

# **An Elementary Account of the Internal Model Principle**

**Matteo Capucci**

**University of Strathclyde / Independent ARIA Creator**

jww Baltieri, Biehl, Virgo  
[arxiv:2503.00511](https://arxiv.org/abs/2503.00511)

# The Classics

*Theorem* : The simplest optimal regulator  $R$  of a reguland  $S$  produces events  $R$  which are related to the events  $S$  by a mapping  $h : S \rightarrow R$ .

Restated somewhat less rigorously, the theorem says that the best regulator of a system is one which is a model of that system in the sense that the regulator's actions are merely the system's actions as seen through a mapping  $h$ . The type of isomorphism here is that expressed (in the form used above) by

$$\exists h : \forall i : \rho(i) = h[\sigma(i)] \quad (8)$$

where  $\rho$  and  $\sigma$  are the mappings that  $R$  and  $S$  impose on their common input  $I$ . This form is essentially that of (5) above.

# The Classics

(Francis and Wonham, 1976)

Our major conclusion of this section can be summarized as: ***The Internal Model Principle: A regulator synthesis is structurally stable only if the controller utilizes feedback of the regulated variable, and incorporates in the feedback path a suitably reduplicated model of the dynamic structure of the exogenous signals which the regulator is required to process.***

Francis, B. A., and W. M. Wonham. 1976. 'The Internal Model Principle of Control Theory'. *Automatica* 12 (5): 457–65.

[https://doi.org/10.1016/0005-1098\(76\)90006-6](https://doi.org/10.1016/0005-1098(76)90006-6).

# The Classics

**We say that a map  $A : \mathcal{X} \rightarrow \mathcal{X}$  incorporates an internal model of  $A_2$  if the minimal polynomial (m.p.) of  $A_2$  divides at least  $d(Z)$  invariant factors (i.f.) of  $A$ . Thus an internal model is an  $\ell$ -fold reduplication in  $A$  of the maximal cyclic component of  $A_2$ , where  $\ell \geq d(Z) =$  the number of independent outputs to be regulated.**

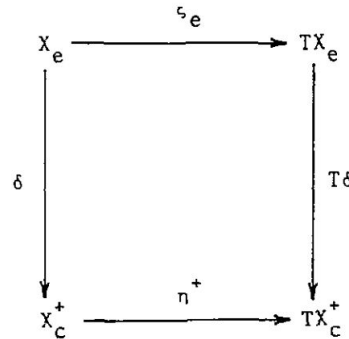
Francis, B. A., and W. M. Wonham. 1976. 'The Internal Model Principle of Control Theory'. *Automatica* 12 (5): 457–65.

[https://doi.org/10.1016/0005-1098\(76\)90006-6](https://doi.org/10.1016/0005-1098(76)90006-6).

# The Classics

## V. EXISTENCE OF INTERNAL MODELS

On the foregoing basis, it can now be established that the controller incorporates an internal model of some kind, corresponding to different observability conditions. The endomorph will be denoted by  $X^-$ , and  $\xi|X^-$  by  $\xi^+$ ,  $pr_1|X^+$  by  $pr_1^+$ ,  $pr_1^+(X^-)$  by  $X_c^+$ . The controller will be said to *incorporate an internal model of the exosystem*, if there are a diffeomorphism  $\delta: X_e \rightarrow X_c^- \subset X_c$  and a vector field  $\eta^+ = \eta|X_c^+$  such that the following diagram commutes.



Hepburn, J., and W. Wonham. 1984. 'Error Feedback and Internal Models on Differentiable Manifolds'. *IEEE Transactions on Automatic Control* 29 (5): 397–403. <https://doi.org/10.1109/TAC.1984.1103563>.

# Ours

- + (light) categorical footing
  - > more generality
  - > prone to compositionality
- + **clearer assumptions** regarding the structure of the systems
  - > we realize some assumptions are not needed!
- + an abstract notion of model
- + a clear link with Bayesian filtering

Still, not the last word & mostly subsumed by older work!

# Systems

**Definition II.1** (System). A *system* (or more precisely, a *fully observable system*)  $\mathsf{X}$  is comprised of a set  $X$  of *states*, a set  $I$  of *inputs* (or *observations*), and an *update* (or *dynamics*) function:

$$\text{upd}_{\mathsf{X}} : X \times I \rightarrow X, \quad (1)$$

The pair  $\left(\frac{I}{X}\right)$  is collectively referred to as the *interface* of the system, and we write  $\mathsf{X} : \mathbf{Sys}\left(\frac{I}{X}\right)$  to mean  $\mathsf{X}$  has such an interface.

Fully observable, discrete-time,  
(deterministic Moore) automata

# Systems

**Definition II.3** (Map of systems). Let  $X : \mathbf{Sys}(I_X)$  and  $X' : \mathbf{Sys}(I_{X'})$  be systems. A *map of systems*  $f : X \rightarrow X'$  is comprised of two parts:

1) a *map on states*, given by a function

$$f_s : X \rightarrow X', \quad (2)$$

2) a *map on inputs*, given by a function

$$f_i : X \times I \rightarrow I', \quad (3)$$

such that the following diagram commutes:

$$\begin{array}{ccc} X \times I & \xrightarrow{(\pi_X \natural f_s, f_i)} & X' \times I' \\ \text{upd}_X \downarrow & & \downarrow \text{upd}_{X'} \\ X & \xrightarrow{f_s} & X' \end{array} \quad (4)$$

meaning that, for every  $x \in X, i \in I$ , the following equation is satisfied:

$$f_s(\text{upd}_X(x, i)) = \text{upd}_{X'}(f_s(x), f_i(x, i)), \quad (5)$$

# Systems

**Construction II.4.** In what follows, we will need to compose maps of systems, so we explain here how that happens. Given maps  $f : X \rightarrow X'$  and  $g : X' \rightarrow X''$ , where  $X : \mathbf{Sys}(I_X)$ ,  $X' : \mathbf{Sys}(I_{X'})$  and  $X'' : \mathbf{Sys}(I_{X''})$ , their composition  $f \circ g : X \rightarrow X''$  is given on states by the composition of the maps on states:

$$(f \circ g)_s = f_s \circ g_s, \quad (6)$$

while on inputs  $(f \circ g)_i : X \times I \rightarrow I''$  is defined as follows:

$$(f \circ g)_i(x, i) = g_i(f_s(x), f_i(x, i)). \quad (7)$$

$$X \xrightarrow{f_s} X' \xrightarrow{g_s} X''$$

$$X \times I \xrightarrow{(f_s, f_i)} X' \times I' \xrightarrow{g_i} I''$$

# Systems

**Definition II.5** (Subsystem). A *subsystem* of  $X$  is a forward-invariant subset of  $X$  together with updates restricted to the subset, thus a map of systems  $\gamma : X' \rightarrow X$  given on states by an inclusion  $\gamma_s : X' \rightarrow X$  and on inputs by projection  $\gamma_i = \pi_I : X' \times I \rightarrow I$ .

**Definition II.6** (Attracting subsystem). An *attracting subsystem* for  $X$  is a subsystem  $X^* \rightarrow X$  such that, for each  $x \in X$ , there exists  $n \in \mathbb{N}$  such that for all  $t \geq n$  and  $\underline{i} \in I^t$ , we have  $\text{upd}_X^t(x, \underline{i}) \in X^*$ .

# IMP setup

**Assumption 1** (Environment, plant, controller). *The following three components are so defined:*

1) *the environment  $E : \mathbf{Sys}(\frac{1}{E})$  is an autonomous system*

$$\text{upd}_E : E \rightarrow E, \quad (8)$$

2) *the plant  $P : \mathbf{Sys}(\frac{E \times C}{P})$  is a system*

$$\text{upd}_P : P \times E \times C \rightarrow P, \quad (9)$$

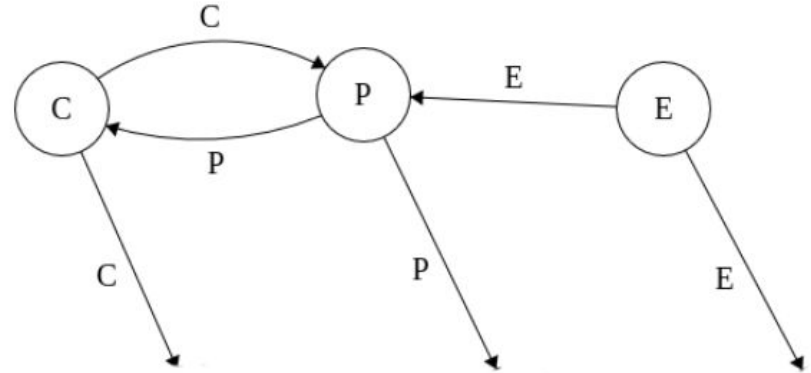
3) *the controller  $C : \mathbf{Sys}(\frac{P}{C})$  is a system*

$$\text{upd}_C : C \times P \rightarrow C. \quad (10)$$

*The full system  $S : \mathbf{Sys}(\frac{1}{E \times P \times C})$  is the following composite autonomous system:*

$$\begin{aligned} \text{upd}_S : E \times P \times C &\longrightarrow E \times P \times C \\ (s_E, s_P, s_C) &\longmapsto (\text{upd}_E(s_E), \text{upd}_P(s_P, s_E, s_C), \\ &\quad \text{upd}_C(s_C, s_P)). \end{aligned} \quad (11)$$

*Let  $S = E \times P \times C$  denote the state space of  $S$ .*



# IMP setup

**Assumption 1** (Environment, plant, controller). *The following three components are so defined:*

1) *the environment  $E : \mathbf{Sys}(\overset{1}{E})$  is an autonomous system*

$$\text{upd}_E : E \rightarrow E, \quad (8)$$

2) *the plant  $P : \mathbf{Sys}(\overset{E \times C}{P})$  is a system*

$$\text{upd}_P : P \times E \times C \rightarrow P, \quad (9)$$

3) *the controller  $C : \mathbf{Sys}(\overset{P}{C})$  is a system*

$$\text{upd}_C : C \times P \rightarrow C. \quad (10)$$

*The full system  $S : \mathbf{Sys}(\overset{1}{E \times P \times C})$  is the following composite autonomous system:*

$$\begin{aligned} \text{upd}_S : E \times P \times C &\longrightarrow E \times P \times C \\ (s_E, s_P, s_C) &\longmapsto (\text{upd}_E(s_E), \text{upd}_P(s_P, s_E, s_C), \\ &\quad \text{upd}_C(s_C, s_P)). \end{aligned} \quad (11)$$

*Let  $S = E \times P \times C$  denote the state space of  $S$ .*

**Remark II.7.** There are maps of systems  $\pi_E : S \rightarrow E$ ,  $\pi_P : S \rightarrow P$  and  $\pi_C : S \rightarrow C$  induced by projecting out states of  $S$ :

$$\begin{array}{ccc} S \times 1 & \xrightarrow{(\pi_E, \text{id}_1)} & E \times 1 & & S \times 1 & \xrightarrow{(\pi_C, \pi_P)} & C \times P \\ \text{upd}_S \downarrow & & \downarrow \text{upd}_E & & \text{upd}_S \downarrow & & \downarrow \text{upd}_C \\ S & \xrightarrow{\pi_E} & E & & S & \xrightarrow{\pi_C} & C \end{array}$$

$$\begin{array}{ccc} S \times 1 & \xrightarrow{(\pi_P, \pi_{E \times C})} & P \times E \times C \\ \text{upd}_S \downarrow & & \downarrow \text{upd}_P \\ S & \xrightarrow{\pi_P} & P \end{array} \quad (12)$$

# IMP setup

**Definition II.8** (Regulation problem). A *regulation problem* (or *reguland* [1]) is a triple  $(E, P, \tilde{K} \subseteq E \times P)$ .

**Assumption 2** (Regulation condition).  $C$  solves its regulation problem, meaning there exists an attracting subsystem (Definition II.6)  $S^* \rightarrow S$  such that, on states,  $S^* \subseteq K$ .

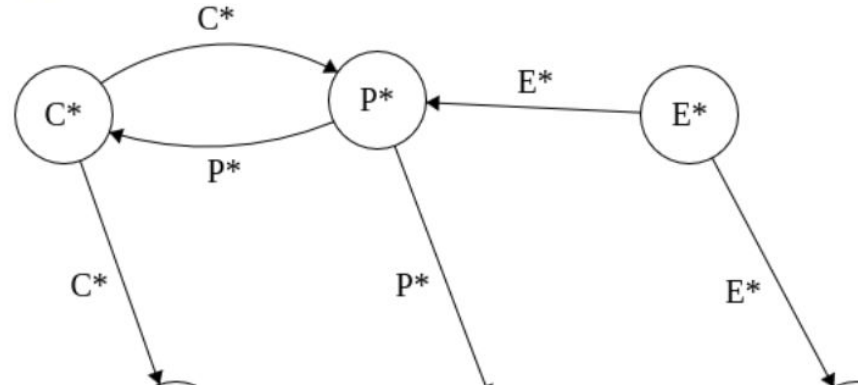
# IMP setup

Even though  $K$  contains all states of  $C$ , this might not be true anymore for  $S^*$ . We can define the set of *attracting control states* as

$$C^* := \{s_C \in C \mid \exists s_E \in E, s_P \in P, (s_E, s_P, s_C) \in S^*\}.$$

It is obtained along with a surjection  $\pi_{C^*} : S^* \rightarrow C^*$  by taking the image of  $S^*$  under the projection map  $\pi_C : S \rightarrow C$ :

$$\begin{array}{ccc} S^* & \xrightarrow{\pi_{C^*}} & C^* \\ \downarrow & & \downarrow \\ S & \xrightarrow{\pi_C} & C \end{array}$$



# Model

**Definition II.9** (Model). A *model* of a system  $X \in \mathbf{Sys}(X^I)$  is:

- a system  $M \in \mathbf{Sys}(M^J)$  (the *archetype*), and
- a map of systems (the *model per se*)

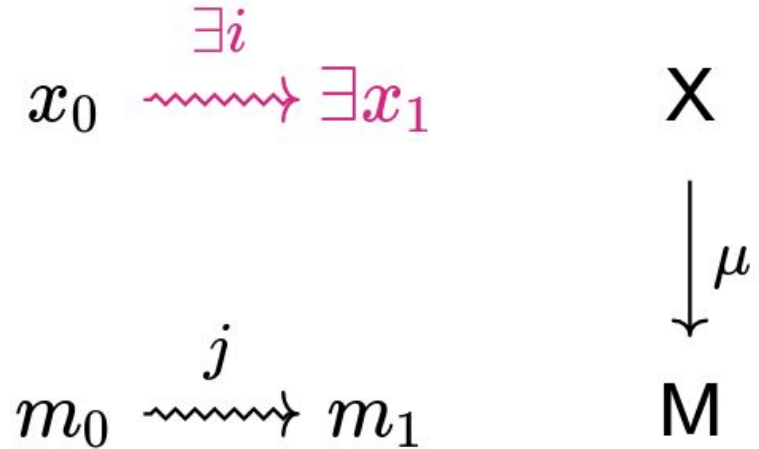
$$X \xrightarrow{\mu} M \quad (15)$$

such that

- 1) its part on states  $\mu_s : X \rightarrow M$  is surjective, and
- 2) its part on inputs  $\mu_i(x, -) : I \rightarrow J$  is surjective for each  $x \in X$ .

# Model

Having a model  $\mu : X \rightarrow M$  means that for each state  $m \in M$  we have a set of states  $\mu^{-1}(m) \in X$ , called the *fibre* of  $m$  and which represents a subset of elements of  $X$  that are *indistinguishable* from the perspective of the simpler system  $M$  as they all map to the same element  $m \in M$  via the surjective function  $\mu$ . As  $m \in M$  varies along the function  $\text{upd}_M$ , this variation is consistent with the variation described by the function  $\text{upd}_X$  for each element  $x$  of the fibre  $\mu^{-1}(m)$  of  $m$ .



everything we do in  $M$  can be reproduced in  $X$

# Model

$$\begin{array}{ccc} M & \xrightarrow{\mu^{-1}} & P^+ X \\ \text{upd}_M \downarrow & \subseteq & \downarrow P^+ \text{upd}_X \\ M & \xrightarrow{\mu^{-1}} & P^+ X \end{array}$$

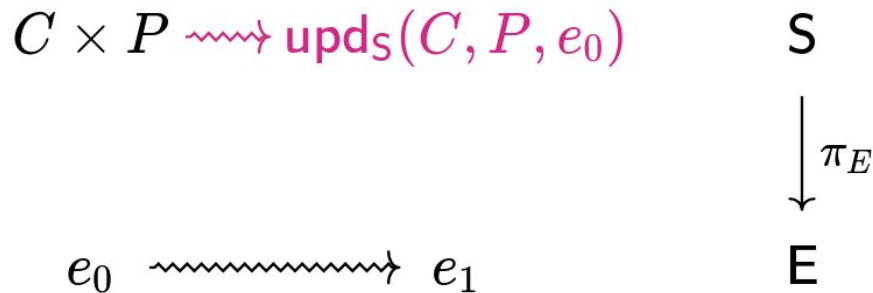
The fibers of models are ‘dynamical beliefs’

# Model

**Remark II.11.** Like any good definition, Definition II.9 admits a trivial instance. A *trivial model* is one where  $\mu : X \rightarrow M$  is a product projection, which means there exists a system  $F \in \mathbf{Sys}_F^H$  such that the state space and inputs of the system  $X$  factor as  $X = M \times F$  and  $I = J \times H$  respectively, and its dynamics decomposes in that of  $M$  and  $F$ :

$$\text{upd}_X((m, f), (j, h)) = (\text{upd}_M(m, j), \text{upd}_F(f, h)). \quad (17)$$

Thus, in a trivial model, knowledge of  $M$  does not afford knowledge about the rest of  $X$ , because  $M$  and  $F$  are not coupled.



# Model

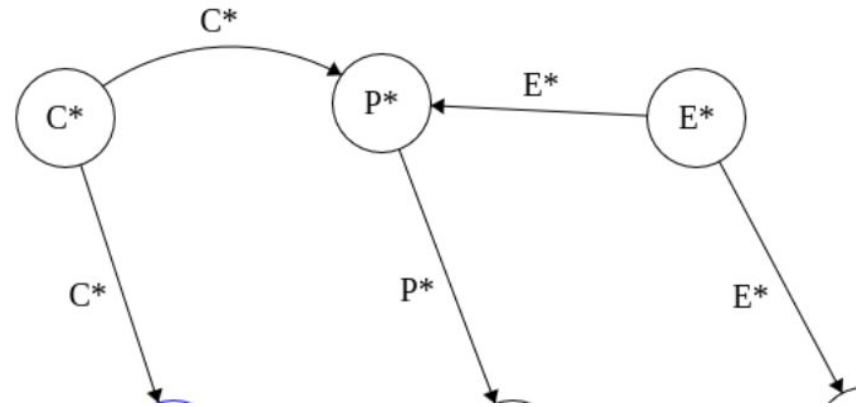
$$\begin{array}{ccc} P^* \times E^* & \xrightarrow{\quad ? \quad} & S^* \\ & & \downarrow \pi_C^* \\ C_0^* & \xrightarrow{p^*} & C_1^* \end{array}$$

Doesn't work!

# IMP setup

**Assumption 3** (Error feedback structure).  $C^*$  supports an autonomous dynamics  $\text{upd}_{C^*} : C^* \rightarrow C^*$ , thus defining a system  $C^* : \text{Sys}(C^*)$ , that we call attracting controller, and making  $\pi_{C^*}$  a full-fledged map of systems:

$$S^* \xrightarrow{\pi_{C^*}} C^*$$



# IMP v1

**Theorem II.12** (*Internal Model Principle*). *Let  $S$  be a system subject to Assumptions 1 to 3. Then  $C^*$  models the attracting full system  $S^*$  via the projection  $\pi_{C^*}$ , as implied by Assumption 3.*

$$\begin{array}{ccc} P^* \times E^* & \rightsquigarrow \text{upd}_S(c_0^*, P^*, E^*) & S^* \\ & & \downarrow \pi_{C^*} \\ c_0^* & \rightsquigarrow c_1^* & C^* \end{array}$$

# IMP v2 – Wonham & Hepburn’s version

**Assumption 4.** *There is an isomorphism of systems  $S^* \cong E^*$ , meaning that for each environment state  $s_E \in E^*$ , there is exactly one  $s \in S^*$  such that  $\pi_E(s) = s_E$ .*

“endomorph”

**Theorem II.14** (*Internal Model Principle (Hepburn and Wonham)*). *Let  $S$  be a system subject to Assumptions 1 to 4. Then  $C^*$  models the attracting environment  $E^*$  via the dashed map below:<sup>2</sup>*

$$\begin{array}{ccc} E^* & \overset{\nu}{\dashrightarrow} & C^* \\ \text{Assumption 4} \curvearrowright \sim & & \nearrow \pi_{C^*} \text{ Assumption 3} \\ & S^* & \end{array} \quad (18)$$

Thanks!

# Bayesian interpretation

We now introduce the main ideas involved in the notion of Bayesian interpretations. To aid intuition, we refer to maps of type  $p : I \multimap X$  for any object  $X$  as *distributions* and to deterministic maps  $x : I \rightarrow X$  of that type as *elements* of  $X$ . We also consider maps  $\psi : \Theta \multimap X$  to correspond to *parametrised families of distributions* (also called “channels” in [34], [65]), with each element  $\theta : I \rightarrow \Theta$  a parameter determining a distribution  $\psi(X | \theta)$  or  $\psi(\theta)$  over  $X$ .

Furthermore, given a map  $f : X \multimap Y$  and an element  $x : I \rightarrow Y$  we call  $x \circ f : I \multimap Y$  the distribution over  $Y$  assigned to the element  $x \in X$  by  $f$ . We write it as  $f(x)$  or  $f(Y | x)$ . In **FinStoch** this terminology and the notation coincide with its common usage in probability theory.

over it. There is a Markov category **FinStoch** whose objects are finite sets and whose morphisms from  $X$  to  $Y$ ,  $X \multimap Y$ , called Markov kernels, are defined as functions  $X \rightarrow D(Y)$ . For a Markov kernel  $f : X \multimap Y$  we write  $f(y | x)$  for  $f(x)(y)$ , the probability assigned to  $y$  when the kernel is given  $x$  as input.

Sequential composition of Markov kernels  $f : X \multimap Y$ ,  $g : Y \multimap Z$  is given by the Markov kernel so defined (i.e. the Chapman–Kolmogorov equation):

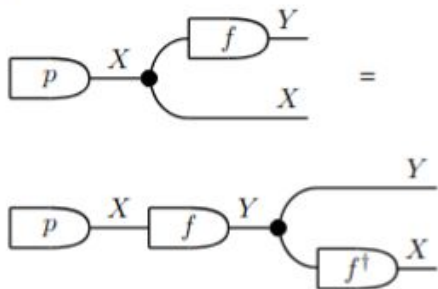
$$f \circ g : X \multimap Z$$
$$(f \circ g)(z | x) := \sum_{y \in Y} g(z | y) f(y | x) \quad (34)$$

There is a Markov category **Rel**<sup>+</sup> whose objects are sets and whose morphisms from  $X$  to  $Y$  are functions  $X \rightarrow P^+(Y)$ , corresponding to left-total relations. Sequential composition of left-total relations  $f : X \multimap Y$ ,  $g : Y \multimap Z$  is given by the left-total relation so defined:

$$f \circ g : X \longrightarrow P^+(Z)$$
$$x \longmapsto \{z \in Z \mid \exists y \in Y, y \in f(x) \text{ and } z \in g(y)\}, \quad (36)$$

# Bayesian interpretation

**Definition III.6** (Bayesian inversion). In a Markov category, a Bayesian inversion [34] of a map  $f : X \rightarrow Y$  with respect to a distribution  $p : I \rightarrow X$  is a map  $f^\dagger : Y \rightarrow X$  satisfying the following equation:



(38)

To get a more concrete feeling for Bayesian inverses, we can examine the form Eq. (38) takes when our Markov category is **FinStoch**, where it corresponds to the equation

$$\begin{aligned} p(x)f(y | x) &= f(y | p(x))f^\dagger(x | y) \\ &= \sum_{x' \in X} p(x')f(y | x')f^\dagger(x | y), \end{aligned} \quad (39)$$

from which we can derive the standard Bayes rule by dividing by  $f(y | p(x))$ , assuming it's positive ( $f(y | p(x)) > 0$ ):

$$\begin{aligned} f^\dagger(x | y) &= \frac{p(x)f(y | x)}{f(y | p(x))} \\ &= \frac{p(x)f(y | x)}{\sum_{x' \in X} p(x')f(y | x')}. \end{aligned} \quad (40)$$

# Bayesian interpretation

**Definition III.7.** In a Markov category, we say that a map  $f^\dagger : Y \otimes \Theta \rightarrow X$  is a Bayesian inversion of a parametrised family  $f : \Theta \rightarrow X$  if

$$\begin{array}{c} \Theta \xrightarrow{\psi} X \begin{array}{l} \nearrow f \rightarrow Y \\ \searrow \rightarrow X \end{array} = \\ \Theta \xrightarrow{\psi} X \xrightarrow{f} Y \begin{array}{l} \nearrow \rightarrow Y \\ \searrow f^\dagger \rightarrow X \end{array} \end{array} \quad (41)$$



# Bayesian interpretation

We need then a form of belief updating that corresponds to *Bayesian filtering*, where the hidden variable also changes. For this we can replace  $\Delta_X \circ f \otimes \text{id}_X : X \rightarrow Y \otimes X$  with a map  $\kappa : X \rightarrow X \otimes Y$  that produces an observation  $Y$  and may change  $X$  instead of just copying it. The analogue of the Bayesian inversion (Definition III.6) with respect to a distribution  $p : I \rightarrow X$  is thus a map  $\kappa^\dagger : Y \rightarrow X$  satisfying

$$(43)$$

$$\kappa^\dagger(x | y) = \frac{\sum_{x' \in X} p(x') \kappa(y, x | x')}{\sum_{x', x'' \in X} p(x') \kappa(y, x'' | x')}.$$

As with ordinary Bayesian inversions (Definition III.6), the map  $\kappa^\dagger$  always exists in any Markov category with conditionals, including thus **FinStoch** and **Rel**<sup>+</sup>. Note that setting  $\kappa = \Delta_X \circ f \otimes \text{id}_X : X \rightarrow Y \otimes X$  recovers Definition III.6.

In Bayesian filtering, the idea is that  $\kappa$  updates  $X$  and generates an observation  $Y$ . Similarly to the Bayesian inference



# Bayesian interpretation

**Definition III.10** (Bayesian filtering interpretation [33]). Given a deterministic map  $c : Y \otimes \Theta \rightarrow \Theta$ , a *Bayesian filtering interpretation* of  $c$  consists of a map  $\psi : \Theta \rightarrow X$  called the *interpretation map*, together with a map  $\kappa : X \rightarrow X \otimes Y$  called *hidden Markov model*, such that Eq. (45) holds. In this context, Eq. (45) is called the *consistency equation*. A map  $c : Y \otimes \Theta \rightarrow \Theta$  together with an interpretation  $(\psi, \kappa)$  is called a *reasoner*.

The diagrammatic equation (45) consists of two parts separated by an equals sign. The top part shows a wire for  $\Theta$  entering a box labeled  $\psi$  with an output wire for  $X$ . This  $X$  wire then enters a box labeled  $\kappa$  with two output wires: one for  $Y$  and one for  $X$ . The bottom part shows a wire for  $\Theta$  entering a box labeled  $\psi$  with an output wire for  $X$ . This  $X$  wire enters a box labeled  $\kappa$  with two output wires: one for  $X$  and one for  $Y$ . The  $Y$  wire from the  $\kappa$  box enters a box labeled  $c$  along with the  $\Theta$  wire. The output of the  $c$  box is a  $\Theta$  wire that loops back to the input of the  $\psi$  box. The output of the  $\psi$  box in the bottom part is a wire for  $X$ .

# Bayesian interpretation

**Theorem IV.5.** Let  $M$  model  $X$  with  $\mu : X \rightarrow M$ , and assume  $M$  and  $X$  are autonomous. Define  $c : X \otimes M \rightarrow M$  as

$$\begin{array}{c} X \\ M \end{array} \begin{array}{|c|} \hline c \\ \hline \end{array} M \quad := \quad \begin{array}{c} X \\ M \end{array} \bullet \begin{array}{|c|} \hline \text{upd}_M \\ \hline \end{array} M \quad (52)$$

and  $\kappa : X \rightarrow X \otimes X$  as:

$$\begin{array}{c} X \end{array} \begin{array}{|c|} \hline \kappa \\ \hline \end{array} \begin{array}{c} X \\ X \end{array} \quad := \quad \begin{array}{c} X \\ X \end{array} \bullet \begin{array}{|c|} \hline \overline{\text{upd}}_X \\ \hline \end{array} \begin{array}{c} X \\ X \end{array} \quad (53)$$

Then  $\kappa$  is the hidden Markov model, and  $\mu_s^{-1} : M \rightarrow X$  the interpretation map of a Bayesian filtering interpretation of  $c$ , i.e. we have:

$$\begin{array}{c} M \\ \mu_s^{-1} \end{array} X \bullet \begin{array}{|c|} \hline \overline{\text{upd}}_X \\ \hline \end{array} \begin{array}{c} X \\ X \end{array} \quad = \quad \begin{array}{c} M \\ \mu_s^{-1} \end{array} X \bullet \begin{array}{|c|} \hline \overline{\text{upd}}_X \\ \hline \end{array} \begin{array}{c} X \\ X \end{array} \bullet \begin{array}{|c|} \hline \text{upd}_M \\ \hline \end{array} M \bullet \begin{array}{|c|} \hline \mu_s^{-1} \\ \hline \end{array} X \quad (54)$$

where the dashed lines show, informally, where we replaced the definitions above in Eq. (45).

Thanks again!