



**A quantitative account of agents
and their actions**
using the causal framework of
Integrated Information Theory (IIT)

Larissa Albantakis

What is an agent?

Agents are open systems that dynamically and informationally interact with the environment.

1. How do we distinguish an agent from its environment?



What is an agent?

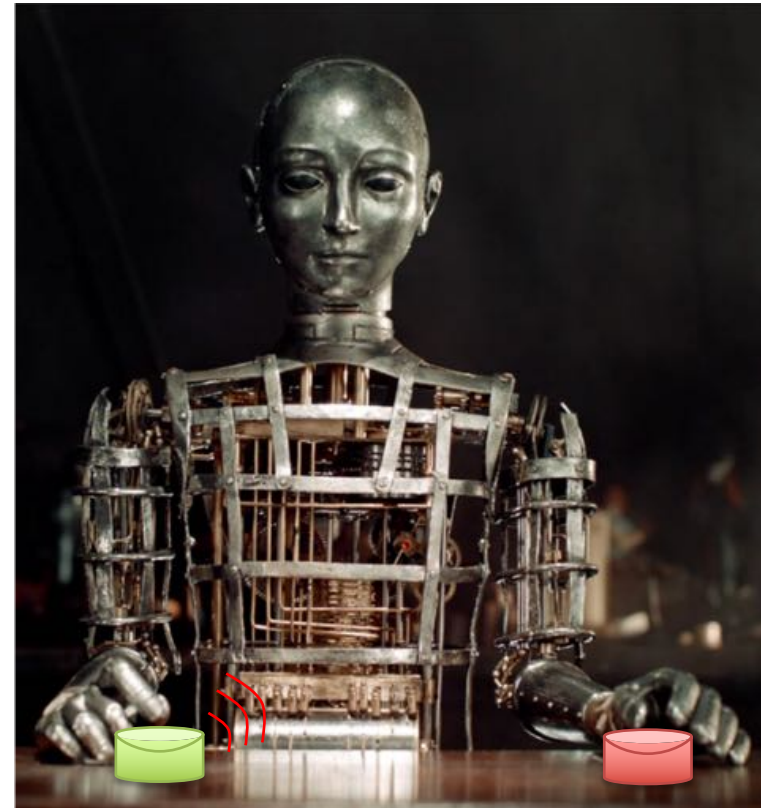
Agents are open systems that dynamically and informationally interact with the environment.

2. How do we distinguish autonomous actions from mere reflexes?

(A) autonomous



(B) automatic



From movie "Hugo" 2011.

What is an agent?

Agents are open systems that dynamically and informationally interact with the environment.

2. How do we distinguish autonomous actions from mere reflexes?

(A) autonomous



(B) automatic



What is an agent?

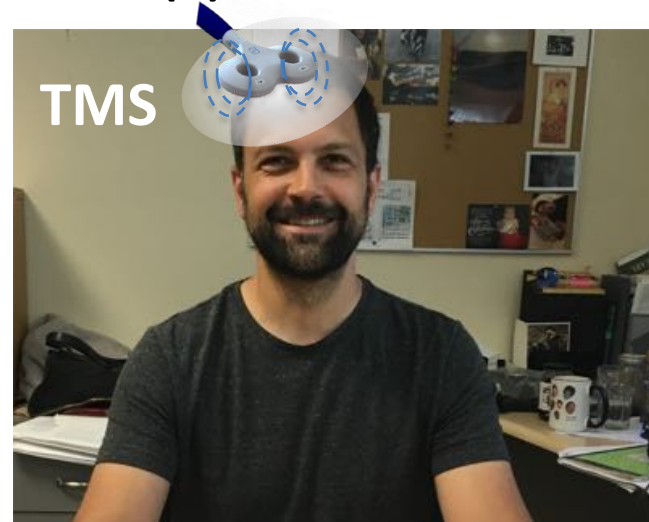
Agents are open systems that dynamically and informationally interact with the environment.

2. How do we distinguish autonomous actions from mere reflexes?

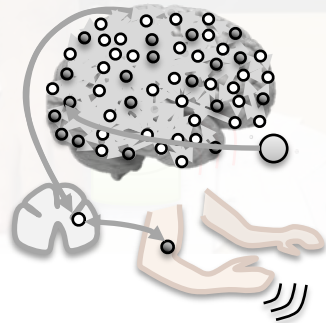
(A) autonomous



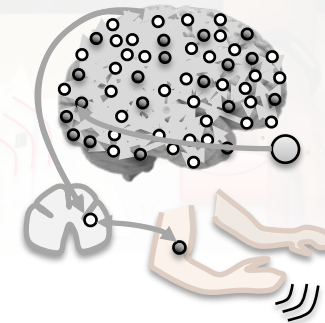
(B) automatic



What was the cause for clicking right?



“Neurons were activated and interacted with other neurons which then triggered a motor response.”



What is an agent?

A minimal working definition:

An autonomous agent is

- (1) an open system with stable, self-defined and self-maintained causal borders,**
- (2) with the capacity to perform actions that are (partially) caused from within**

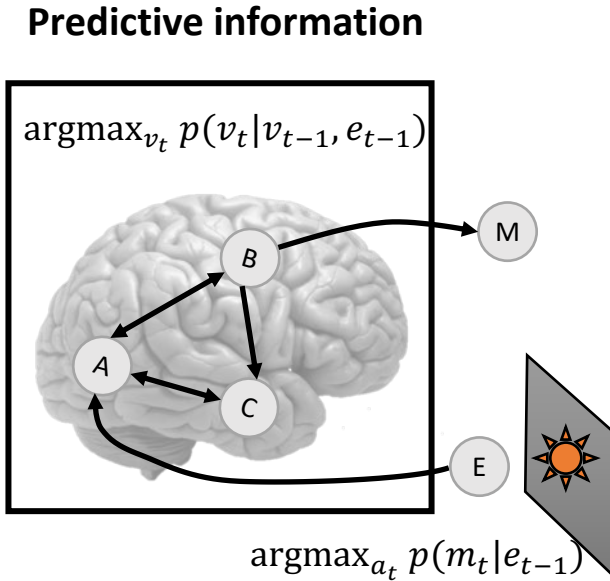
Can we quantify this?

Recent related work:

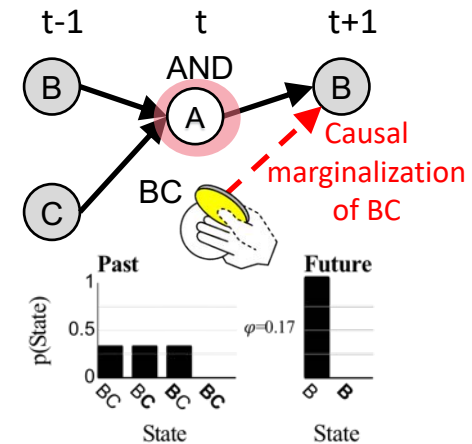
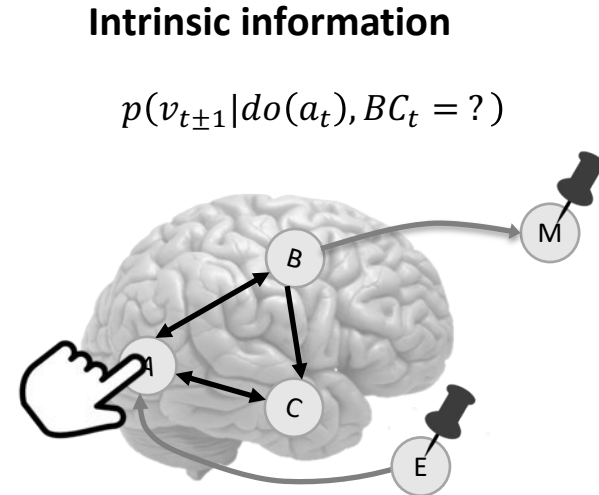
- Bertschinger N, Olbrich E, Ay N, Jost J (2008) Autonomy: An information theoretic perspective. *Biosystems* 91:331–345.
- Krakauer D, Bertschinger N, Olbrich E, Ay N, Flack JC (2014) The Information Theory of Individuality. *Arxiv* 1412.2447.
- Biehl M, Ikegami T, Polani D (2016) Towards information based spatiotemporal patterns as a foundation for agent representation in dynamical systems. In: *Artificial Life Conference 2016*, pp 722.
- Kolchinsky A, Wolpert DH (2018) Semantic information, autonomous agency and non-equilibrium statistical physics. *Interface Focus* 8:20180041.
- Kirchhoff M, Parr T, Palacios E, Friston K, Kiverstein J (2018) The Markov blankets of life: autonomy, active inference and the free energy principle. *J R Soc Interface* 15:20170792.
- Walker SI, Davies PCW (2013) The algorithmic origins of life. *J R Soc Interface* 10:20120869.

Extrinsic vs. intrinsic information

extrinsic and correlational



intrinsic and causal



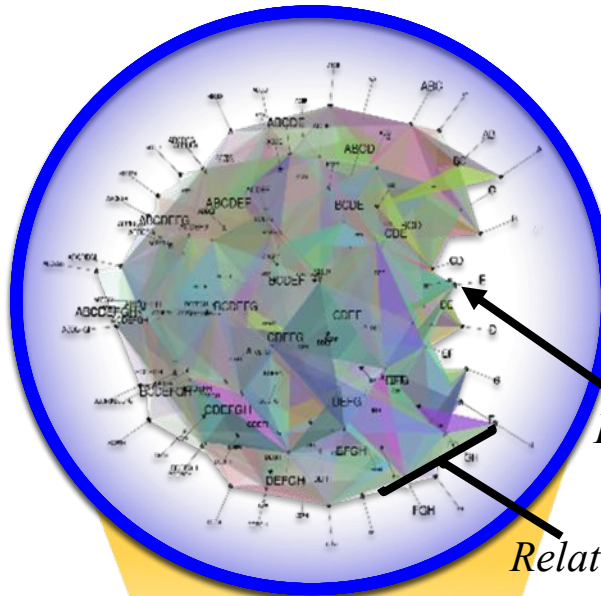
Integrated Information Theory (IIT) of Consciousness

An experience is a maximum (Φ) of cause-effect power (a cause-effect structure)



Experience

=



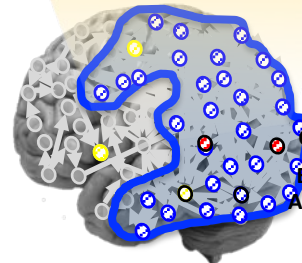
Cause-effect structure

Quantity:
irreducibility Φ
of the cause-effect structure

Quality:
'form'
of the cause-effect structure

Distinctions

Relations

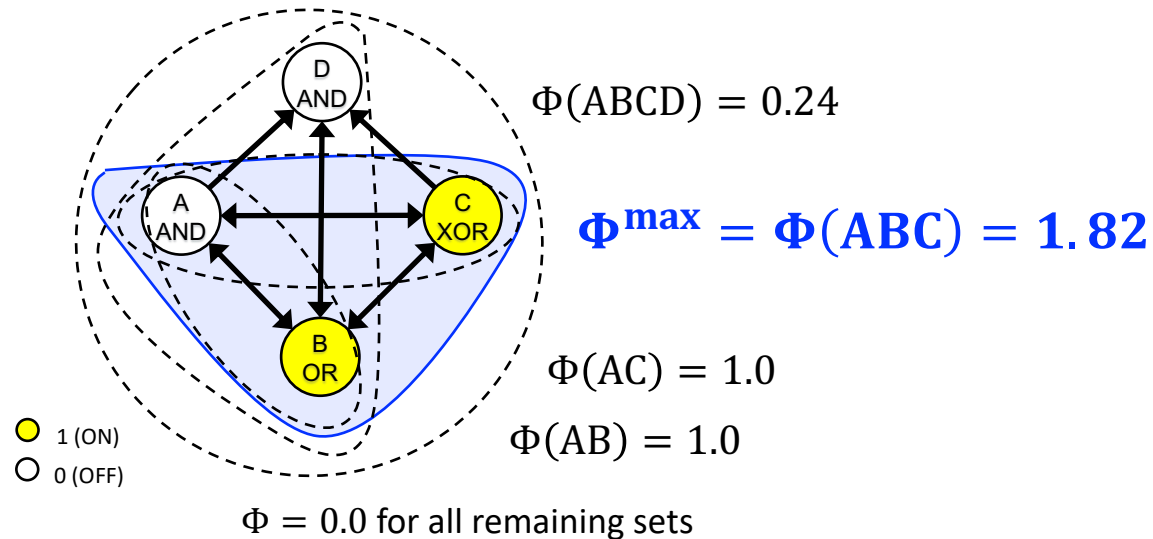


Φ^{\max}

Physical substrate

Causal framework of IIT

1. Identifying self-defined causal borders – maxima of Φ



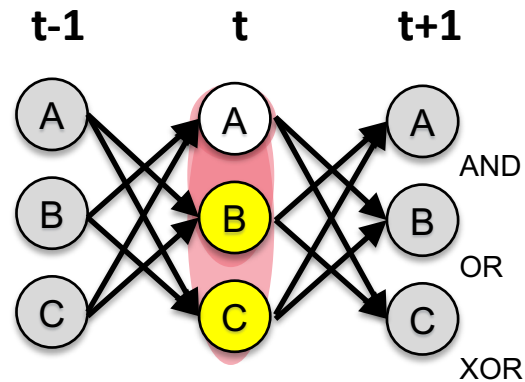
An autonomous agent is:

1. an open system with stable, self-defined and self-maintained causal borders,
2. with the capacity to perform actions that are (partially) caused from within

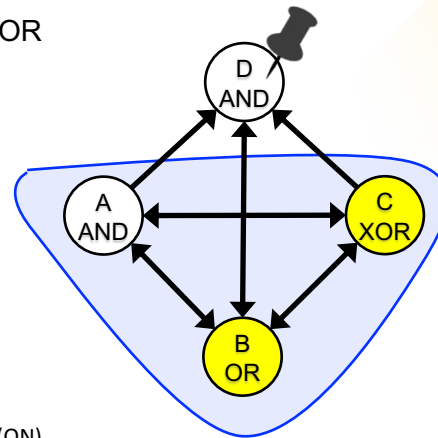
Causal framework of IIT

1. Identifying self-defined causal borders – maxima of Φ

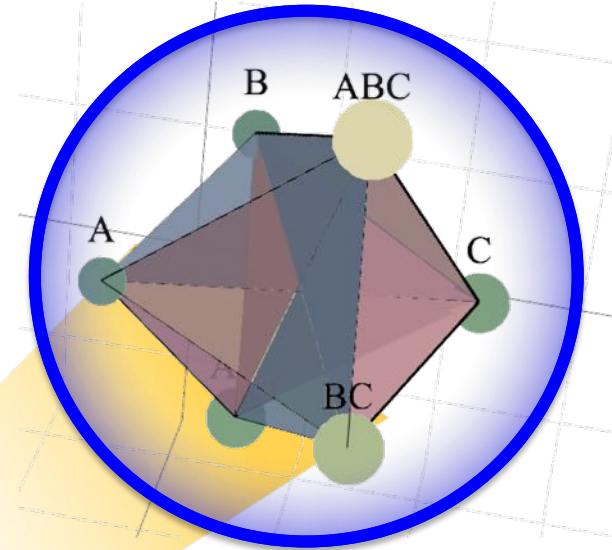
Unfolding the causes and effects of a system in a state



$$\Phi(ABC) = 1.82$$

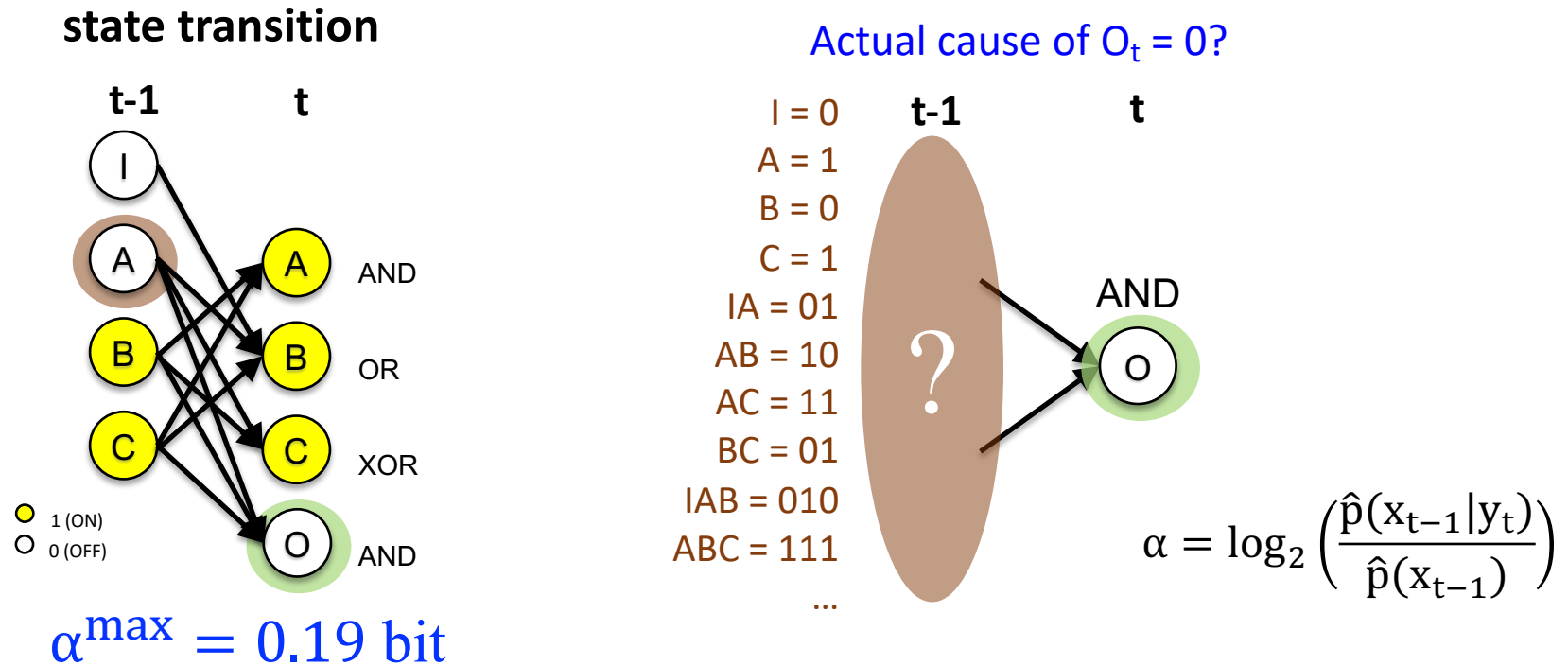


● 1 (ON)
○ 0 (OFF)



Causal framework of IIT

2. Identifying actual causes (of actions) – maxima of α



An autonomous agent is:

1. an open system with stable, self-defined and self-maintained causal borders,
2. **with the capacity to perform actions that are (partially) caused from within**

Albantakis L, Marshall W, Hoel E, Tononi G (2019) What caused what? A quantitative account of actual causation using dynamical causal networks. Entropy 21:459.

An autonomous agent is:

1. an open system with stable, self-defined and self-maintained causal borders,
2. with the capacity to perform actions that are (partially) caused from within

1. Identifying self-defined causal borders

Example 1: Fission yeast cell cycle model

Marshall W, Kim H, Walker SI, Tononi G, Albantakis L (2017)

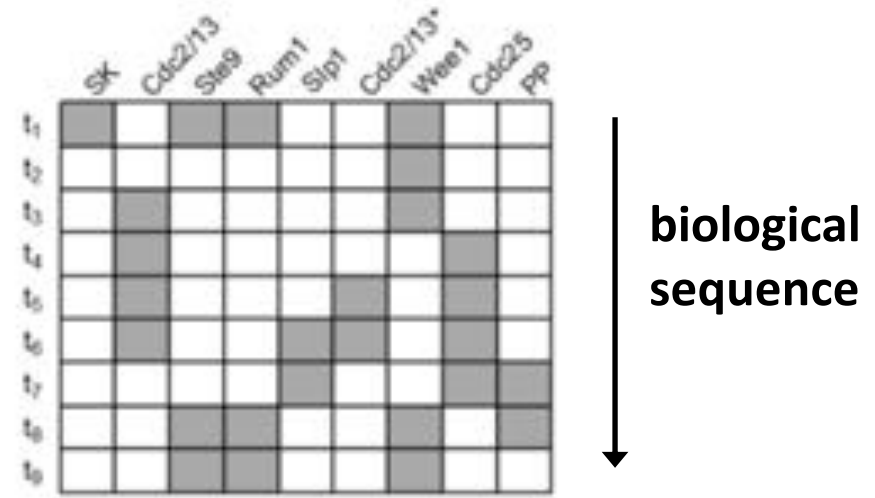
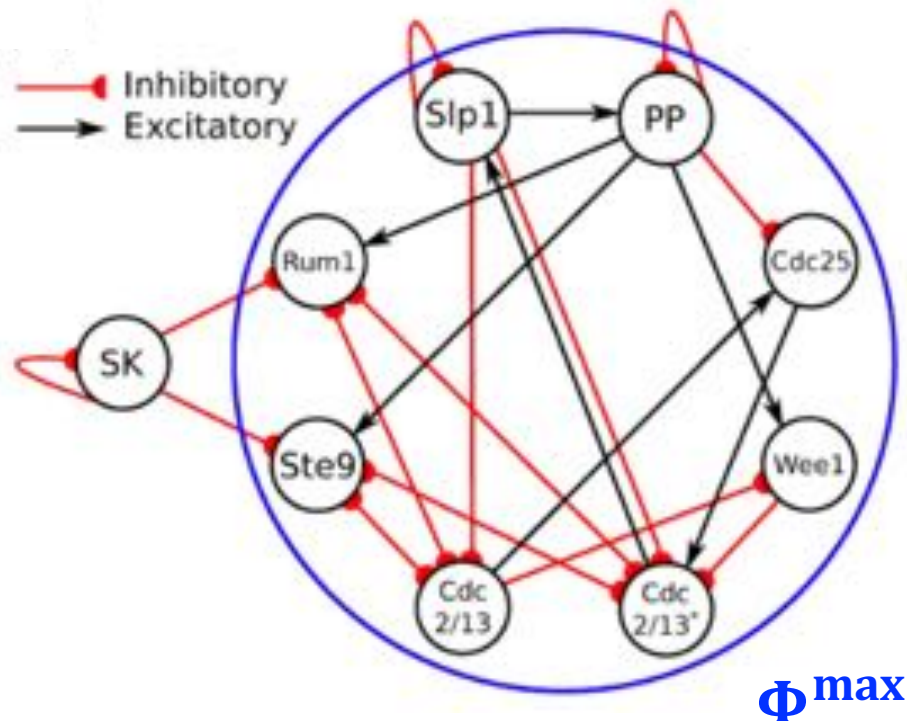
Example 2: Evolved artificial agents (animats)

Albantakis L, Hintze A, Koch C, Adami C, Tononi G (2014)



The yeast cell-cycle model has stable causal borders throughout its biological sequence

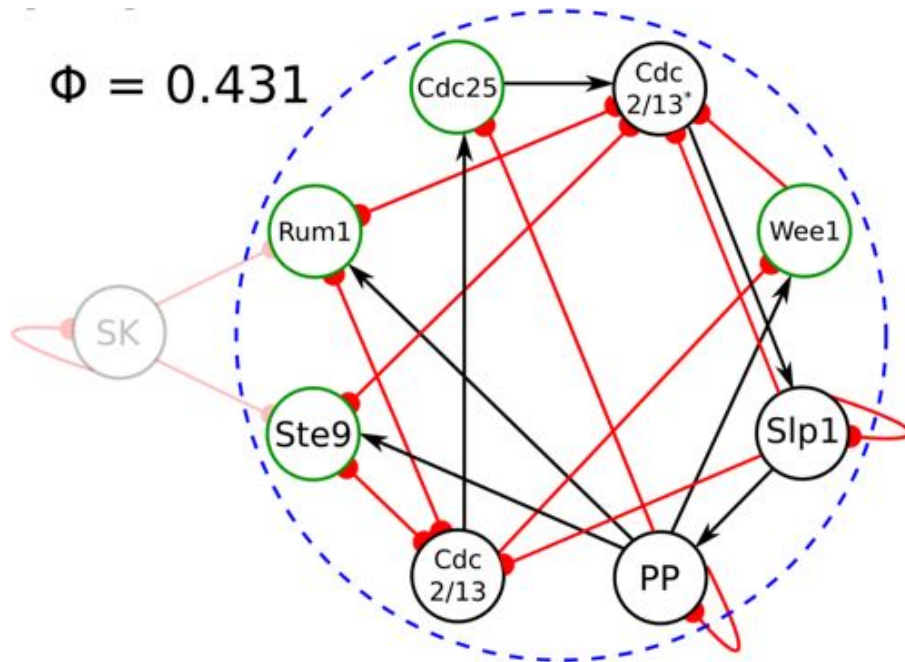
Boolean network model of the fission-yeast cell-cycle



Φ^{\max}

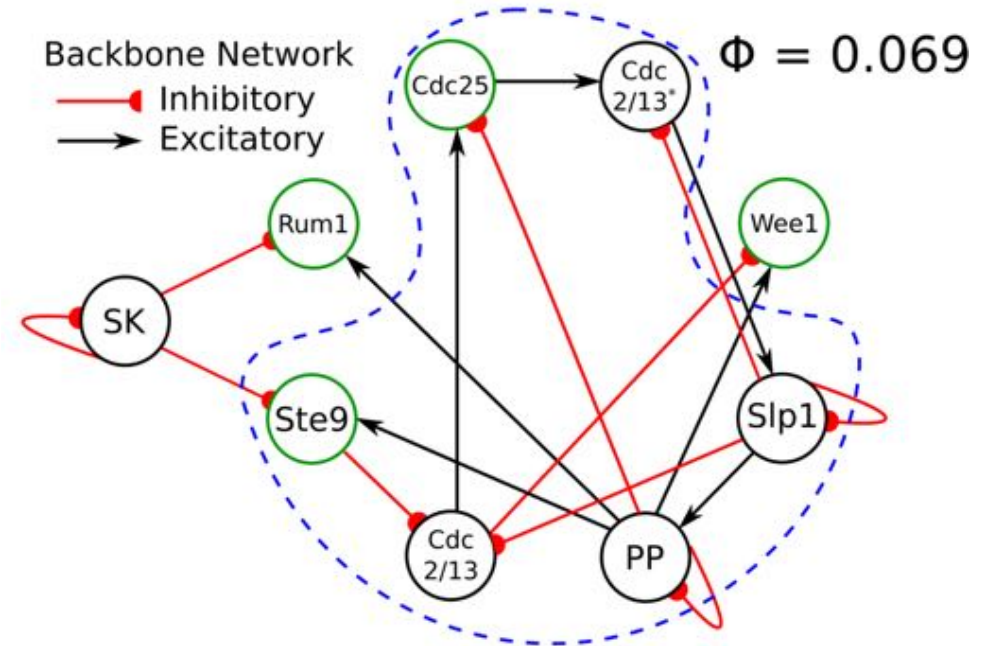
A functionally equivalent reduced version of the network does not form a stable entity

fission yeast cell cycle model

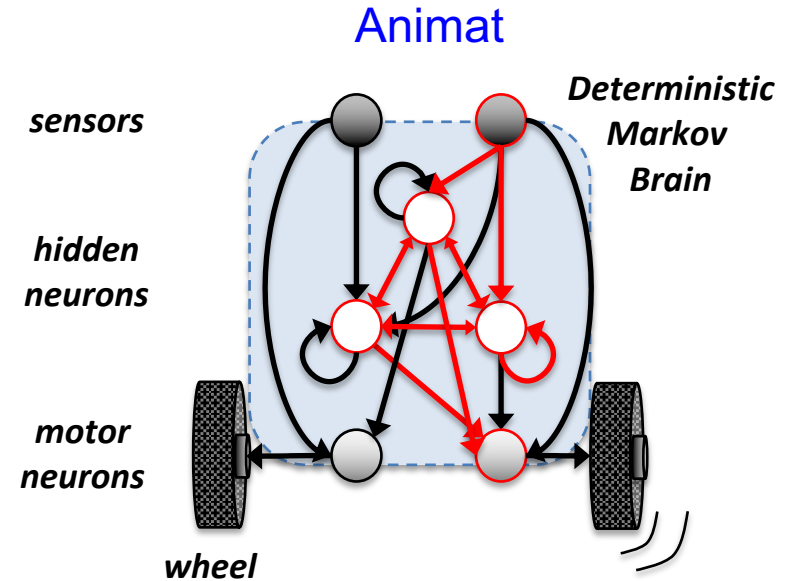
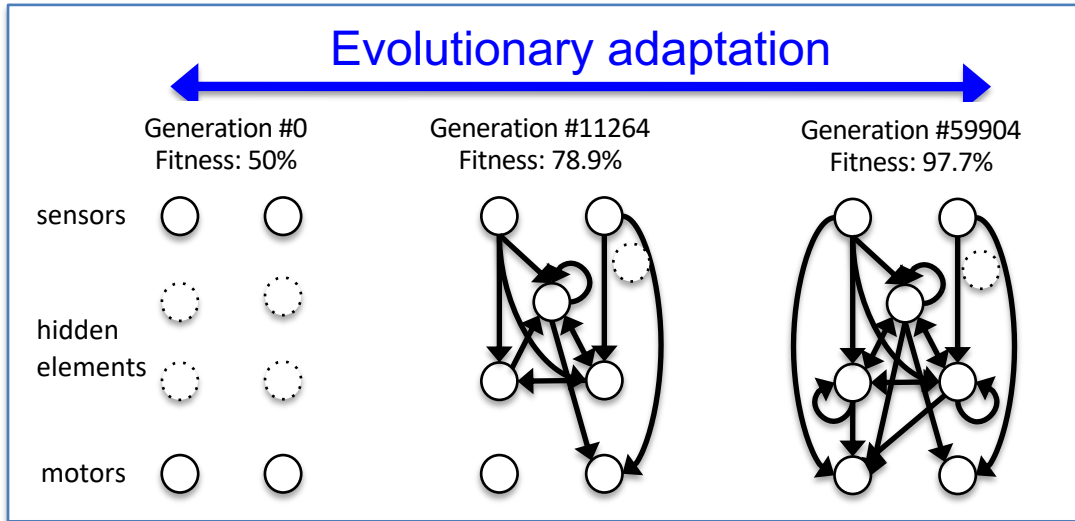


maximum of Φ in every state in the cycle

functionally identical
'backbone' network



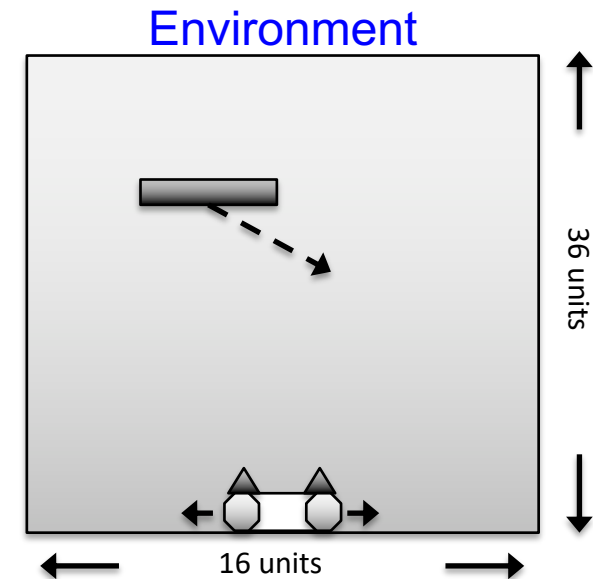
Evolved artificial agents



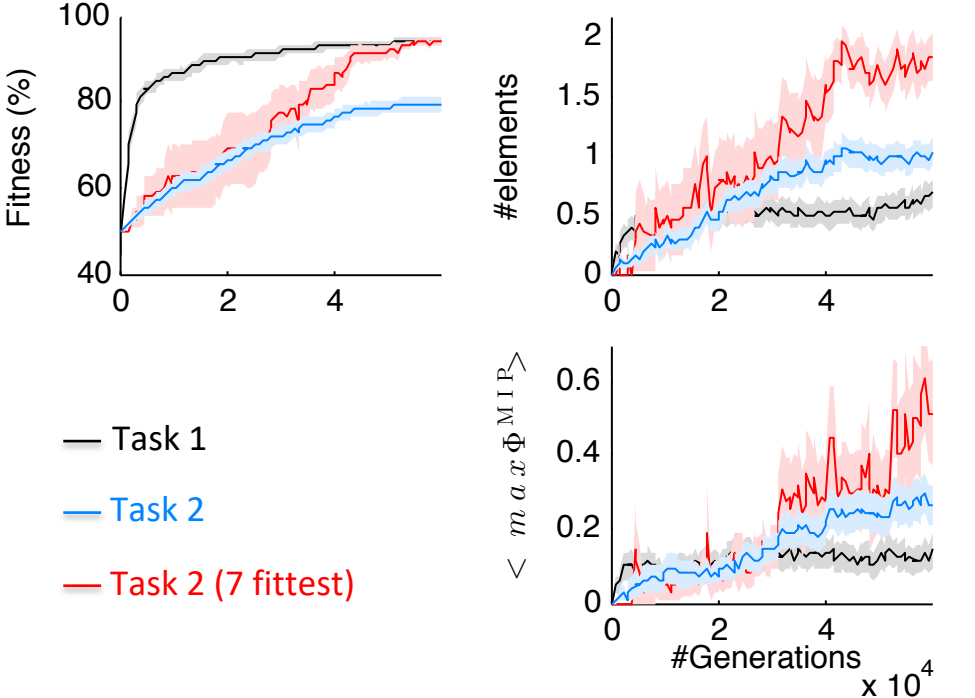
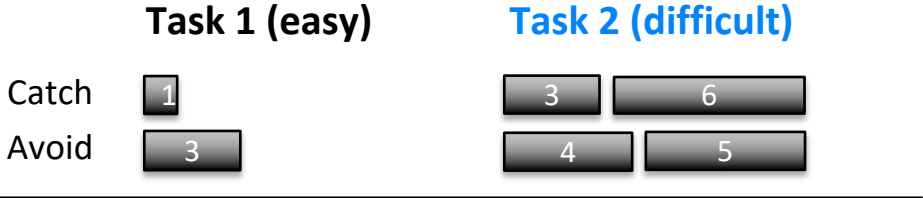
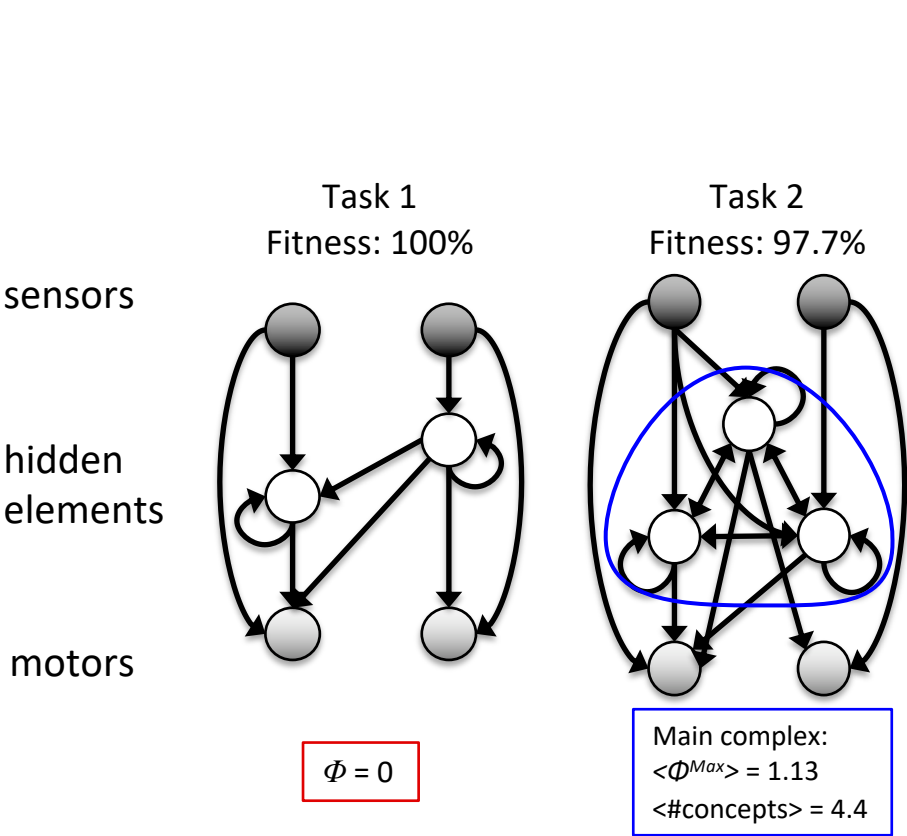
Selection & Mutation

Fitness: % of successfully caught and avoided blocks (out of 128 trials)

Point mutations, deletions, and duplications in the genome after each selection

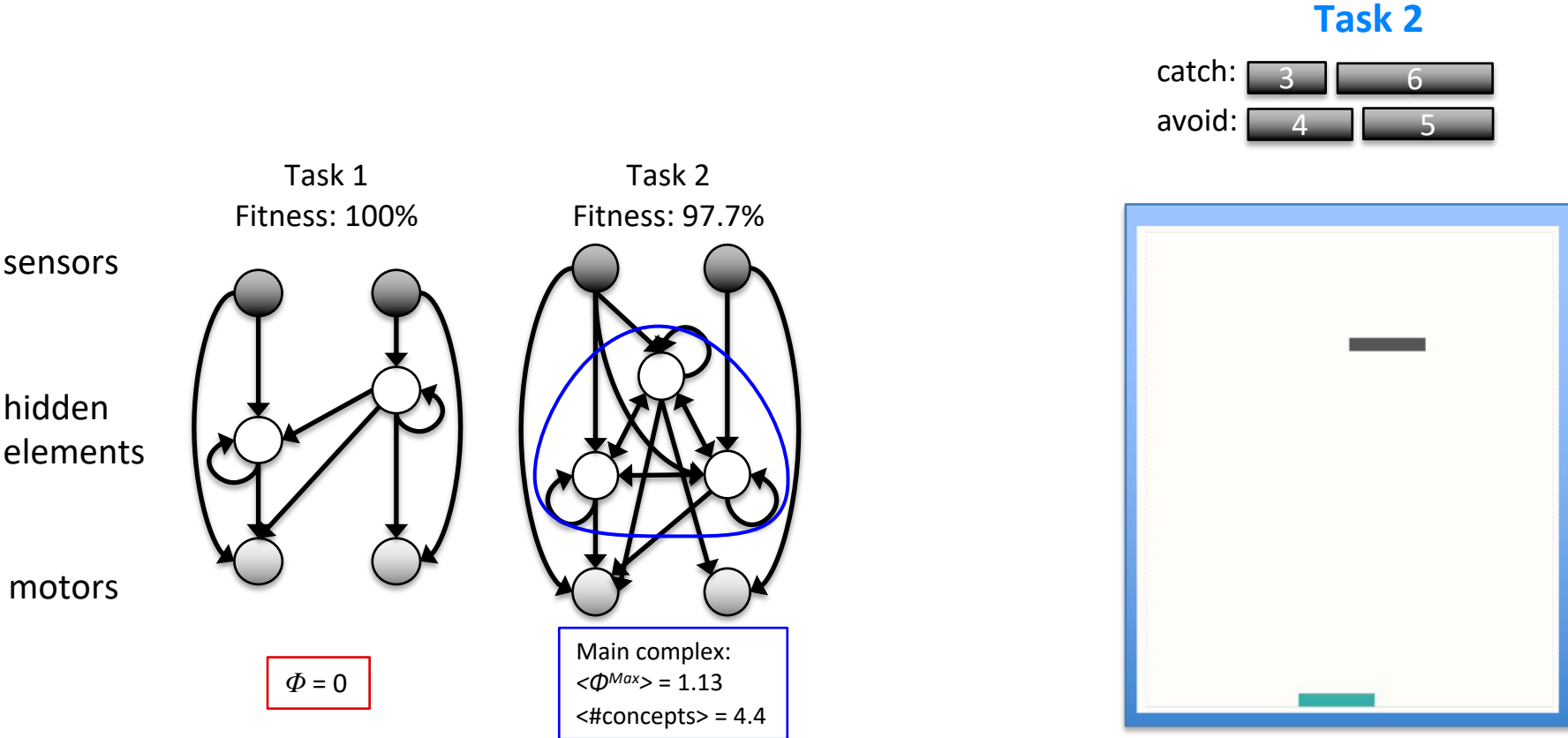


Larger and more integrated complexes evolved in tasks that require more memory and context-dependent behavior



Albantakis L, Hintze A, Koch C, Adami C, Tononi G (2014) Evolution of Integrated Causal Structures in Animats Exposed to Environments of Increasing Complexity. PLoS Comput Biol 10:e1003966.

Larger and more integrated complexes evolved in tasks that require more memory and context-dependent behavior



Animat is partially autonomous due to memory.

An autonomous agent is:

1. an open system with stable, self-defined and self-maintained causal borders,
2. **with the capacity to perform actions that are (partially) caused from within**

2. Identifying the actual causes of actions

Example 1: Tracing back the causal chain

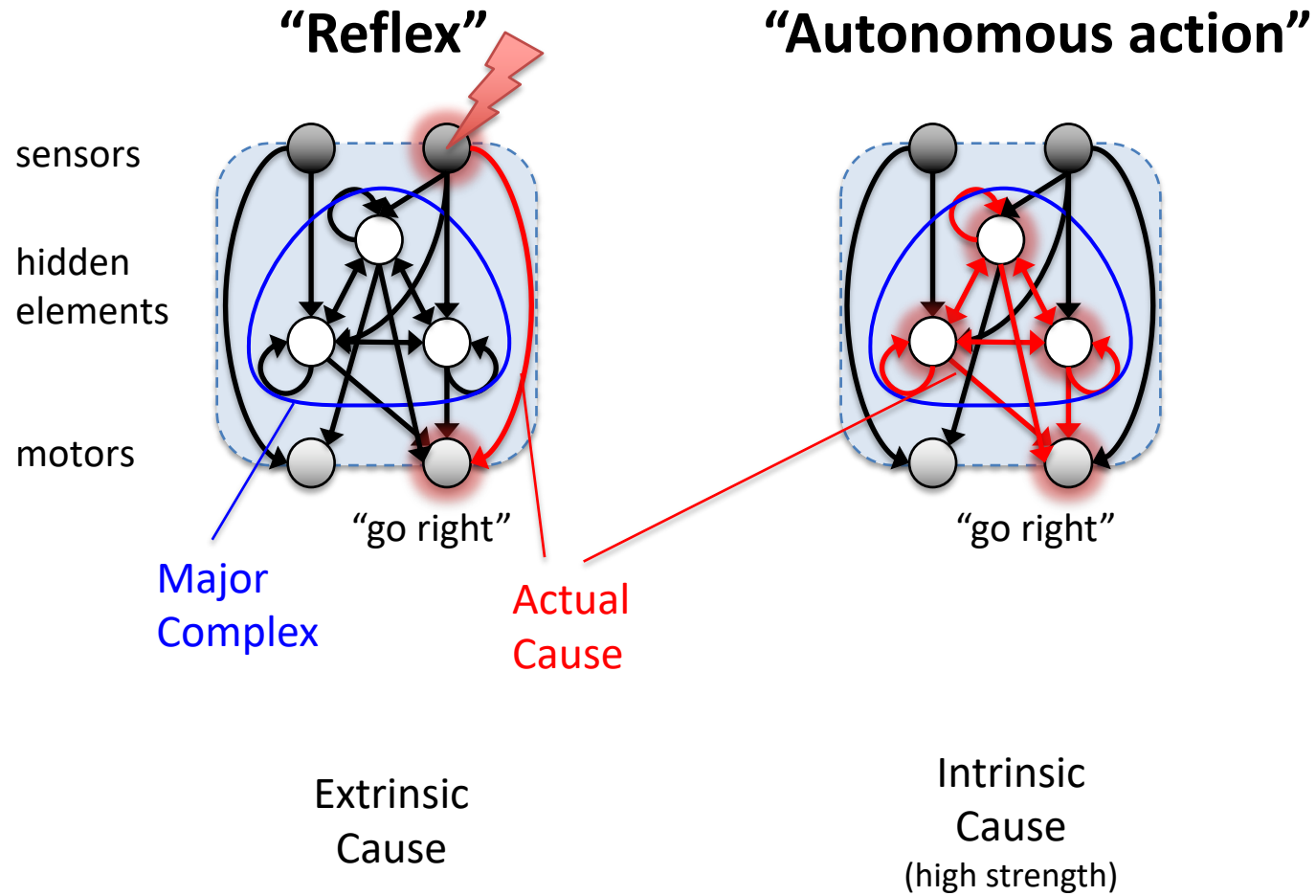
Albantakis L, Massari F, Tononi G (in prep)

Example 2: Causes of actions across various evolutionary environments

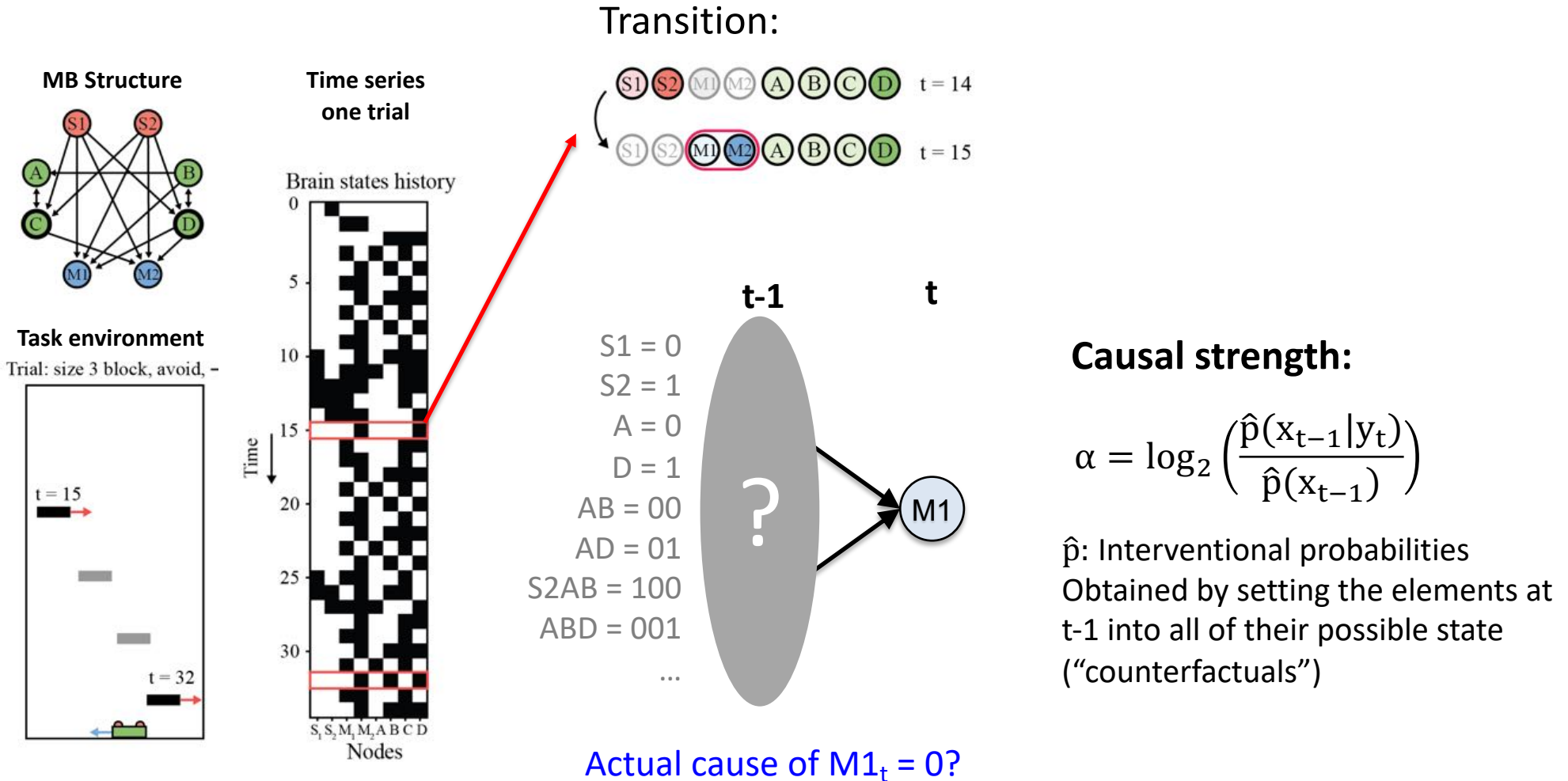
Juel BE, Comolatti R, Tononi G, Albantakis L (2019)



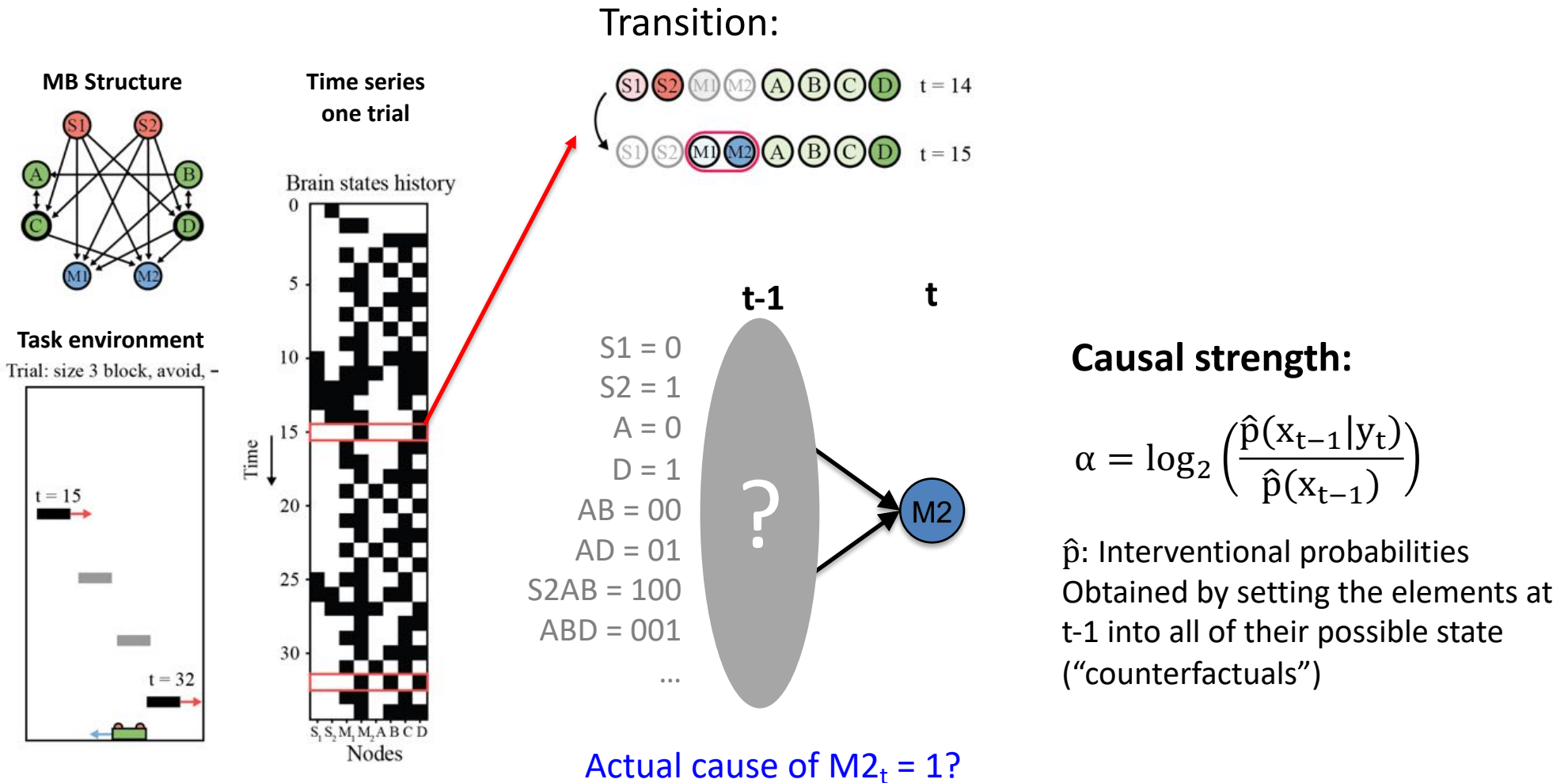
Extrinsic vs. intrinsic actual causes



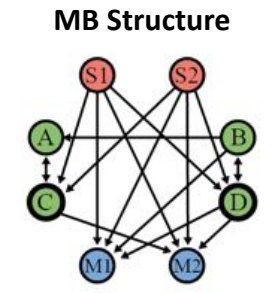
Causal analysis: direct actual causes



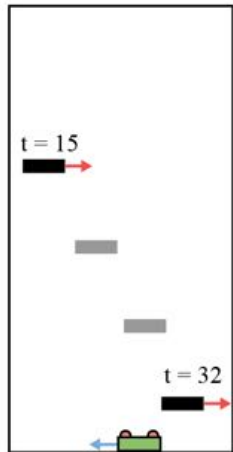
Causal analysis: direct actual causes



Causal analysis: direct actual causes

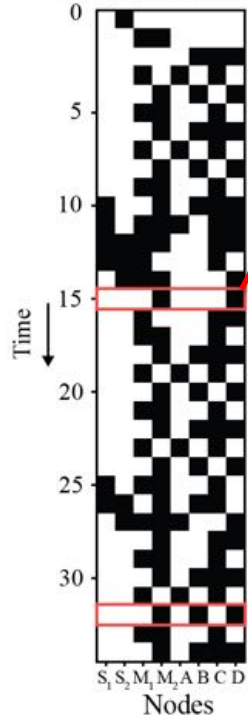


Task environment
Trial: size 3 block, avoid, -

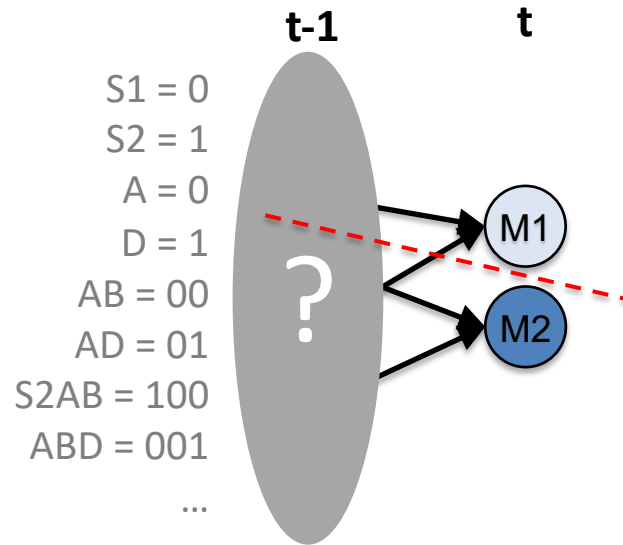
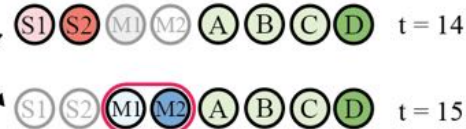


Time series one trial

Brain states history



Transition:



Actual cause of $M1M2_t = 01$?

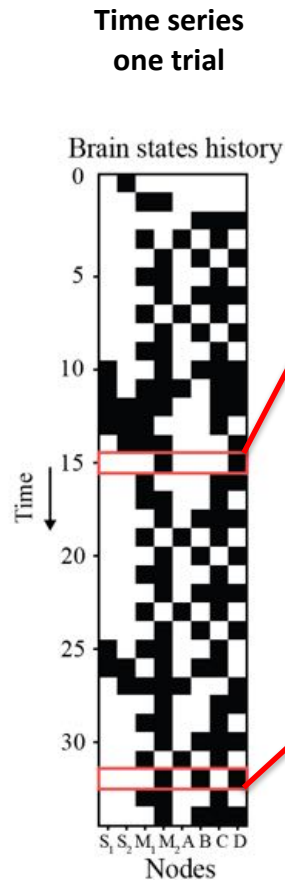
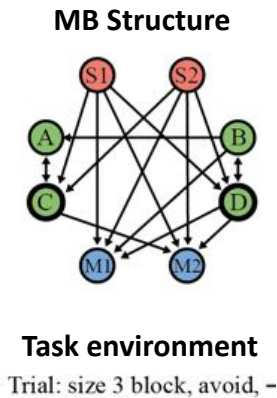
Causal strength:

$$\alpha = \min_{\Psi} \left(\log_2 \left(\frac{\hat{p}(x_{t-1}|y_t)}{\Psi(\hat{p}(x_{t-1}|y_t))} \right) \right)$$

Ψ : partition

\hat{p} : Interventional probabilities
Obtained by setting the elements at t-1 into all of their possible state ("counterfactuals")

Same action, same agent, different causes



Transition:



$\alpha = 1.415$ [S2, D] ← [M1]
 $\alpha = 0.415$ [S2, B, C] ← [M2]
 $\alpha = 0.415$ [S2, B, C] ← [M1, M2]

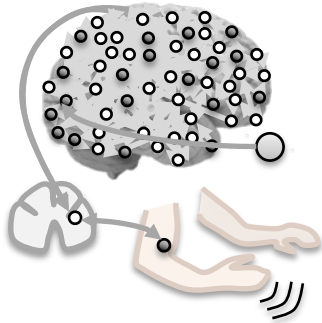
$\alpha_{TOTAL} = 2.24$ <purview length> = 2.67
 $\alpha_{SENSOR} = 0.98$ $\alpha_{HIDDEN} = 1.26$



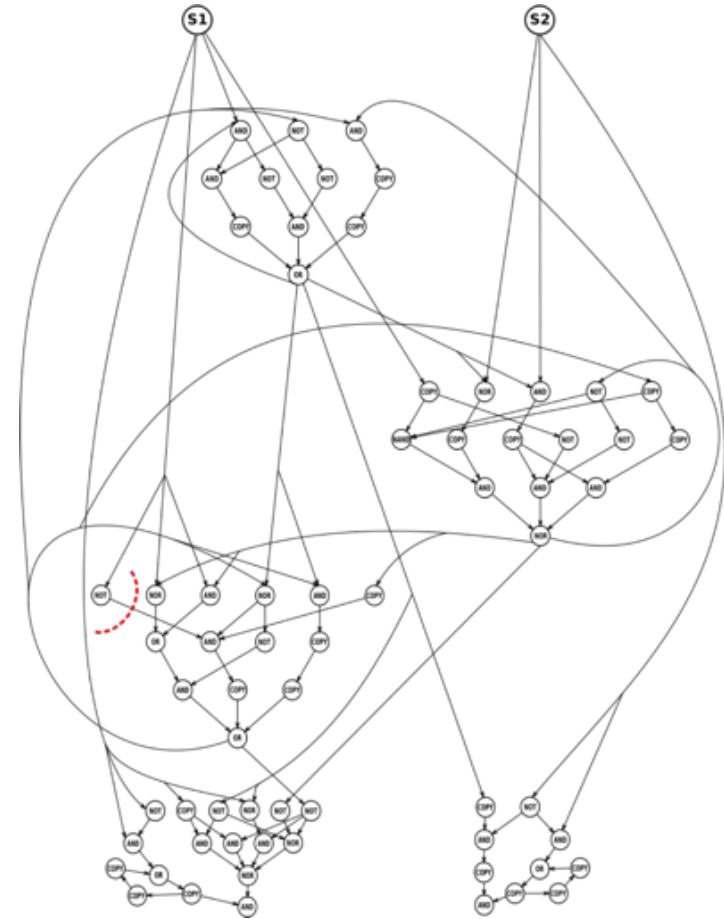
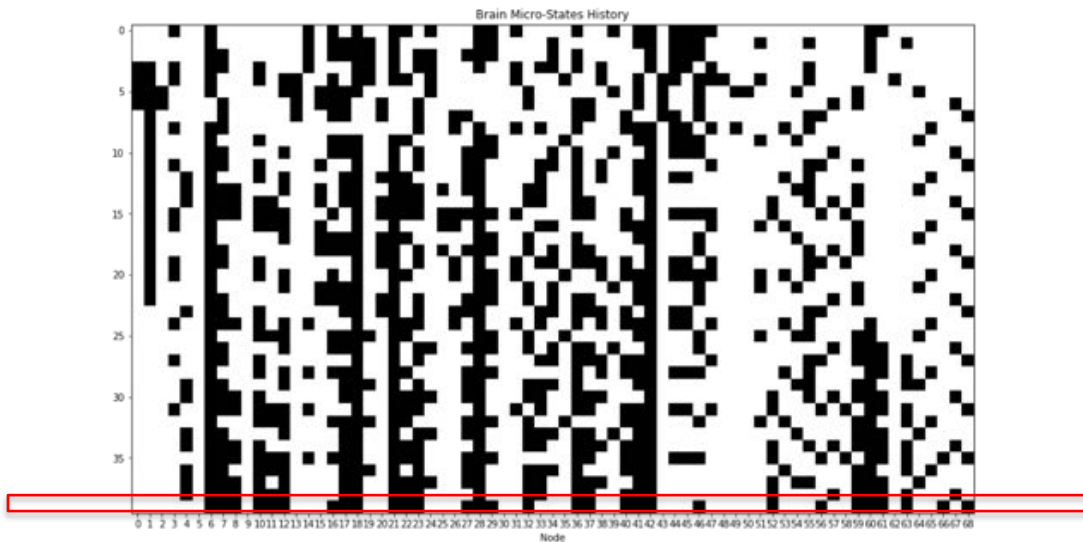
$\alpha = 1.415$ [S1, S2, D] ← [M1]
 $\alpha = 0.415$ [[S2, C], [S1, S2, D]] ← [M2]
 $\alpha = 0.263$ [S1, S2, D] ← [M1, M2]

$\alpha_{TOTAL} = 2.09$ <purview length> = 2.83
 $\alpha_{SENSOR} = 1.36$ $\alpha_{HIDDEN} = 0.73$

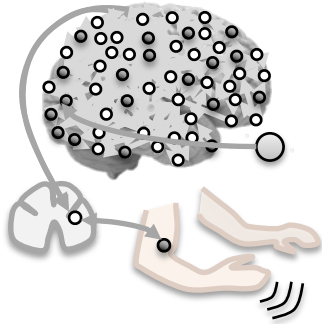
Tracing back the causal chain



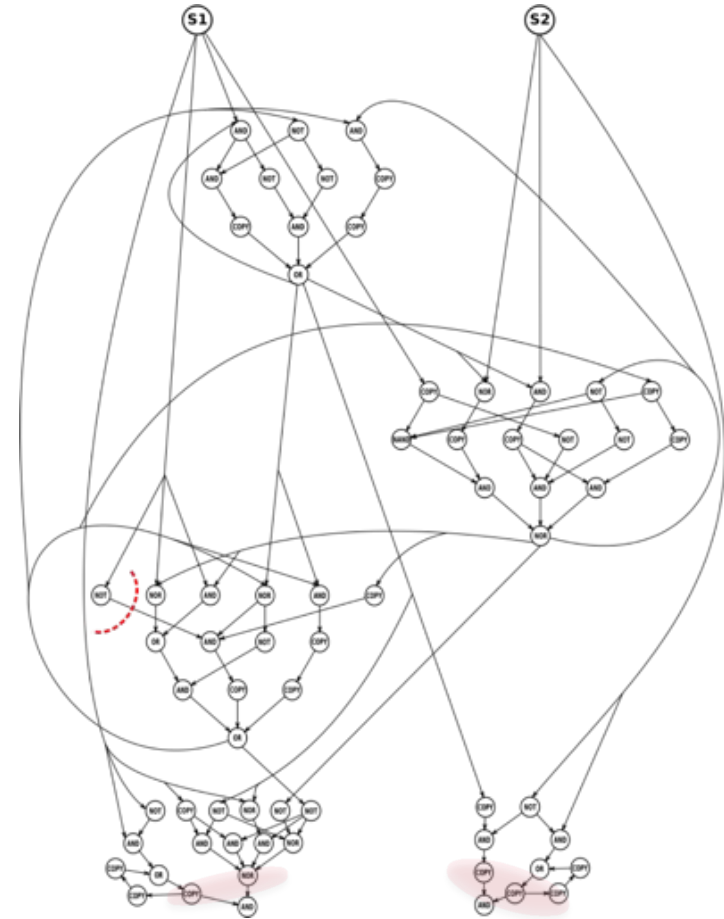
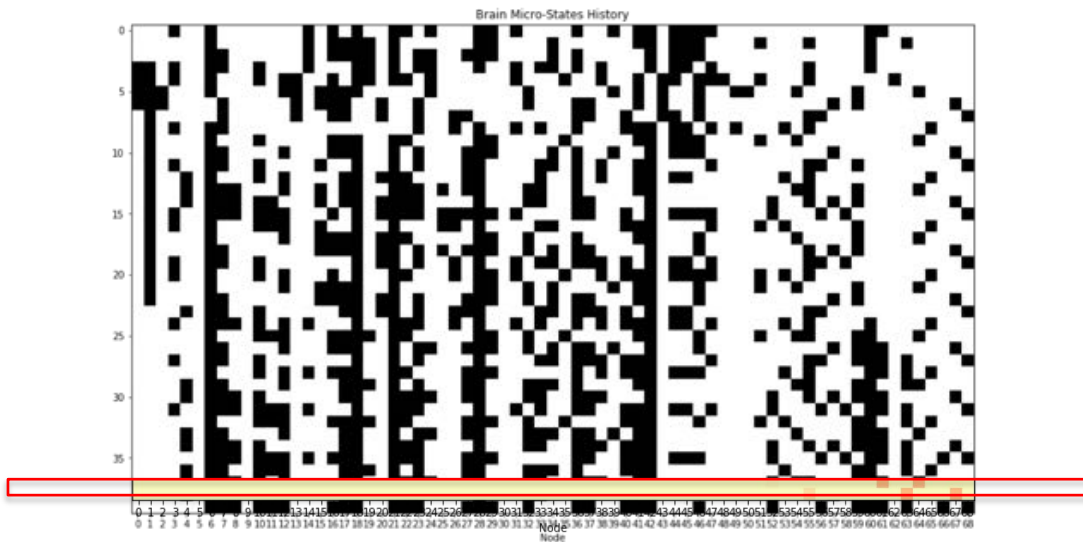
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



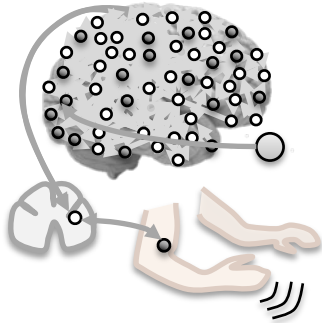
Tracing back the causal chain



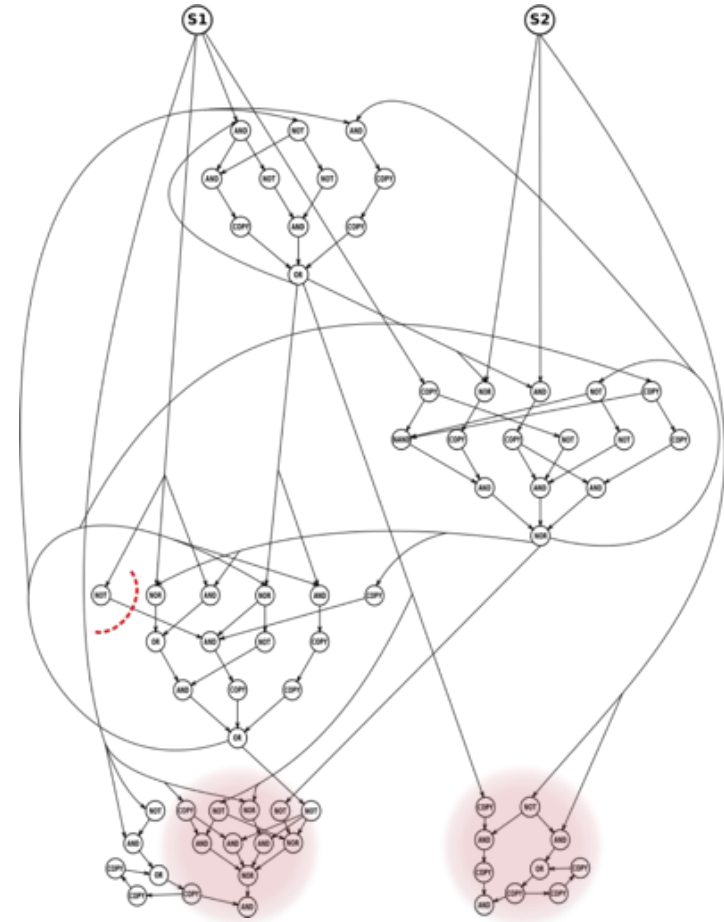
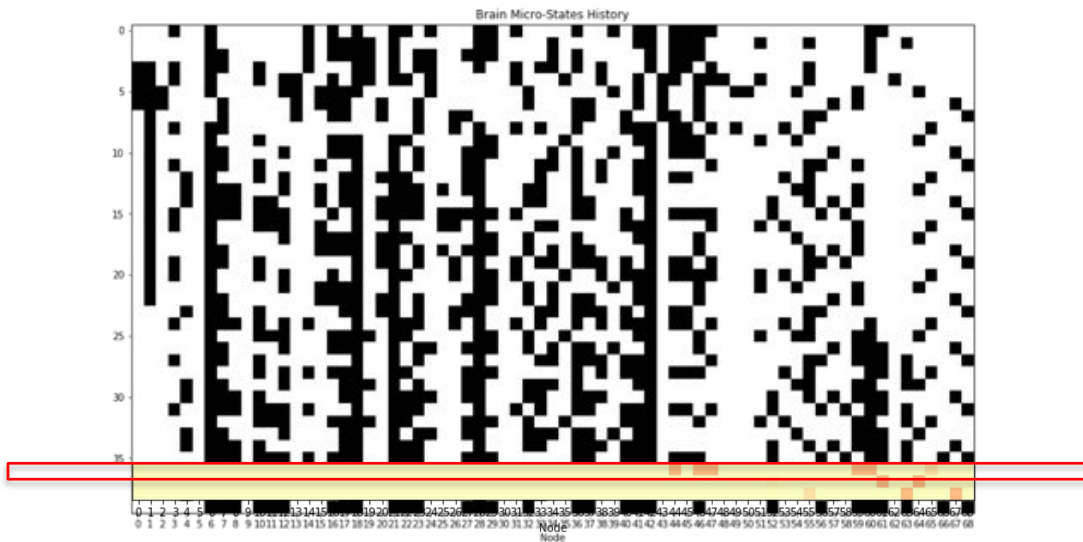
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



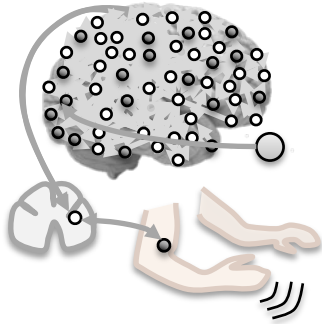
Tracing back the causal chain



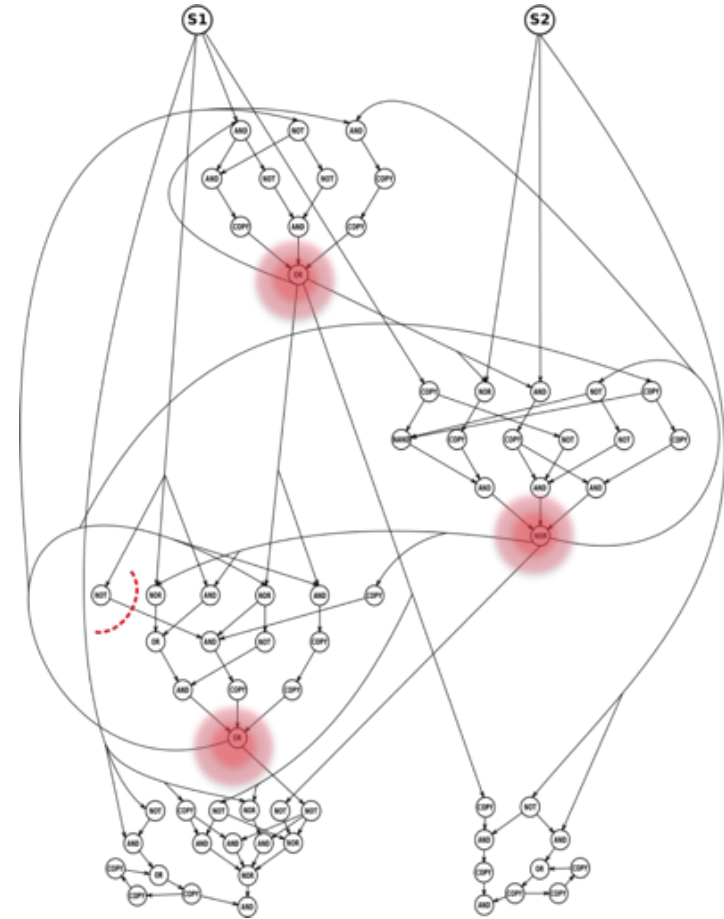
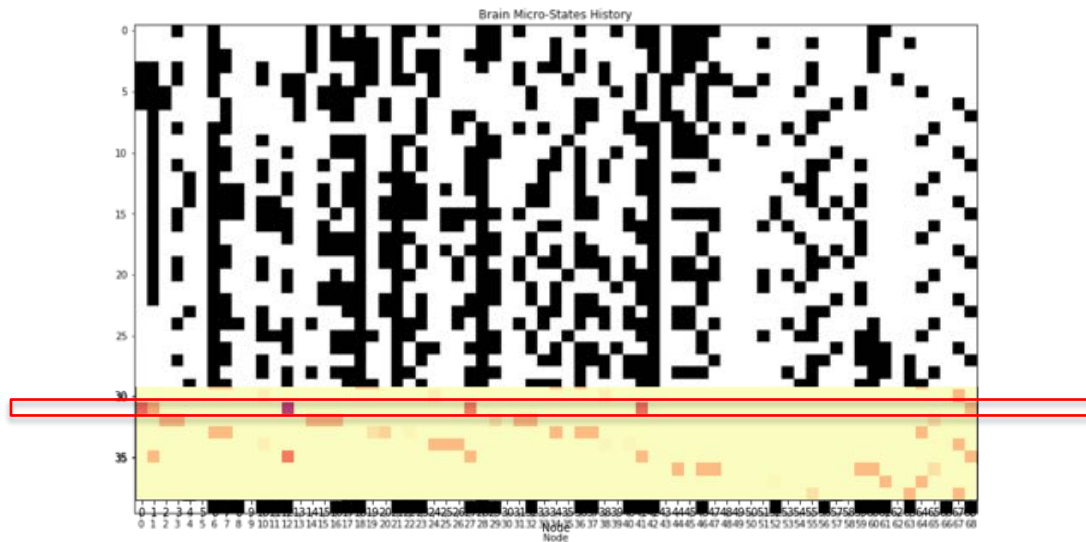
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



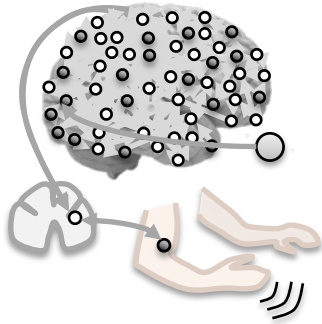
Tracing back the causal chain



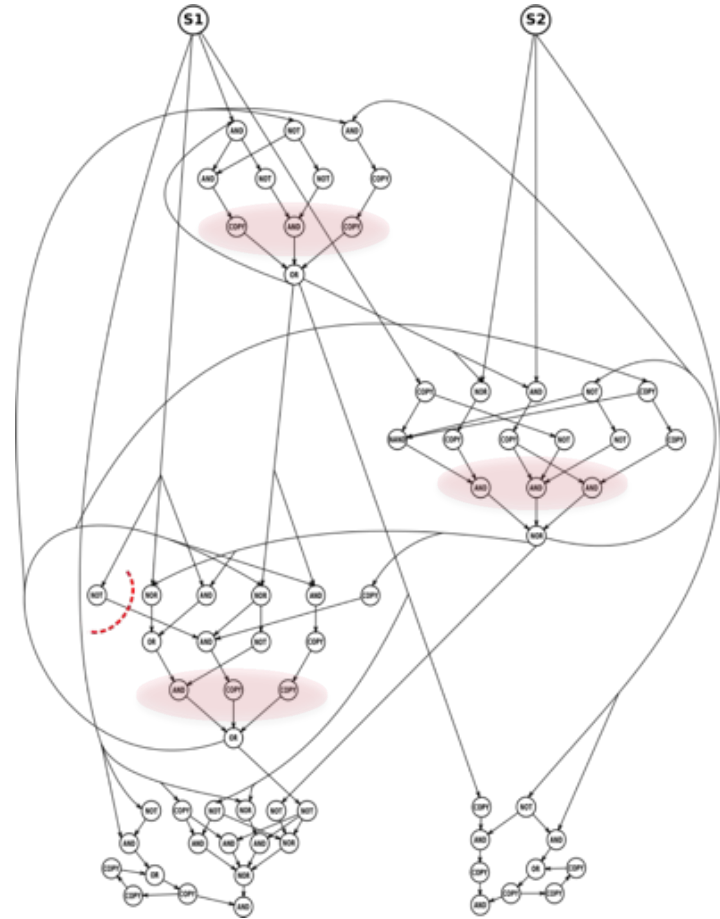
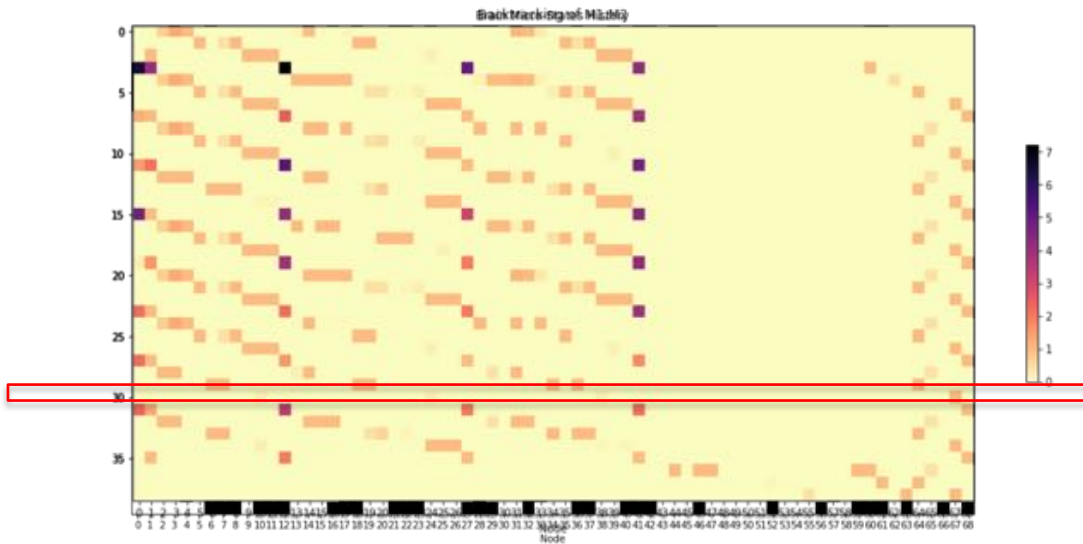
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



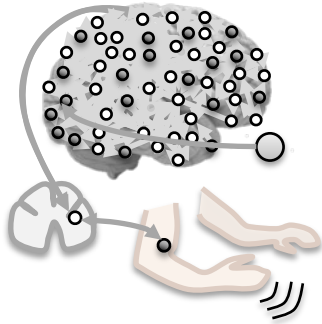
Tracing back the causal chain



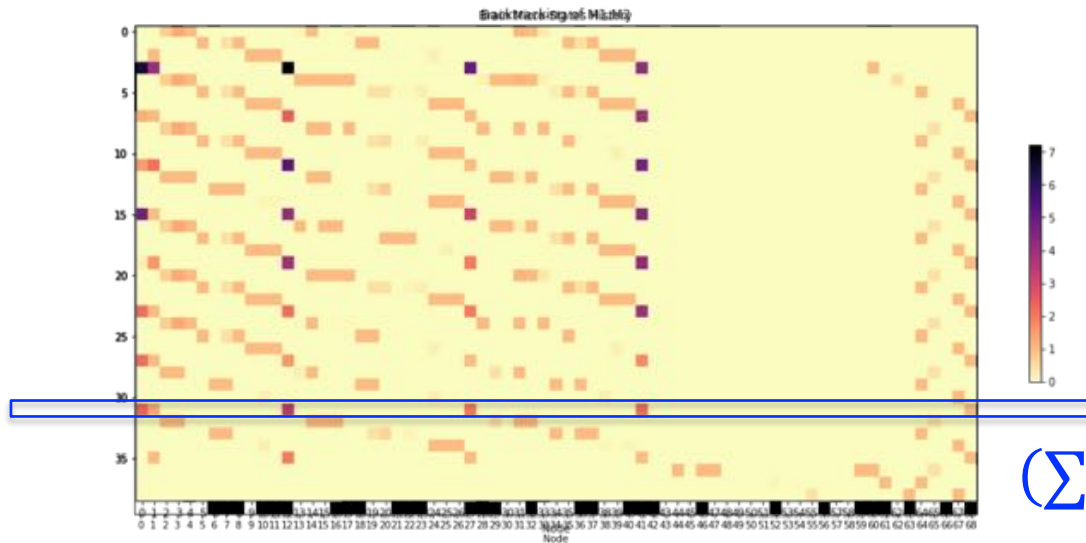
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



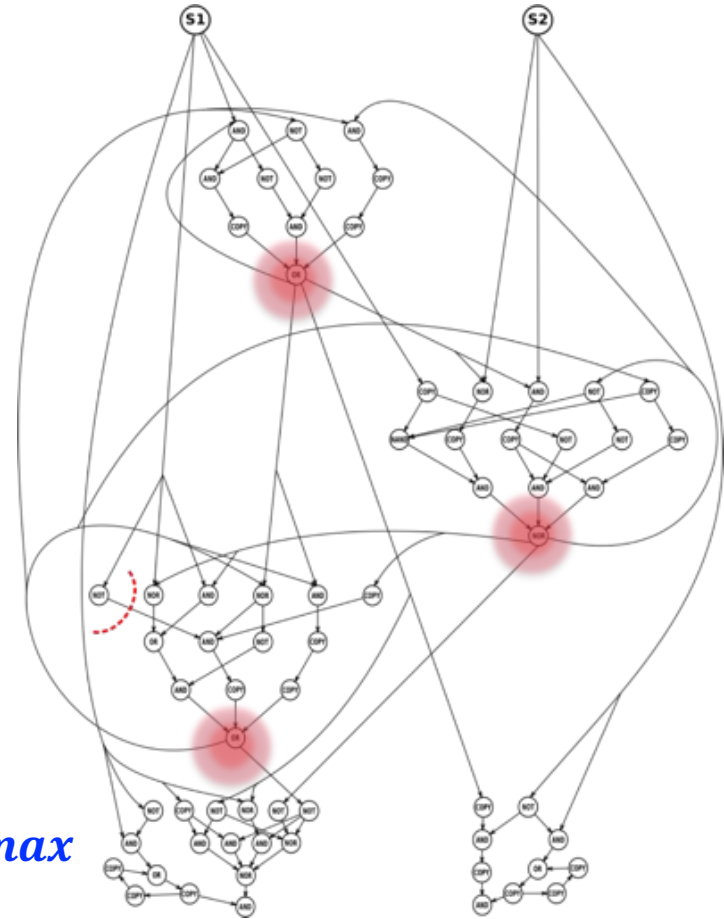
Tracing back the causal chain



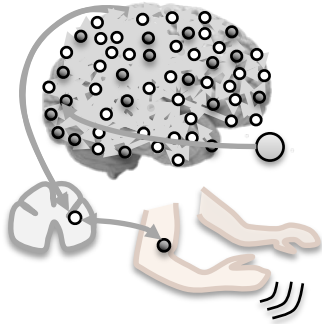
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



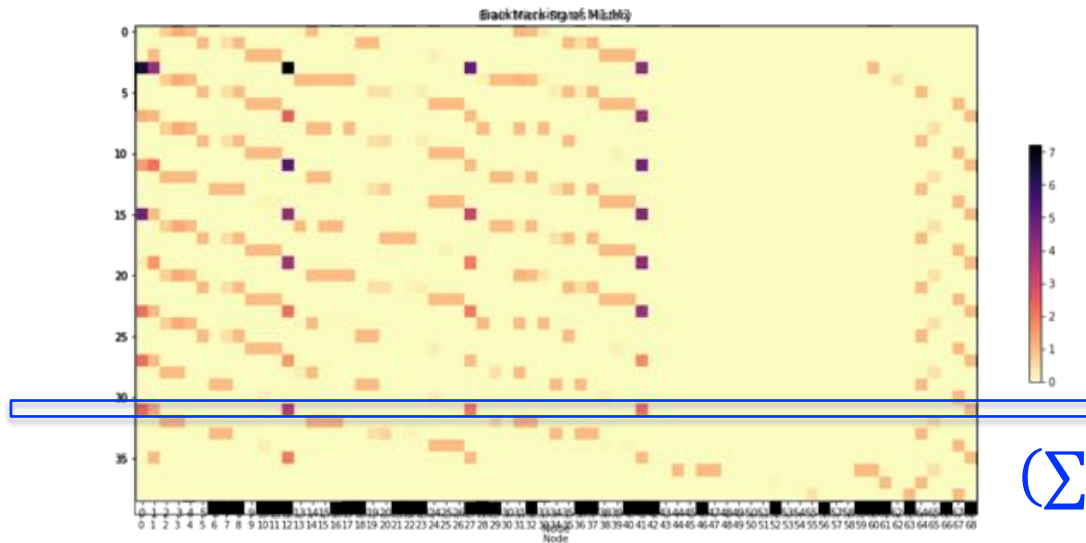
$$(\sum \alpha)^{max}$$



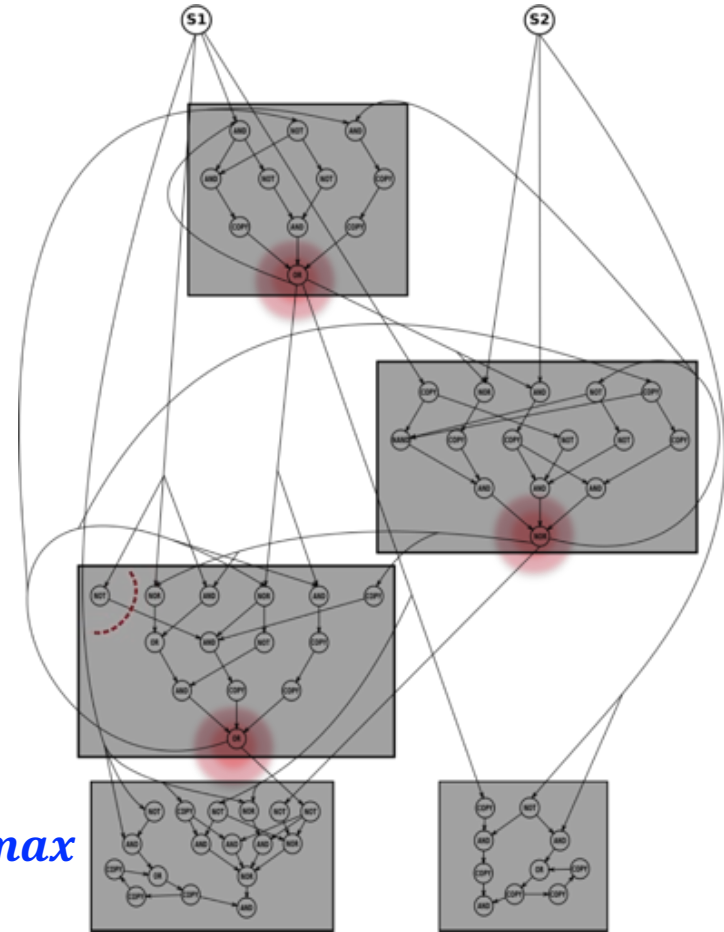
Tracing back the causal chain



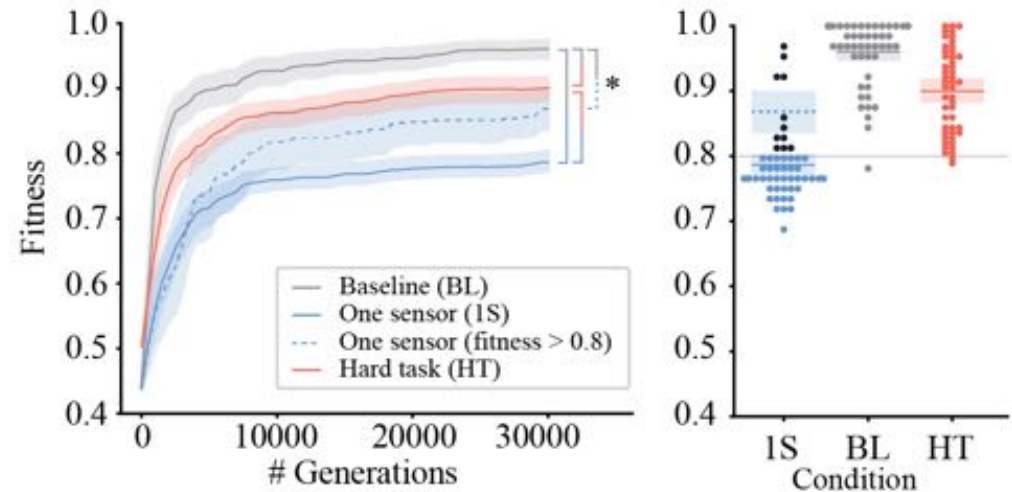
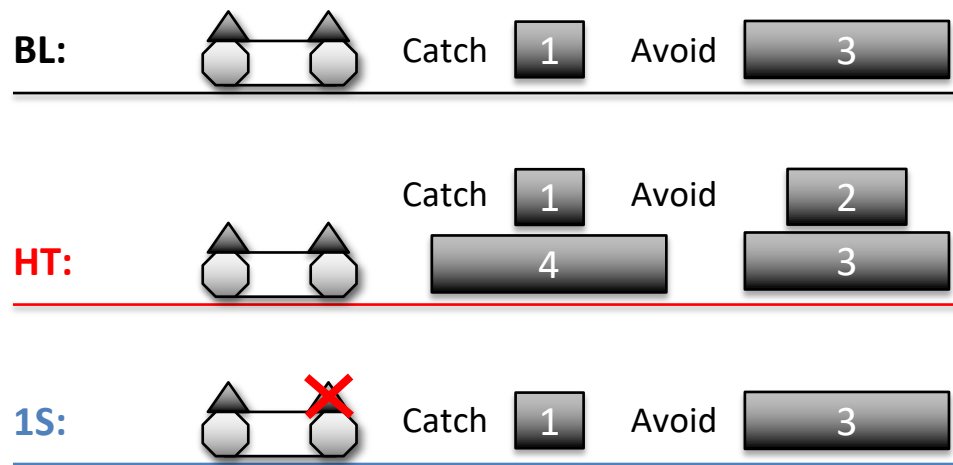
“Neurons were activated and interacted with other neurons which then triggered a motor response.”



$$(\sum \alpha)^{max}$$



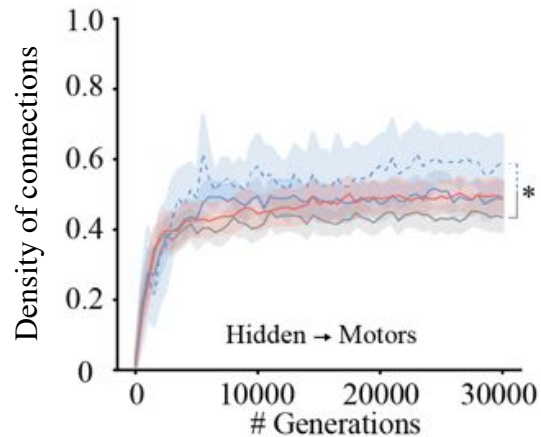
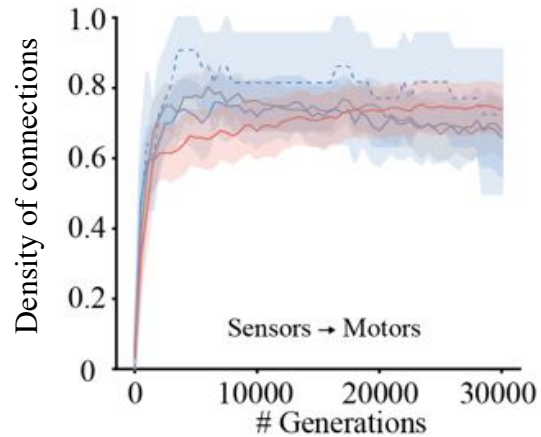
Three different evolutionary environments



Relative task difficulty: $BL < HT < 1S$

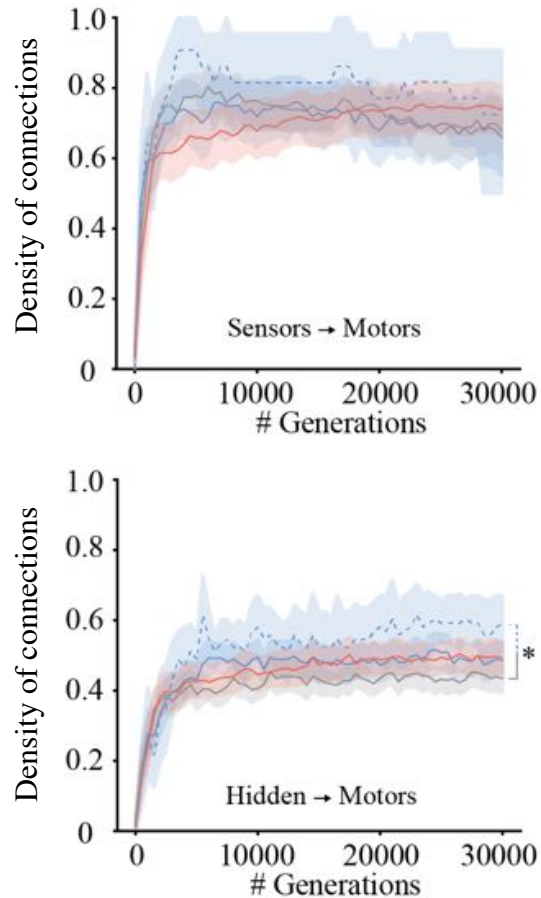
No difference in motor connectivity across conditions

Structural properties:

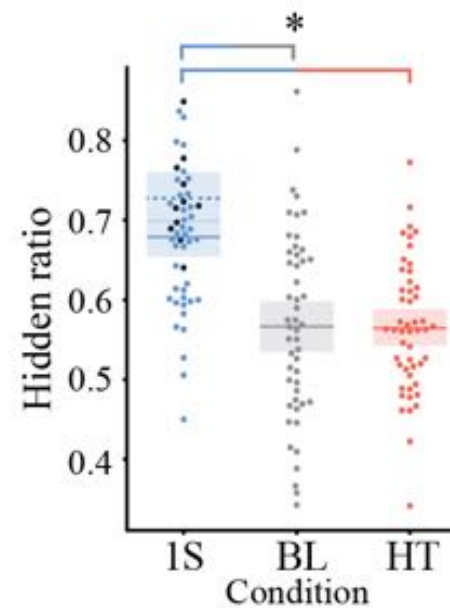


Backtracking analysis reveals differences between evolutionary conditions

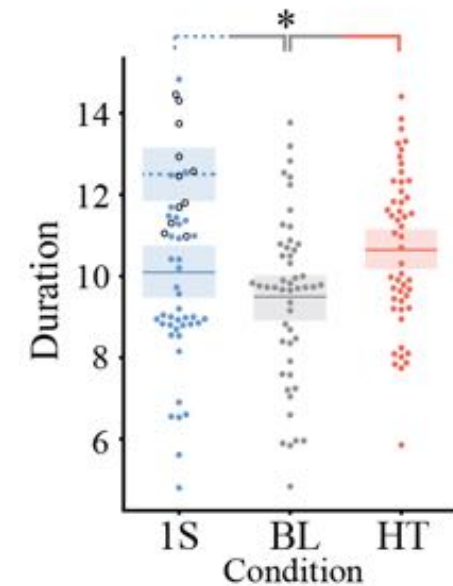
Structural properties:



Backtracking:



1S: More intrinsic



More memory and involvement of hidden states in 1S and HT.

An autonomous agent is:

1. an open system with stable, self-defined and self-maintained causal borders,
2. with the capacity to perform actions that are (partially) caused from within

1. Identifying self-defined causal borders

Example 1: Fission yeast cell cycle model

Marshall W, Kim H, Walker SI, Tononi G, Albantakis L (2017)

Example 2: Evolved artificial agents (animats)

Albantakis L, Hintze A, Koch C, Adami C, Tononi G (2014)

2. Identifying the actual causes of actions

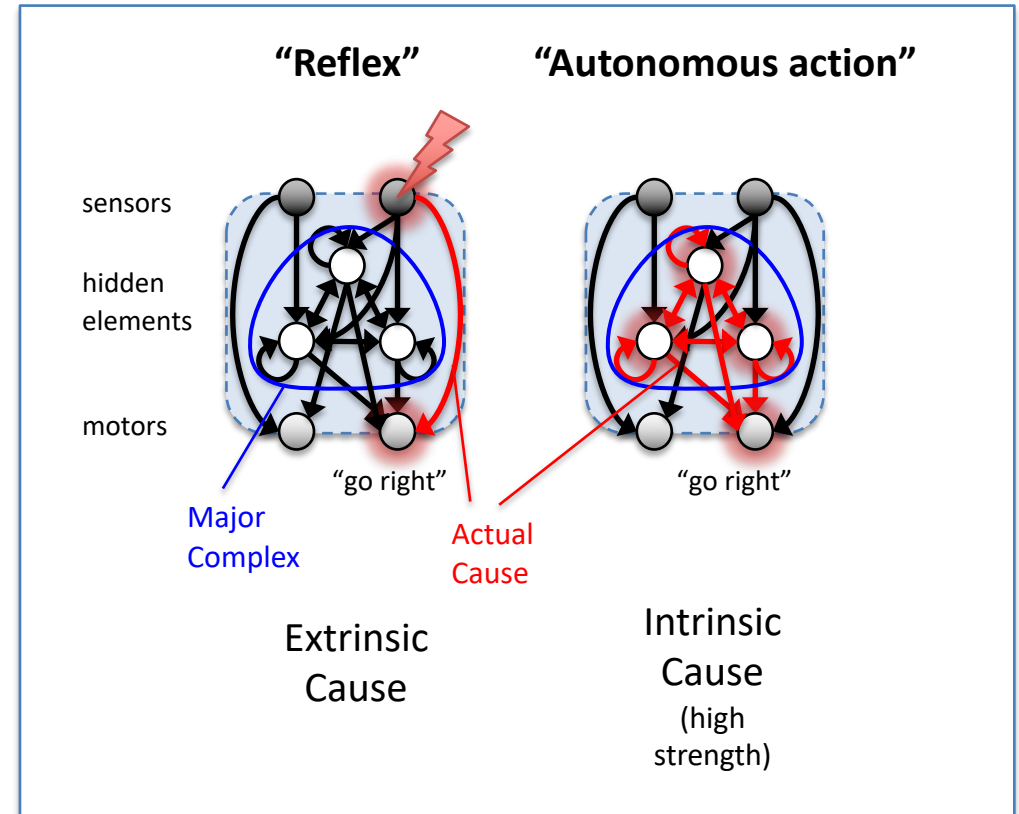
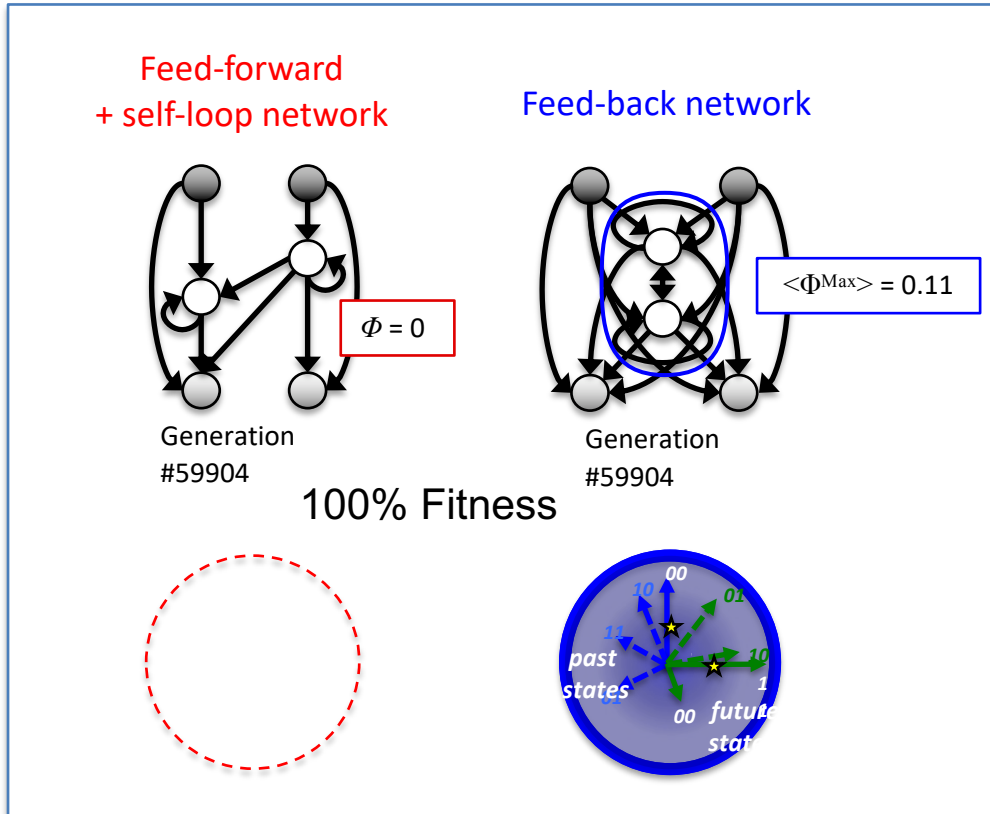
Example 1: Tracing back the causal chain

Albantakis L, Massari F, Tononi G (in prep)

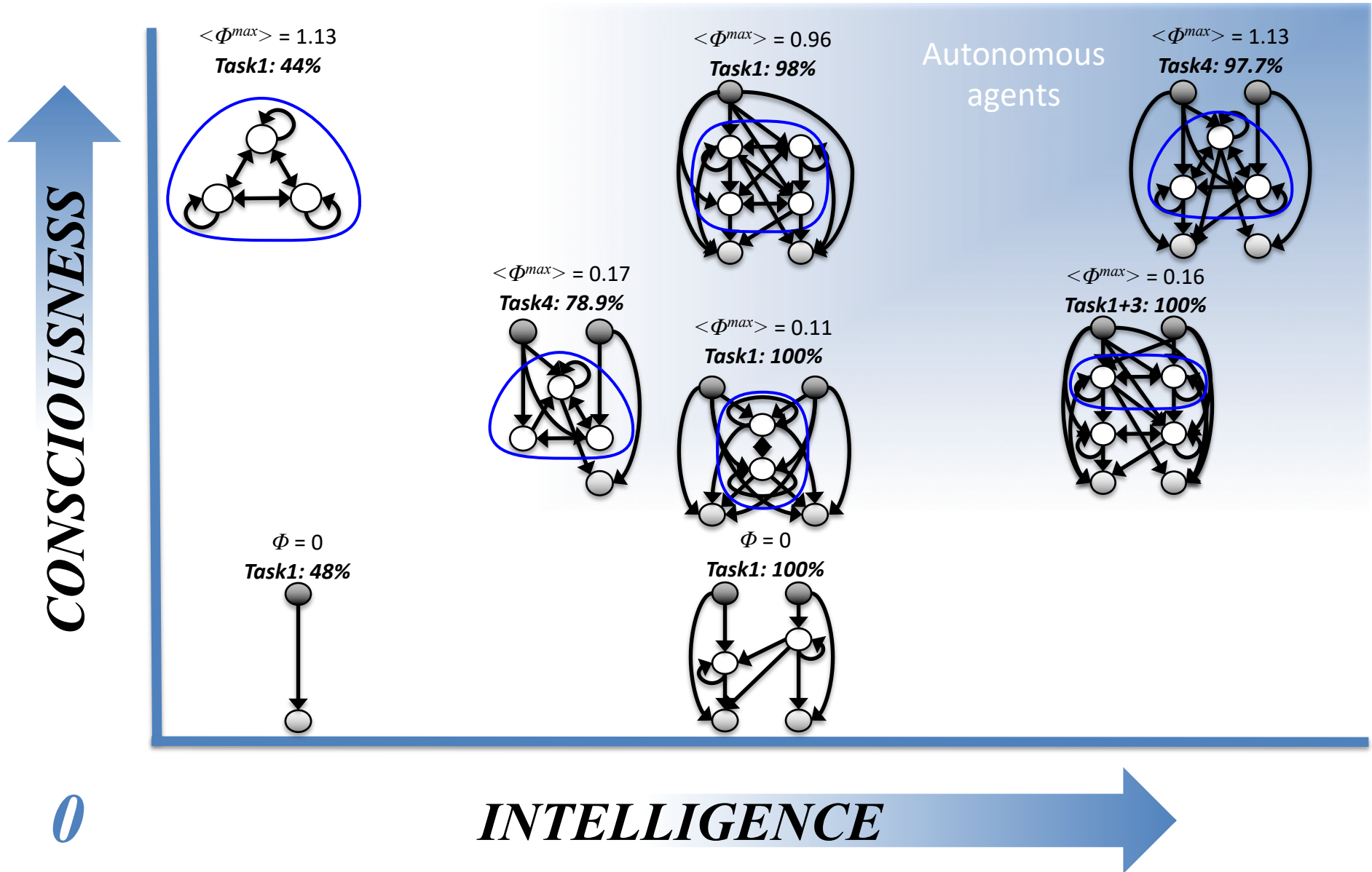
Example 2: Causes of actions across various evolutionary environments

Juel BE, Comolatti R, Tononi G, Albantakis L (2019)

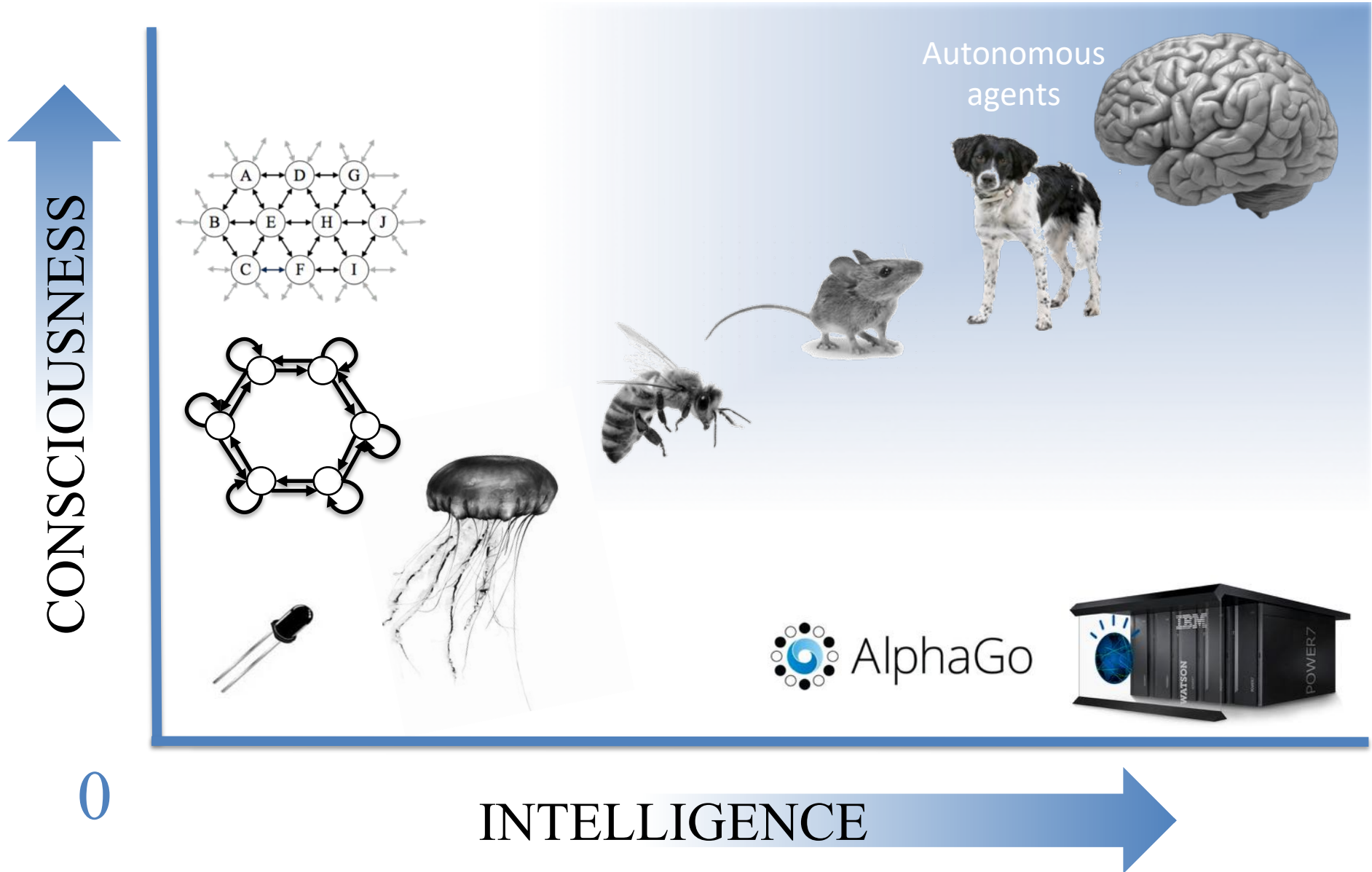
Two dissociations



Intelligence vs. Consciousness



Intelligence vs. Consciousness



Thank you!

<http://integratedinformationtheory.org>



Wisconsin Institute for
Sleep and Consciousness
UNIVERSITY OF WISCONSIN-MADISON



TEMPLETON WORLD
CHARITY FOUNDATION

Giulio Tononi
William Marshall
Bjørn Juel
Renzo Comolatti
Will Mayner

THE
TINY BLUE DOT
FOUNDATION