# (variational) Bayesian inference

**J. Daunizeau**

*ICM, Paris, France*
*TNU, Zurich, Switzerland*

# Overview of the talk

✓ Introduction to Bayesian inference

✓ The variational approach to approximate Bayesian inference

✓ VBA toolbox

# Overview of the talk

✓ **Introduction to Bayesian inference**

✓ The variational approach to approximate Bayesian inference
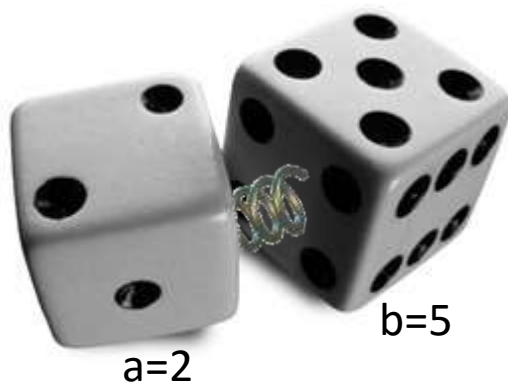
✓ VBA toolbox

# Probability theory: basics

***Degree of*** *plausibility desiderata*:

❑ should be represented using real numbers        (D1)

❑ should conform with intuition        (D2)

❑ should be consistent        (D3)

a=2

→ normalization:

$$\sum_a P(a) = 1$$

a=2    b=5

→ marginalization:

$$P(b) = \sum_a P(a,b)$$

→ conditioning :

(*Bayes rule*)

$$P(a,b) = P(a|b)\,P(b)$$
$$= P(b|a)\,P(a)$$

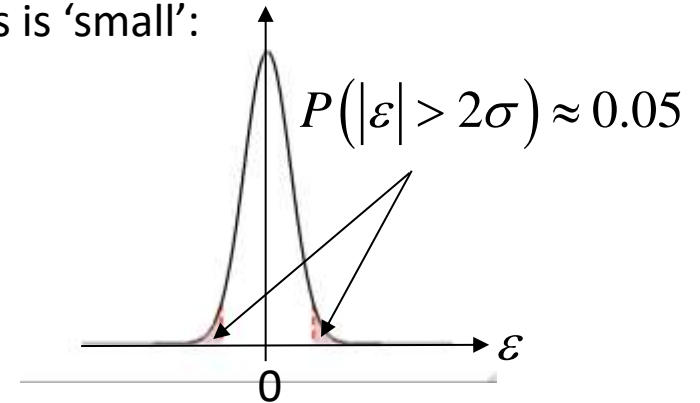# Deriving the likelihood function



$\theta$

$f$

$y$

- Model of data with unknown parameters:

$$y = f(\theta) \qquad \text{e.g., GLM:} \qquad f(\theta) = X\theta$$

- But data is noisy:
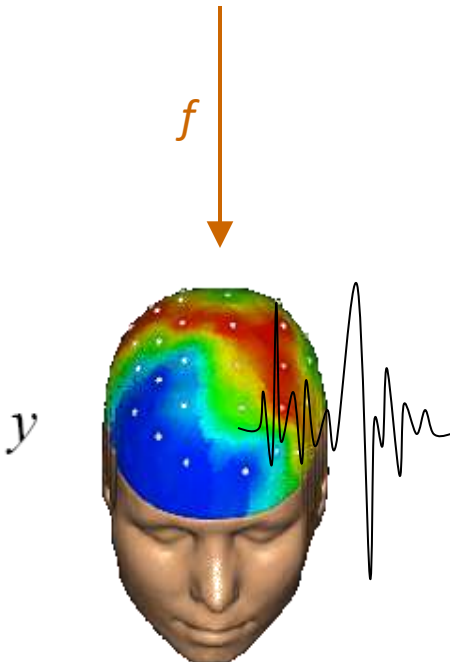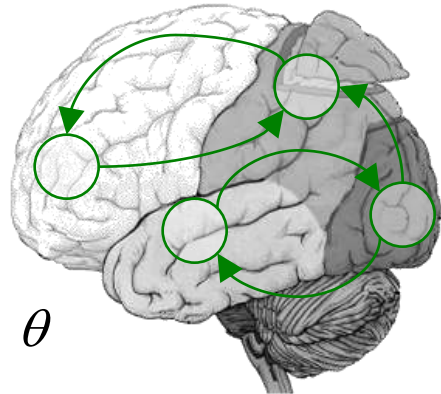
$$y = f(\theta) + \varepsilon$$

- Assume noise/residuals is 'small':

$$P(|\varepsilon| > 2\sigma) \approx 0.05$$

$$p(\varepsilon) \propto \exp\left(-\frac{1}{2\sigma^2}\varepsilon^2\right)$$
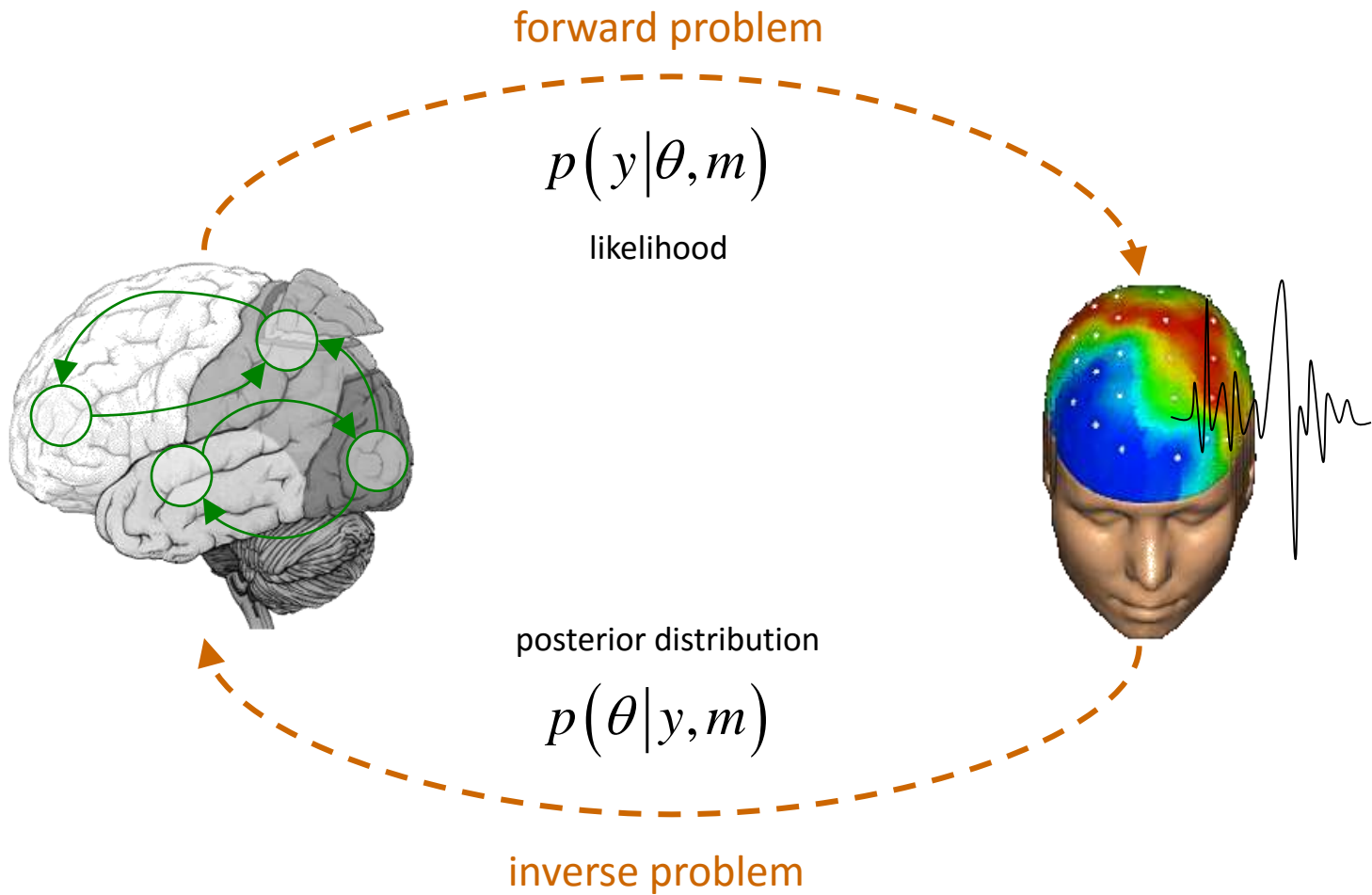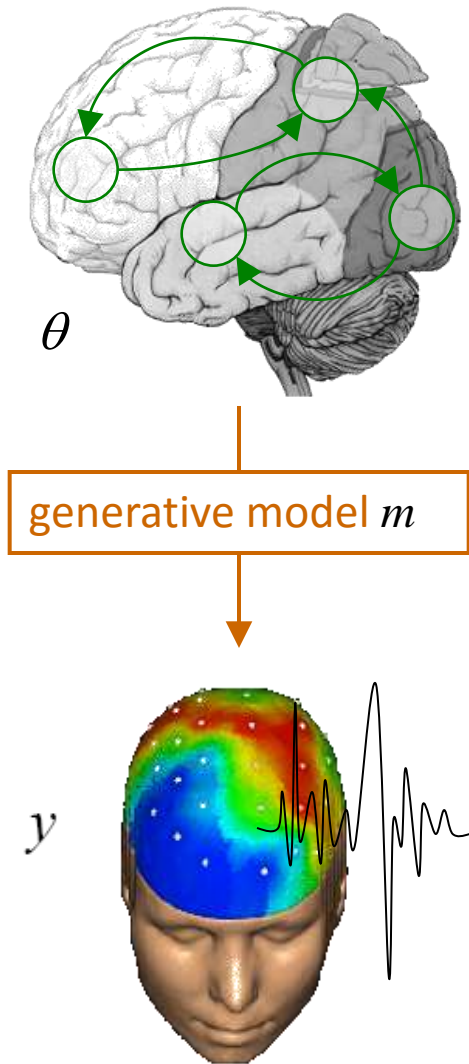


→ Distribution of data, *given fixed parameters*:

$$p(y|\theta) \propto \exp\left(-\frac{1}{2\sigma^2}(y - f(\theta))^2\right)$$
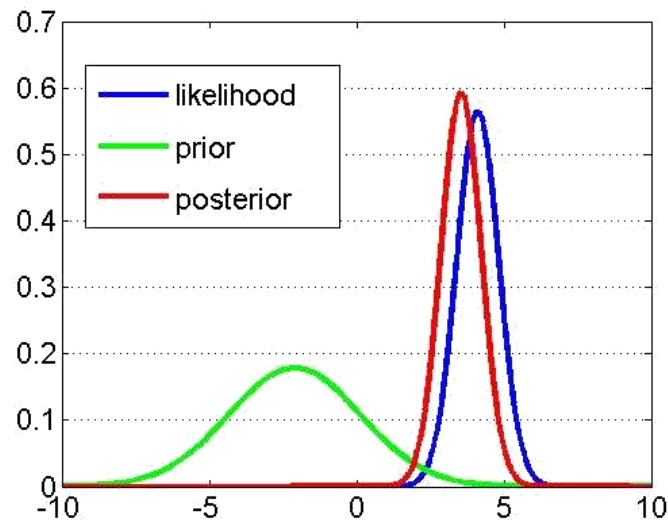
# Probabilistic model inversion



forward problem

$$p\left(y|\theta,m\right)$$

likelihood

posterior distribution

$$p\left(\theta|y,m\right)$$

inverse problem

# Posterior inference on model parameters



Likelihood: $p(y|\theta, m)$
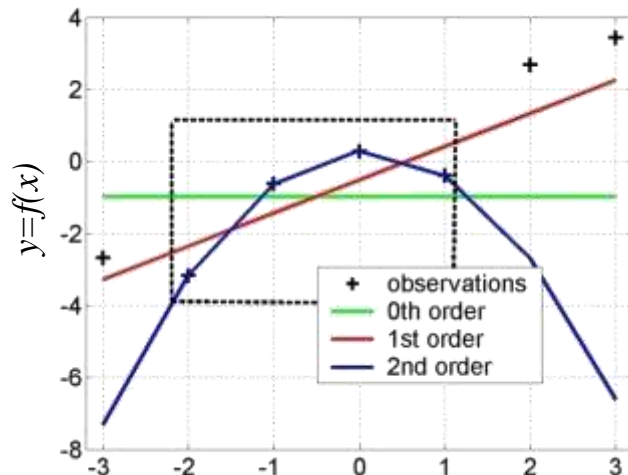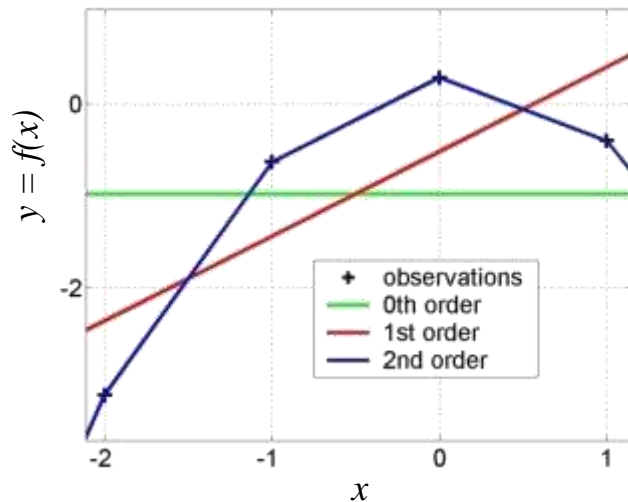
Prior: $p(\theta|m)$

Bayes rule: $p(\theta|y, m) = \dfrac{p(y|\theta, m)\, p(\theta|m)}{p(y|m)}$
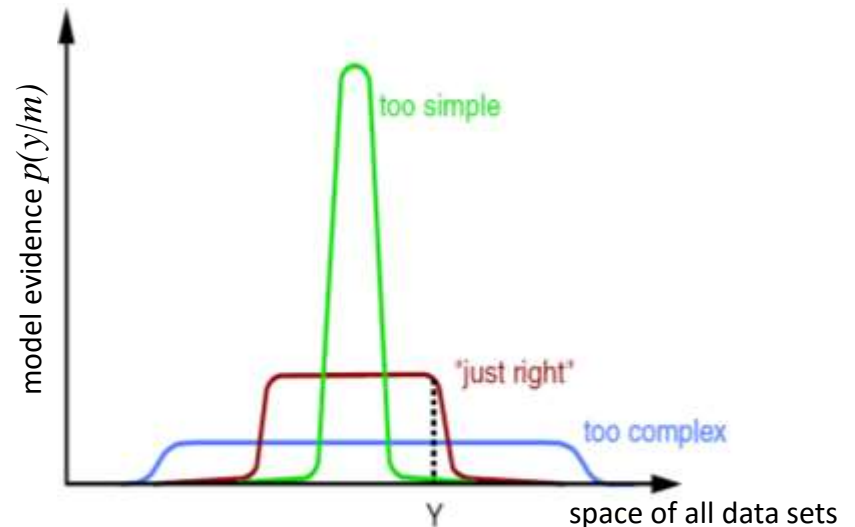
# Bayesian model comparison

*Principle of parsimony* :
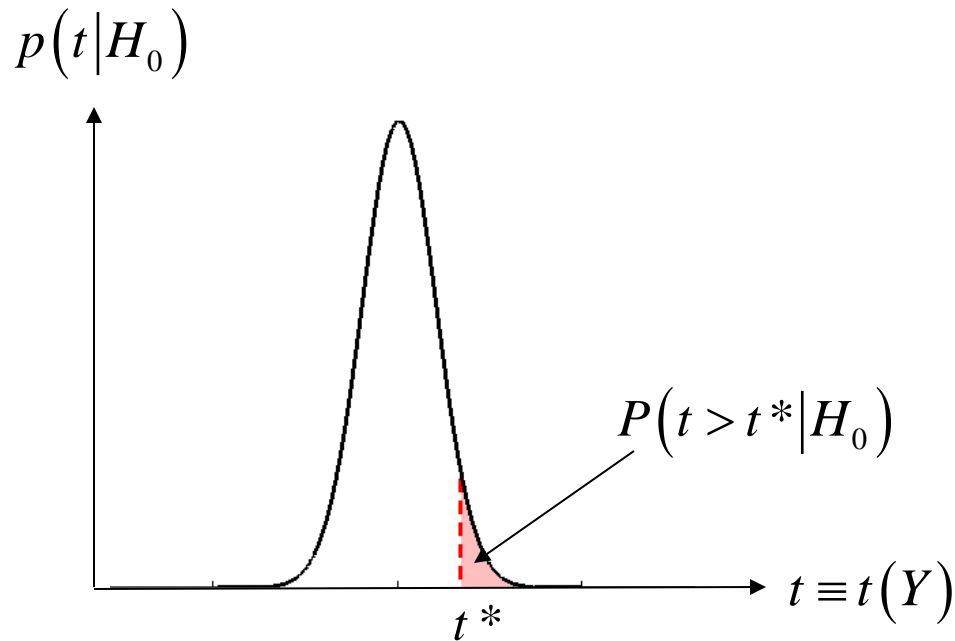« plurality should not be assumed without necessity »



Model evidence:

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

"Occam's razor" :

# Bayesian versus frequentist hypothesis testing

• define the null, e.g.: $H_0 : \theta = 0$
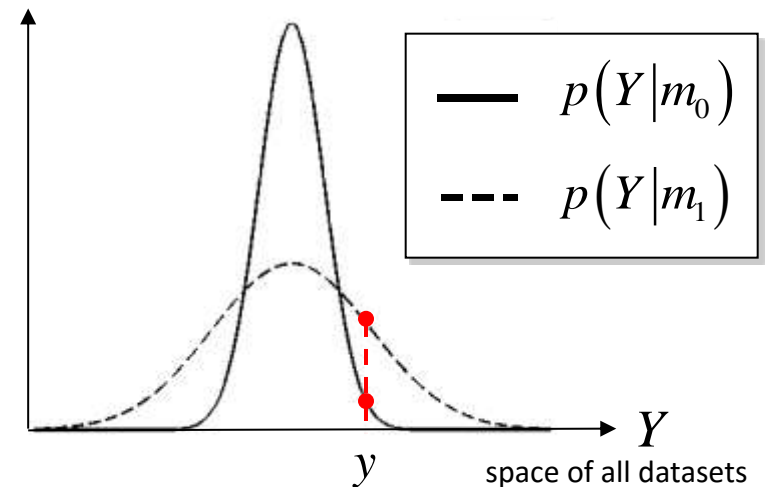
$p(t|H_0)$



$P(t > t^*|H_0)$

$t \equiv t(Y)$

$t^*$

• estimate parameters (obtain test stat.)

• apply decision rule, i.e.:

if $P(t > t^*|H_0) \leq \alpha$ then reject H0

classical (null) hypothesis testing

• define two alternative models, e.g.:

$$m_0 : p(\theta|m_0) = \begin{cases} 1 & \text{if } \theta = 0 \\ 0 & \text{otherwise} \end{cases}$$
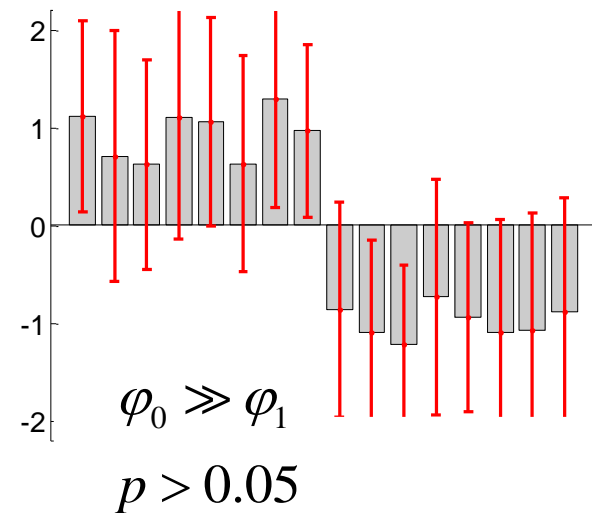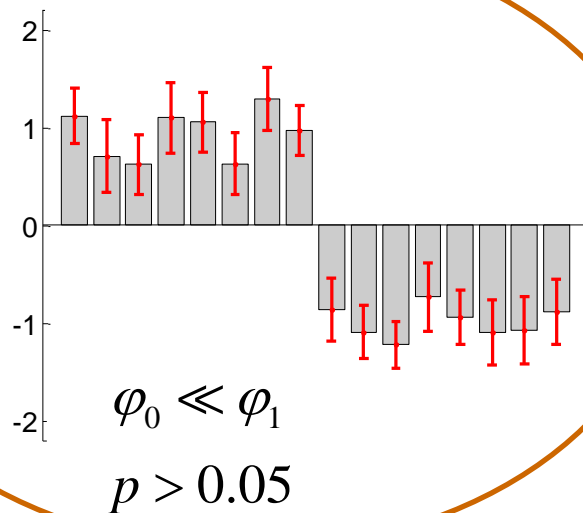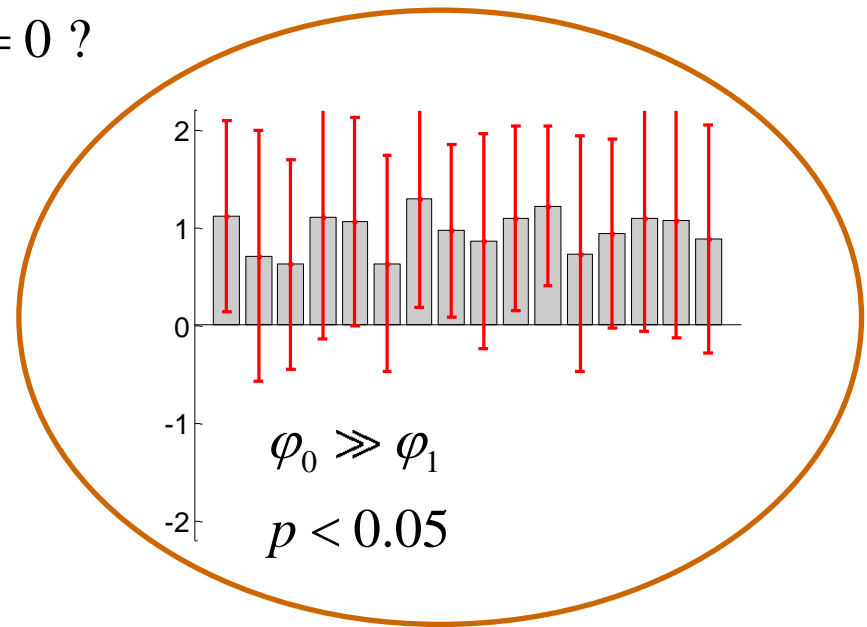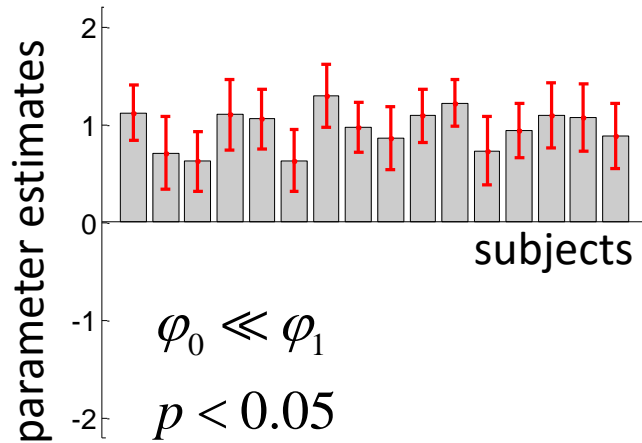
$$m_1 : p(\theta|m_1) = N(0, \Sigma)$$



$p(Y|m_0)$

$p(Y|m_1)$

$y$    space of all datasets

$Y$

• apply decision rule, e.g.:

if $\dfrac{P(m_0|y)}{P(m_1|y)} \geq \alpha$ then accept $m_0$

Bayesian model comparison

# Group-level model selection

# Family-level inference
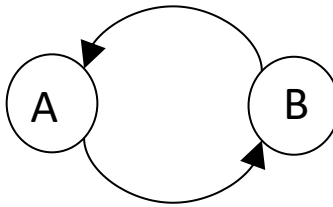
model selection error risk:

P(m$_1$|y) = 0.04



P(m$_2$|y) = 0.25



$$P\left(e=1\middle|y\right)=1-\max_{m}P\left(m\middle|y\right)$$

$$=0.3$$

P(m$_2$|y) = 0.01



u

P(m$_2$|y) = 0.7



u

# Family-level inference



$$P(e=1|y) = 1 - \max_m P(m|y)$$

$$= 0.3$$

**family inference**
(pool statistical evidence)

$$P(f|y) = \sum_{m \in f} P(m|y)$$

$$P(e=1|y) = 1 - \max_f P(f|y)$$

$$= 0.05$$

# Priors and the bias-variance trade-off



correct model

wrong model

# Type, role and impact of priors

- Types of priors:

  - ✓ Explicit priors on *model parameters* (e.g., Gaussian)

  - ✓ Implicit priors on *model functional form* (e.g., evolution & observation functions)

  - ✓ Choice of "interesting" *data features* (e.g., response magnitude vs response profile)
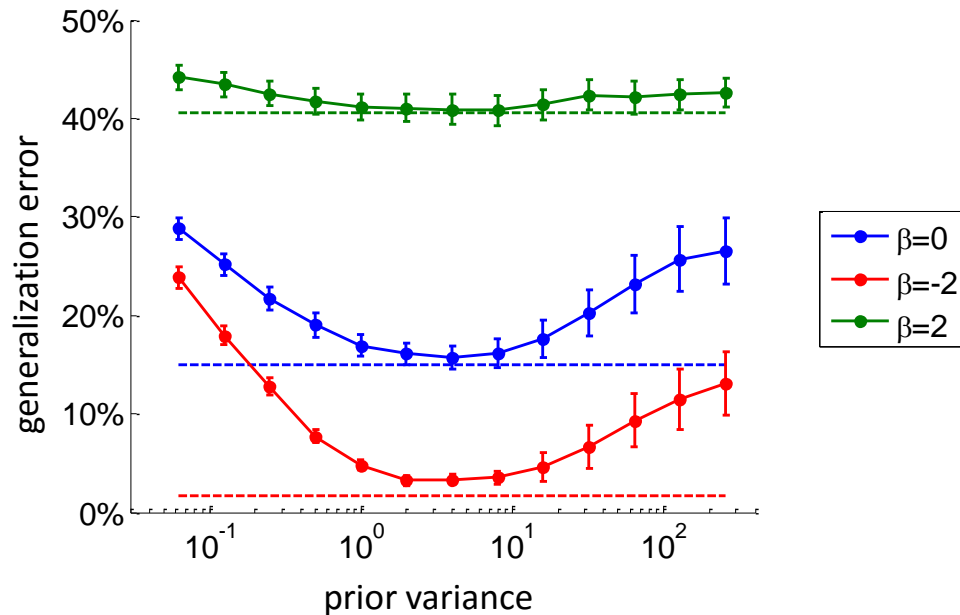
- Role of explicit priors (on model parameters):

  - ✓ Resolving the *ill-posedness* of the inverse problem

  - ✓ Avoiding *overfitting* (cf. generalization error)

- Impact of priors:

  - ✓ On parameter posterior distributions (cf. "shrinkage to the mean" effect)

  - ✓ On model evidence (cf. "Occam's razor")

# Overview of the talk

✓ Introduction to Bayesian inference

✓ The variational approach to approximate Bayesian inference

✓ VBA toolbox

# Why do we need approximations?



$p\left(\theta_1, \theta_2 \middle| y, m\right)$

$p\left(\theta_{1 \text{ or } 2} \middle| y, m\right)$

$q\left(\theta_{1 \text{ or } 2}\right)$

# The Laplace approximation

$$t(\theta) = \ln p(y|\theta, m) + \ln p(\theta|m)$$

$$\approx t(\hat{\theta}) + \underbrace{(\theta - \hat{\theta})^{\mathrm{T}} \left. \frac{\partial t}{\partial \theta} \right|_{\hat{\theta}}}_{0} + \frac{1}{2} (\theta - \hat{\theta})^{\mathrm{T}} \underbrace{\left. \frac{\partial^2 t}{\partial \theta^2} \right|_{\hat{\theta}}}_{-H(\hat{\theta})} (\theta - \hat{\theta})$$

$$\ln p(y|m) = \ln \int \exp(t(\theta)) d\theta$$

$$\approx \underbrace{t(\hat{\theta}) + \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln \left| H(\hat{\theta}) \right|}_{F_{\mathrm{Laplace}}}$$
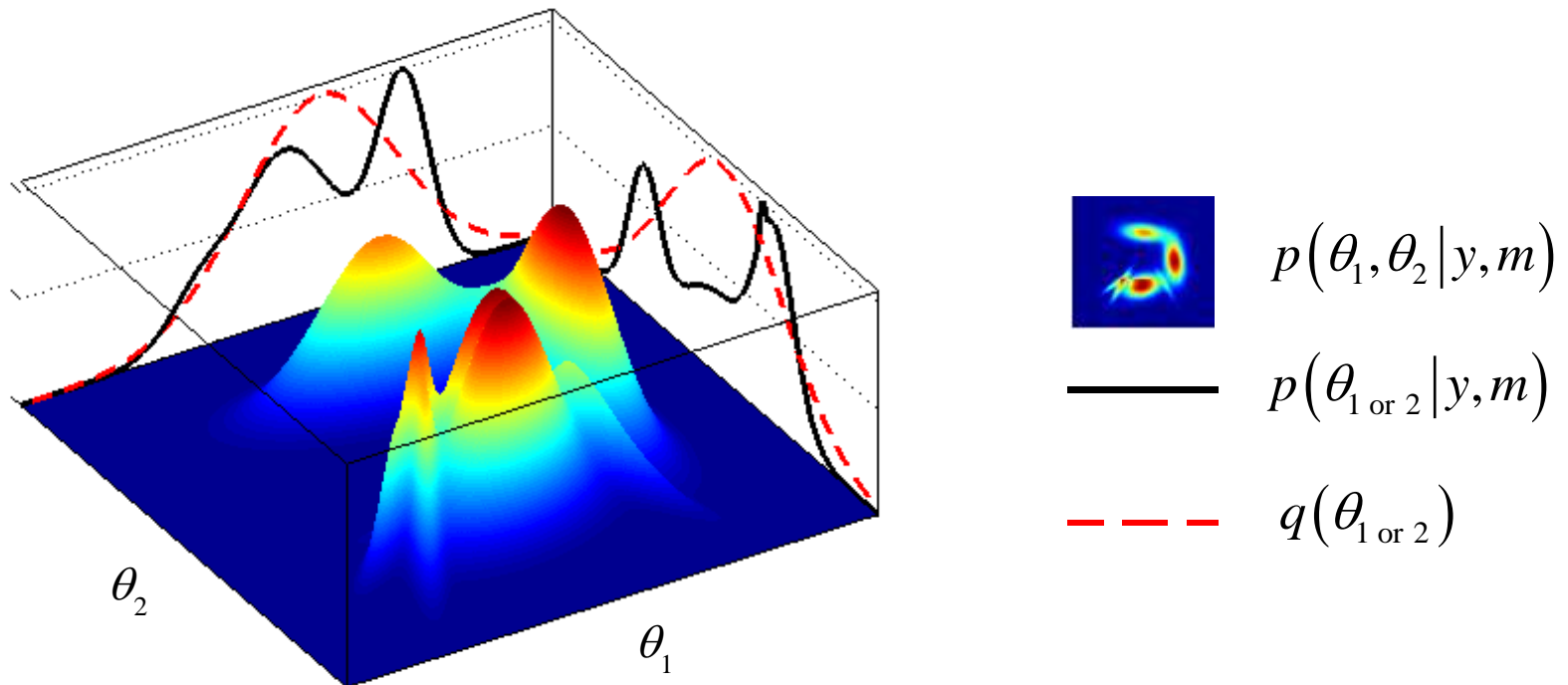
# The Free energy lower bound

$$F = \left\langle \ln p\left(y|\theta,m\right) + \ln p\left(\theta|m\right) \right\rangle_q + S\left(q\right)$$

$$= \left\langle \ln p\left(y|\theta,m\right) \right\rangle_q - KL\left( p\left(\theta|m\right) \; ; q\left(\theta\right) \right)$$

$$= \ln p\left(y|m\right) - KL\left( p\left(\theta|y,m\right) \; ; q\left(\theta\right) \right)$$

# VB and the Free Energy

$$\ln p(y \mid m) = F(q) + KL\big(p(\theta \mid y, m); q(\theta)\big)$$

→ **VB** :    maximize the free energy $F(q)$ w.r.t. the approximate posterior $q(\theta)$

under some (e.g., *mean field, Laplace*) simplifying constraint



$p(\theta_1, \theta_2 \mid y, m)$

$p(\theta_{1\ or\ 2} \mid y, m)$

$q(\theta_{1\ or\ 2})$

# The mean-field approximation

$$F = \left\langle \ln p\left(y|\theta,m\right) + \ln p\left(\theta|m\right) \right\rangle_q + S\left(q\right)$$

$$q\left(\theta\right) \approx q_1\left(\theta_1\right) q_2\left(\theta_2\right)$$

$$\frac{\delta F}{\delta q_2} = 0 \Rightarrow q_2\left(\theta_2\right) \propto \exp\left\langle \ln p\left(y|\theta,m\right) + \ln p\left(\theta|m\right) \right\rangle_{q_1}$$

# The frequentist limit to the model evidence

$$F = \left\langle \ln p\left(y\middle|\theta, m\right) + \ln p\left(\theta\middle|m\right) \right\rangle_q + S\left(q\right)$$

$$\xrightarrow[\text{flat priors}]{p(\theta)\to 1} \left\langle \ln p\left(y\middle|\theta, m\right) \right\rangle_q + S\left(q\right)$$

$$\xrightarrow[\substack{\text{point mass}\\\text{approximation}}]{q(\theta)\to\delta(\hat{\theta})} \underbrace{\ln p\left(y\middle|\hat{\theta}, m\right)}_{\text{frequentist log-likelihood}}$$

# BIC and AIC

→ BIC: Laplace approximation at the asymptotic limit

$$\Sigma \xrightarrow{\;n\to\infty\;} \frac{1}{n} I_p$$

$$F_{\text{Laplace}} \xrightarrow{\;n\to\infty\;} \underbrace{\ln p\left(y\middle|\hat{\theta},m\right) - \frac{p}{2}\ln n}_{\text{BIC}}$$
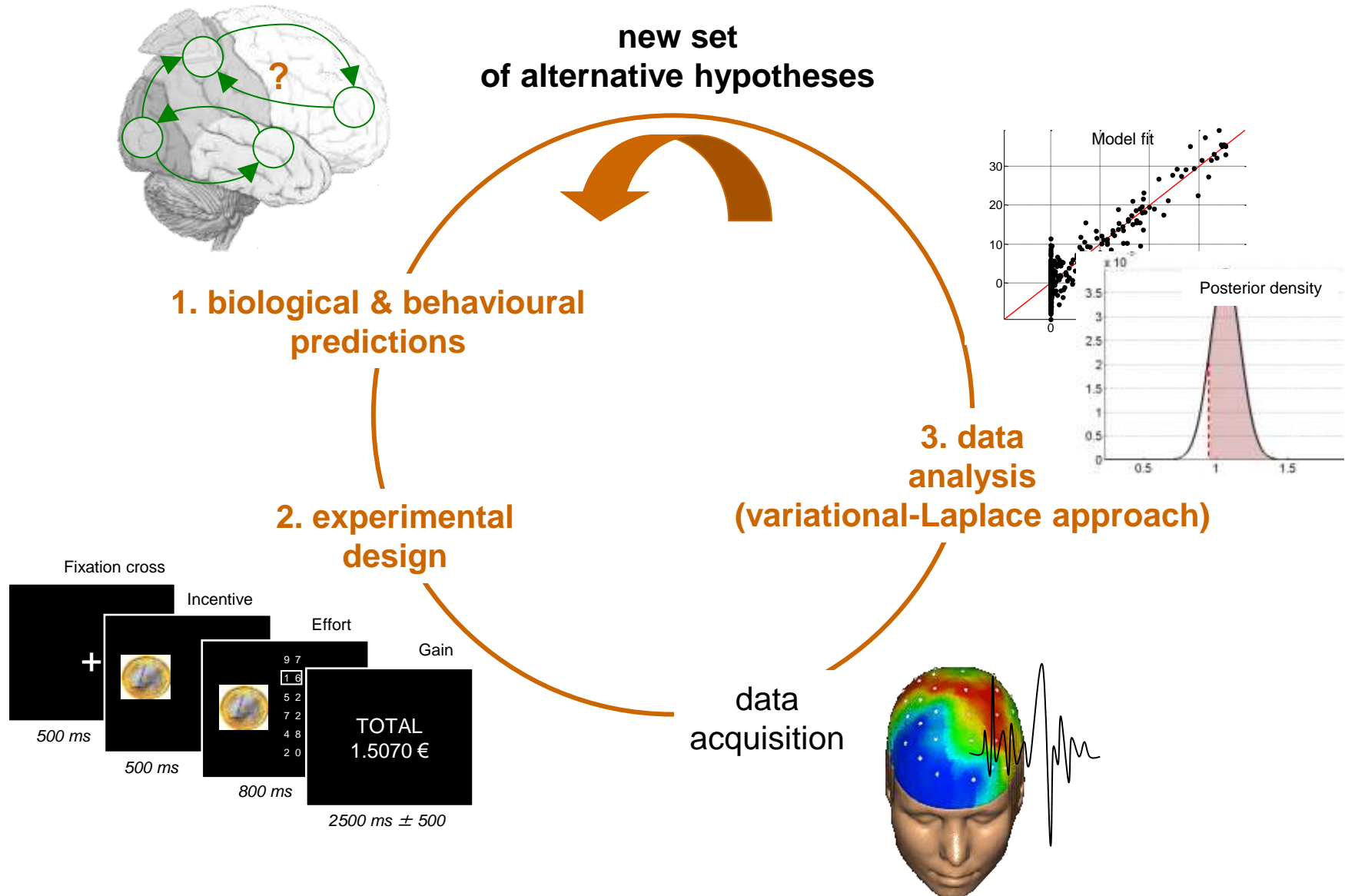
→ AIC: approximation to a frequentist KL-divergence risk!

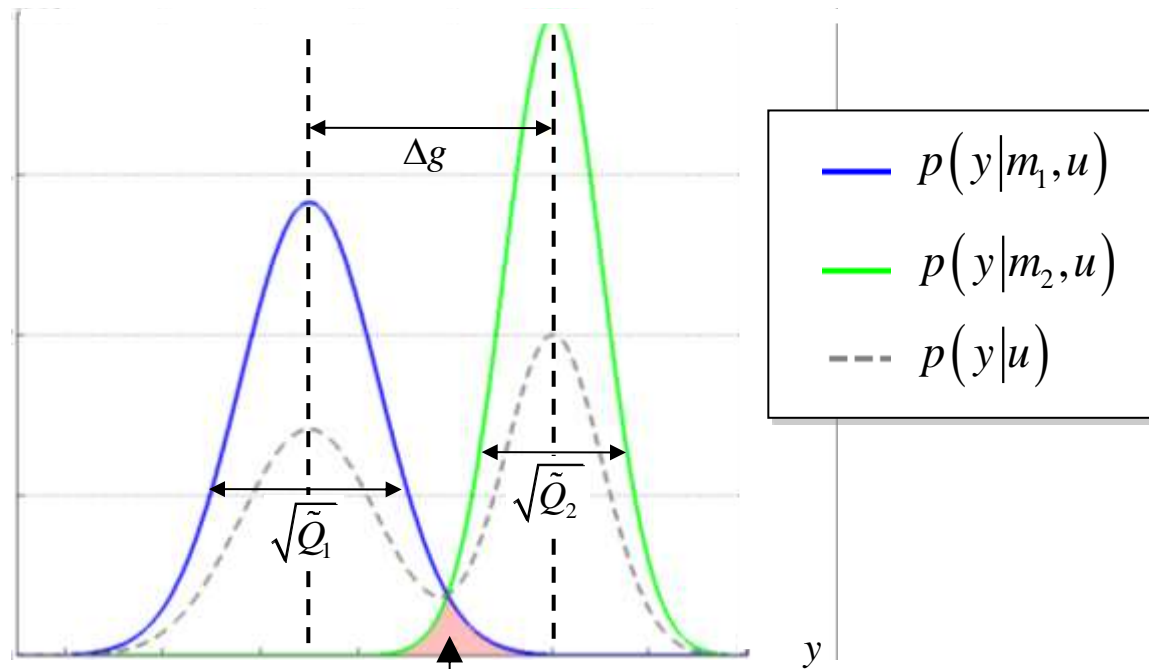$$AIC = \ln p\left(y\middle|\hat{\theta},m\right) - p$$

# Overview of the talk

✓ Introduction: Bayesian inference

✓ The variational approach to approximate Bayesian inference

✓ VBA toolbox

# The three core mathematical problems of VBA



new set
of alternative hypotheses

1. biological & behavioural
predictions

2. experimental
design

3. data
analysis
(variational-Laplace approach)

data
acquisition

Model fit

Posterior density

Fixation cross

Incentive

Effort

Gain

TOTAL
1.5070 €

500 ms

500 ms

800 ms

2500 ms ± 500

# Selection error rate and the Laplace-Chernoff risk

$$b_{LC}(u) = 1 - \frac{1}{2}\log\left(\frac{\Delta g(u)^2}{4\tilde{Q}(u)} + 1\right) \quad \text{if} \quad \tilde{Q}_1(u) \approx \tilde{Q}_2(u) \equiv \tilde{Q}(u)$$



$$p(\hat{e} = 1|u) = 1 - \int_Y \max_m\left[p(m)\,p(y|m,u)\right]dy$$

# VBA: model structure



- evolution parameters
- stochastic innovations variance

$\theta, \alpha$

inputs

$u_t$

- observation parameters
- measurement noise variance

initial conditions

$x_0$

$\varphi, \sigma$

$x_{t-1}$  $f$  $x_t$

system's states

$g$

observations  $y_t$

$t = 1, ..., T$

Daunizeau, Friston et al., Physica D, 2009

# VBA: how to?

✓ You need to provide:

- The data

- The system's inputs (can be left blank)

- The observation and evolution functions (the latter can be left blank)

- The model dimensions

✓ You can specify (otherwise VBA uses defaults):

- The priors (mean and covariances)

- Inversion options  (free-energy tolerance, display flag, etc…)

# Model inversion diagnostics (I)
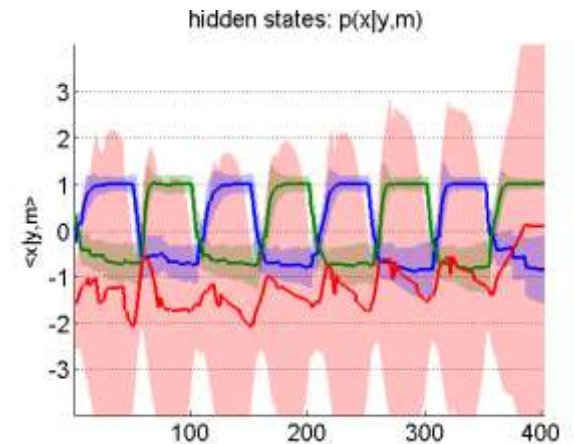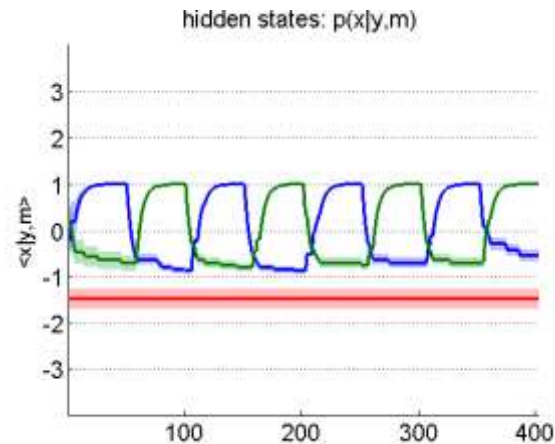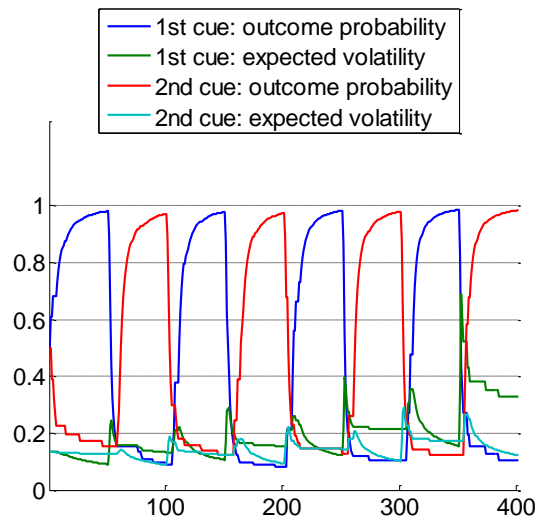
# Model inversion diagnostics (II)

# Model improvement

# Statistical features

- ✓ Mediation analysis (Sobel test, Monte-Carlo sampling)

- ✓ Cross-validation (balanced accuracy, PRESS, etc…)

- ✓ Missing data handling

- ✓ Sparse estimation (approximate lasso)

- ✓ Autoregressive and/or state-dependent noise

- ✓ Arbitrary likelihood functions (e.g., reaction time distributions…)

- ✓ Mixed-effects modelling (empirical priors based upon group statistics)

- ✓ Clustering and Dirichlet processes

- ✓ On-line adaptive experimental designs

- ✓ …

# VBA's (growing) library of models

- ✓ Learning models (Q-learning, hierarchical Bayesian learning, etc...)

- ✓ Decision models (delay/risk/effort discounting, etc...)

- ✓ Neural systems models (DCM, Hodgkin-Huxley, Fitz-Hugh-Nagumo, etc...)

- ✓ Dynamical systems models (double-well, Lorenz attractor, Henon map, etc...)

- ✓ ...

- ✓ Recursive Theory of Mind

- ✓ Spiking models for calcium imaging

- ✓ Race models and other variants of diffusion-drift models

- ✓ ...

# Acknowledgements