



# Integrated Information Theory (IIT) with Simple Maths (slides)

David Rudrauf, Tonglin Yan, Nils Ruet, Kenneth Williford, Grégoire Sergeant-Perthuis

## ► To cite this version:

David Rudrauf, Tonglin Yan, Nils Ruet, Kenneth Williford, Grégoire Sergeant-Perthuis. Integrated Information Theory (IIT) with Simple Maths (slides). 27th annual meeting of the Association for the Scientific Study of Consciousness (ASSC 27), Ryota Kanai, Jul 2024, Tokyo, Japan. hal-04636522

**HAL Id: hal-04636522**

**<https://hal.sorbonne-universite.fr/hal-04636522v1>**

Submitted on 5 Jul 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# Integrated Information Theory (IIT) with Simple Maths

David Rudrauf, **Tonglin Yan**<sup>1</sup>, Nils Ruet, Kenneth Williford,  
**Grégoire Sergeant-Perthuis**<sup>2</sup>

[1] CIAMS, Université Paris-Saclay, [2] LCQB, Sorbonne Université

ASSC 27  
Tokyo, 4 July, 2024

For the 'preprint' go here: [here](#)

- What we do:
  - **Reinforcement learning** + **phenomenological** aspects of consciousness

## Reproduce the experience of *space* in 'robots'

*Consciousness involves a subjective perspective, characterized by viewpoint-structured organization, a sense of unity (holistic world), embodiment, and an internal representation of the world in perspective from a specific standpoint.*

### ↪ the **Projective Consciousness Model**

- Initiated by D. Rudrauf, K. Williford, D. Bennequin, K. Friston [RBG<sup>+</sup>17, WBFR18, RBW20]

→ See K. Williford's MoC4 presentation of phenomenological motivation  
<https://www.youtube.com/watch?v=eHyVZWZMqzg&t=853s>

- What do I do?
  - Background: mathematical physics (PhD)
  - Machine learning  $\cap$  geometry → computational biology
  - ↪ Embarked on the *consciousness* adventure
    - ↔ Geometric structure of the world model [RSPB<sup>+</sup>20, SPRR<sup>+</sup>23]
- Why this work?
  - Disseminate formal and computational models of consciousness
  - Simplify entering into the mathematical details of the models: IIT, active inference, etc.

Based on a seminar given at:

- PMMC2: Paris Mathematical Models of Cognition and Consciousness (PMMC2, [link to the seminar page](#))
  - On mathematical and computational models of consciousness and cognition
  - Longer version on YouTube: [link to the talk](#)
  - ↪ Come and give a talk

Today's presentation is a simplification of the excellent presentation:

- *The mathematical structure of integrated information theory*, Johannes Kleiner and Sean Tull. In Frontiers in Applied Mathematics and Statistics, 2020.

## Definition (Probability measures)

Let  $E$  be a finite set, denote by  $\mathbb{P}(E)$  the set of probability measures on  $E$ ,

$$p \in \mathbb{P}(E) \iff \forall x \in E, p(x) \geq 0 \text{ and } \sum_{x \in E} p(x) = 1 \quad (0.1)$$

## Definition (Markov kernels (stochastic maps))

A Markov kernel or stochastic map  $T$  from  $E$  to  $F$ , denoted as  $T : E \rightarrow \mathbb{P}(F)$ , sends any point  $x \in E$  to a probability measure  $T_x \in \mathbb{P}(F)$ . For  $x \in E$  and  $y \in F$ , we will denote  $T_x(y)$  as  $T(y|x)$ .

- We will be considering multiple variables ( $X_i, i \in S$ ) denotes  $X_S$
- When  $S$  is a finite set:
  - enumerate  $S$ ,  $S \simeq [1, \dots, N]$
- We will consider sub-collections of variable:
  - $a \subseteq S$
  - $\bar{a}$  the complement of  $a$  in  $S$ ,

$$a \cup \bar{a} = S \qquad a \cap \bar{a} = \emptyset$$

- a subset  $a \subseteq S$  is associates to a random variable ( $X_i, i \in a$ )
  - $X_i$  takes values in  $E_{X_i}$
  - Denote  $(X_i, i \in a)$  as  $X_a$
  - $X_a$  takes values in  $\prod_{i \in a} E_{X_i}$  denoted as  $E_{X_a}$

## Definition (Conditional expectation)

Let  $(Y_i, i \in S_1)$  be a collection of variables taking values in  $F = \prod_{i \in S_1} F_{Y_i}$ . Each  $F_{Y_i}$  is a finite set. Let  $p \in \mathbb{P}(F)$ . For any  $a \subseteq S_1$ , and any function  $f : F \rightarrow \mathbb{R}$ , one defines the conditional expectation with respect to  $Y_a$  as;

$$\forall y_a \in F_a, \quad \mathbb{E}[f | Y_a](y_a) = \sum_{y_{\bar{a}} \in Y_{\bar{a}}} \frac{f(y_{\bar{a}}, y_a) p(y_{\bar{a}}, y_a)}{\sum_{y_{\bar{a}} \in Y_{\bar{a}}} p(y_{\bar{a}}, y_a)}$$

Example with two random variables  $X, Y$ , assume:

- $X \in E, Y \in F$ , where  $E, F$  are finite sets
- $x \in E, y \in F$

then,

- $\mathbb{E}[Y = y | X = x] := P(y|x)$
- and

$$P(y|x) = \frac{P(X = x, Y = y)}{P(X = x)} \quad (\text{Bayes' Rule})$$



We want to see the effect of the stochastic dynamic  $T$  on sub-collections of variables.

- We start with  $T : E_{X_1} \times \dots \times E_{X_N} \rightarrow \mathbb{P}(F_{Y_1} \times \dots \times F_{Y_M})$
- Choose  $a = (1, \dots, n_1)$  and  $b = (1, \dots, m_1)$  with  $n_1 \leq N$  and  $m_1 \leq M$
- How to deduce a transition

$$T^{a,b} : E_{X_1} \times \dots \times E_{X_{n_1}} \rightarrow \mathbb{P}(F_{Y_1} \times \dots \times F_{Y_{m_1}})$$

**Definition** (Building a Markov kernel  $E_{X_a} \rightarrow \mathbb{P}(F_{Y_b})$  from a prior  $Q$ )

Let  $X_S := (X_i, i \in S)$ ,  $Y_{S_1} := (Y_i, i \in S_1)$  and  $E = \prod_{i \in S} E_{X_i}$ ,  
 $F = \prod_{i \in S_1} F_{Y_i}$ .

Let  $T = E \rightarrow \mathbb{P}(F)$  be a Markov kernel. For any  $a \subseteq S$  and  $b \subseteq S_1$ , a choice of  $Q \in \mathbb{P}(E)$  allows us to derive from  $T$  the kernel denoted  $T^{Q,a,b} : E_{X_a} \rightarrow \mathbb{P}(F_{Y_b})$ , which encodes the effect of the variables  $X_a$  on  $Y_b$ . It is defined as,

$$\forall y_b \in F_{Y_b}, \forall x_a \in E_{X_a} \quad T^{Q,a,b}(y_b|x_a) := \mathbb{E}[Y_b = y_b | X_a = x_a]$$

– How?  $\hookrightarrow$  Joint distribution,

$$P(Y_{S_1} = y, X_S = x) := T(y|x) \times Q(x)$$

$\hookrightarrow$  Sum out (*marginalize*)  $X_{\bar{a}}, Y_{\bar{b}}$

## ***‘Cutting’ interactions: the central operation***

$$X_1 \dots X_{n_1}$$

$$X_{n_1+1} \dots X_N$$

$$X_1 \dots X_N$$

$$\begin{array}{c} a \\ \Downarrow \tau^{a,b} \\ b \end{array}$$

$$\begin{array}{c} \bar{a} \\ \Downarrow \tau^{\bar{a},\bar{b}} \\ \bar{b} \end{array}$$

$$\begin{array}{c} S \\ \Downarrow T \\ S_1 \end{array}$$

$$Y_1 \dots Y_{m_1}$$

$$Y_{m_1+1} \dots Y_M$$

$$Y_1 \dots Y_M$$

Left: Cutting the interactions. Right: Overall interaction.

- We quantify the effect of ‘cutting’ interactions between variables
- $(T^{a,b}, T^{\bar{a},\bar{b}})$  should be in the same space as  $T$

### Definition (Product of local kernels)

For any two probability kernels,  $T^{a,b} : E_{X_a} \rightarrow \mathbb{P}(F_{Y_b})$  and  $T^{\bar{a},\bar{b}} : E_{X_{\bar{a}}} \rightarrow \mathbb{P}(F_{Y_{\bar{b}}})$  posit,

$$(T^{a,b} \otimes T^{\bar{a},\bar{b}})(y|x) := T^{a,b}(y_b|x_a) \cdot T^{\bar{a},\bar{b}}(y_{\bar{b}}|x_{\bar{a}}) \quad (0.2)$$

→ We want to quantify how far  $T^{a,b} \otimes T^{\bar{a},\bar{b}}$  is from  $T$

### Definition (Informal definition of divergence)

For a finite space  $Y$ , we define a divergence  $D$  on  $\mathbb{P}(Y)$  as a function  $D : \mathbb{P}(Y) \times \mathbb{P}(Y) \rightarrow \mathbb{R}_{\geq 0}$  such that, for any two probability distributions  $P$  and  $P_1$  in  $\mathbb{P}(Y)$ ,  $D(P, P_1)$  decreases as the two distributions  $P$  and  $P_1$  get 'closer'; and it reaches its minimum value of 0 when and only when  $P = P_1$ .

- The dissimilarity is on *probability distributions*.
  - Fix  $x \in E$  then,  $T^{a,b} \otimes T^{\bar{a},\bar{b}}(.|x)$  is a probability distribution.
  - Similarly  $T(.|x) \in \mathbb{P}(F)$
- For a fixed  $x$
- ↪ denote  $T^{a,b} \otimes T^{\bar{a},\bar{b}}(.|x)$  as  $T_{x_a}^{a,b} \otimes T_{x_{\bar{a}}}^{\bar{a},\bar{b}}$
  - ↪ denote  $T(.|x)$  as  $T_x$

- Little  $\varphi_e$  focusing on effects.
  - $S \simeq [1, \dots, N]$  and  $X_1 \dots X_N$
  - $M \subseteq S$ , 'old'  $X_S$  is now  $X_M$
  - $P \subseteq S$  'old'  $Y_{S_1}$  is now  $X_P$
  - To remember that we start with  $X_M$  and we go to  $X_P$  denote  $T_M^P$  the associated kernel

## Definition

For any  $M, P \subseteq S$  and  $x_M \in X_M$ ,

$$\varphi_{M, x_M}^P := \inf_{\substack{a \subseteq M \\ b \subseteq P}} D(T_{M, x_M}^P | T_{M, x_a}^{P, (a, b)} \otimes T_{M, x_{\bar{a}}}^{P, (\bar{a}, \bar{b})}) \quad (0.3)$$

And

$$\varphi_{M, x_M}^* := \max_{P \subseteq S} \varphi_{M, x_M}^P \quad (0.4)$$

→ One more step to compute *big*  $\Phi$ , focusing on effects (see [here](#))

$$\Phi_{M,x,b} := \sum_{a \subseteq M} \varphi_{a,x_a}^* D(T_{a,x_a}^{\psi(a,x_a)} | T_{a,x_a}^{b_a} \otimes T_{a,x_a}^{\bar{b}_a})$$

with  $\psi(M, x) := \mathbf{argmax} \varphi_M^P$

$$\Phi_{M,x} = \mathbf{argmin}_{b \subseteq \psi(M, x_M)} \Phi_{M,x,b}$$

→ IIT 4.0 introduces an important difference in how information is quantified  $\rightsquigarrow$  dissimilarity function.

→ Go check their paper [ABF<sup>+</sup>23].

# Thank you very much for your attention

Thank you very much for your attention!



# References I



Larissa Albantakis, Leonardo Barbosa, Graham Findlay, Matteo Grasso, Andrew M. Haun, William Marshall, William G. P. Mayner, Alireza Zaeemzadeh, Melanie Boly, Bjørn E. Juel, Shuntaro Sasai, Keiko Fujii, Isaac David, Jeremiah Hendren, Jonathan P. Lang, and Giulio Tononi, *Integrated information theory (iit) 4.0: Formulating the properties of phenomenal existence in physical terms*, PLOS Computational Biology **19** (2023), no. 10, 1–45.



D. Rudrauf, D. Bennequin, I. Granic, G. Landini, , K. Friston, and K. Williford, *A mathematical model of embodied consciousness*, Journal of Theoretical Biology (2017).



David Rudrauf, Daniel Bennequin, and Kenneth Williford, *The Moon illusion explained by the projective consciousness model*, Journal of Theoretical Biology (2020).

# References II



David Rudrauf, Grégoire Sergeant-Perthuis, Olivier Belli, Yvain Tisserand, and Giovanna Di Marzo Serugendo, *Modeling the subjective perspective of consciousness and its role in the control of behaviours*, arXiv:2012.12963, 2020.



Grégoire Sergeant-Perthuis, Nils Ruet, David Rudrauf, Dimitri Ognibene, and Yvain Tisserand, *Influence of the geometry of the feature space on curiosity based exploration*, NeurIPS 2023 workshop: Information-Theoretic Principles in Cognitive Systems, 2023.



Kenneth Williford, Daniel Bennequin, Karl Friston, and David Rudrauf, *The projective consciousness model and phenomenal selfhood*, Frontiers in Psychology (2018).