

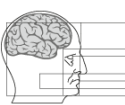
Inference in a hierarchically structured, dynamic environment

Christoph Mathys

Wellcome Trust Centre for Neuroimaging at UCL, London, UK

Max Planck UCL Centre for Computational Psychiatry and Ageing Research, London, UK

Glencoe, June 1, 2015



Uncertainty: a shamelessly artificial example

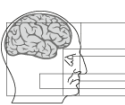
Imagine the following situation:

You're on a boat, you're lost in a storm and trying to get back to shore. A lighthouse has just appeared on the horizon, but you can only see it when you're at the peak of a wave. Your GPS etc., has all been washed overboard, but what you can still do to get an idea of your position is to measure the angle between north and the lighthouse. These are your measurements (in degrees):

76, 73, 75, 72, 77

What number are you going to base your calculation on?

Right. The mean: 74.6. How do you calculate that?



Updates to the mean

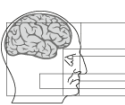
The usual way to calculate the mean \bar{x} of x_1, x_2, \dots, x_n is to take

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

This requires you to remember all x_i , which can become inefficient. Since the measurements arrive sequentially, we would like to update \bar{x} sequentially as the x_i come in – without having to remember them.

It turns out that this is possible. After some algebra (see next slide), we get

$$\bar{x}_{n+1} = \bar{x}_n + \frac{1}{n+1} (x_{n+1} - \bar{x}_n)$$



Updates to the mean

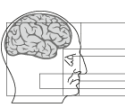
Proof of sequential update formula:

$$\bar{x}_{n+1} = \frac{1}{n+1} \sum_{i=1}^{n+1} x_i = \frac{x_{n+1}}{n+1} + \frac{1}{n+1} \sum_{i=1}^n x_i = \frac{x_{n+1}}{n+1} + \frac{n}{n+1} \underbrace{\frac{1}{n} \sum_{i=1}^n x_i}_{=\bar{x}_n} =$$

$$= \frac{x_{n+1}}{n+1} + \frac{n}{n+1} \bar{x}_n = \bar{x}_n + \frac{x_{n+1}}{n+1} + \frac{n}{n+1} \bar{x}_n - \frac{n+1}{n+1} \bar{x}_n =$$

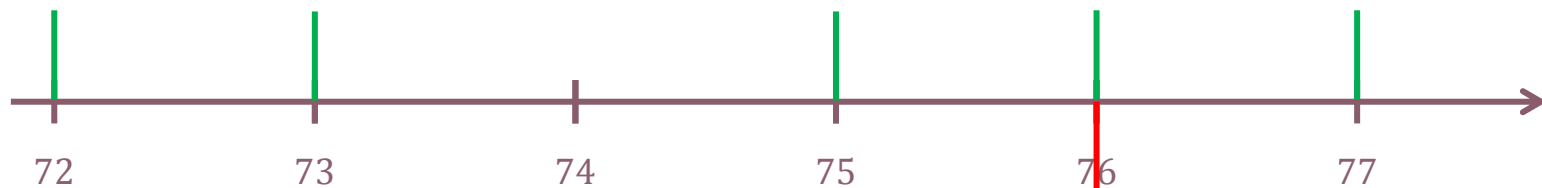
$$= \bar{x}_n + \frac{1}{n+1} (x_{n+1} + (n - n - 1)\bar{x}_n) = \bar{x}_n + \frac{1}{n+1} (x_{n+1} - \bar{x}_n)$$

q.e.d.



Uncertainty: updates to the mean

The sequential updates in our example now look like this:



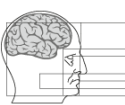
$$\bar{x}_1 = 76$$

$$\bar{x}_2 = 76 + \frac{1}{2}(73 - 76) = 74.5$$

$$\bar{x}_3 = 74.5 + \frac{1}{3}(75 - 74.5) = 74.\bar{6}$$

$$\bar{x}_4 = 74.\bar{6} + \frac{1}{4}(72 - 74.\bar{6}) = 74$$

$$\bar{x}_5 = 74 + \frac{1}{5}(77 - 74) = 74.6$$

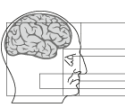


What are the building blocks of the updates we've just seen?

$$\bar{x}_{n+1} = \bar{x}_n + \frac{1}{n+1} (x_{n+1} - \bar{x}_n)$$

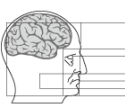
Diagram illustrating the components of the update equation:

- new input**: Points to x_{n+1} in the term $(x_{n+1} - \bar{x}_n)$.
- prediction error**: Points to the term $(x_{n+1} - \bar{x}_n)$.
- weight (learning rate)**: Points to the fraction $\frac{1}{n+1}$.
- prediction**: Points to \bar{x}_n .



Is this a general pattern?

- More specifically, does it generalize to Bayesian inference?
- «Bayesian inference» simply means inference on uncertain quantities according to the rules of probability theory (i.e., according to logic).
- Agents who use Bayesian inference will make better predictions (provided they have a good model of their environment), which will give them an evolutionary advantage.
- We may therefore assume that evolved biological agents use Bayesian inference, or a close approximation to it.
- So is Bayesian inference based on predictions that are updated using uncertainty-weighted prediction errors?



Updates in a simple Gaussian model

- Think boat, lighthouse, etc., again, but now we're doing Bayesian inference.
- Before we make the next observation, our belief about the true angle ϑ can be described by a Gaussian prior:

$$p(\vartheta) \sim \mathcal{N}(\mu_{\vartheta}, \pi_{\vartheta}^{-1})$$

- The likelihood of our observation is also Gaussian, with precision π_{ε} :

$$p(x|\vartheta) \sim \mathcal{N}(\vartheta, \pi_{\varepsilon}^{-1})$$

- Bayes' rule now tells us that the posterior is Gaussian again:

$$p(\vartheta|x) = \frac{p(x|\vartheta)p(\vartheta)}{\int p(x|\vartheta')p(\vartheta')d\vartheta'} \sim \mathcal{N}(\mu_{\vartheta|x}, \pi_{\vartheta|x}^{-1})$$

Updates in a simple Gaussian model

- Here's how the updates to the sufficient statistics μ and π describing our belief look like:

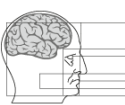
$$\pi_{\vartheta|x} = \pi_{\vartheta} + \pi_{\varepsilon}$$

$$\mu_{\vartheta|x} = \mu_{\vartheta} + \frac{\pi_{\varepsilon}}{\pi_{\vartheta|x}} (x - \mu_{\vartheta})$$

Diagram illustrating the update of the mean $\mu_{\vartheta|x}$:

- μ_{ϑ} is labeled "prediction" (red arrow).
- $\frac{\pi_{\varepsilon}}{\pi_{\vartheta|x}}$ is labeled "weight (learning rate)" (blue arrow).
- $(x - \mu_{\vartheta})$ is labeled "prediction error" (purple oval and arrow).
- The weight is further defined as: $\text{weight (learning rate)} = \frac{\text{how much we're learning here}}{\text{how much we already know}}$

- So it's the same story all over again: the mean is updated by an uncertainty-weighted (more specifically: precision-weighted) prediction error.
- The size of the update is proportional to the likelihood precision and inversely proportional to the posterior precision.
- This pattern is not specific to the univariate Gaussian case, but generalizes to Bayesian updates for all exponential families of likelihood distributions with conjugate priors (i.e., to all formal descriptions of inference you are ever likely to need).



The analogy with simple mean updating goes further

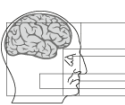
- Reminder (Gaussian update):

$$\mu_{\vartheta|x} = \mu_{\vartheta} + \frac{\pi_{\varepsilon}}{\pi_{\vartheta|x}}(x - \mu_{\vartheta}) = \mu_{\vartheta} + \frac{\pi_{\varepsilon}}{\pi_{\vartheta} + \pi_{\varepsilon}}(x - \mu_{\vartheta})$$

- Reducing by π_{ε} the fraction of precisions that make the learning rate, we get

$$\mu_{\vartheta|x} = \mu_{\vartheta} + \frac{1}{\frac{\pi_{\vartheta}}{\pi_{\varepsilon}} + 1}(x - \mu_{\vartheta})$$

- This is again our equation for updating an arithmetic mean, but with n replaced by $\frac{\pi_{\vartheta}}{\pi_{\varepsilon}}$.
- This shows that Bayesian inference on the mean of a Gaussian distribution entails nothing more than updating the arithmetic mean of observations with $\frac{\pi_{\vartheta}}{\pi_{\varepsilon}} =: \nu$ as a proxy for the number of prior observations, i.e. for the **weight of the prior relative to the observation**.



Generalization to all exponential families of distributions

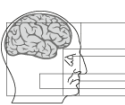
- Many of the most widely used probability distributions are families of exponential distributions.
- For example, the Gaussian distribution is an exponential family of distributions (and so are the beta, gamma, binomial, Bernoulli, multinomial, categorical, Dirichlet, Wishart, Gaussian-gamma, log-Gaussian, multivariate Gaussian, Poisson, and exponential distributions, among others). This means it can be written the following way:

$$p(\mathbf{x}|\boldsymbol{\vartheta}) = h(\mathbf{x}) \exp(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \mathbf{T}(\mathbf{x}) - A(\boldsymbol{\vartheta})) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma}\right)$$

with

$$\mathbf{x} = x, \quad \boldsymbol{\vartheta} = (\mu, \sigma)^T, \quad h(\mathbf{x}) = \frac{1}{\sqrt{2\pi}}, \quad \boldsymbol{\eta}(\boldsymbol{\vartheta}) = \left(\frac{\mu}{\sigma}, -\frac{1}{2\sigma}\right)^T, \quad \mathbf{T}(\mathbf{x}) = (x, x^2)^T, \quad A(\boldsymbol{\vartheta}) = \frac{\mu^2}{\sigma} + \frac{\ln \sigma}{2}$$

- This allows us to look at Bayesian belief updating in a very general way for all exponential families of distributions.



Generalization to all exponential families of distributions

- Our likelihood is an exponential family in its general form:

$$p(\mathbf{x}|\boldsymbol{\vartheta}) = h(\mathbf{x}) \exp(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \mathbf{T}(\mathbf{x}) - A(\boldsymbol{\vartheta}))$$

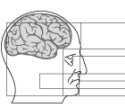
- The vector $\mathbf{T}(\mathbf{x})$ (a function of the observation \mathbf{x}) is called the sufficient statistic.
- For the prior, we may assume that we have made ν observations with sufficient statistic $\boldsymbol{\xi}$:

$$p(\boldsymbol{\vartheta}|\boldsymbol{\xi}, \nu) = z(\boldsymbol{\xi}, \nu) \exp(\nu(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \boldsymbol{\xi} - A(\boldsymbol{\vartheta}))) \quad (\text{where } z(\boldsymbol{\xi}, \nu) \text{ is a normalization constant})$$

- It then turns out that the posterior has the same form, but with an updated $\boldsymbol{\xi}$ and ν replaced with $\nu + 1$:

$$p(\boldsymbol{\vartheta}|\mathbf{x}, \boldsymbol{\xi}, \nu) = z(\boldsymbol{\xi}', \nu + 1) \exp((\nu + 1)(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \boldsymbol{\xi}' - A(\boldsymbol{\vartheta})))$$

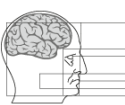
$$\boldsymbol{\xi}' = \boldsymbol{\xi} + \frac{1}{\nu + 1} (\mathbf{T}(\mathbf{x}) - \boldsymbol{\xi})$$



Proof of the update equation

$$\begin{aligned} \overbrace{p(\boldsymbol{\vartheta}|\mathbf{x}, \boldsymbol{\xi}, \nu)}^{\text{posterior}} &\propto \overbrace{p(\mathbf{x}|\boldsymbol{\vartheta})}^{\text{likelihood}} \overbrace{p(\boldsymbol{\vartheta}|\boldsymbol{\xi}, \nu)}^{\text{prior}} \\ &= h(\mathbf{x}) \exp(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \mathbf{T}(\mathbf{x}) - A(\boldsymbol{\vartheta})) z(\boldsymbol{\xi}, \nu) \exp(\nu(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \boldsymbol{\xi} - A(\boldsymbol{\vartheta}))) \\ &\propto \exp(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot (\mathbf{T}(\mathbf{x}) + \nu\boldsymbol{\xi}) - (\nu + 1)A(\boldsymbol{\vartheta})) \\ &= \exp\left((\nu + 1) \left(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \frac{1}{\nu + 1} (\mathbf{T}(\mathbf{x}) + \nu\boldsymbol{\xi}) - A(\boldsymbol{\vartheta}) \right)\right) \\ &= \exp\left((\nu + 1) \left(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \left(\boldsymbol{\xi} + \frac{1}{\nu + 1} (\mathbf{T}(\mathbf{x}) + \nu\boldsymbol{\xi} - (\nu + 1)\boldsymbol{\xi}) \right) - A(\boldsymbol{\vartheta}) \right)\right) \\ &= \exp\left((\nu + 1) \left(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \underbrace{\left(\boldsymbol{\xi} + \frac{1}{\nu + 1} (\mathbf{T}(\mathbf{x}) - \boldsymbol{\xi}) \right)}_{=:\boldsymbol{\xi}'} - A(\boldsymbol{\vartheta}) \right)\right) \\ &\Rightarrow p(\boldsymbol{\vartheta}|\mathbf{x}, \boldsymbol{\xi}, \nu) = z(\boldsymbol{\xi}', \nu') \exp(\nu'(\boldsymbol{\eta}(\boldsymbol{\vartheta}) \cdot \boldsymbol{\xi}' - A(\boldsymbol{\vartheta}))) \\ &\quad \text{with } \nu' := \nu + 1, \quad \boldsymbol{\xi}' := \boldsymbol{\xi} + \frac{1}{\nu + 1} (\mathbf{T}(\mathbf{x}) - \boldsymbol{\xi}) \end{aligned}$$

q.e.d.



Some examples

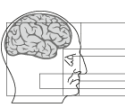
- **Univariate Gaussian** model with unknown mean but **known precision** (our example from the beginning):

$$T(x) = x$$

- This means updating beliefs about the mean simply requires tracking the mean of observations
- **Univariate Gaussian** model with unknown mean and unknown precision:

$$T(x) = (x, x^2)^T$$

- Updating beliefs about both mean and precision of a Gaussian requires tracking the means of observations and squared observations; this amounts to the first and second moments by which a Gaussian distribution is fully characterized.
- In the **multivariate Gaussian** case we have $T(x) = (x, xx^T)^T$



Some examples

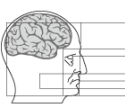
- **Bernoulli** model (one out of two possible outcomes, coded as 0 and 1; e.g., coin flipping):

$$T(x) = x$$

- The prior here turns out to be a **beta distribution** corresponding to ν pseudo-observations with mean ξ . All we need to do to get the posterior (i.e., to update our belief) is to update the mean as new observations come in.
- **Categorical** model (one out of several possible outcomes, with the observed outcome coded as 1, the rest as 0)

$$T(x) = x$$

- The prior and posterior here are **Dirichlet distributions**, and again, all we need to do to update beliefs that have a Dirichlet form is to track the means of observed successes (1) and failures (0).



Some examples

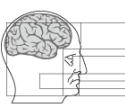
- **Beta** model (an outcome bounded between 0 and 1):

$$\mathbf{T}(x) = (\ln x, \ln(1 - x))^T$$

- **Gamma** model (an outcome bounded below at 0):

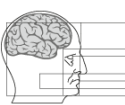
$$\mathbf{T}(x) = (\ln x, x)^T$$

- Now that we have dealt with beliefs about states that are binary (Bernoulli), categorical, bounded on both sides (beta), bounded on one side (gamma), and unbounded (Gaussian), we have most kinds of states we can have beliefs about.
- **All Bayesian (i.e., probabilistic, rational) updates of such beliefs take the form of precision-weighted prediction errors.**



Limitations

- Examples of distributions that are not exponential families: Student's t , Cauchy
- These distributions are popular because of their «fat tails». However, fat tails can also be achieved with appropriate hierarchies of Gaussians (cf. the hierarchical Gaussian filter, HGF)
- A further kind of distributions that are not exponential families are found in mixture models.
- Such models are popular because of they provide multimodal distributions. But again, appropriate hierarchies of distributions may save the day.



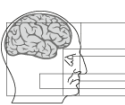
“Every good regulator of a system must be a model of that system” (Conant & Ashby, 1970)

Abstract:

«The design of a complex regulator often includes the making of a model of the system to be regulated. The making of such a model has hitherto been regarded as optional, as merely one of many possible ways.

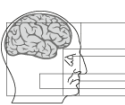
In this paper a theorem is presented which shows, under very broad conditions, that any regulator that is maximally both successful and simple must be isomorphic with the system being regulated. (The exact assumptions are given.) Making a model is thus necessary.

*The theorem has the interesting corollary that **the living brain**, so far as it is to be successful and efficient as a regulator for survival, **must proceed**, in learning, **by the formation of a model (or models) of its environment.**»*

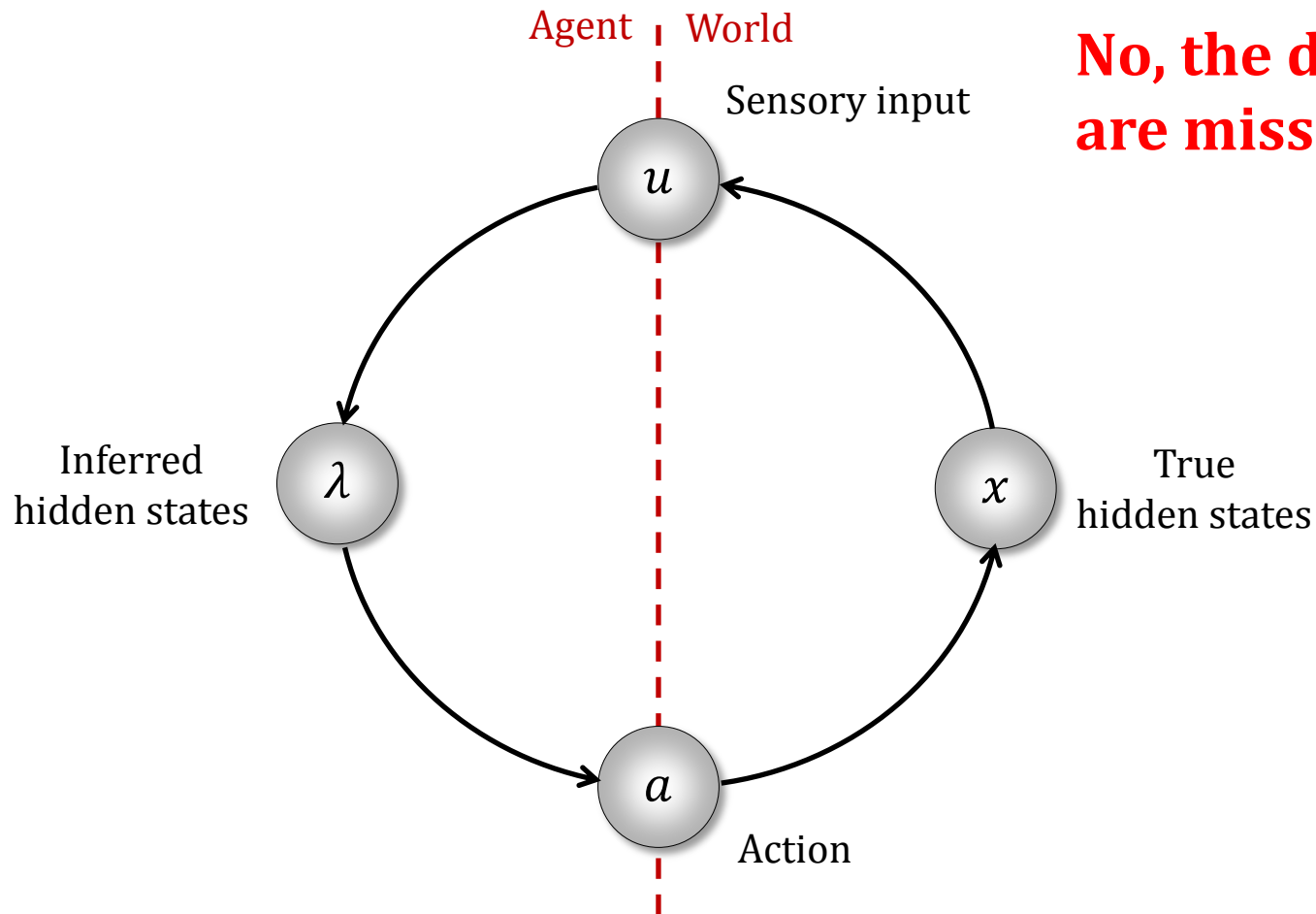


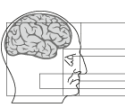
Systems theory as the conceptual bridge between clinical phenomena and neuronal pathophysiology

- Belief updating by precision-weighted prediction errors provides a conceptual framework in which both clinical phenomena and neurobiological findings can be interpreted.
- For examples of this approach, see Adams et al. (2013) (psychosis), or Lawson et al. (2014), Quattrocki & Friston (2014) (autism).
- Summary: the mind needs to be a model of its environment \Rightarrow needs to perform Bayesian inference \Rightarrow needs to use precision-weighting of prediction errors \Rightarrow if that's all the mind does, it's also all that can go wrong \Rightarrow both clinical manifestations and the neurobiology of psychiatric disorders must be interpretable in these terms.
- Now that we have this conceptual framework, we can start filling it with content.



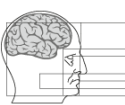
Does inference as we've described it adequately describe the situation of actual biological agents?





What about dynamics?

- Up to now, we've only looked at inference on static quantities, but biological agents live in a continually changing world.
- In our example, the boat's position changes and with it the angle to the lighthouse.
- How can we take into account that old information becomes obsolete? If we don't, our learning rate becomes smaller and smaller because our equations were derived under the assumption that we're accumulating information about a stable quantity.



What's the simplest way to keep the learning rate from going too low?

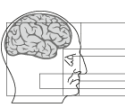
- Keep it constant!
- So, taking the update equation for the mean of our observations as our point of departure...

$$\bar{x}_n = \bar{x}_{n-1} + \frac{1}{n}(x_n - \bar{x}_{n-1}),$$

- ... we simply replace $\frac{1}{n}$ with a constant α :

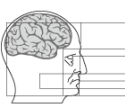
$$\mu_n = \mu_{n-1} + \alpha(x_n - \mu_{n-1}).$$

- This is called *Rescorla-Wagner learning* [although it wasn't this line of reasoning that led Rescorla & Wagner (1972) to their formulation].



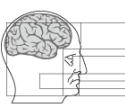
Does a constant learning rate solve our problems?

- Partly: it implies a certain rate of forgetting because it amounts to taking only the $n = \frac{1}{\alpha}$ last data points into account. But...
- ... if the learning rate is supposed to reflect uncertainty in Bayesian inference, then how do we
 - (a) know that α reflects the right level of uncertainty at any one time, and
 - (b) account for changes in uncertainty if α is constant?
- What we really need is an adaptive learning that accurately reflects uncertainty.



Needed: an adaptive learning rate that accurately reflects uncertainty

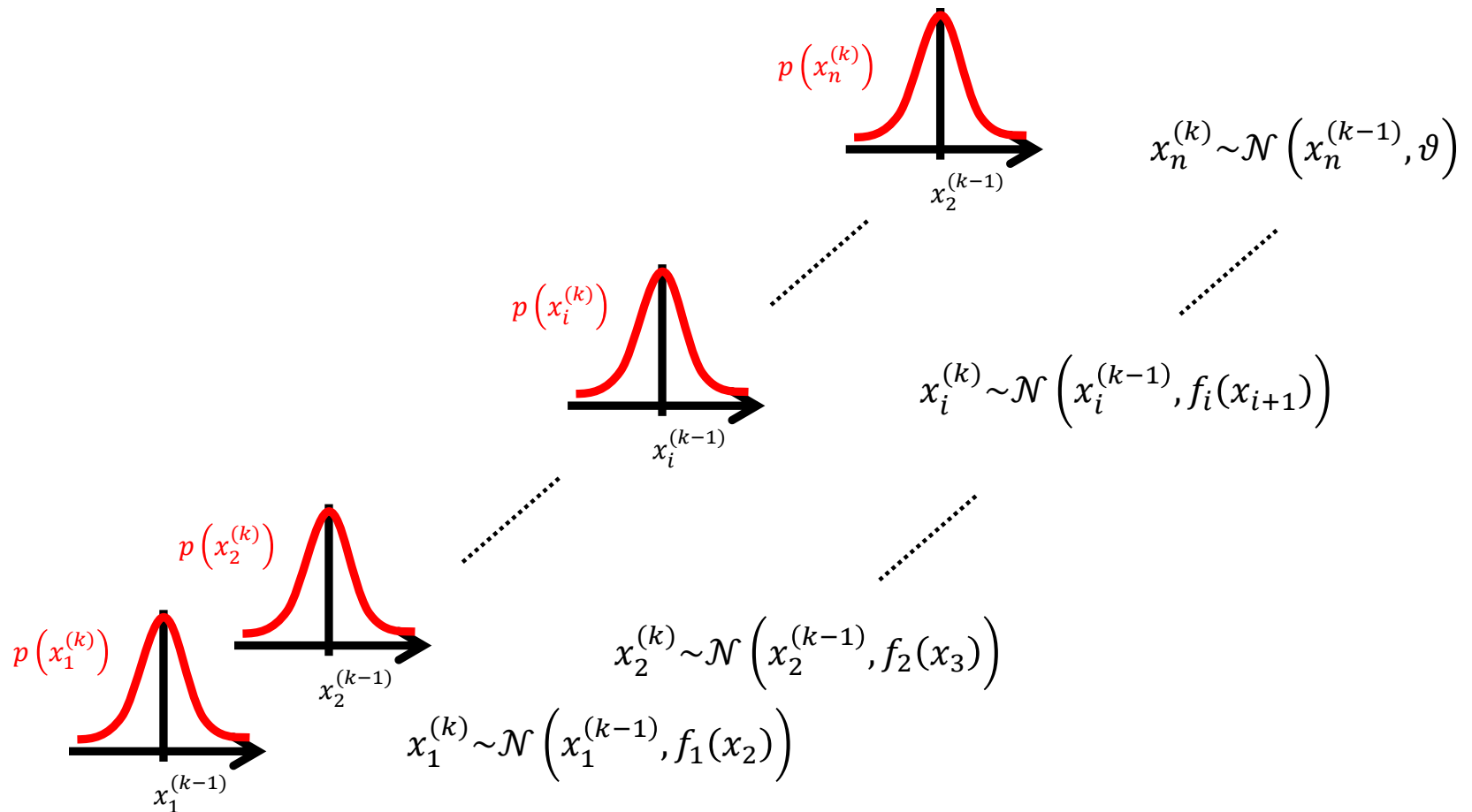
- This requires us to think a bit more about what kinds of uncertainty we are dealing with.
- A possible taxonomy of uncertainty is (cf. Yu & Dayan, 2003; Payzan-LeNestour & Bossaerts, 2011):
 - (a) **outcome uncertainty** that remains unaccounted for by the model, called *risk* by economists (π_ε in our Bayesian example); this uncertainty remains even when we know all parameters exactly,
 - (b) **informational** or *expected* uncertainty about the value of model parameters ($\pi_{\vartheta|x}$ in the Bayesian example),
 - (c) **environmental** or *unexpected* uncertainty owing to changes in model parameters (not accounted for in our Bayesian example, hence unexpected).

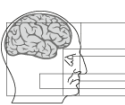


An adaptive learning rate that accurately reflects uncertainty

- Various efforts have been made to come up with an adaptive learning rate:
 - Kalman (1960)
 - Sutton (1992)
 - Nassar et al. (2010)
 - Payzan-LeNestour & Bossaerts (2011)
 - Mathys et al. (2011)
 - Wilson et al. (2013)
- The Kalman filter is optimal for linear dynamical systems, but realistic data usually require non-linear models.
- Mathys et al. use a generic non-linear hierarchical Bayesian model that allows us to derive update equations that are optimal in the sense that they minimize surprise.

The hierarchical Gaussian filter (HGF)





The hierarchical Gaussian filter (HGF)

- At the outcome level (i.e., at the very bottom of the hierarchy), we have

$$u^{(k)} \sim \mathcal{N} \left(x_1^{(k)}, \hat{\pi}_u^{-1} \right)$$

- This gives us the following update for our belief on x_1 (our quantity of interest):

$$\pi_1^{(k)} = \hat{\pi}_1^{(k)} + \hat{\pi}_u$$

$$\mu_1^{(k)} = \mu_1^{(k-1)} + \frac{\hat{\pi}_u}{\pi_1^{(k)}} \left(u^{(k)} - \mu_1^{(k-1)} \right)$$

- The familiar structure again – but now with a learning rate that is responsive to all kinds of uncertainty, including environmental (unexpected) uncertainty.

The learning rate in the HGF

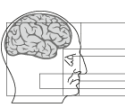
Unpacking the learning rate, we see:

$$\frac{\hat{\pi}_u}{\pi_1^{(k)}} = \frac{\hat{\pi}_u}{\hat{\pi}_1^{(k)} + \hat{\pi}_u} = \frac{\hat{\pi}_u}{\frac{1}{\sigma_1^{(k-1)} + \exp(\kappa_1 \mu_2^{(k-1)} + \omega_1)} + \hat{\pi}_u}$$

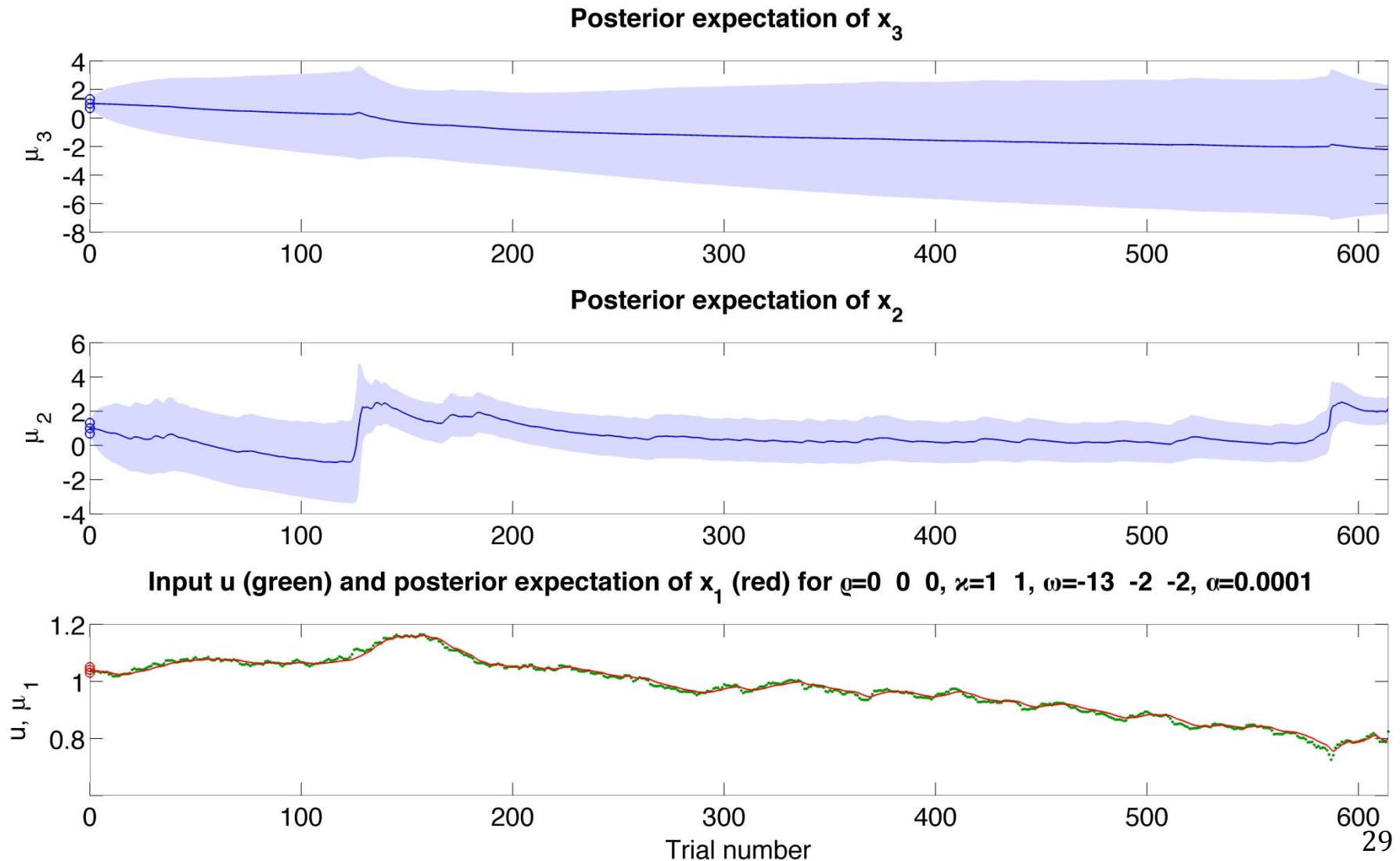
outcome uncertainty

informational uncertainty

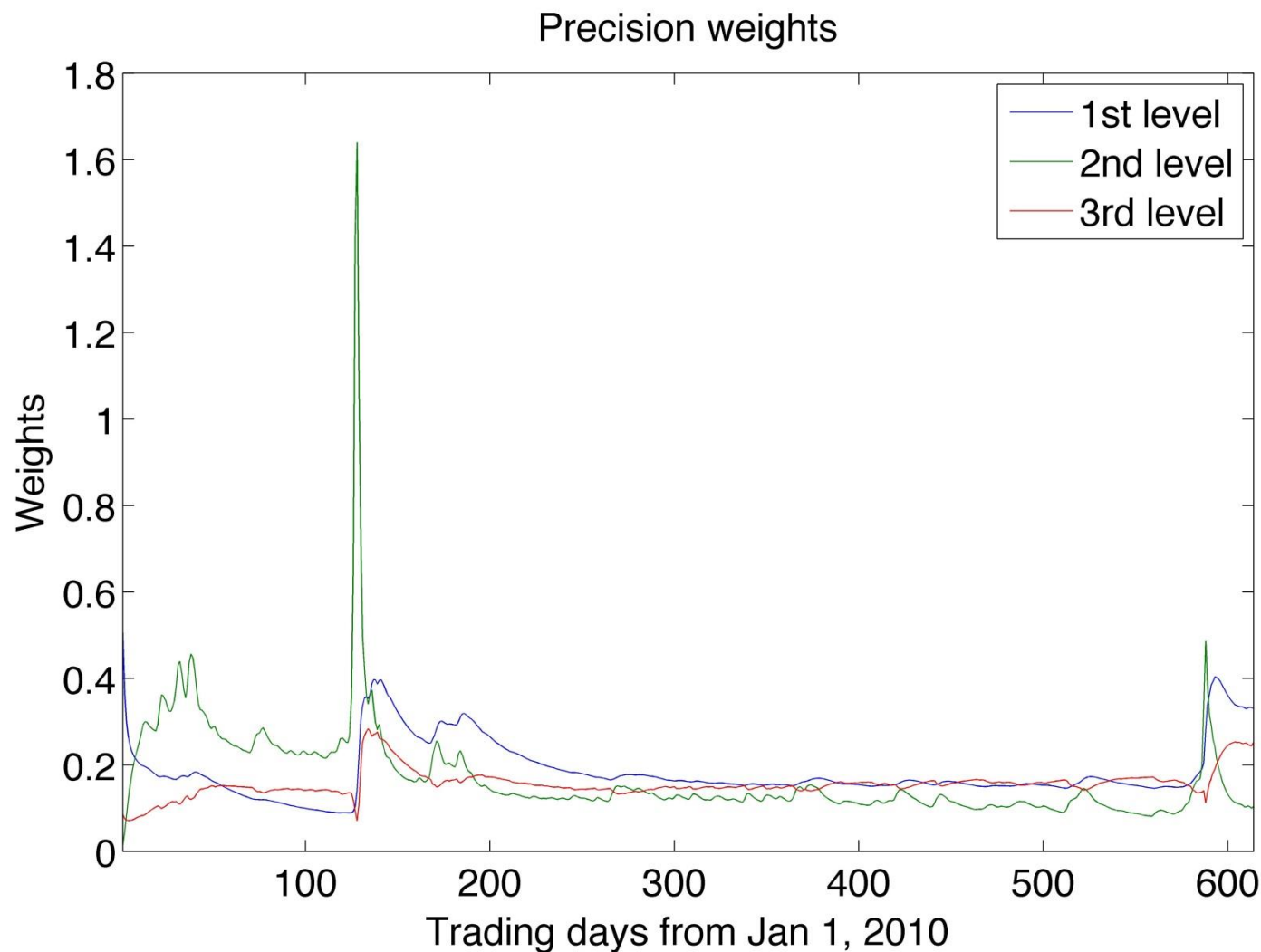
environmental uncertainty

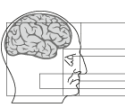


3-level HGF for continuous observations



3-level HGF for continuous observations





VAPes and VOPEs

The updates of the belief on x_1 are driven by value prediction errors (VAPes)

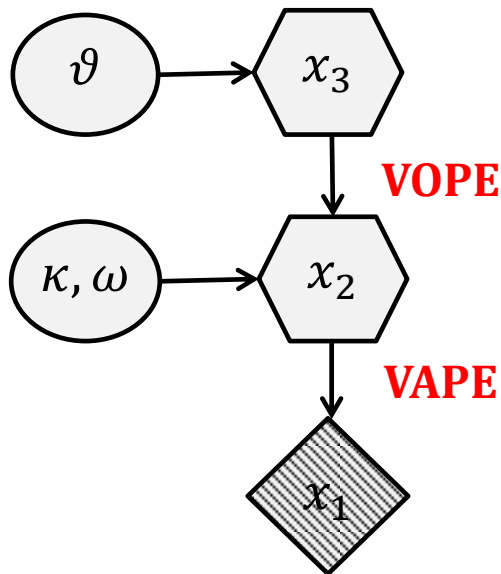
$$\mu_1^{(k)} = \mu_1^{(k-1)} + \frac{\hat{\pi}_u}{\pi_1^{(k)}} (u^{(k)} - \mu_1^{(k-1)}), \text{ VAPE}$$

while the x_2 -updates are driven by volatility prediction errors (VOPEs)

$$\mu_2^{(k)} = \mu_2^{(k-1)} + \frac{1}{2} \kappa_1 v_1^{(k)} \frac{\hat{\pi}_1^{(k)}}{\pi_2^{(k)}} \delta_1^{(k)} \text{ VOPE}$$

$$\delta_1^{(k)} \stackrel{\text{def}}{=} \frac{\sigma_1^{(k)} + \left(\mu_1^{(k)} - \mu_1^{(k-1)}\right)^2}{\sigma_1^{(k-1)} + \exp\left(\kappa_1 \mu_2^{(k-1)} + \omega_1\right)} - 1$$

3-level HGF for binary observations



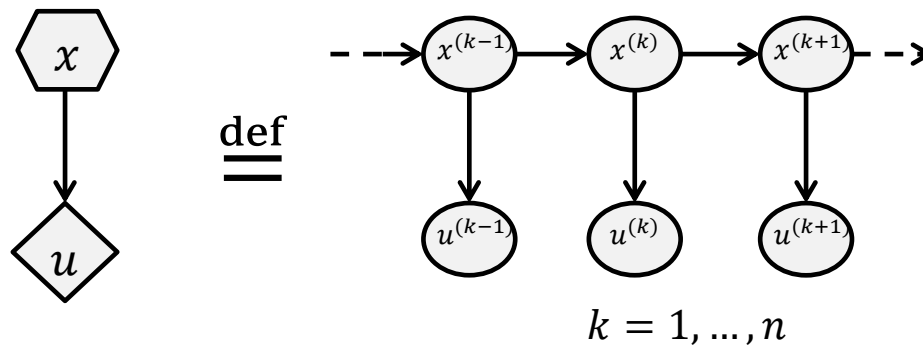
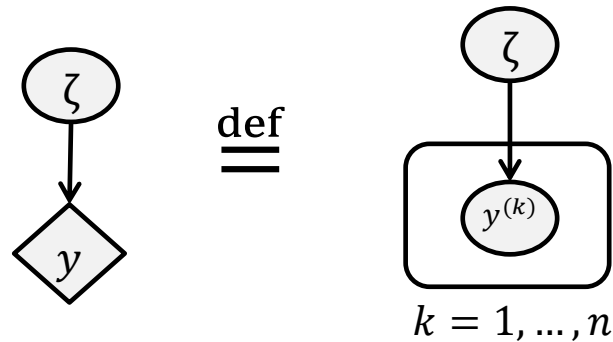
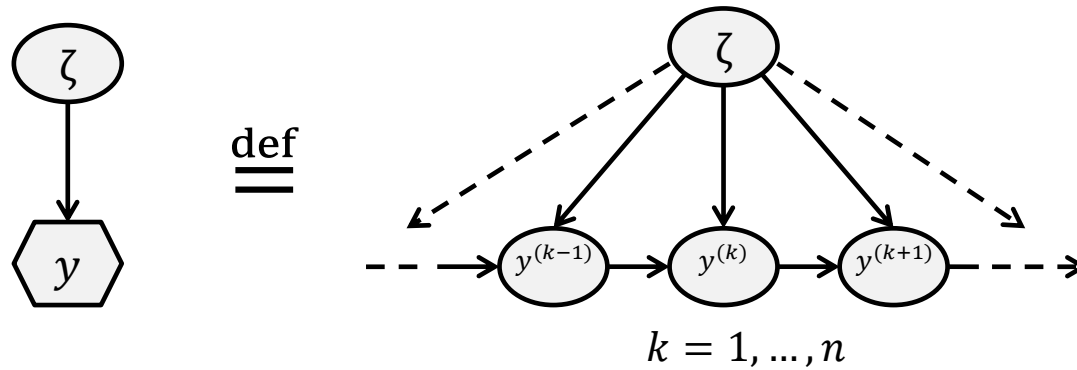
$$x_3^{(k)} \sim \mathcal{N}(x_3^{(k-1)}, \vartheta)$$

$$x_2^{(k)} \sim \mathcal{N}(x_2^{(k-1)}, \exp(\kappa x_3^{(k)} + \omega))$$

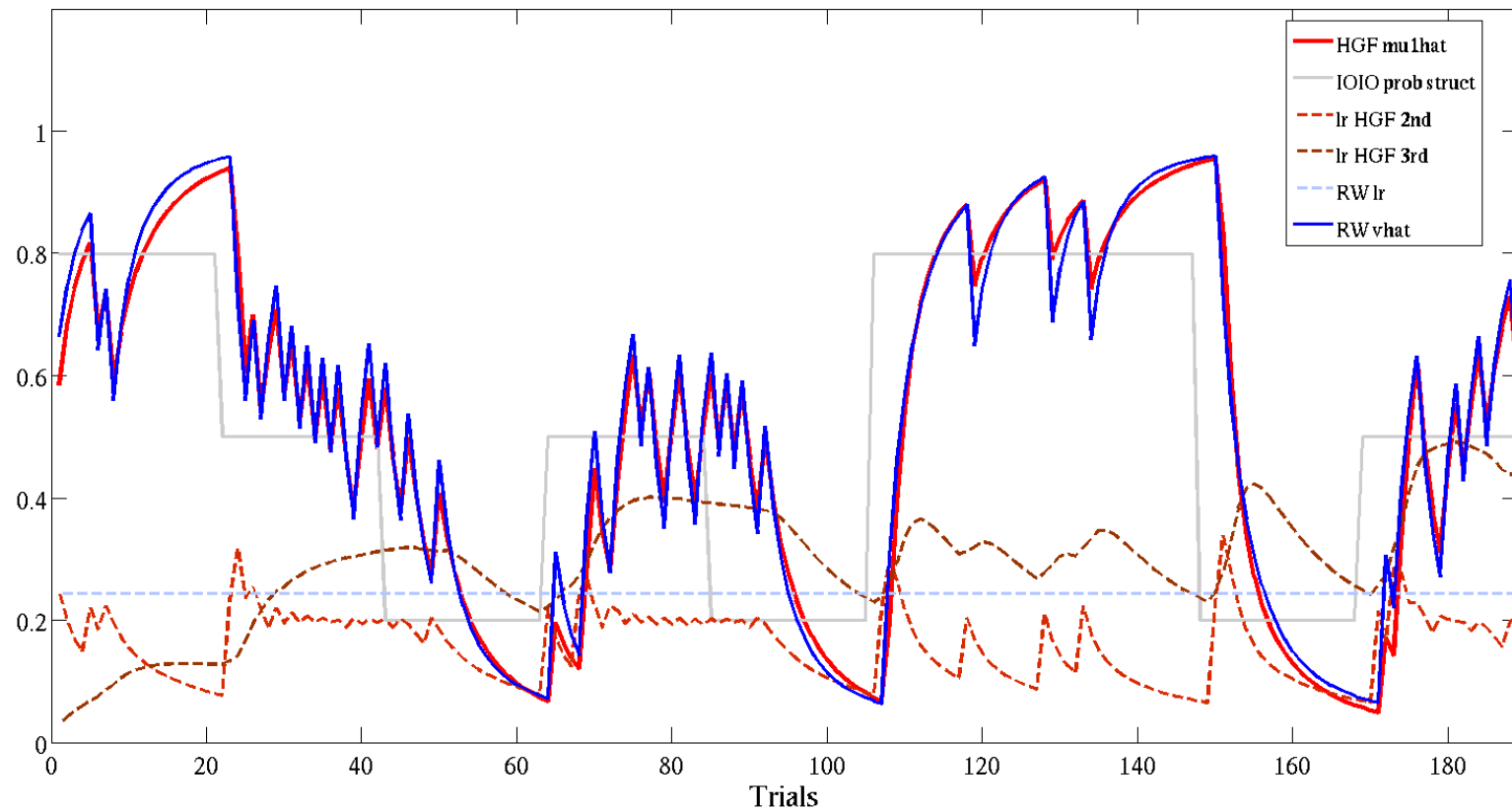
$$x_1^{(k)} \sim \text{Bern}(x_2^{(k)})$$

Mathys et al., 2011; Iglesias et al., 2013; Vossel et al., 2014a; Hauser et al., 2014; Diaconescu et al., 2014; Vossel et al., 2014b; ...

Taking it all together: notation

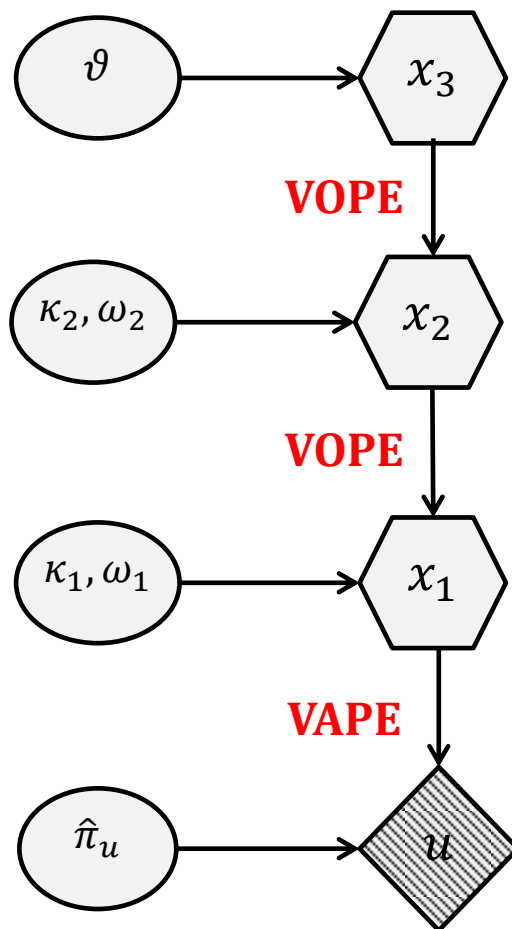


The learning rate in the HGF



Andreea Diaconescu

3-level HGF for continuous observations



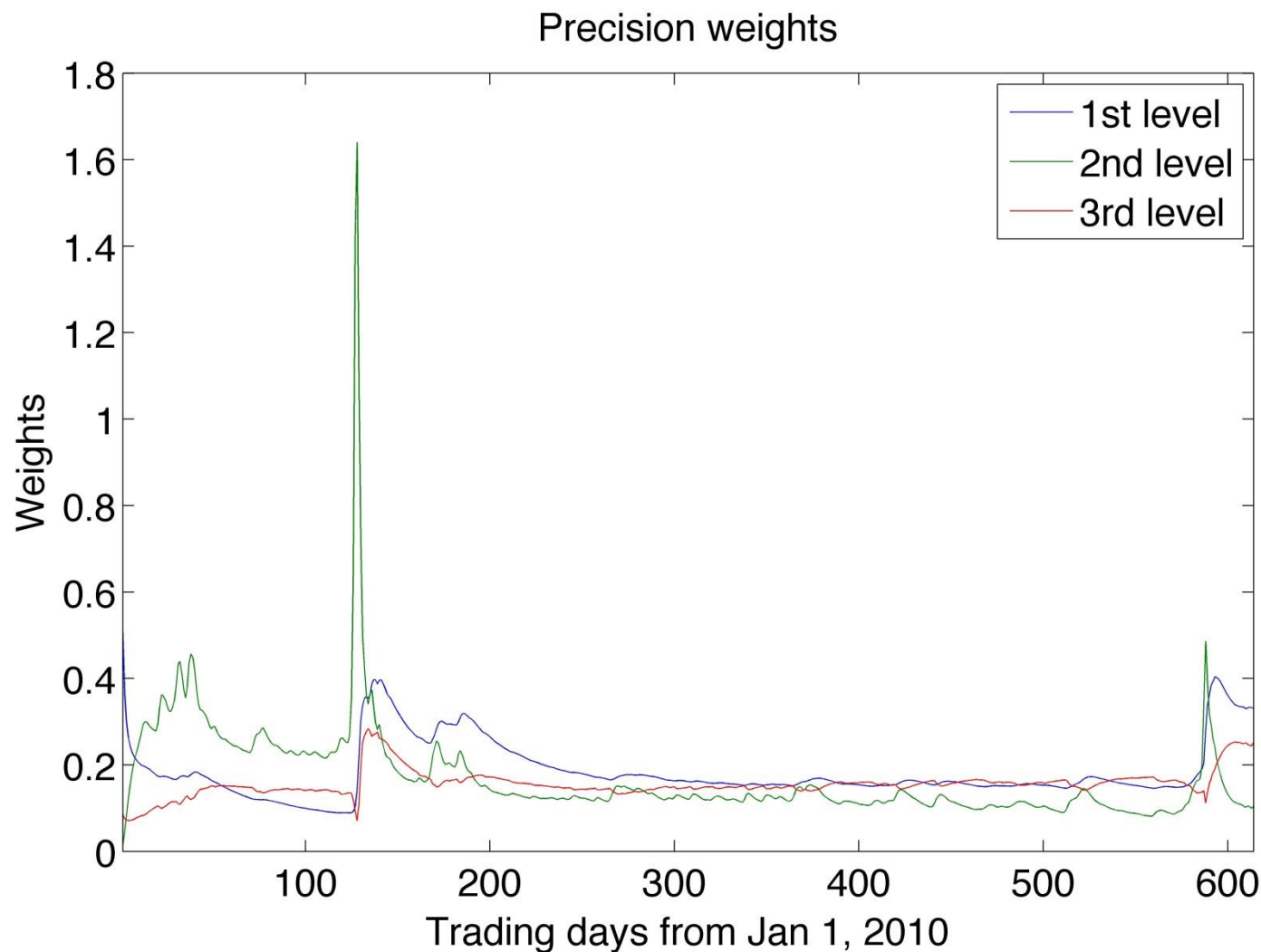
$$x_3^{(k)} \sim \mathcal{N} \left(x_3^{(k-1)}, \vartheta \right)$$

$$x_2^{(k)} \sim \mathcal{N} \left(x_2^{(k-1)}, \exp \left(\kappa_2 x_3^{(k)} + \omega_2 \right) \right)$$

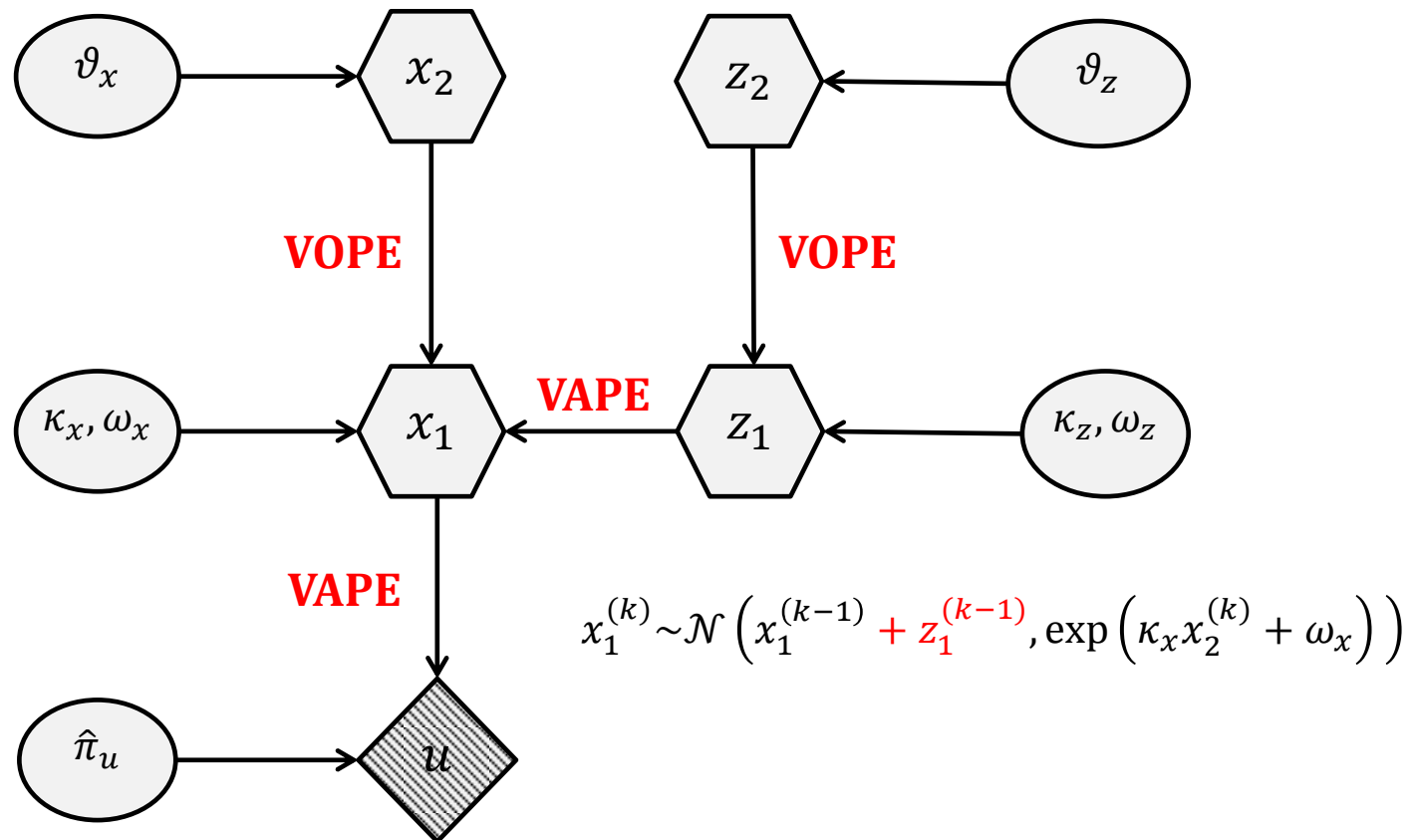
$$x_1^{(k)} \sim \mathcal{N} \left(x_1^{(k-1)}, \exp \left(\kappa_1 x_2^{(k)} + \omega_1 \right) \right)$$

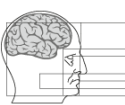
$$u^{(k)} \sim \mathcal{N} \left(x_1^{(k)}, \hat{\pi}_u^{-1} \right)$$

3-level HGF for continuous observations

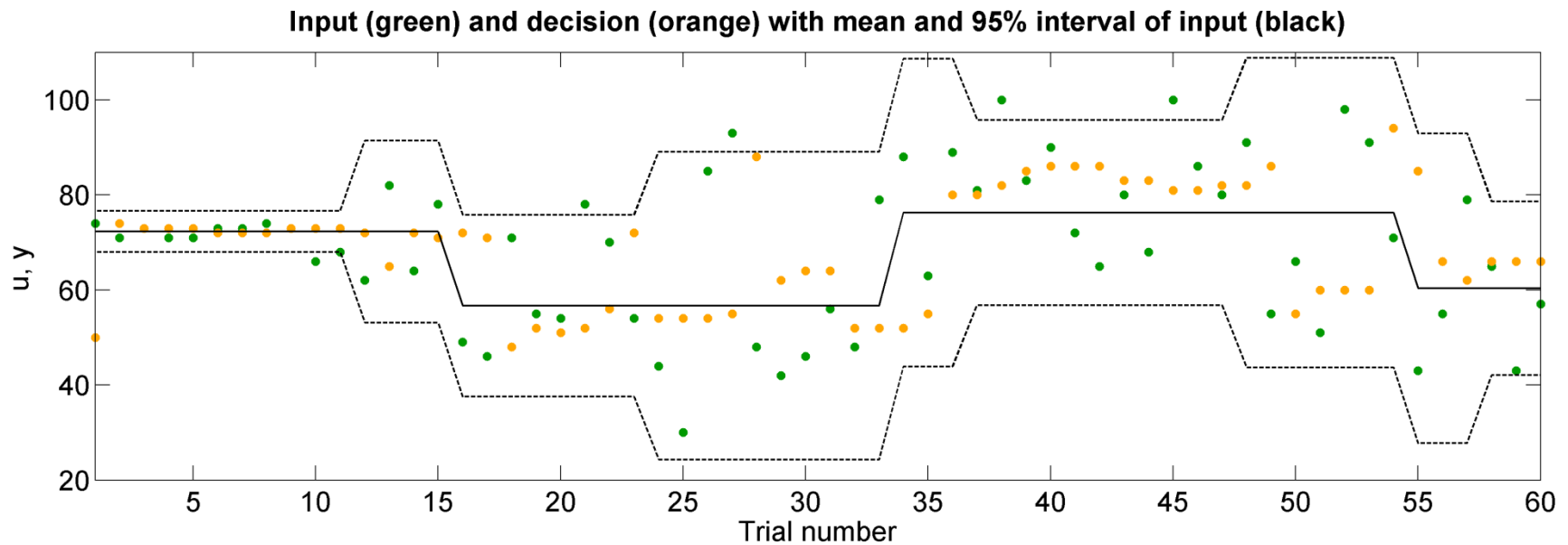


Variable drift



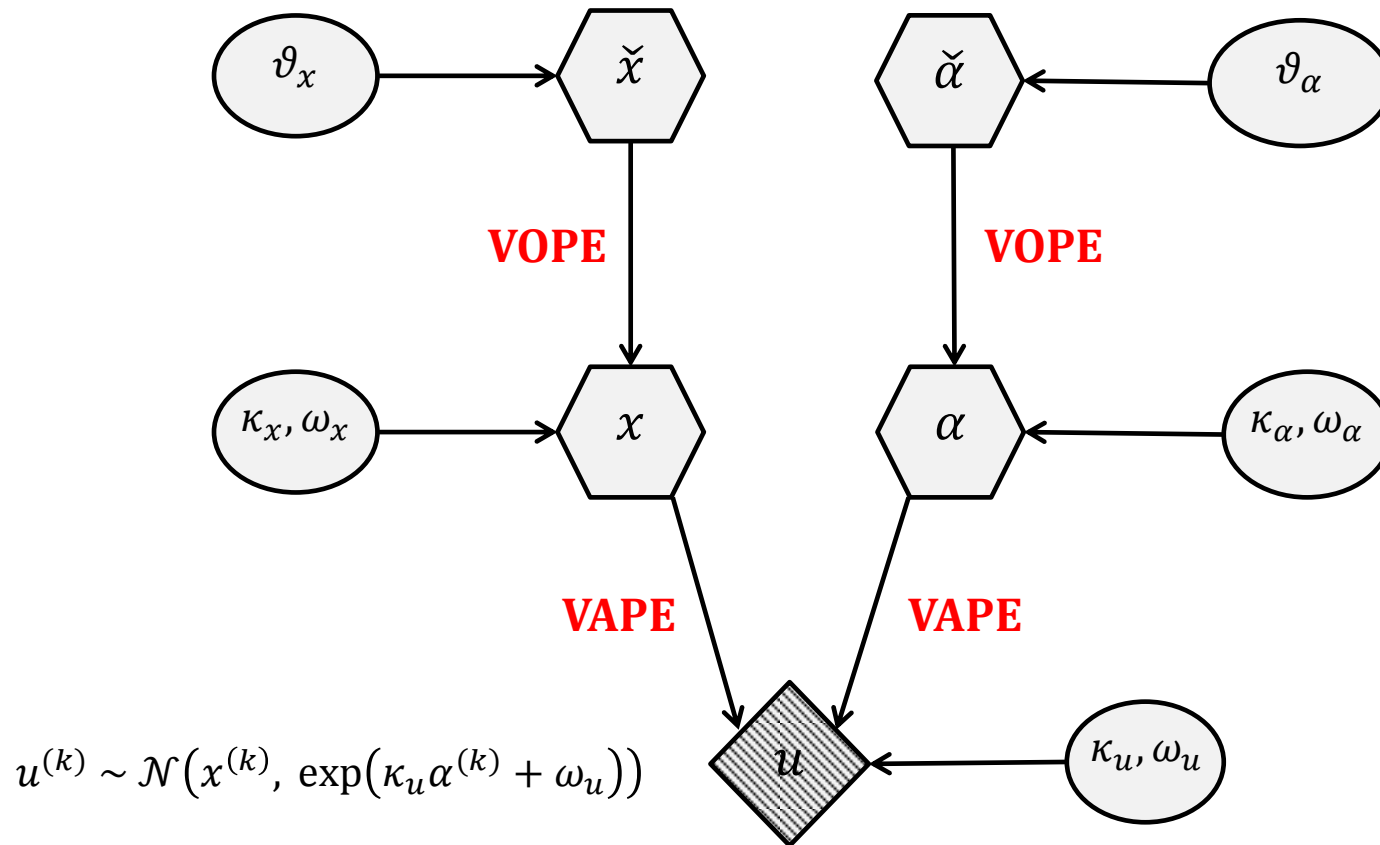


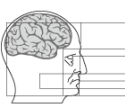
Jumping Gaussian estimation task



Chaohui Guo

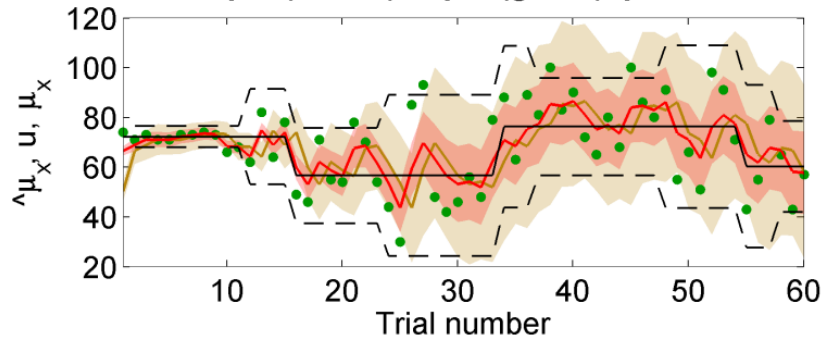
Independent mean and variance model



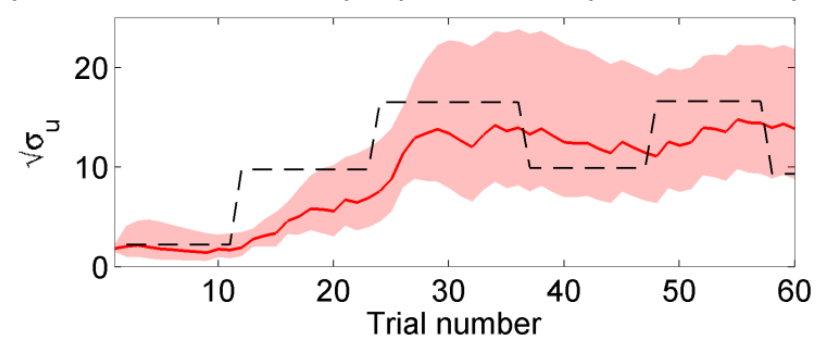


Jumping Gaussian estimation task

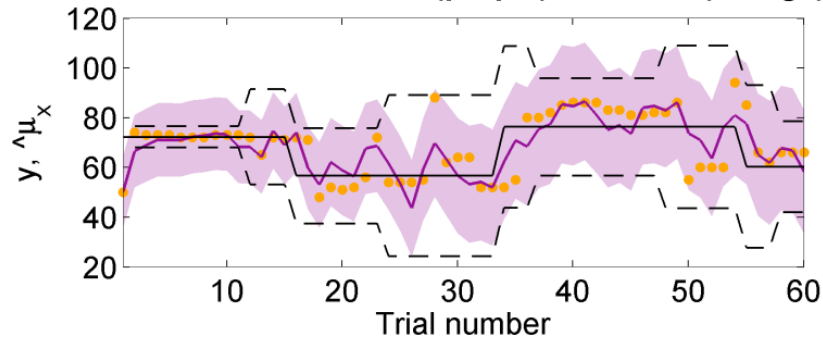
Prediction of input (brown), input (green), posterior belief (red)



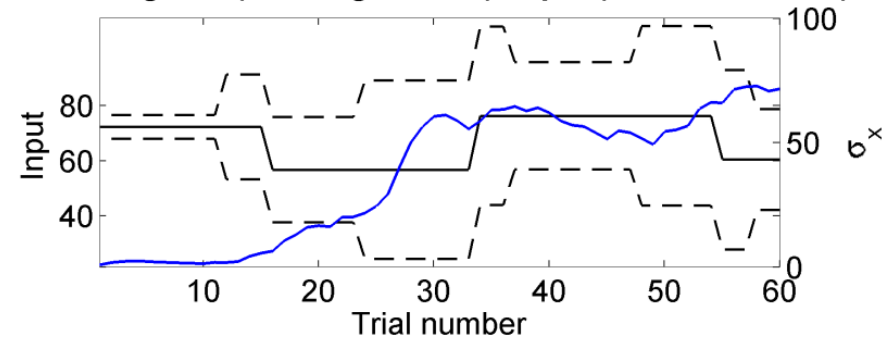
Belief on noise (red), true noise (dashed black)

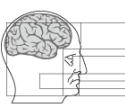


Prediction of decision (purple), decision (orange)



Learning rate (blue; right scale), input (black; left scale)

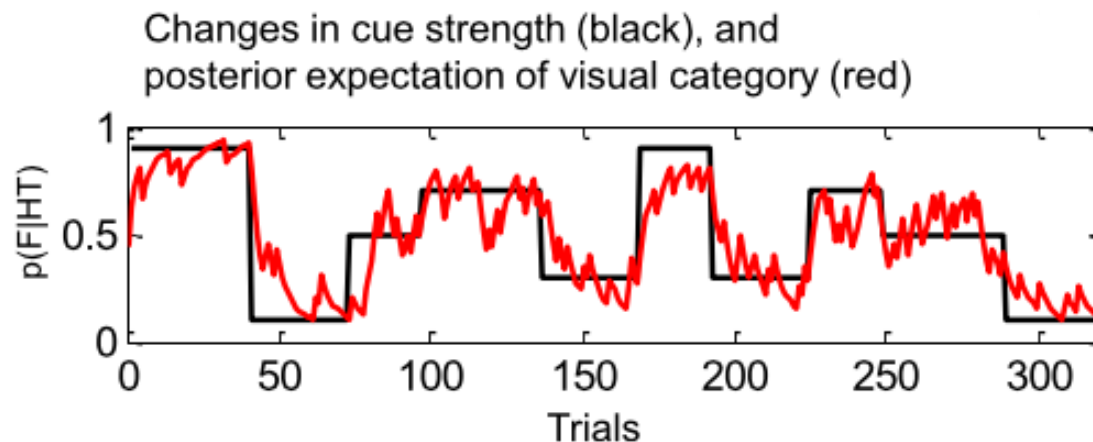
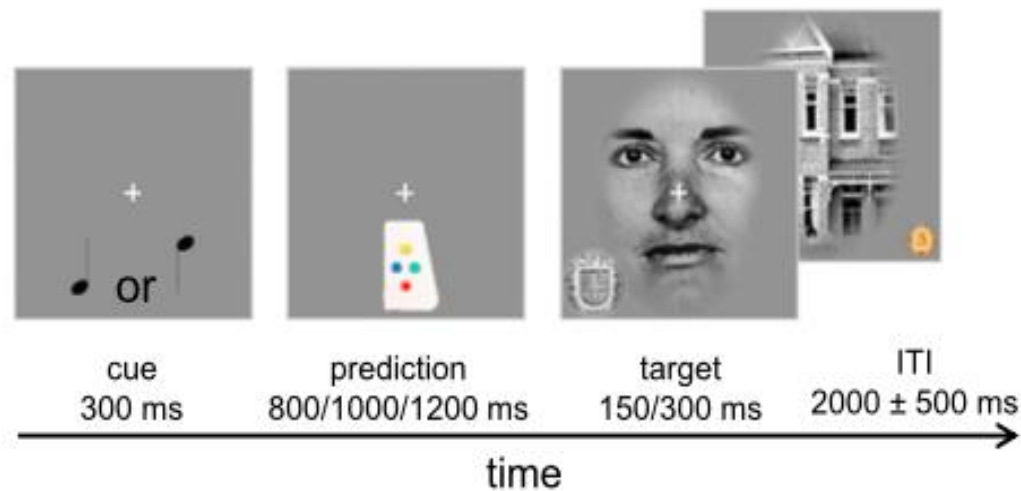


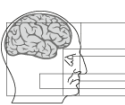


How to reveal the precision-weighting of prediction errors when simple exponential-family likelihoods will not do

- Formulate the problem hierarchically (i.e., imitate evolution: when it built a brain that supports a mind which is a model of its environment, it came up with a (largely) hierarchical solution)
- Separate levels using a mean-field approximation
- Derive update equations

HGF: empirical evidence (Iglesias et al., 2013)

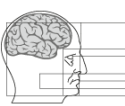




HGF: empirical evidence (Iglesias et al., 2013)

Model comparison:

BMS results	Behavioral study		fMRI study 1		fMRI study 2	
	PP	XP	PP	XP	PP	XP
HGF1	0.8435	1	0.7422	1	0.7166	1
HGF2	0.0259	0	0.0200	0	-	-
HGF3	0.0361	0	0.1404	0	0.1304	0
Sutton	0.0685	0	0.0710	0	0.0761	0
RW	0.0260	0	0.0264	0	0.0769	0



HGF: empirical evidence (Iglesias et al., 2013)

Model comparison:

BMS results	Behavioral study		fMRI study 1		fMRI study 2	
	PP	XP	PP	XP	PP	XP
HGF1	0.8435	1	0.7422	1	0.7166	1
HGF2	0.0259	0	0.0200	0	-	-
HGF3	0.0361	0	0.1404	0	0.1304	0
Sutton	0.0685	0	0.0710	0	0.0761	0
RW	0.0260	0	0.0264	0	0.0769	0

HGF: empirical evidence (Iglesias et al., 2013)

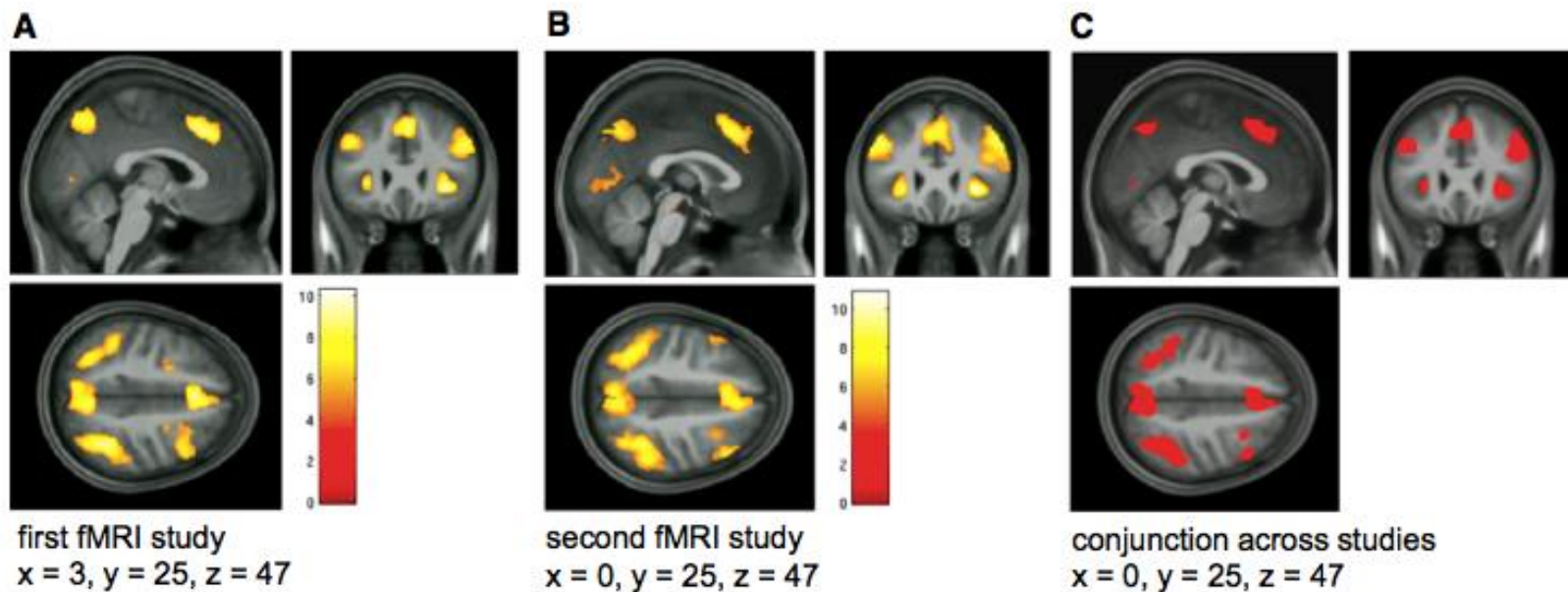


Figure 2. Whole-Brain Activations by ε_2

Activations by precision-weighted prediction error about visual stimulus outcome, ε_2 , in the first fMRI study (A) and the second fMRI study (B). Both activation maps are shown at a threshold of $p < 0.05$, FWE corrected for multiple comparisons across the whole brain. To highlight replication across studies, (C) shows the results of a “logical AND” conjunction, illustrating voxels that were significantly activated in both studies.

HGF: empirical evidence (Iglesias et al., 2013)

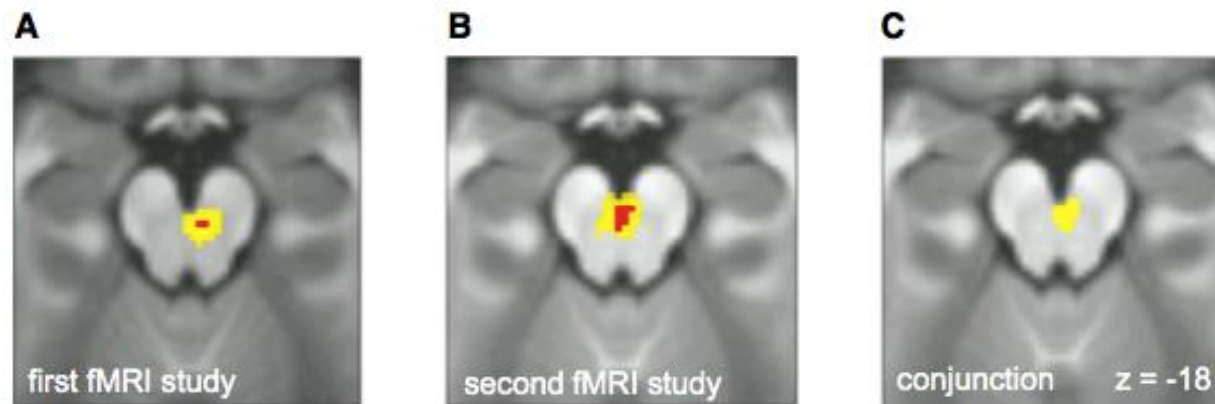


Figure 3. Midbrain Activation by ϵ_2

Activation of the dopaminergic VTA/SN associated with precision-weighted prediction error about stimulus category, ϵ_2 . This activation is shown both at $p < 0.05$ FWE whole-brain corrected (red) and $p < 0.05$ FWE corrected for the volume of our anatomical mask comprising both dopaminergic and cholinergic nuclei (yellow).

(A) Results from the first fMRI study.

(B) Second fMRI study.

(C) Conjunction (logical AND) across both studies.

HGF: empirical evidence (Iglesias et al., 2013)

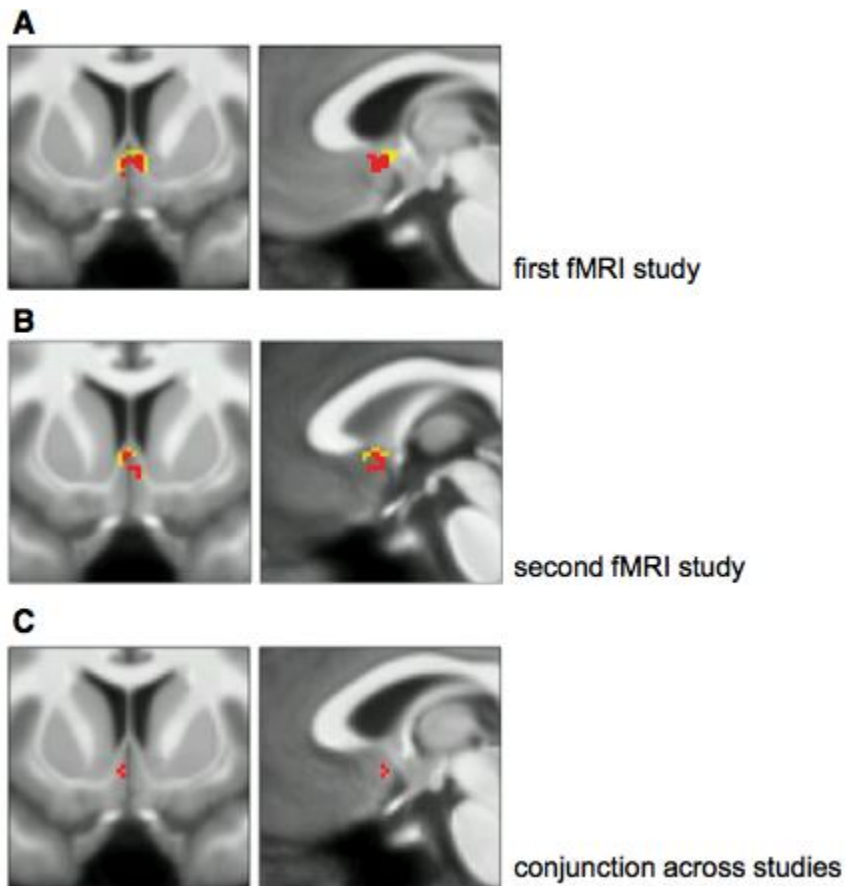


Figure 6. Basal Forebrain Activations by ε_3

Activation of the cholinergic basal forebrain associated with precision-weighted prediction error about stimulus probabilities ε_3 within the anatomically defined mask. For visualization of the activation area we overlay the results thresholded at $p < 0.05$ FWE corrected for the entire anatomical mask (red) on the results thresholded at $p < 0.001$ uncorrected (yellow) in the first (A: $x = 3, y = 9, z = -8$) and the second fMRI study (B: $x = 0, y = 10, z = -8$). (C) The conjunction analysis ("logical AND") across both studies ($x = 2, y = 11, z = -8$).

HGF: empirical evidence (Diaconescu et al., in preparation)

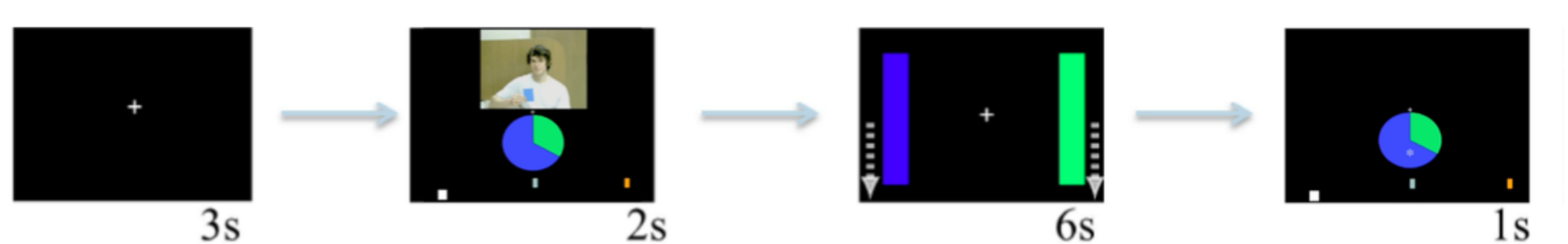
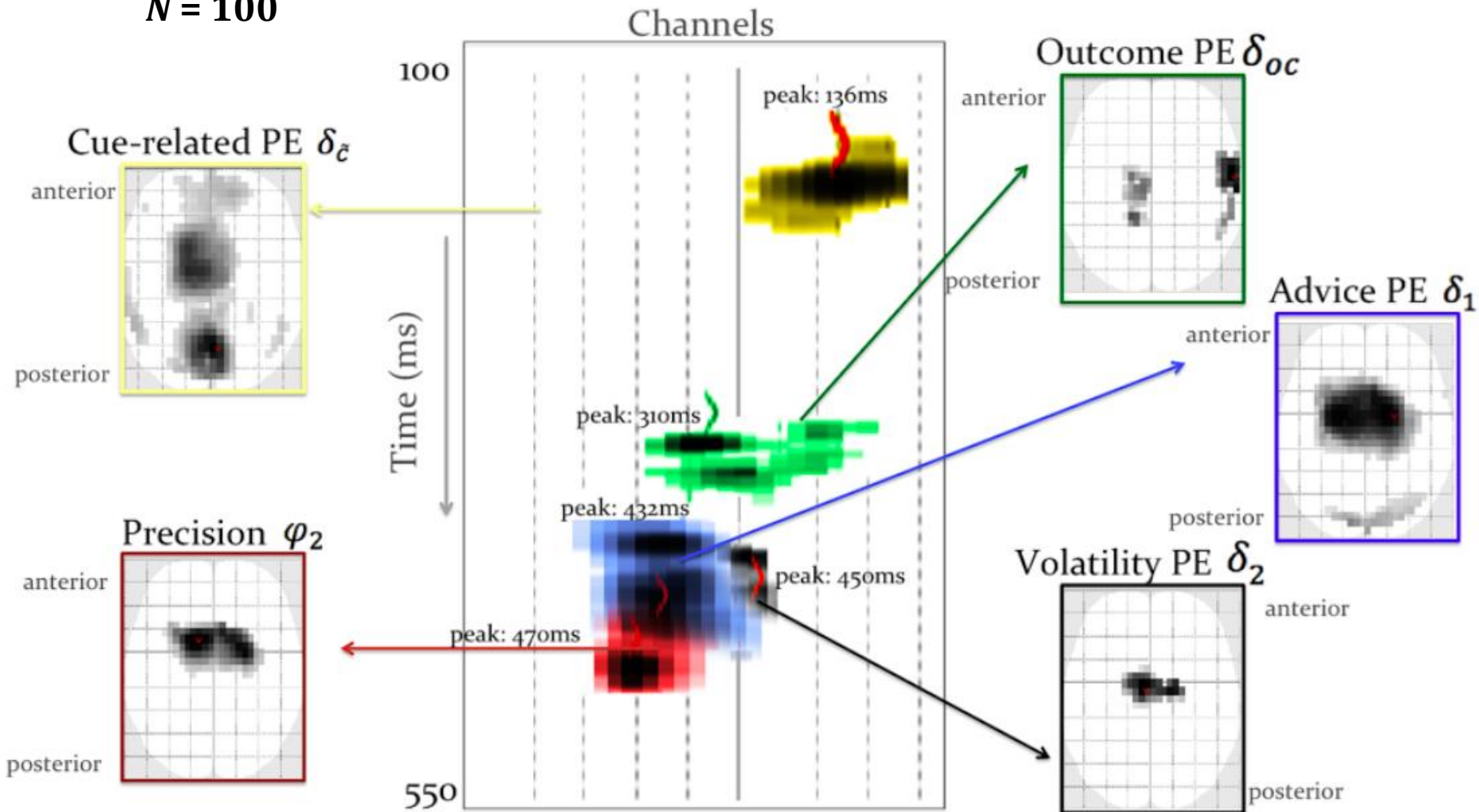
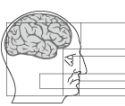


Fig. 1. Experimental Paradigm: 100 male volunteers played a binary lottery task and received advice about which option to choose from a more informed agent who was also incentivized to influence the participants' choices. To decide whether to take his advice into account, participants also inferred on the other's intentions and how they changed in time.

HGF: empirical evidence (Diaconescu et al., in preparation)

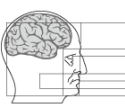
$N = 100$





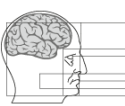
How to estimate and compare models: the HGF Toolbox

- Available at
<http://www.translationalneuromodeling.org/tapas>
- Start with README, manual, and interactive demo
- Modular, extensible
- Matlab-based



Summary

- We have to make good predictions to avoid surprise and survive, that is we have to use probabilistic (i.e., Bayesian) inference based on a good model of our environment.
- Bayesian inference means updating beliefs by uncertainty- (i.e., precision-) weighted prediction errors.
- Precision-weighting has to take account of all forms of uncertainty.
- A breakdown in this may be the root of many psychopathological phenomena.



Thanks

- Rick Adams
- Kay Brodersen
- Jean Daunizeau
- Andreea Diaconescu
- Chaohui Guo
- Karl Friston
- Sandra Iglesias
- Lars Kasper
- Ekaterina Lomakina
- Saeed Paliwal
- Klaas Enno Stephan
- Simone Vossel
- Lilian Weber

