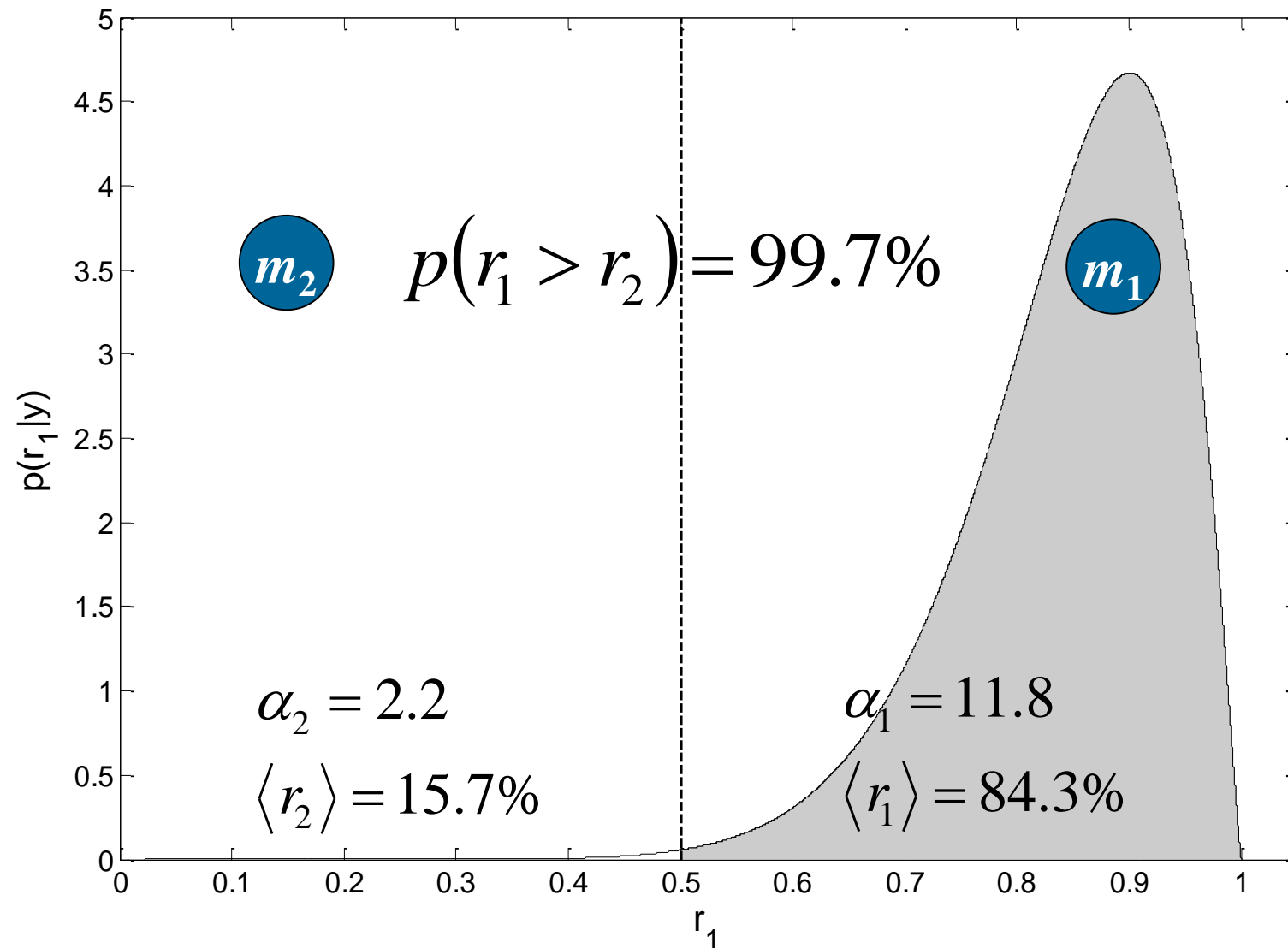
our (normalized)  
posterior belief that  
model  $k$  generated the  
data from subject  $n$

$$\beta_k = \sum_n g_{nk}$$

expected number of  
subjects whose data we  
believe were generated  
by model  $k$



**Data:** Stephan et al. 2003, *Science*  
**Models:** Stephan et al. 2007, *J. Neurosci.*





# How can we report the results of random effects BMS?

## 1. Dirichlet parameter estimates

$$\alpha$$

## 2. **expected posterior probability** of obtaining the $k$ -th model for any randomly selected subject

$$\langle r_k \rangle_q = \alpha_k / (\alpha_1 + \dots + \alpha_K)$$

## 3. **exceedance probability** that a particular model $k$ is more likely than any other model (of the $K$ models tested), given the group data

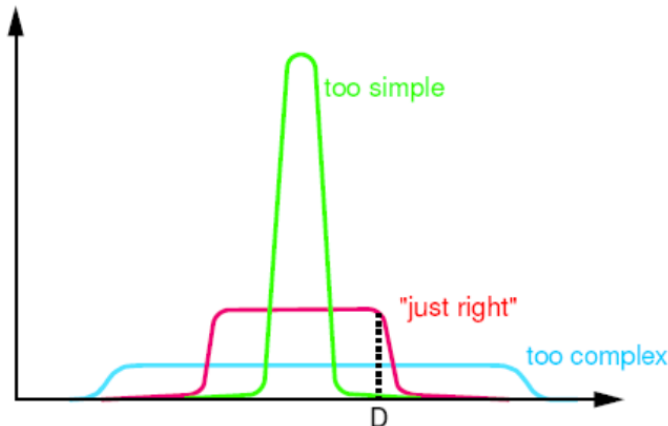
$$\exists k \in \{1 \dots K\}, \forall j \in \{1 \dots K \mid j \neq k\} :$$

$$\varphi_k = p(r_k > r_j \mid y; \alpha)$$

## 4. **protected exceedance probability**: see below

# Overfitting at the level of models

- $\uparrow \# \text{models} \Rightarrow \uparrow \text{risk of overfitting}$
- solutions:
  - regularisation: definition of model space = choosing priors  $p(m)$
  - family-level BMS
  - Bayesian model averaging (BMA)



posterior model probability:

$$p(m | y) = \frac{p(y | m) p(m)}{\sum_m p(y | m) p(m)}$$

BMA:

$$p(\theta | y) = \sum_m p(\theta | y, m) p(m | y)$$

# Model space partitioning or: Comparing model families

- partitioning model space into K subsets or families:

$$M = \{f_1, \dots, f_K\}$$

- pooling information over all models in these subsets allows one to compute the probability of a model family, given the data

$$p(f_k)$$

- effectively removes uncertainty about any aspect of model structure, other than the attribute of interest (which defines the partition)

# Family-level inference: fixed effects

- We wish to have a uniform prior at the family level:
- This is related to the model level via the sum of the priors on models:
- Hence the uniform prior at the family level is:
- The probability of each family is then obtained by summing the posterior probabilities of the models it includes:

$$p(f_k) = \frac{1}{K}$$

$$p(f_k) = \sum_{m \in f_k} p(m)$$

$$\forall m \in f_k : p(m) = \frac{1}{K|f_k|}$$

$$p(f_k | y_{1..N}) = \sum_{m \in f_k} p(m | y_{1..N})$$

# Family-level inference: random effects

- The frequency of a family in the population is given by:
- In RFX-BMS, this follows a Dirichlet distribution, with a uniform prior on the parameters  $\alpha$  (see above).
- A uniform prior over family probabilities can be obtained by setting:

$$s_k = \sum_{m \in f_k} r_m$$

$$p(s) = \text{Dir}(\alpha)$$

$$\forall m \in f_k : \alpha_{\text{prior}}(m) = \frac{1}{|f_k|}$$

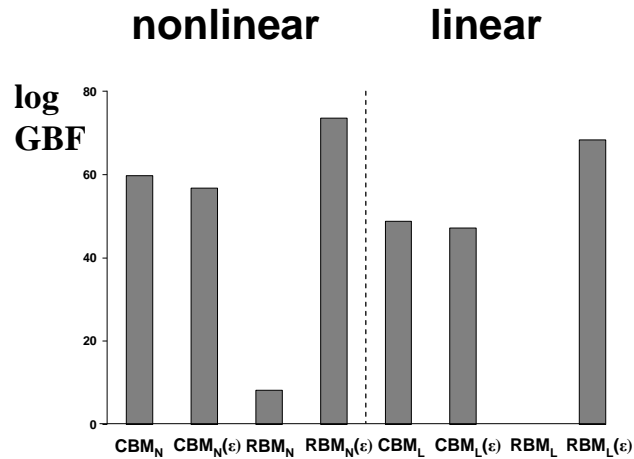
# Family-level inference: random effects – a special case

- When the families are of equal size, one can simply sum the posterior model probabilities within families by exploiting the agglomerative property of the Dirichlet distribution:

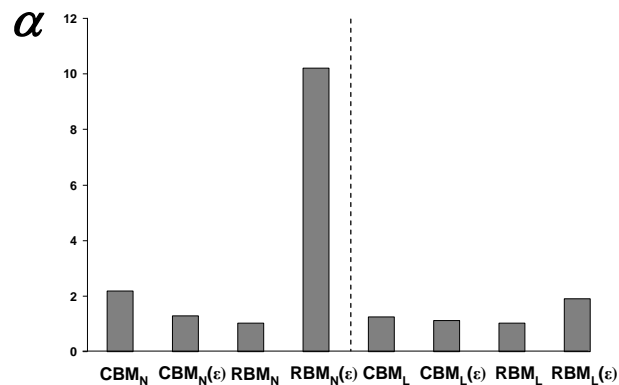
$$\begin{aligned} (r_1, r_2, \dots, r_K) &\sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_K) \\ \Rightarrow r_1^* &= \sum_{k \in N_1} r_k, r_2^* = \sum_{k \in N_2} r_k, \dots, r_J^* = \sum_{k \in N_J} r_k \\ &\sim \text{Dir}\left(\alpha_1^* = \sum_{k \in N_1} \alpha_k, \alpha_2^* = \sum_{k \in N_2} \alpha_k, \dots, \alpha_J^* = \sum_{k \in N_J} \alpha_k\right) \end{aligned}$$

# Model space partitioning: comparing model families

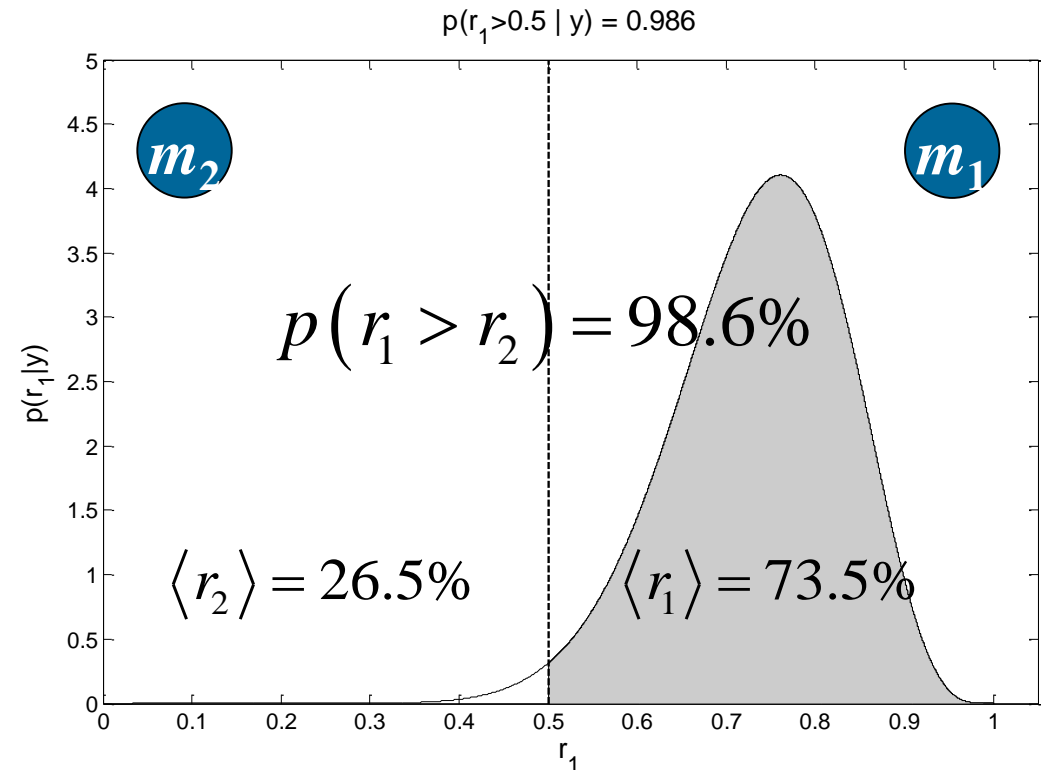
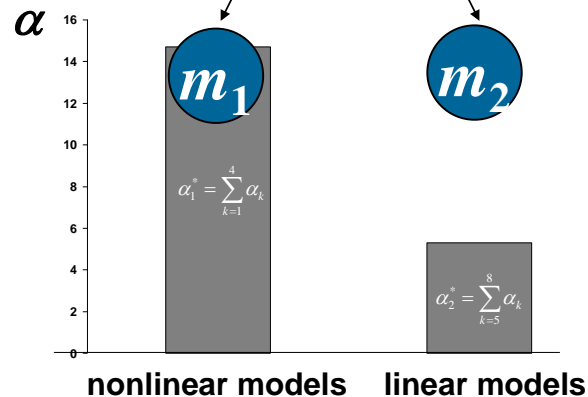
FFX



RFX

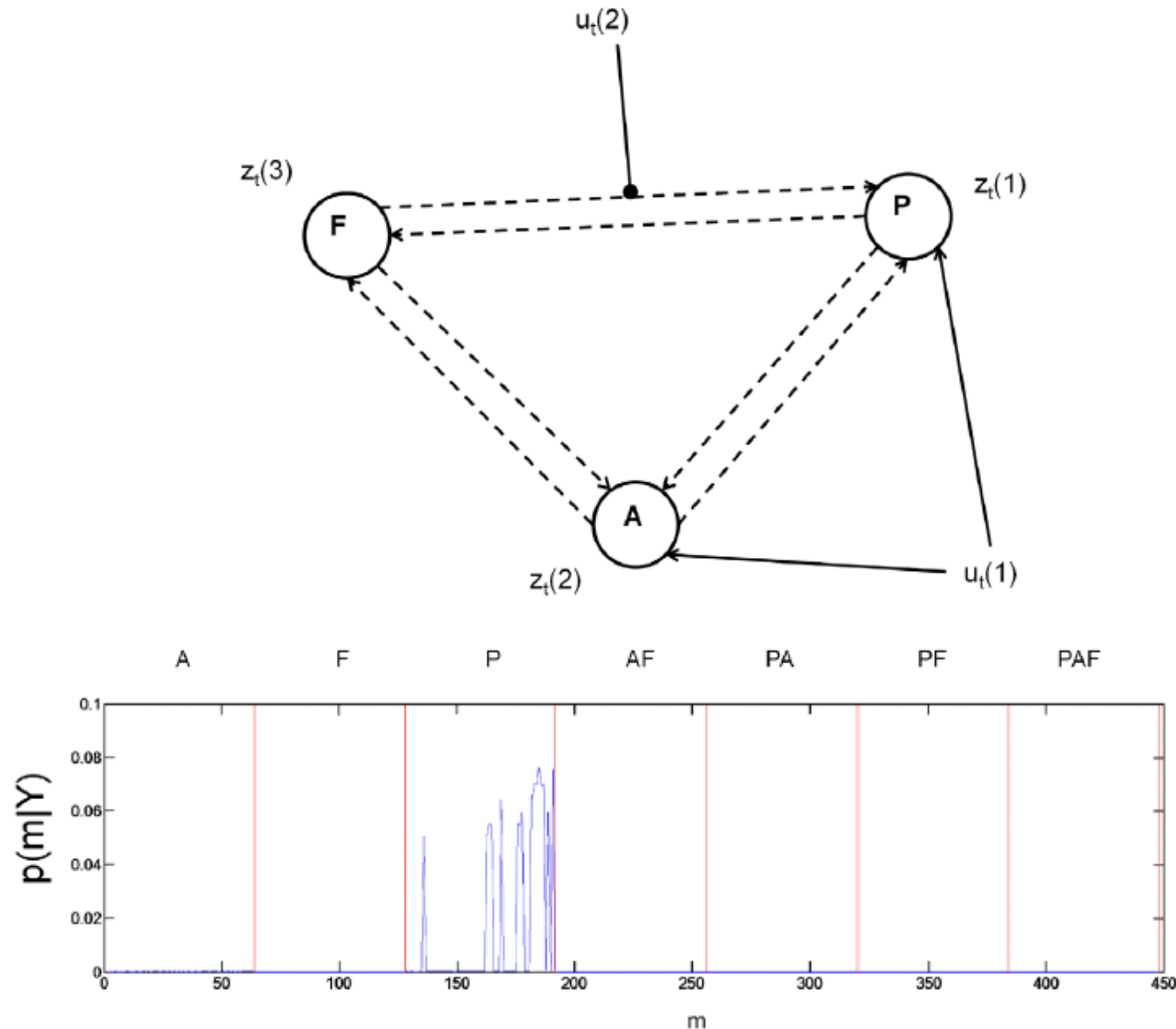


Model  
space  
partitioning



# Comparing model families – a second example

- data from Leff et al. 2008, J. Neurosci
- one driving input, one modulatory input
- $2^6 = 64$  possible modulations
- $2^3 - 1$  input patterns
- $7 \times 64 = 448$  models
- integrate out uncertainty about modulatory patterns and ask where auditory input enters





# Bayesian Model Averaging (BMA)

- abandons dependence of parameter inference on a single model and takes into account model uncertainty
- uses the entire model space considered (or an optimal family of models)
- averages parameter estimates, weighted by posterior model probabilities
- represents a particularly useful alternative
  - when none of the models (or model subspaces) considered clearly outperforms all others
  - when comparing groups for which the optimal model differs

## single-subject BMA:

$$p(\theta | y) \\ = \sum_m p(\theta | y, m) p(m | y)$$

## group-level BMA:

$$p(\theta_n | y_{1..N}) \\ = \sum_m p(\theta_n | y_n, m) p(m | y_{1..N})$$

NB:  $p(m|y_{1..N})$  can be obtained by either FFX or RFX BMS

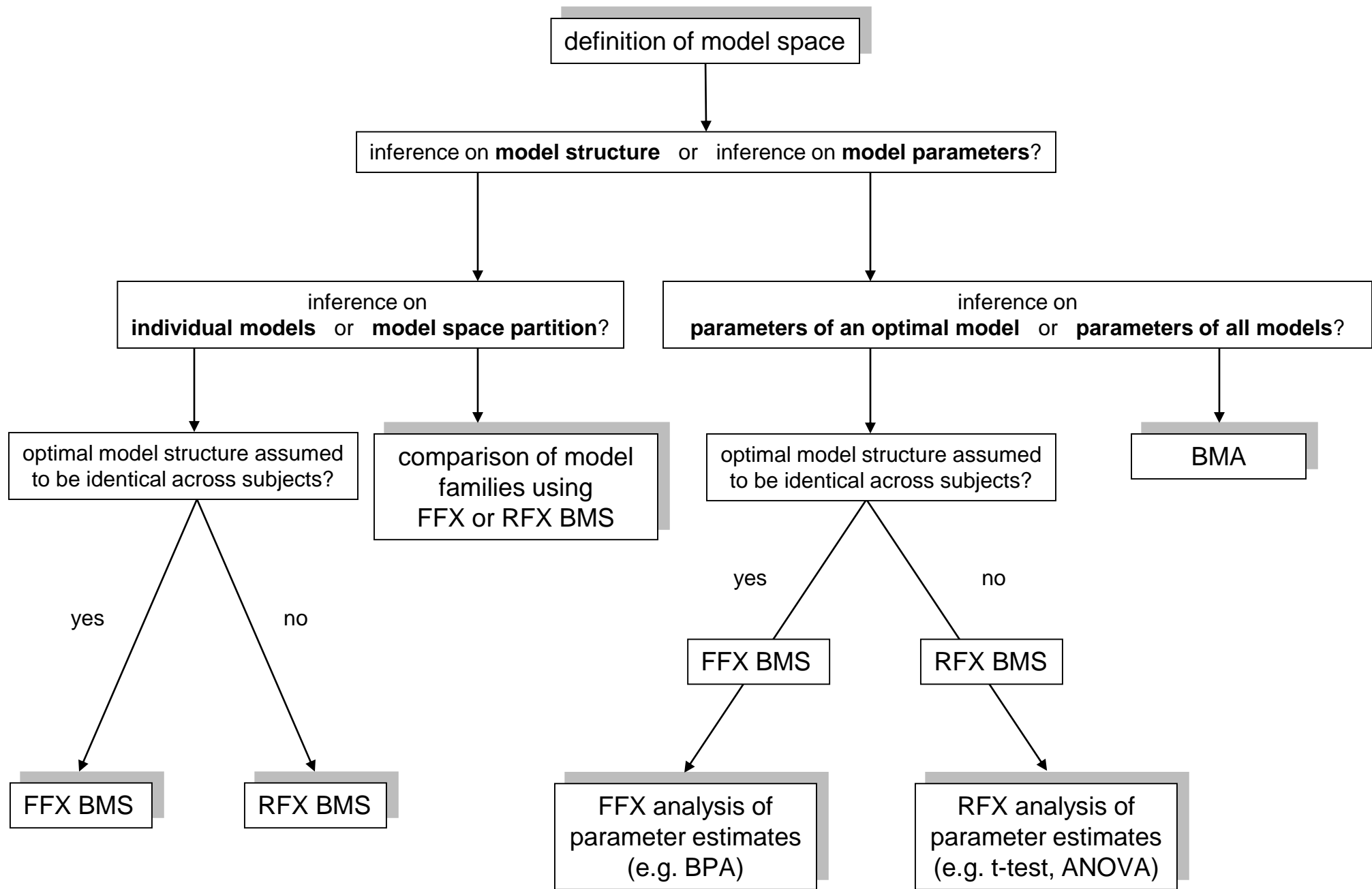
# Protected exceedance probability: Using BMA to protect against chance findings

- EPs express our confidence that the posterior probabilities of models are different – under the hypothesis  $H_1$  that models differ in probability:  $r_k \neq 1/K$
- does not account for possibility "null hypothesis"  $H_0$ :  $r_k = 1/K$
- **Bayesian omnibus risk (BOR)** of wrongly accepting  $H_1$  over  $H_0$ :

$$P_o = \frac{1}{1 + \frac{p(m|H_1)}{p(m|H_0)}}.$$

- **protected EP**: Bayesian model averaging over  $H_0$  and  $H_1$ :

$$\begin{aligned}\tilde{\varphi}_k &= P(r_k \geq r_{k' \neq k} | y) \\ &= P(r_k \geq r_{k' \neq k} | y, H_1)P(H_1 | y) + P(r_k \geq r_{k' \neq k} | y, H_0)P(H_0 | y) \\ &= \varphi_k(1 - P_o) + \frac{1}{K}P_o\end{aligned}$$



# Some examples of empirical BMS/BMA applications

Behavioral/Systems/Cognitive

## Effective Connectivity Determines the Nature of Subjective Experience in Grapheme-Color Synesthesia

Tessa M. van Leeuwen,<sup>1</sup> Hanneke E. M. den Ouden,<sup>1</sup> and Peter Hagoort<sup>1,2</sup>

<sup>1</sup>Donders Institute for Brain, Cognition and Behaviour, Centre for Cognitive Neuroimaging, Radboud University Nijmegen, 6500 HB, Nijmegen, the Netherlands, and <sup>2</sup>Max Planck Institute for Psycholinguistics, 6500 AH, Nijmegen, the Netherlands

van Leeuwen et al. 2011,  
*J. Neurosci.*

doi:10.1093/brain/awv261

BRAIN 2015: Page 1 of 13 | 1

**BRAIN**  
A JOURNAL OF NEUROLOGY

## Network dysfunction of emotional and cognitive processes in those at genetic risk of bipolar disorder

Michael Breakspear,<sup>1,2,3,\*</sup> Gloria Roberts,<sup>3,4,\*</sup> Melissa J. Green,<sup>3,4,5,6</sup> Vinh T. Nguyen,<sup>1</sup> Andrew Frankland,<sup>3,4</sup> Florence Levy,<sup>3</sup> Rhoshel Lenroot<sup>3,6</sup> and Philip B. Mitchell<sup>3,4</sup>

Breakspear et al. 2015,  
*Brain*

Original Investigation

## Brain Connectivity Abnormalities Predating the Onset of Psychosis Correlation With the Effect of Medication

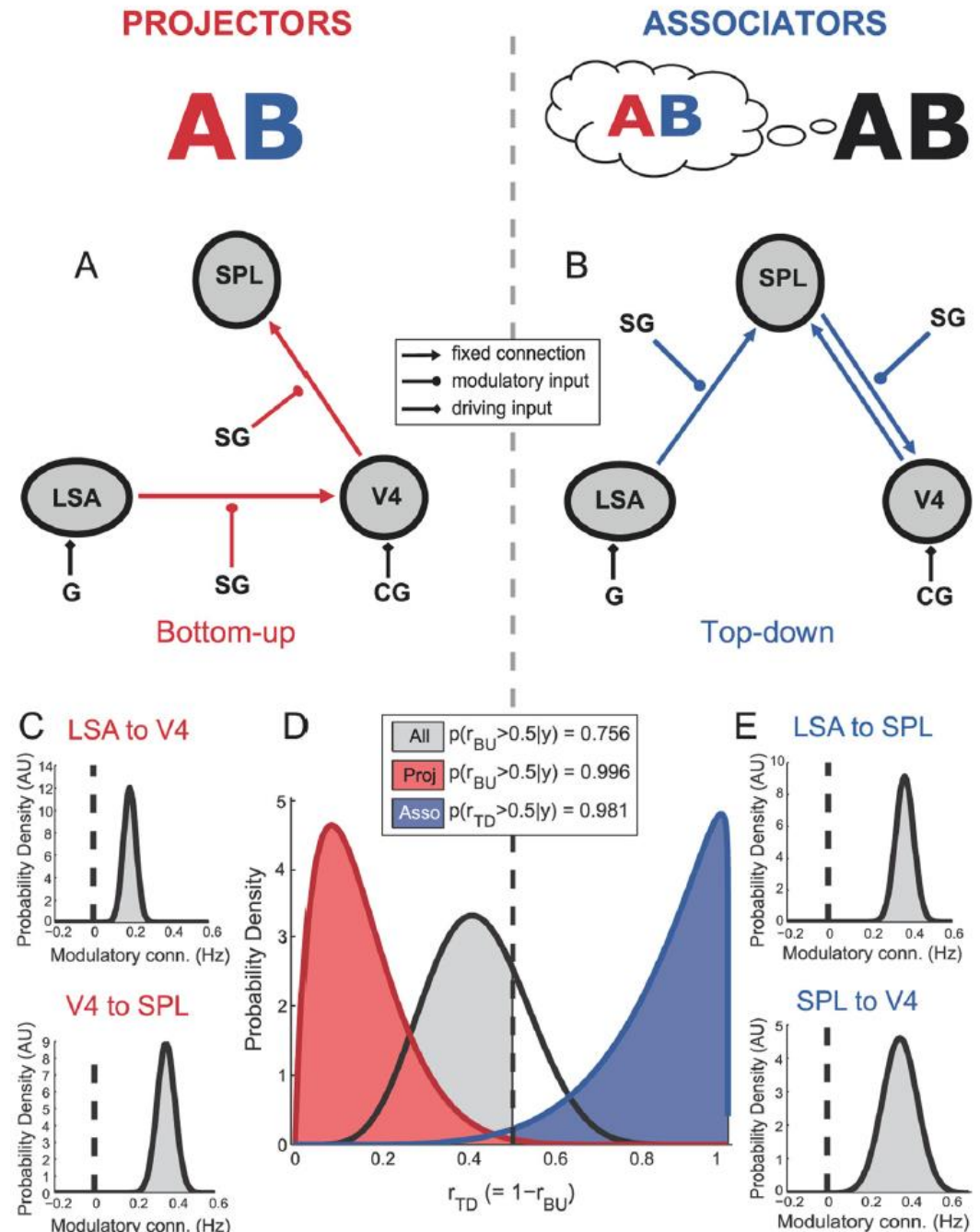
André Schmidt, PhD; Renata Smieskova, PhD; Jacqueline Aston, MD; Andor Simon, MD; Paul Allen, PhD; Paolo Fusar-Poli, MD, PhD; Philip K. McGuire, MD, PhD; Anita Riecher-Rössler, MD, PhD; Klaas E. Stephan, MD, PhD; Stefan Borgwardt, MD, PhD

Schmidt et al. 2013,  
*JAMA Psychiatry*

# Application: Synaesthesia

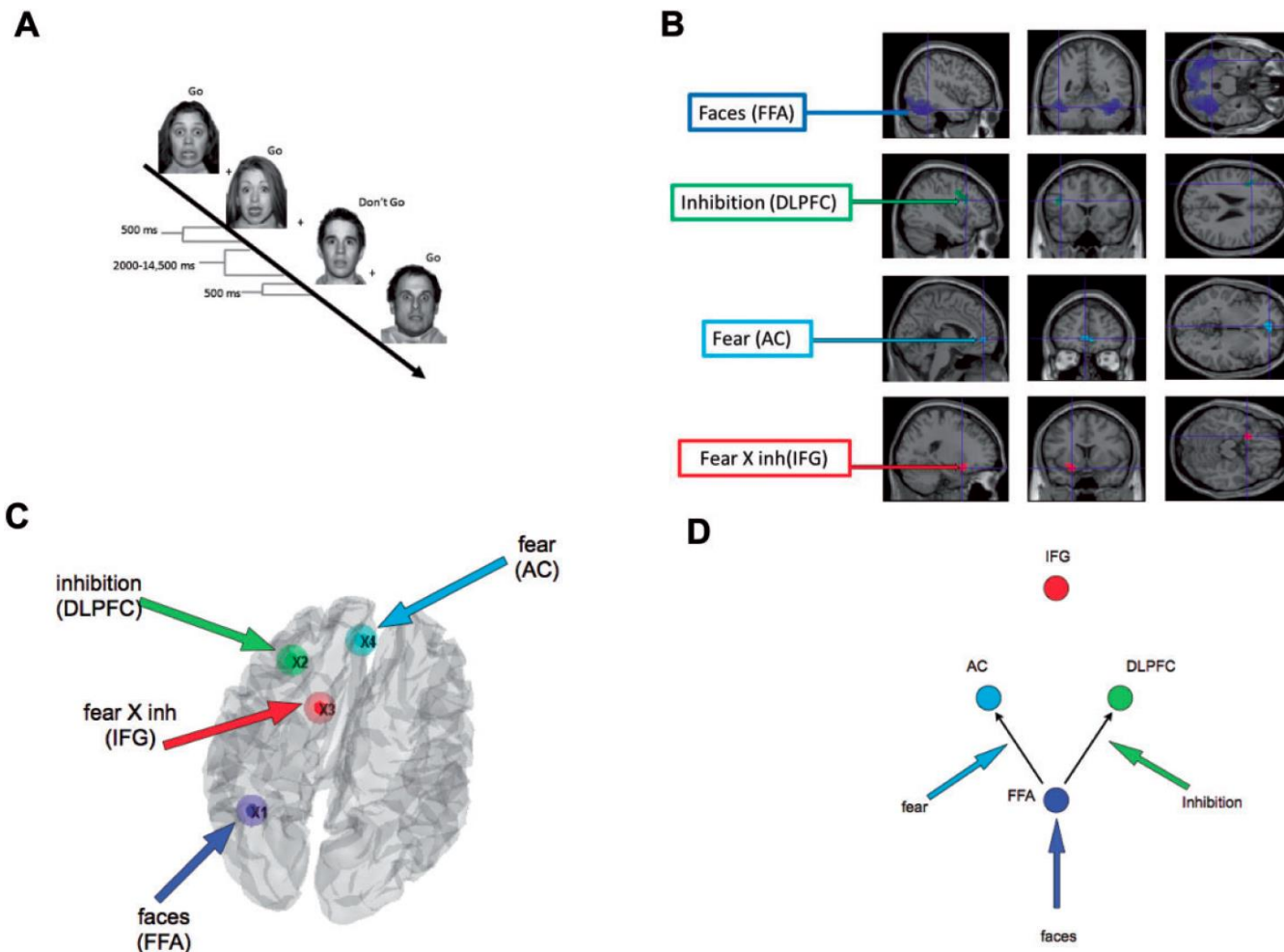
- “projectors” experience color externally colocalized with a presented grapheme
- “associators” report an internally evoked association
- across all subjects: no evidence for either model
- but BMS results map precisely onto projectors (bottom-up mechanisms) and associators (top-down)

van Leeuwen et al. 2011, *J. Neurosci.*



# Go/No-Go task to emotional faces (bipolar patients, at-risk individuals, controls)

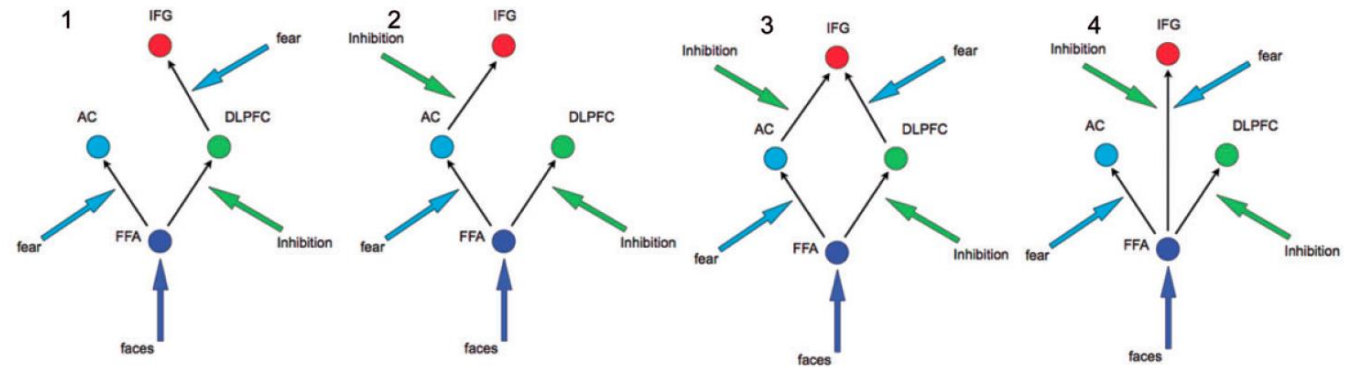
- interaction of motor inhibition and fear perception
- hypoactivation of left IFG in the at-risk group during fearful distractor trials
- What is the most likely circuit mechanism explaining the fear x inhibition interaction in IFG?



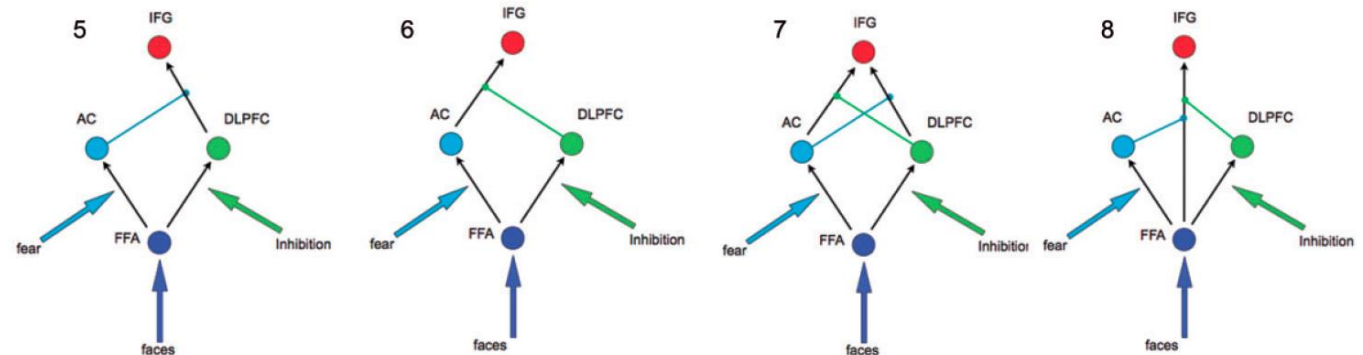
# Model space

- models of serial (1-3), parallel (4) and hierarchical (5-8) processes

**A: Bilinear models**

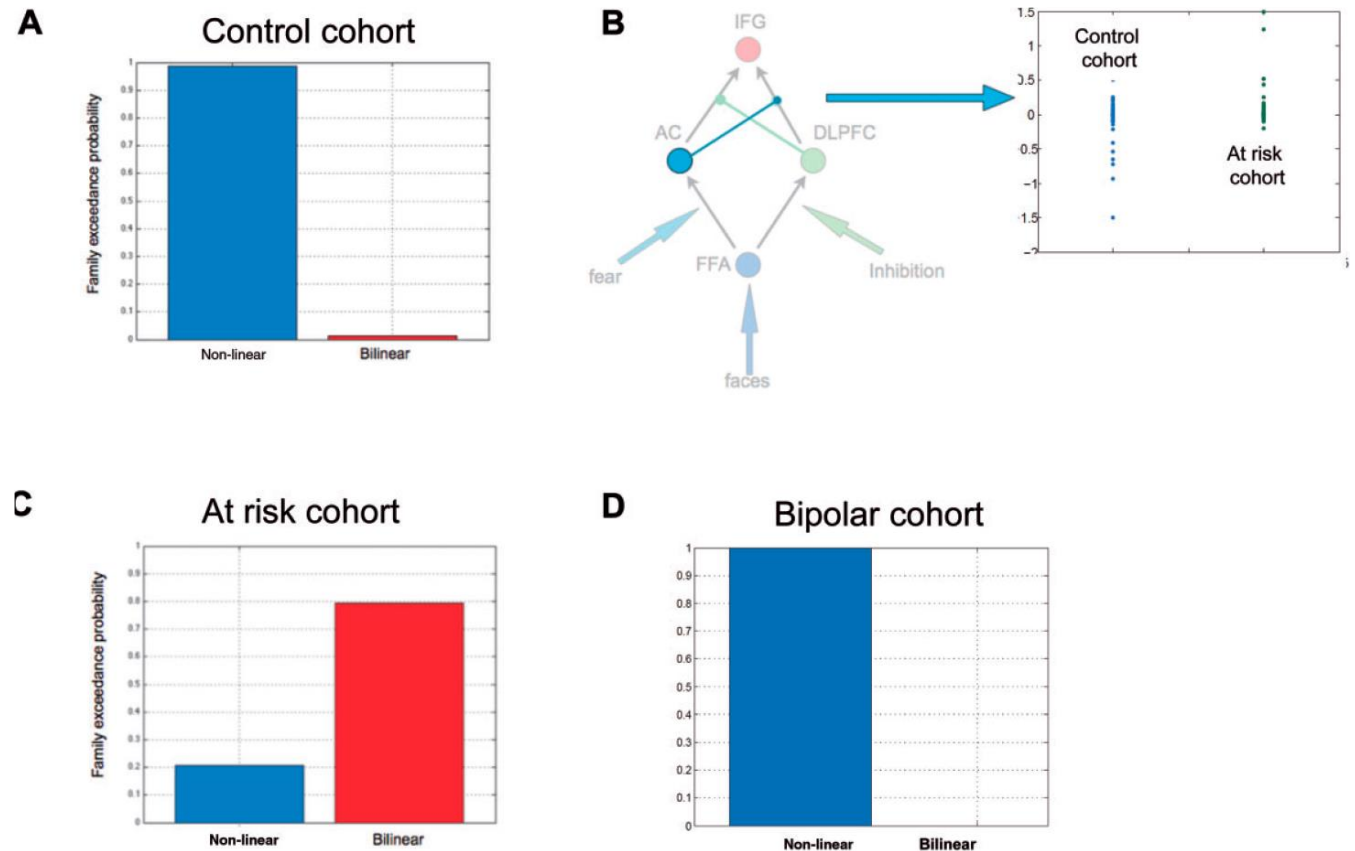


**B: Non-linear models**



# Family-level BMS

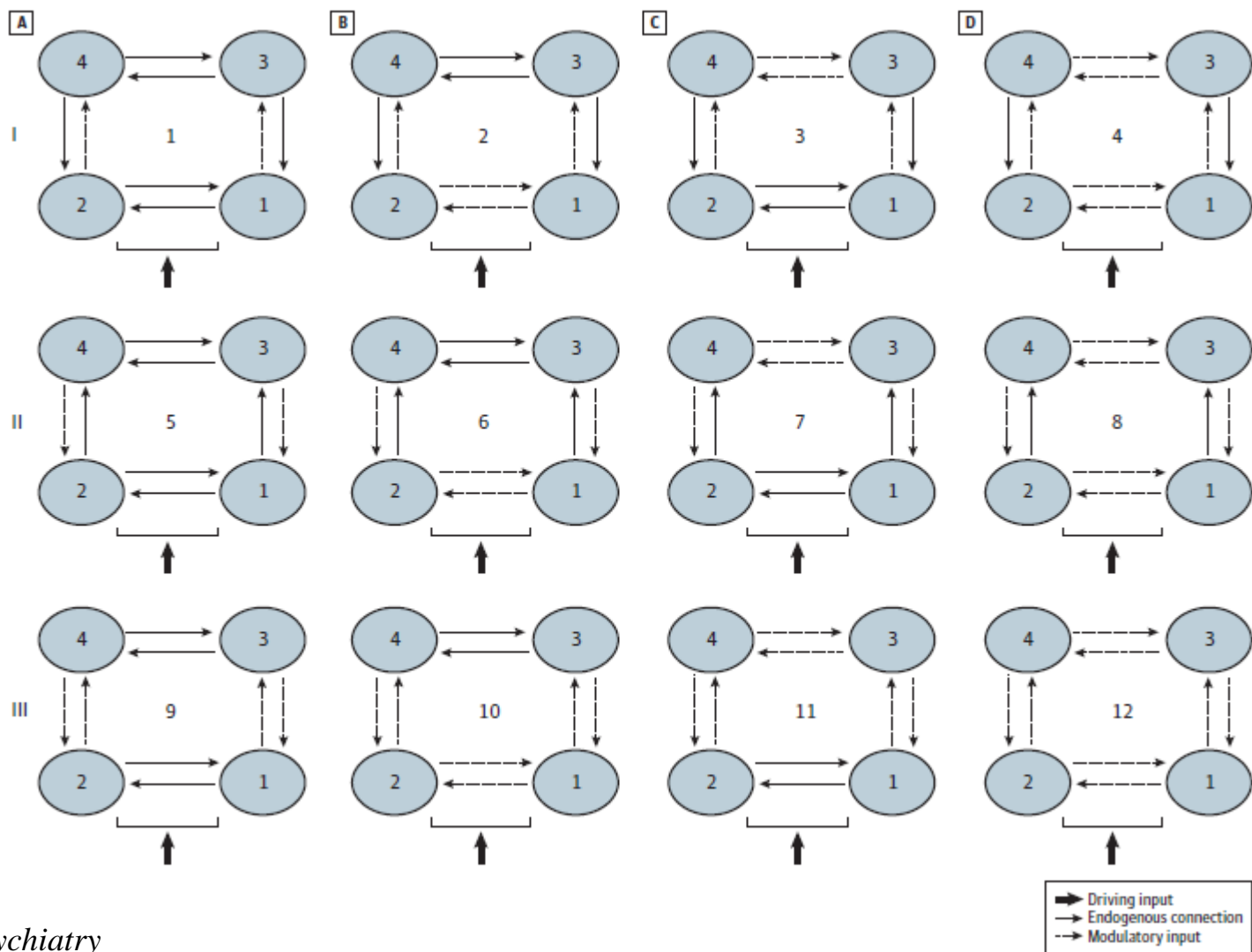
- family-level comparison: nonlinear models more likely than bilinear ones in both healthy controls and bipolar patients
- at-risk group: bilinear models more likely
- significant group difference in ACC modulation of DLPFC→IFG interaction





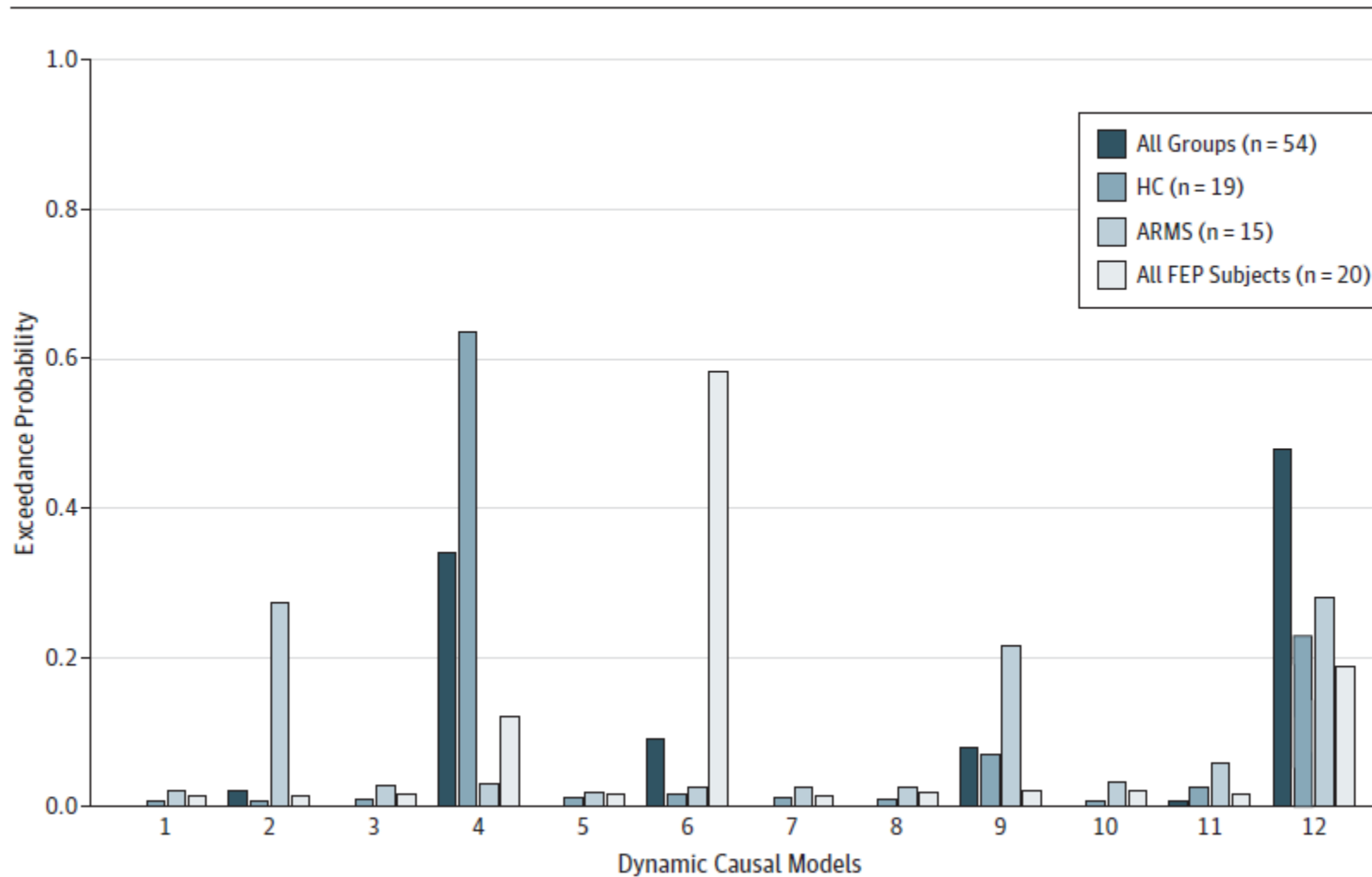


# Prefrontal-parietal connectivity during working memory in schizophrenia

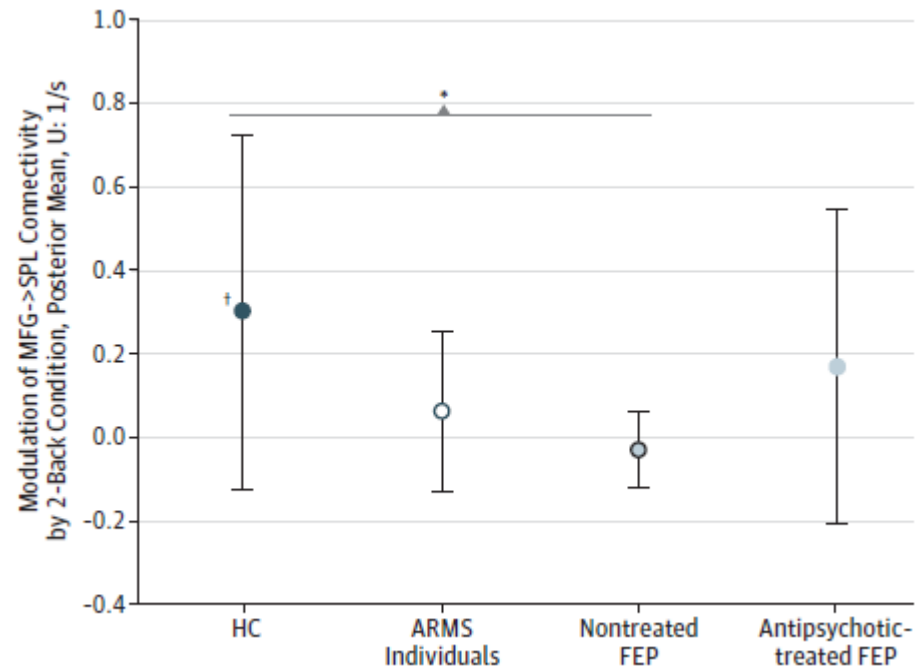
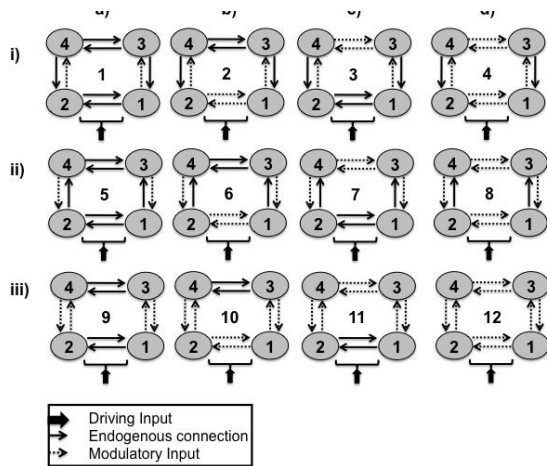
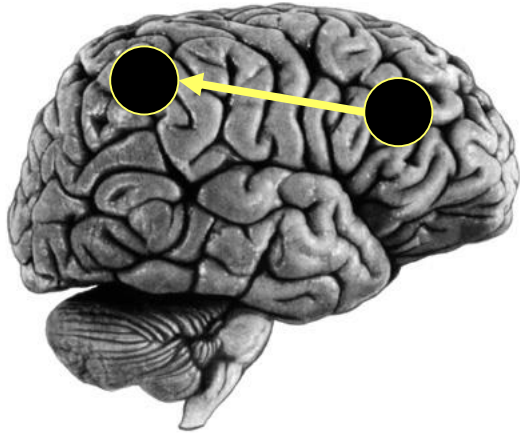


- 17 at-risk mental state (ARMS) individuals
- 21 first-episode patients (13 non-treated)
- 20 controls

# BMS results for all groups



# BMA results: PFC → PPC connectivity



17 ARMS, 21 first-episode (13 non-treated),  
20 controls

## Further reading on BMS

- Penny WD, Stephan KE, Mechelli A, Friston KJ (2004a) Comparing dynamic causal models. *NeuroImage* 22:1157-1172.
- Penny WD, Stephan KE, Daunizeau J, Joao M, Friston K, Schofield T, Leff AP (2010) Comparing Families of Dynamic Causal Models. *PLoS Computational Biology* 6: e1000709.
- Penny WD (2012) Comparing dynamic causal models using AIC, BIC and free energy. *Neuroimage* 59: 319-330.
- Rigoux L, Stephan KE, Friston KJ, Daunizeau J (2014) Bayesian model selection for group studies – revisited. *NeuroImage* 84: 971-985.
- Stephan KE, Weiskopf N, Drysdale PM, Robinson PA, Friston KJ (2007) Comparing hemodynamic models with DCM. *NeuroImage* 38:387-401.
- Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ (2009) Bayesian model selection for group studies. *NeuroImage* 46:1004-1017.
- Stephan KE, Penny WD, Moran RJ, den Ouden HEM, Daunizeau J, Friston KJ (2010) Ten simple rules for Dynamic Causal Modelling. *NeuroImage* 49: 3099-3109.

**Thank you**