

# Integrated information theory (IIT) 4.0: Formulating the properties of phenomenal existence in physical terms

Larissa Albantakis<sup>1¶</sup>, Leonardo Barbosa<sup>1,2¶</sup>, Graham Findlay<sup>1,3¶</sup>, Matteo Grasso<sup>1¶</sup>, Andrew M Haun<sup>1¶</sup>, William Marshall<sup>1,4¶</sup>, William GP Mayner<sup>1,3¶</sup>, Alireza Zaeemzadeh<sup>1¶</sup>, Melanie Boly<sup>1,5</sup>, Bjørn E Juel<sup>1,6</sup>, Shuntaro Sasai<sup>1,7</sup>, Keiko Fujii<sup>1</sup>, Isaac David<sup>1</sup>, Jeremiah Hendren<sup>1,8</sup>, Jonathan P Lang<sup>1</sup>, Giulio Tononi<sup>1\*</sup>

**1** Department of Psychiatry, University of Wisconsin, Madison, Wisconsin, USA

**2** Fralin Biomedical Research Institute at VTC, Virginia Tech, Roanoke, Virginia, USA

**3** Neuroscience Training Program, University of Wisconsin, Madison, Wisconsin, USA

**4** Department of Mathematics and Statistics, Brock University, St. Catharines, Ontario, Canada

**5** Department of Neurology, University of Wisconsin, Madison, Wisconsin, USA

**6** Institute of Basic Medical Sciences, University of Oslo, Oslo, Norway

**7** Araya Inc., Tokyo, Japan

**8** Graduate School Language & Literature, Ludwig Maximilian University of Munich, Munich, Germany

¶These authors contributed equally to this work.

\* gtononi@wisc.edu

Let the physical substrate  $U$  be a stochastic system of interacting units  $\{U_1, U_2, \dots, U_n\}$ .

Let  $u$  be its current state, and  $u \rightarrow \bar{u}$  an update within state space  $\Omega_U = \prod_i \Omega_{U_i}$ .

Let  $\mathcal{T}_U \equiv p(\bar{u} | do(u)) = \prod_{i=1}^n p(\bar{u}_i | do(u))$  be its interventional\* transition probability function.  
 (\*): Impose all possible current states uniformly.

For each candidate system  $S \subseteq U$  in state  $s$  and background units  $W = U \setminus S$  in state  $w$ :

Compute the effect TPM  $\mathcal{T}_e \equiv p_e(\bar{s} | s) = p(\bar{s} | s, w)$

Compute the cause TPM  $\mathcal{T}_c \equiv p_c(s | \bar{s}) = \prod_{k=1}^{|\Omega_S|} \sum_{\bar{s}_k} p(s_k | \bar{s}, \bar{w}) \left( \frac{\sum_{s'} p(s' | \bar{s}, \bar{w})}{\sum_{s'} p(s' | \bar{s})} \right)$

<p>Compute the unconstrained cause probability</p> $p_c(s) =  \Omega_S ^{-1} \sum_{\bar{s} \in \Omega_{\bar{S}}} p_c(s   \bar{s})$	<p>Compute the unconstrained effect probability</p> $p_e(s) =  \Omega_S ^{-1} \sum_{\bar{s} \in \Omega_{\bar{S}}} p_e(\bar{s}   s)$
<p>Compute the probability over <math>\bar{s}</math> using Bayes' rule</p> $p_c(\bar{s}   s) = \frac{p_c(s   \bar{s}) \cdot  \Omega_{\bar{S}} ^{-1}}{p_c(s)}$	<p>Compute the probability over <math>s</math> using Bayes' rule</p> $p_e(s   \bar{s}) = \frac{p_e(\bar{s}   s) \cdot  \Omega_S ^{-1}}{p_e(\bar{s})}$
<p>For each candidate cause state <math>\bar{s}</math>:</p> <p>Compute intrinsic cause information</p> $ii_c(s, \bar{s}) = p_c^*(\bar{s}   s) \log \left( \frac{p_c(s   \bar{s})}{p_c(s)} \right)$	<p>For each candidate effect state <math>\bar{s}</math>:</p> <p>Compute intrinsic effect information</p> $ii_e(s, \bar{s}) = p_e(\bar{s}   s) \log \left( \frac{p_e(\bar{s}   s)}{p_e(\bar{s})} \right)$
<p>Find the maximal cause state</p> $s'_c(\mathcal{T}_c, s) = \operatorname{argmax}_{\bar{s} \in \Omega_{\bar{S}}} ii_c(s, \bar{s})$	<p>Find the maximal effect state</p> $s'_e(\mathcal{T}_e, s) = \operatorname{argmax}_{\bar{s} \in \Omega_{\bar{S}}} ii_e(s, \bar{s})$

For each directional system partition  $\theta$ :

Compute the partitioned transition probability functions  $\mathcal{T}_c^\theta, \mathcal{T}_e^\theta$

<p>Compute the integrated cause information</p> $\varphi_c(\mathcal{T}_c, s, \theta) = p_c(s'_c   s) \log \left( \frac{p_c(s   s'_c)}{p_c(s'_c   s)} \right)$	<p>Compute the integrated effect information</p> $\varphi_e(\mathcal{T}_e, s, \theta) = p_e(s'_e   s) \log \left( \frac{p_e(s   s'_e)}{p_e(s'_e   s)} \right)$
<p>Compute the (candidate) system integrated information <math>\varphi_s(\mathcal{T}_c, \mathcal{T}_e, s, \theta) = \min\{\varphi_c(\mathcal{T}_c, s, \theta), \varphi_e(\mathcal{T}_e, s, \theta)\}</math></p>	
<p>Find the minimum partition (MIP) <math>\theta' = \operatorname{argmin}_{\theta \in \Theta(S)} \frac{\varphi_s(\mathcal{T}_c, \mathcal{T}_e, s, \theta)}{\max_{\theta \in \Theta(S)} \varphi_s(\mathcal{T}_c, \mathcal{T}_e, s, \theta)}</math></p>	
<p>Identify system integrated information <math>\varphi_s(\mathcal{T}_c, \mathcal{T}_e, s) := \varphi_s(\mathcal{T}_c, \mathcal{T}_e, s, \theta')</math></p>	
<p>Find the first complex <math>S^* = \operatorname{argmax}_{S \subseteq U} \varphi_s(\mathcal{T}_c, \mathcal{T}_e, s)</math>. This is the PSC*</p>	

(\*) In principle, not only sets of units, but also the grain of units, updates, and states should be considered

Unfold the cause-effect structure of the complex:

For each candidate mechanism  $M \subseteq S^*$  in state  $m$ :

<p>For each candidate purview <math>Z \subseteq S^*</math>:</p> <p>Compute the probability over <math>M</math>, marginalizing out external influences <math>Y = S^* \setminus Z</math></p> $\pi_c(m   z) = \prod_{i=1}^{ \Omega_Z }  \Omega_i ^{-1} \sum_{y \in \Omega_Y} p(m_i   z, y)$	<p>For each candidate purview <math>Z \subseteq S^*</math>:</p> <p>Compute the probability over <math>Z</math>, marginalizing out external influences <math>X = S^* \setminus M</math></p> $\pi_e(z   m) = \prod_{i=1}^{ \Omega_Z }  \Omega_i ^{-1} \sum_{x \in \Omega_X} p(z_i   m, x)$
<p>Compute the unconstrained cause probability</p> $\pi_c(m; Z) =  \Omega_Z ^{-1} \sum_{z \in \Omega_Z} \pi_c(m   z)$	<p>Compute the unconstrained effect probability</p> $\pi_e(z; M) =  \Omega_M ^{-1} \sum_{m \in \Omega_M} \pi_e(z   m)$
<p>Compute the probability over <math>Z</math> using Bayes' rule</p> $\pi_c^*(z   m) = \frac{\pi_c(m   z) \cdot  \Omega_Z ^{-1}}{\pi_c(m; Z)}$	<p>Compute the probability over <math>M</math> using Bayes' rule</p> $\pi_e^*(m   z) = \frac{\pi_e(z   m) \cdot  \Omega_M ^{-1}}{\pi_e(z; M)}$
<p>For each candidate cause purview state <math>z</math>:</p> <p>Compute the intrinsic cause information</p> $ii_c(m, z) = \pi_c^*(z   m) \log \left( \frac{\pi_c(m   z)}{\pi_c(m; Z)} \right)$	<p>For each candidate effect purview state <math>m</math>:</p> <p>Compute the intrinsic effect information</p> $ii_e(z, m) = \pi_e^*(m   z) \log \left( \frac{\pi_e(z   m)}{\pi_e(z; M)} \right)$
<p>Find the maximal cause state</p> $z'_c(m, Z) = \operatorname{argmax}_{z \in \Omega_Z} ii_c(m, z)$	<p>Find the maximal effect state</p> $m'_e(z, M) = \operatorname{argmax}_{m \in \Omega_M} ii_e(z, m)$
<p>For each disintegrating mechanism partition <math>\theta</math>:</p> <p>Compute the partitioned probability</p> $\pi_c^*(m   z) = \prod_{i=1}^{ \Omega_Z } \pi_c^*(m_i   z_i)$	<p>For each disintegrating mechanism partition <math>\theta</math>:</p> <p>Compute the partitioned probability</p> $\pi_e^*(z   m) = \prod_{i=1}^{ \Omega_M } \pi_e^*(z_i   m_i)$
<p>Compute the integrated cause information</p> $\varphi_c(m, Z, \theta) = \pi_c^*(z'_c   m) \log \left( \frac{\pi_c(m   z'_c)}{\pi_c(z'_c   m)} \right)$	<p>Compute the integrated effect information</p> $\varphi_e(m, Z, \theta) = \pi_e^*(m'_e   z) \log \left( \frac{\pi_e(z   m'_e)}{\pi_e(m'_e   z)} \right)$
<p>Find the minimum partition (MIP)</p> $\theta' = \operatorname{argmin}_{\theta \in \Theta(M, Z)} \frac{\varphi_c(m, Z, \theta)}{\max_{\theta \in \Theta(M, Z)} \varphi_c(m, Z, \theta)}$	<p>Find the minimum partition (MIP)</p> $\theta' = \operatorname{argmin}_{\theta \in \Theta(M, Z)} \frac{\varphi_e(m, Z, \theta)}{\max_{\theta \in \Theta(M, Z)} \varphi_e(m, Z, \theta)}$
<p>Identify integrated information</p> $\varphi_c(m, Z) := \varphi_c(m, Z, \theta')$	<p>Identify integrated information</p> $\varphi_e(m, Z) := \varphi_e(m, Z, \theta')$
<p>Find the maximally irreducible cause</p> $z^*(m) = \operatorname{argmax}_{z \in \Omega_Z} \varphi_c(m, z^*(m, Z))$	<p>Find the maximally irreducible effect</p> $m^*(z) = \operatorname{argmax}_{m \in \Omega_M} \varphi_e(m, z^*(m, Z))$
<p>Identify the integrated cause information</p> $\varphi_c(m) := \max_{z \in \Omega_Z} \varphi_c(m, z^*(m, Z))$	<p>Identify the integrated effect information</p> $\varphi_e(m) := \max_{z \in \Omega_Z} \varphi_e(m, z^*(m, Z))$
<p>Compute the candidate distinction's integrated information <math>\varphi_d(m) = \min\{\varphi_c(m), \varphi_e(m)\}</math></p> <p>The candidate distinction is <math>d(m) = (m, z^* = \{z^*_c, z^*_e\}, \varphi_d)</math></p>	

Compute the set of congruent distinctions  $D(\mathcal{T}_c, \mathcal{T}_e, s^*) = \{d : \varphi_d > 0, z^*_c \subseteq s^*, z^*_e \subseteq s^*\}$

For each candidate set of distinctions  $d \in D(\mathcal{T}_c, \mathcal{T}_e, s^*)$ :

For each set of causes and/or effects  $z$  such that:

$$z : z \cap \{z^*_c(d), z^*_e(d)\} \neq \emptyset \forall d \in d, \prod_{z^* \in z} z^* \neq \emptyset, |z| > 1$$

Compute the maximal overlap  $\alpha^*(z) = \prod_{z^* \in z} \alpha^*(z^*)$

The relation face is  $f(z) = (z, \alpha^*(z))$

The set of relation faces is  $f(d) = \{f(z)\}_d$

Compute the integrated information of the candidate relation

$$\varphi_s(d) = \min_{f \in f(d)} \bigcup_{\alpha^*} \left| \frac{\varphi_d}{z^*_c(d) \cup z^*_e(d)} \right|$$

The candidate relation is  $r(d) = (d, f(d), \varphi_s)$

Compute the set of relations  $R(D) = \{r(d) : \varphi_s(d) > 0\}$

The cause-effect structure of the complex  $S^*$  (its  $\Phi$ -structure) is

$$C(\mathcal{T}_c, \mathcal{T}_e, s^*) = D \cup R(D) = \{D(\mathcal{T}_c, \mathcal{T}_e, s^*) \cup R(D(\mathcal{T}_c, \mathcal{T}_e, s^*))\}$$

Compute  $\Phi(\mathcal{T}_c, \mathcal{T}_e, s^*) = \sum_{C(\mathcal{T}_c, \mathcal{T}_e, s^*)} \varphi_s$

