

# The Cost of Sharing: From Topological Multitasking Limits to Semantic Horizons

G. Petri  
NSIA Workshop, Netsci 2026  
1/6/2026

**Network Science Institute**  
at **Northeastern University**

[ **NPL** ]  
RESEARCH



# The Cost of Sharing: From Topological Multitasking Limits to Semantic Horizons

G. Petri  
NSIA Workshop, Netsci 2026  
1/6/2026

**Network Science Institute**  
at **Northeastern University**

[ **NPL** ]  
RESEARCH



# Experiment #1

Name the color of the following stimulus  
and, *at the same time*, point to where it is...

# Experiment #1

Name the color of the following stimulus  
and, *at the same time*, point to where it is...

**BROWN**

# Experiment #1

Name the color of the following stimulus  
and, *at the same time*, point to where it is...

**BLUE**



**YELLOW**



# Experiment #2

point left if the written word is

 **RED**

point right if the written word is

**GREEN** 

# Experiment #2

point left if the written word is

← RED

point right if the written word is

GREEN →

RED

# Experiment #2

point left if the written word is

 **RED**

point right if the written word is

**GREEN** 

**GREEN**



**RED**



# Experiment #3

Name the color of the following stimulus  
and, *at the same time*:

point left if the written word is

← **RED**

point right if the written word is

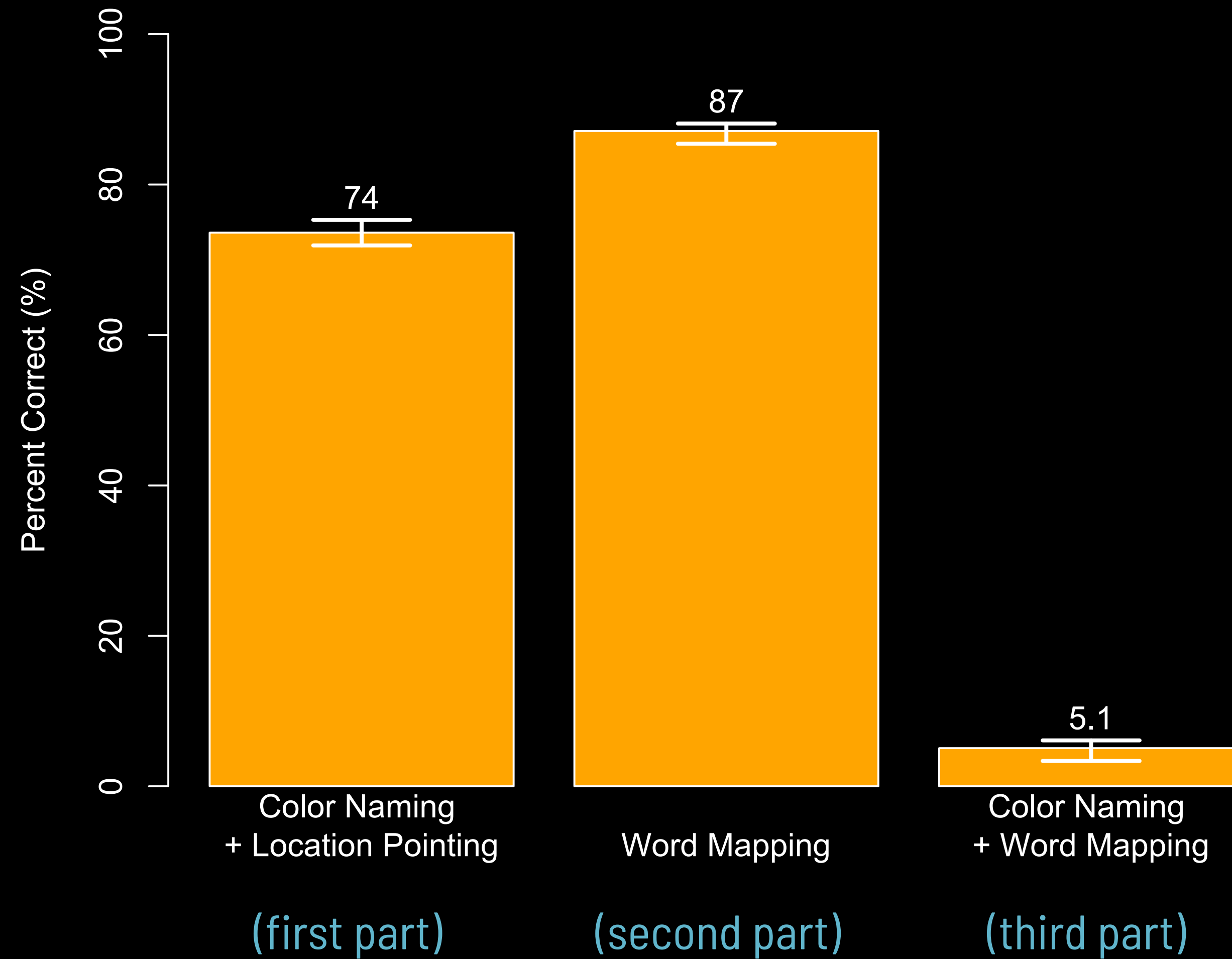
**GREEN** →

**GREEN**

**RED**

# Experiment

## Accuracy Results



Under which conditions can we multitask?

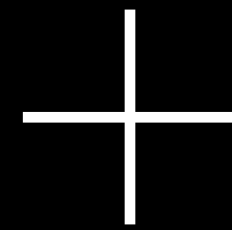
+

Under which conditions can we multitask?

+

Understanding general computation

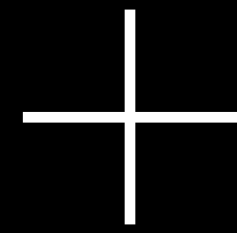
Under which conditions can we multitask?



Understanding general computation

Practical application in design of sensitive systems

Under which conditions can we multitask?

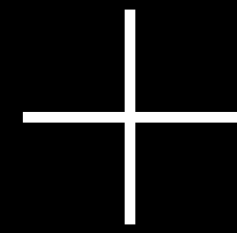


Understanding general computation

Practical application in design of sensitive systems

Automated AI agents

Under which conditions can we multitask?



**What is the closest object to the red cross?**

**What is the closest object to the red cross?**















These two tasks share a  
fundamental similarity

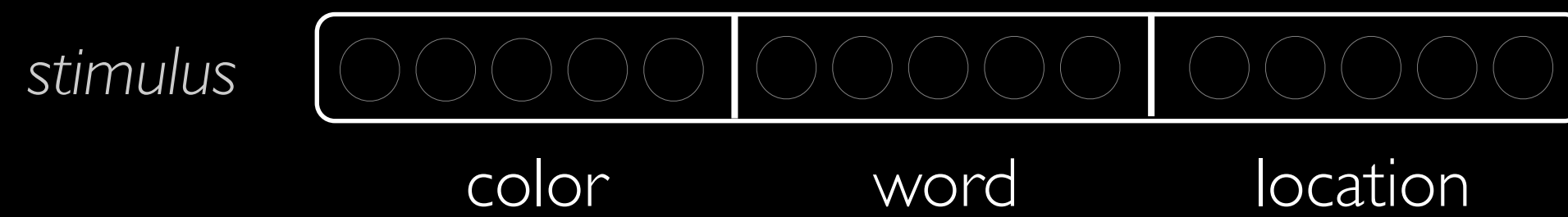
These two tasks share a  
fundamental similarity

**Limits emerging from  
information encoding**

**Ok... Why?**

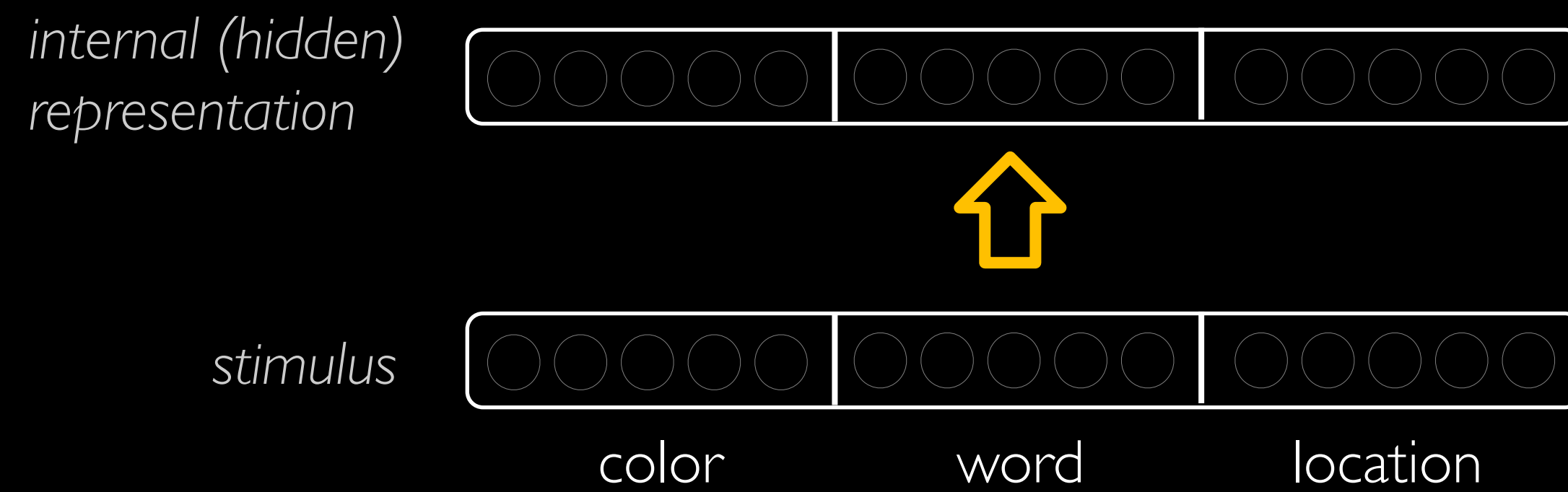
# Neural Network Model of Controlled Processing

# Neural Network Model of Controlled Processing

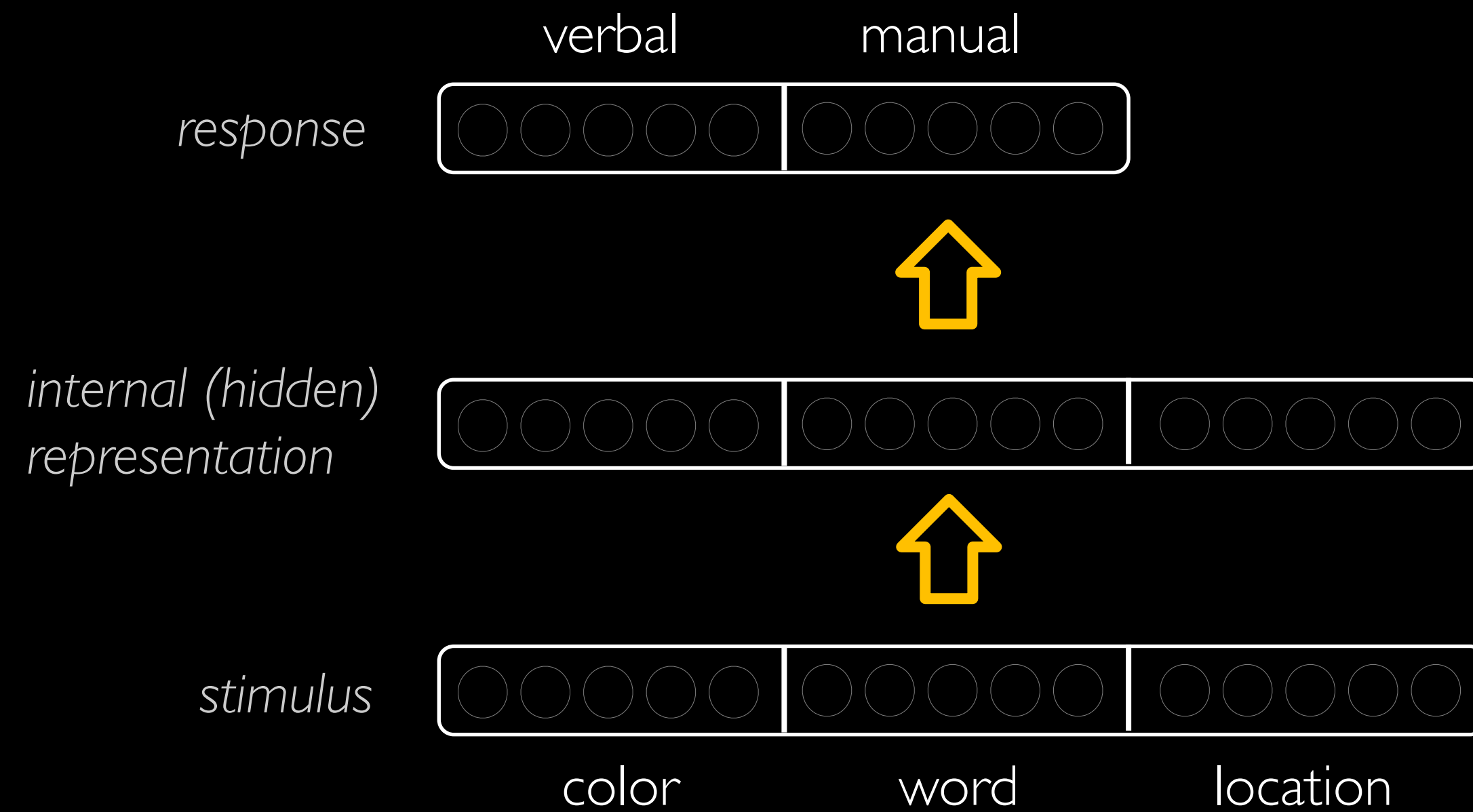


(Cohen et al., 1990 ; Feng et al., 2014; Musslick et al., 2016)

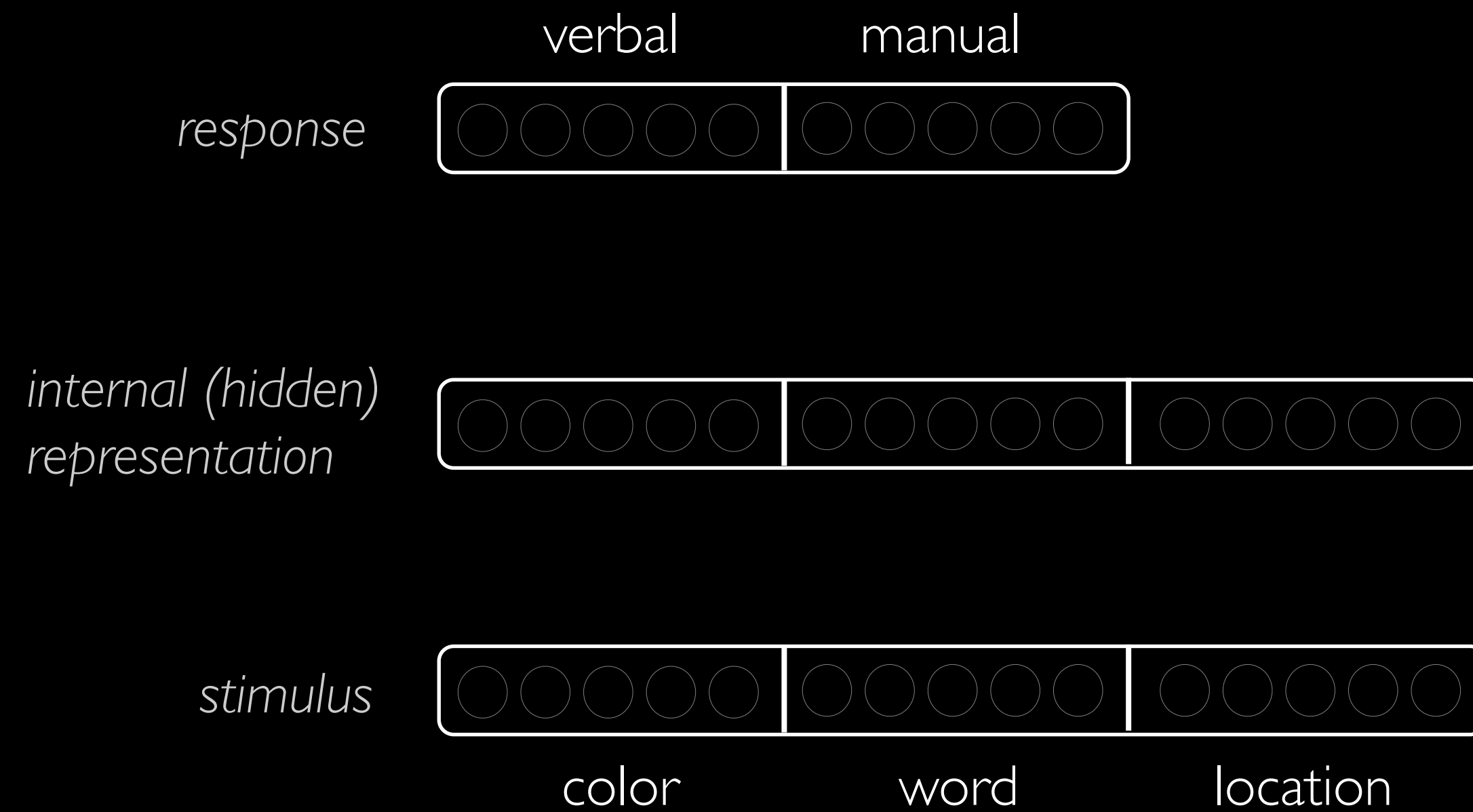
# Neural Network Model of Controlled Processing



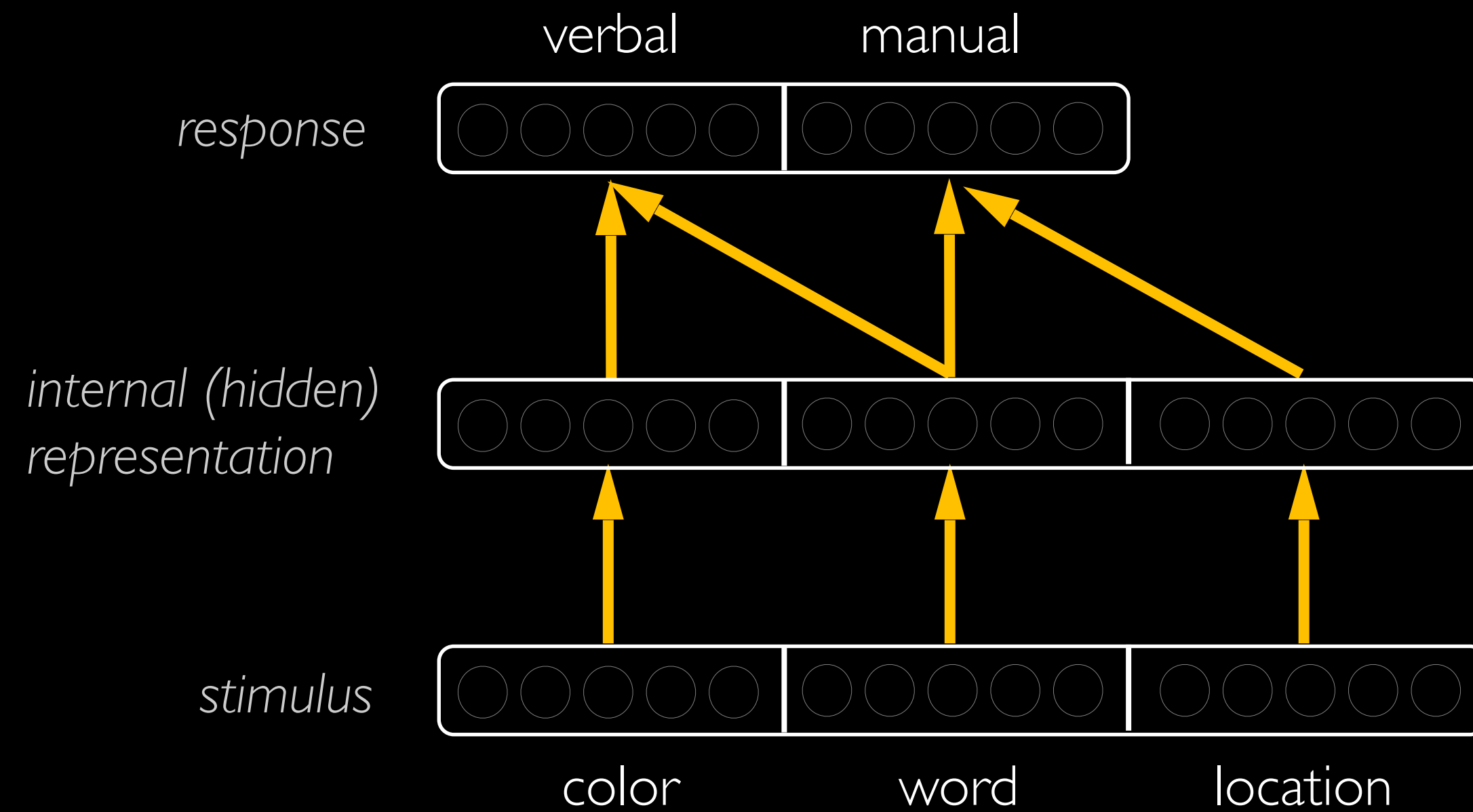
# Neural Network Model of Controlled Processing



# Neural Network Model of Controlled Processing

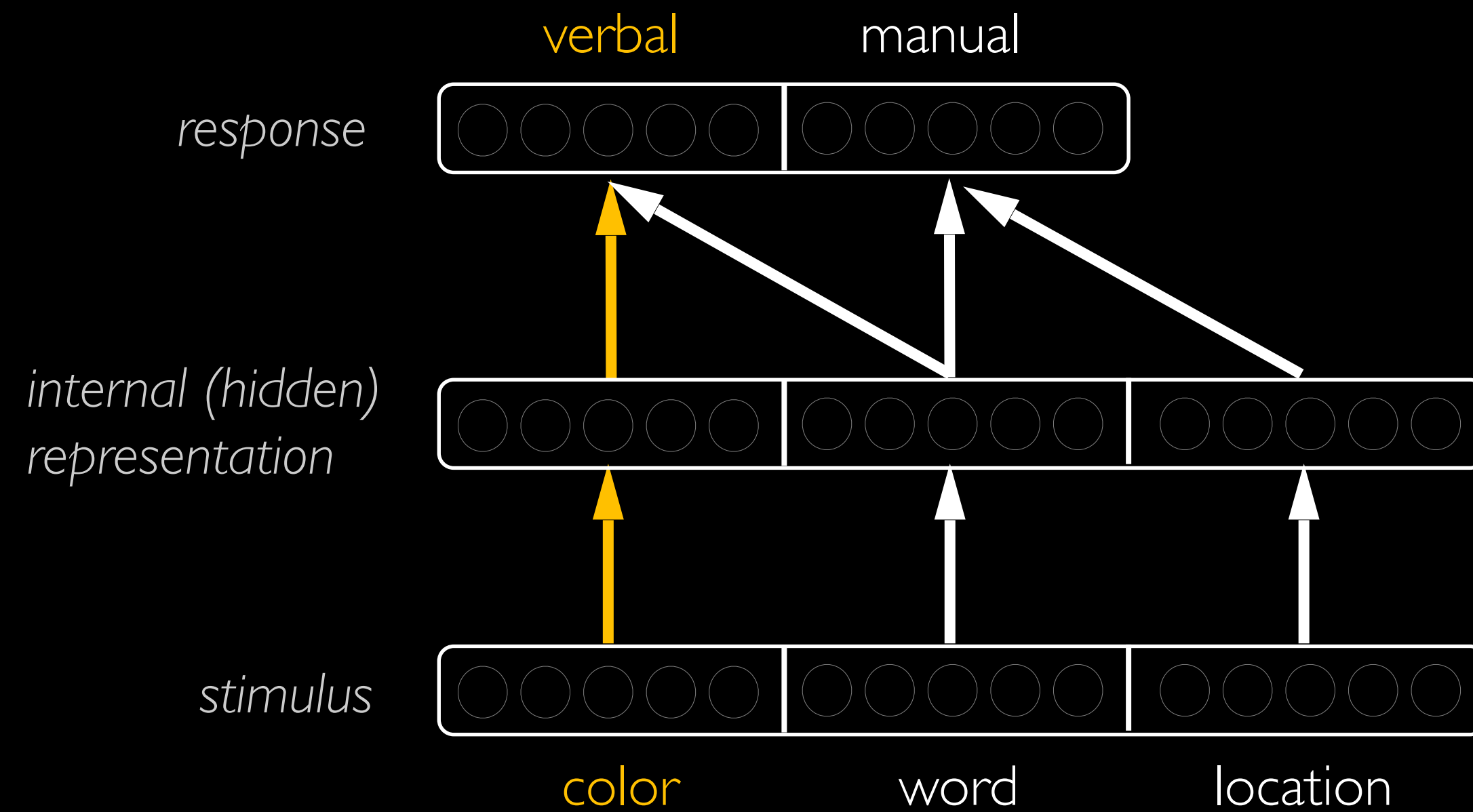


# Neural Network Model of Controlled Processing

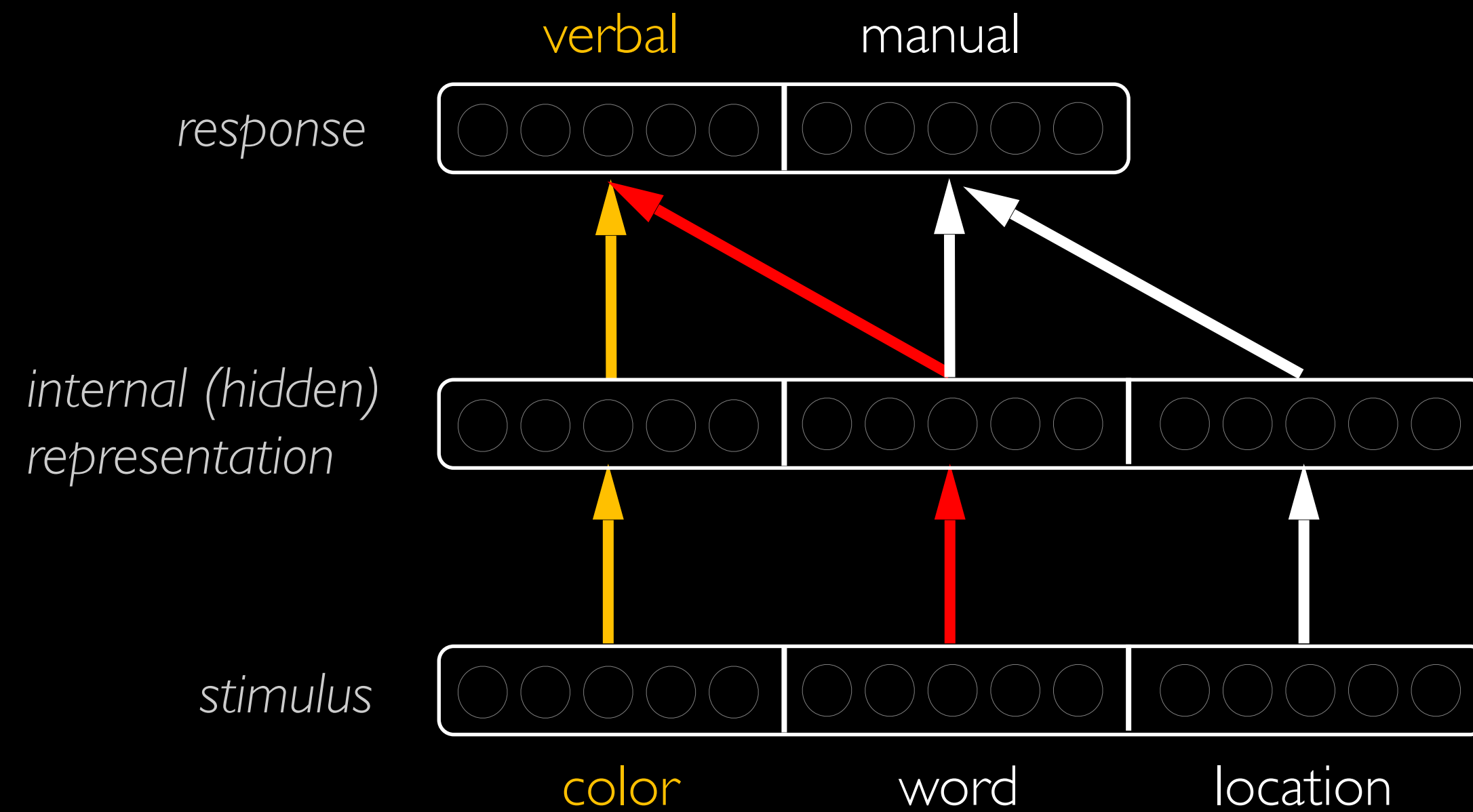


(Cohen et al., 1990 ; Feng et al., 2014; Musslick et al., 2016)

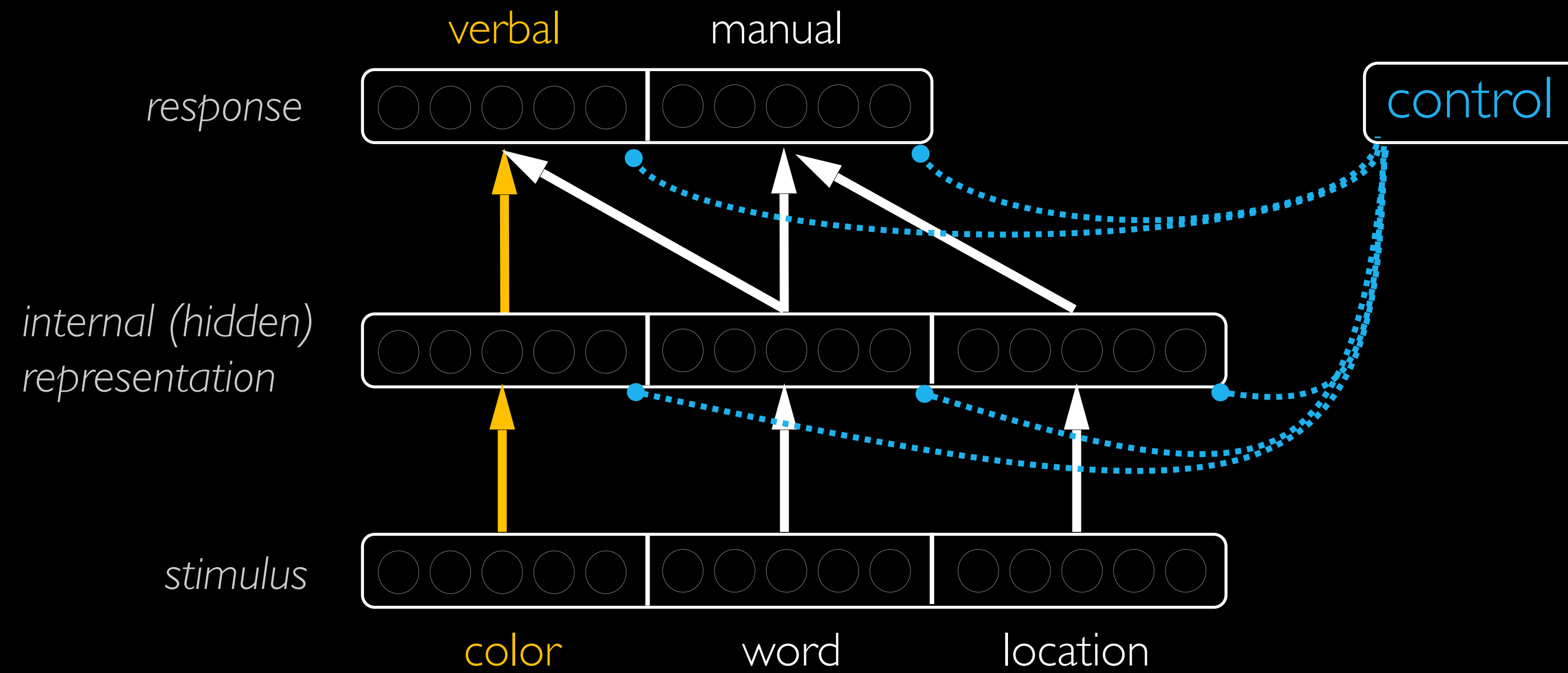
# Neural Network Model of Controlled Processing



# Neural Network Model of Controlled Processing

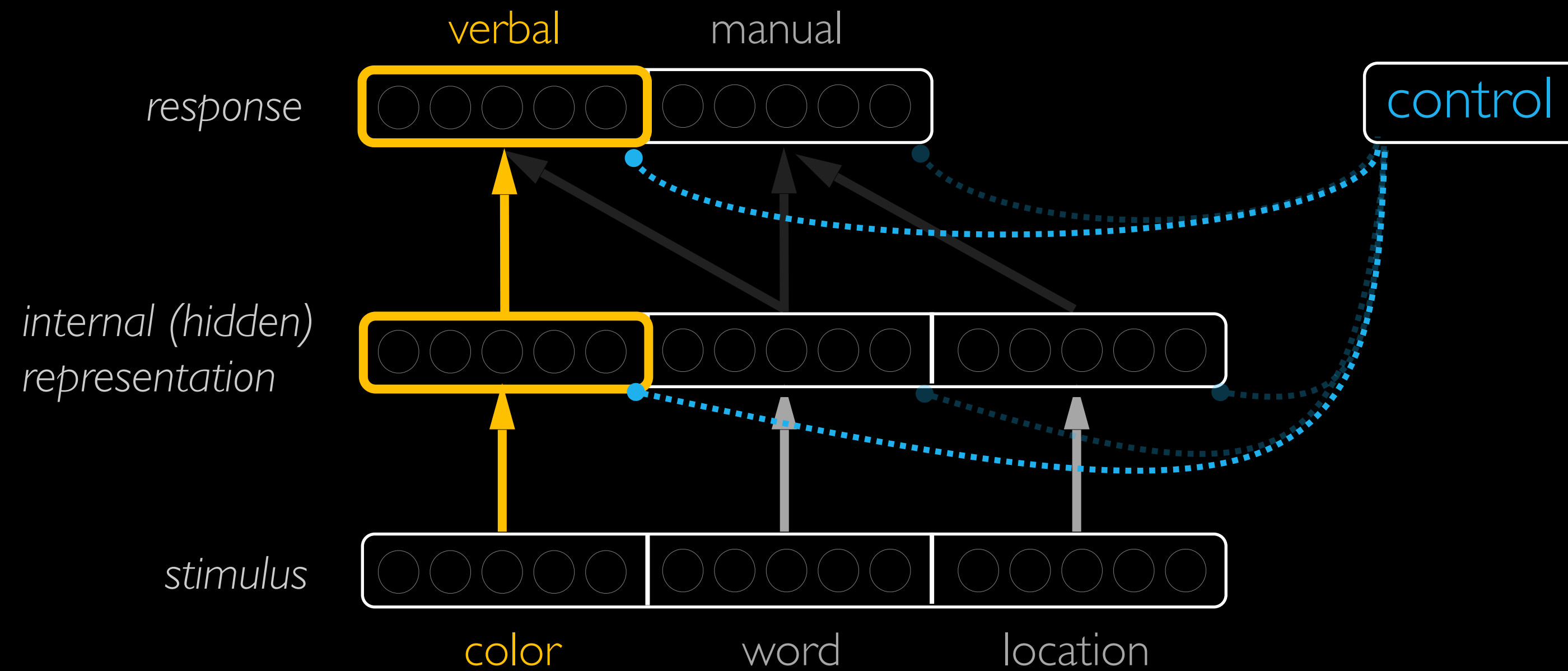


# Neural Network Model of Controlled Processing



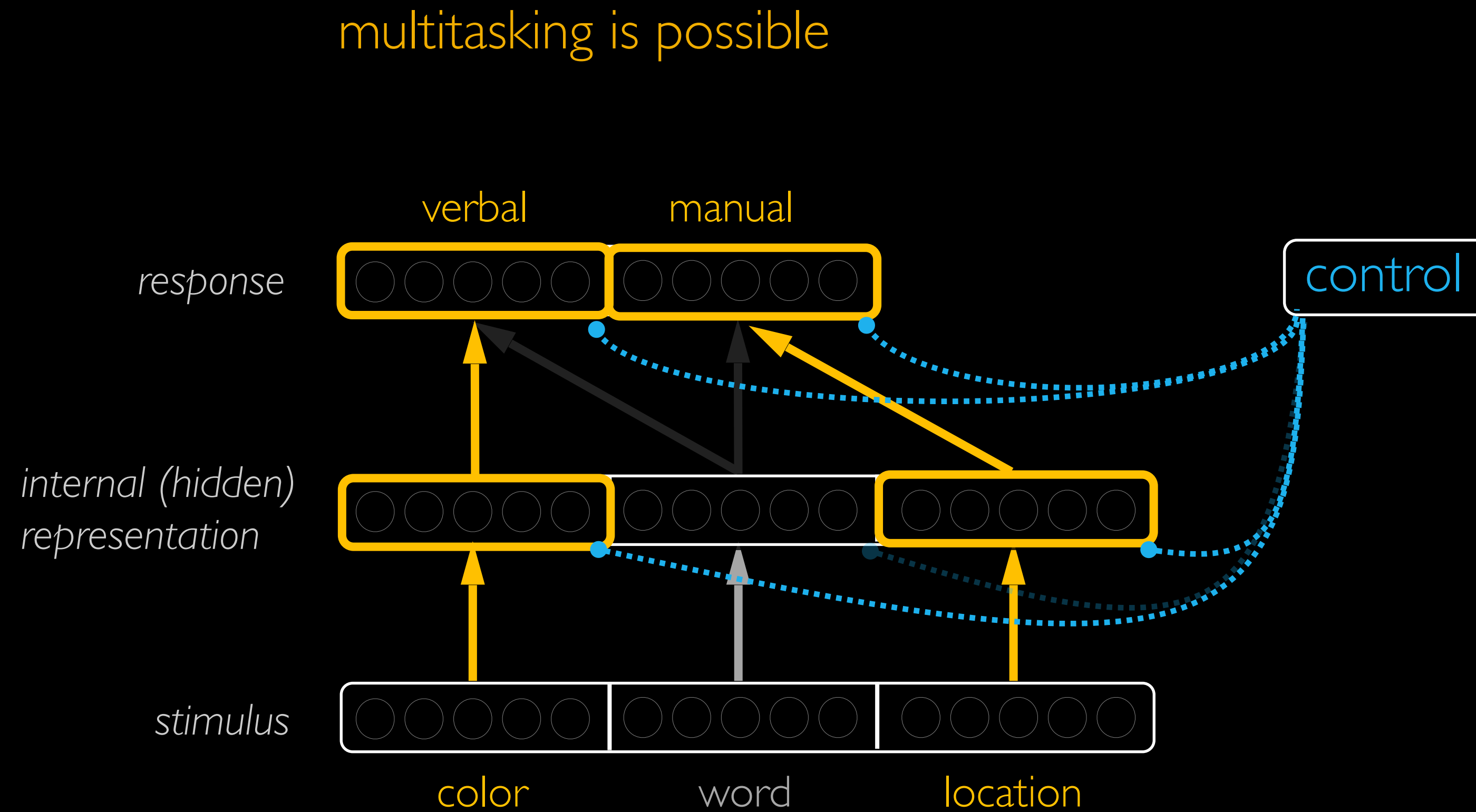
(Cohen et al., 1990 ; Feng et al., 2014; Musslick et al., 2016)

# Neural Network Model of Controlled Processing



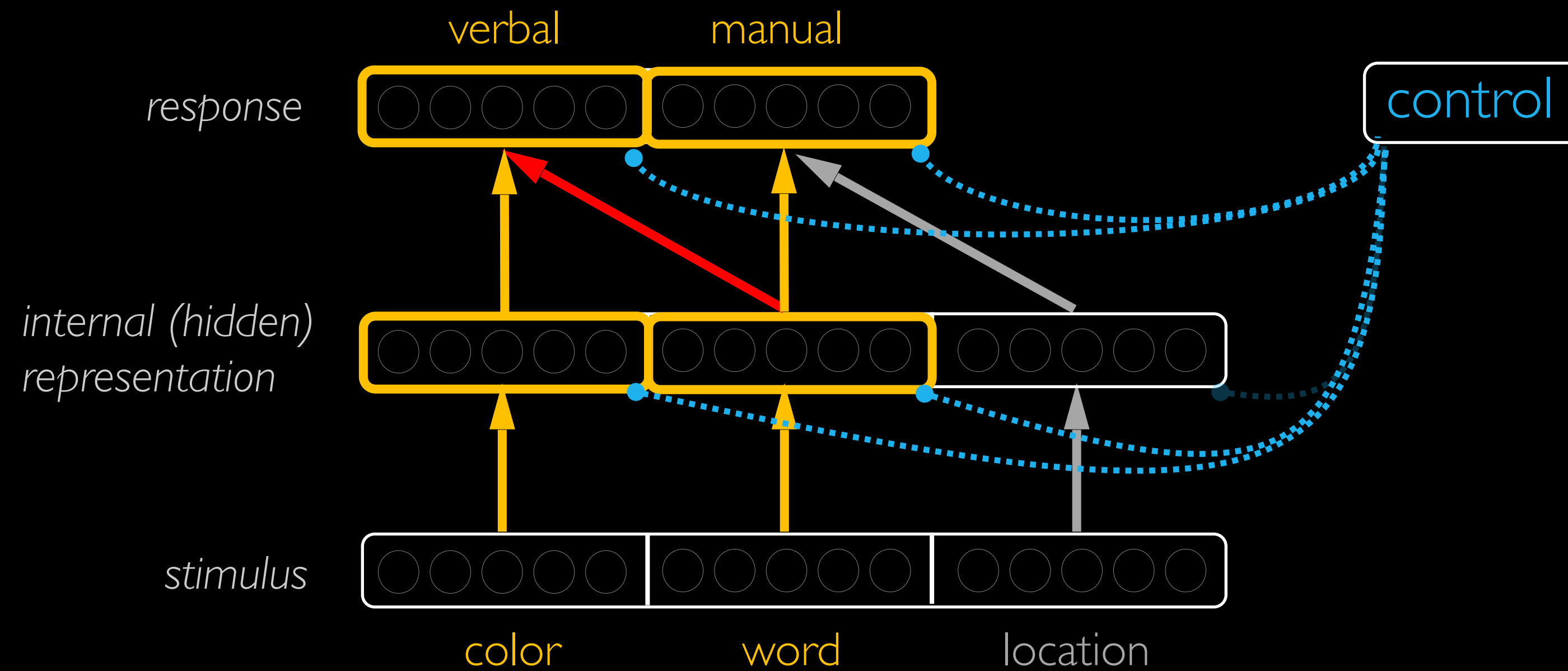
(Cohen et al., 1990 ; Feng et al., 2014; Musslick et al., 2016)

# Neural Network Model of Controlled Processing



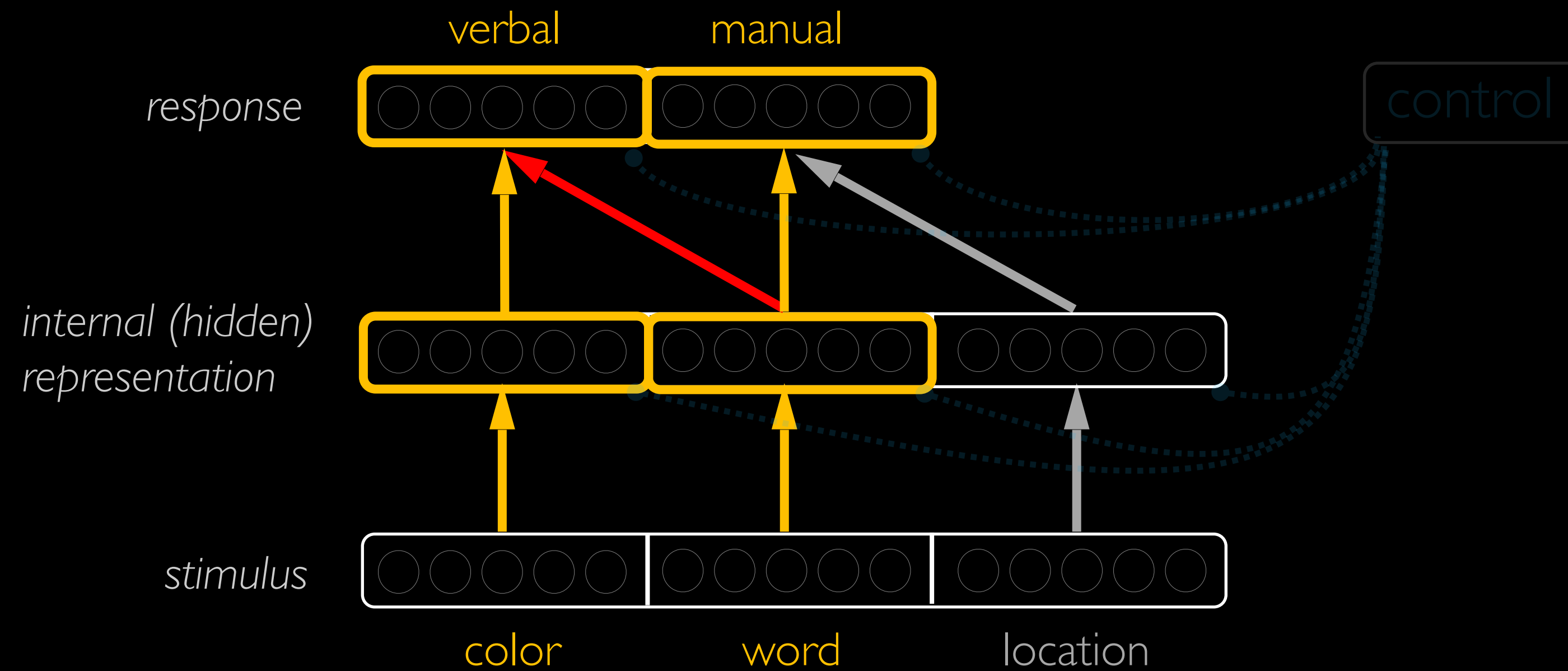
# Neural Network Model of Controlled Processing

multitasking is **not** possible



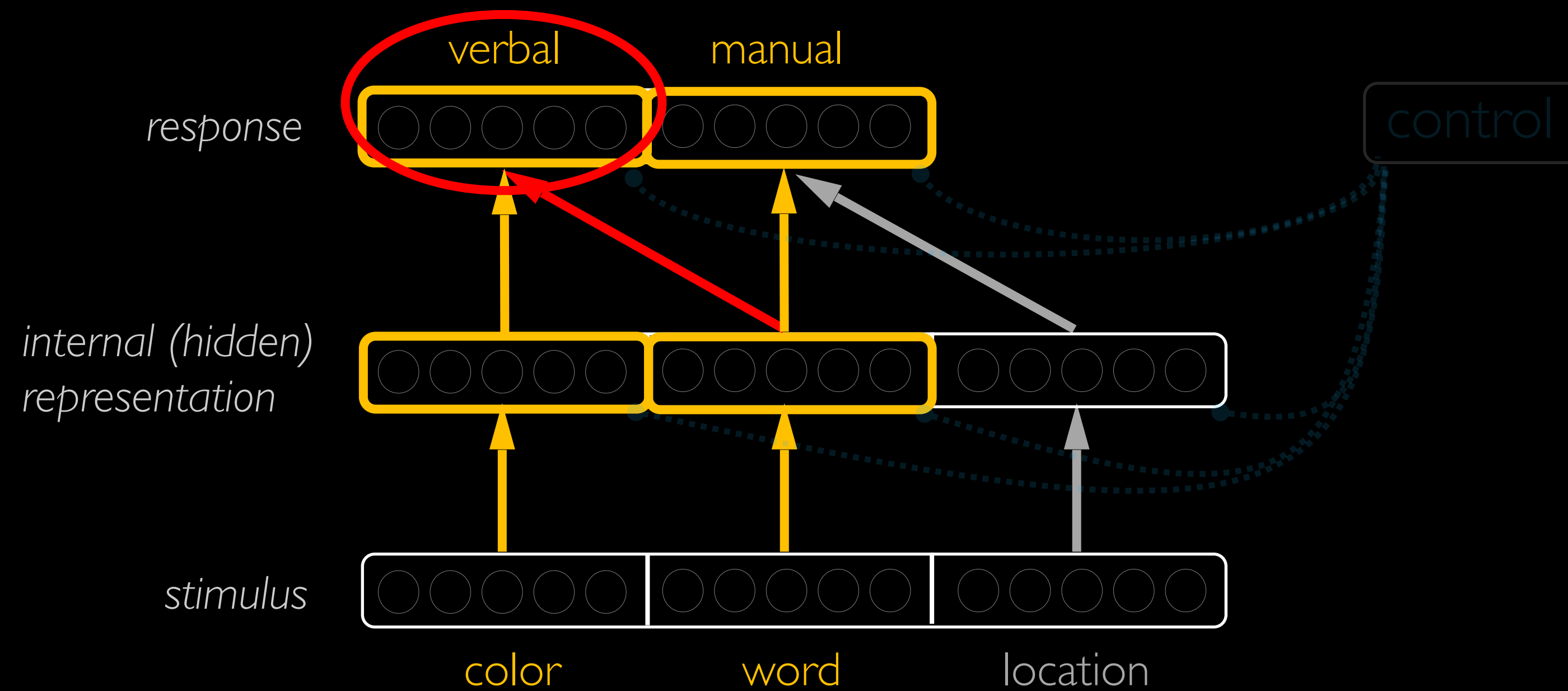
# Neural Network Model of Controlled Processing

multitasking is **not** possible



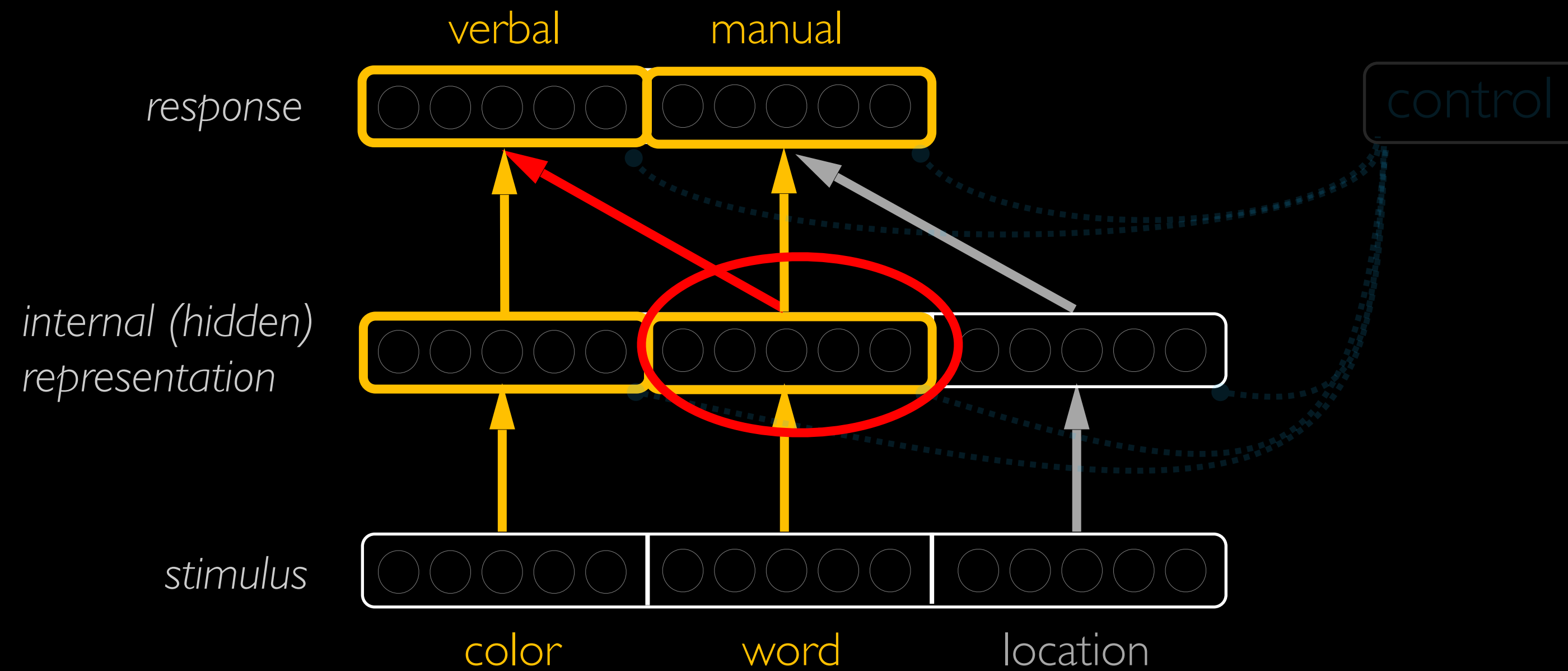
# Neural Network Model of Controlled Processing

multitasking is **not** possible

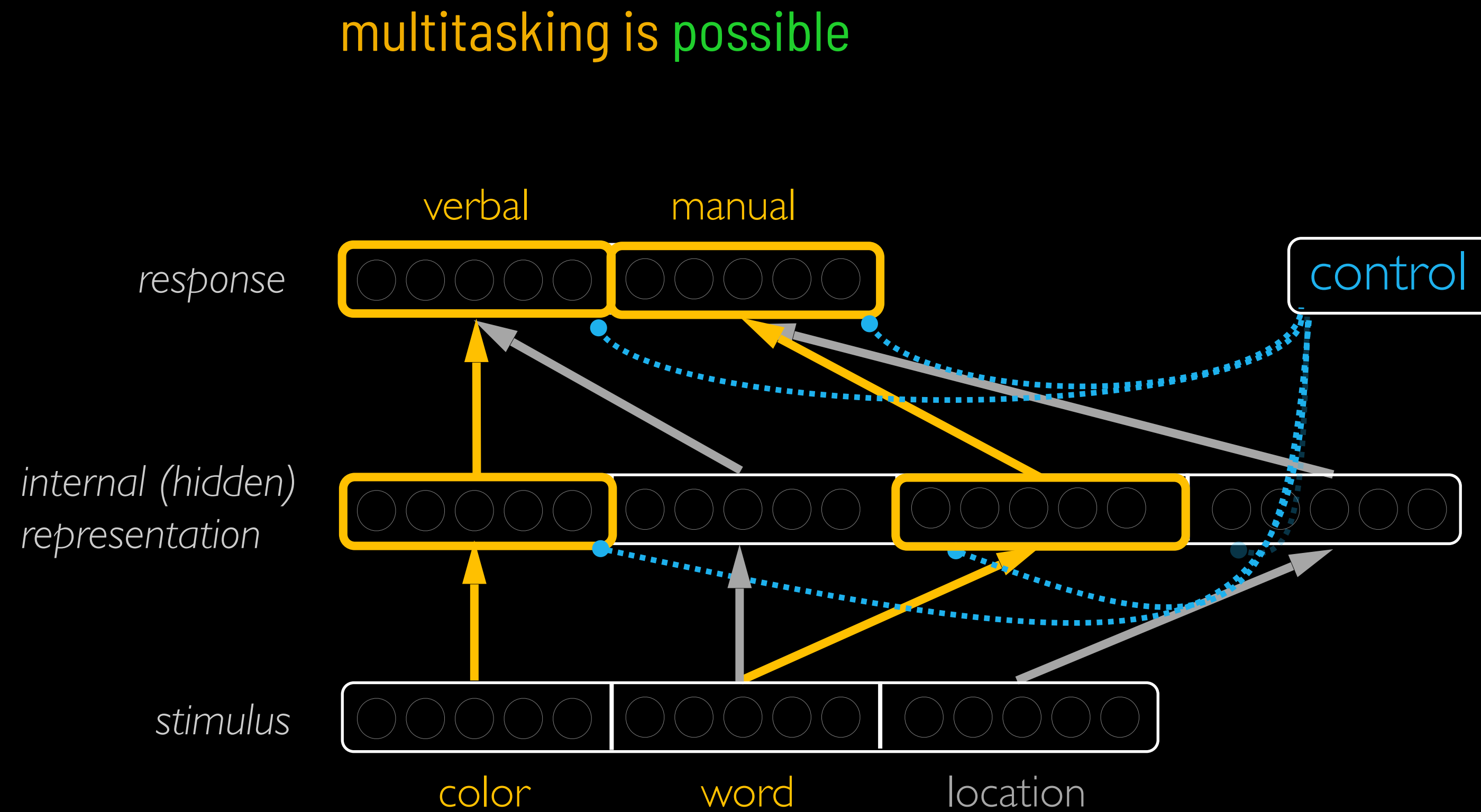


# Neural Network Model of Controlled Processing

multitasking is **not** possible

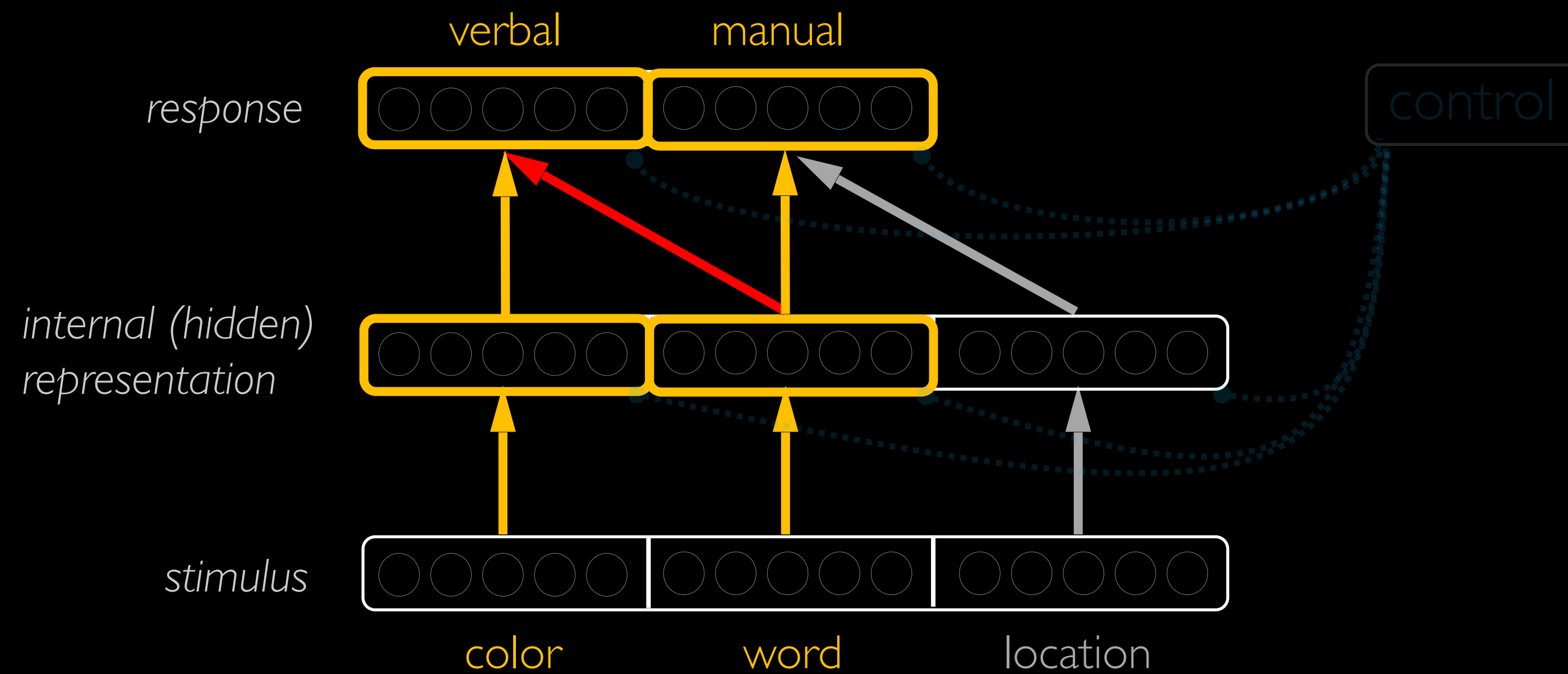


# Neural Network Model of Controlled Processing



(Cohen et al., 1990 ; Feng et al., 2014; Musslick et al., 2016)

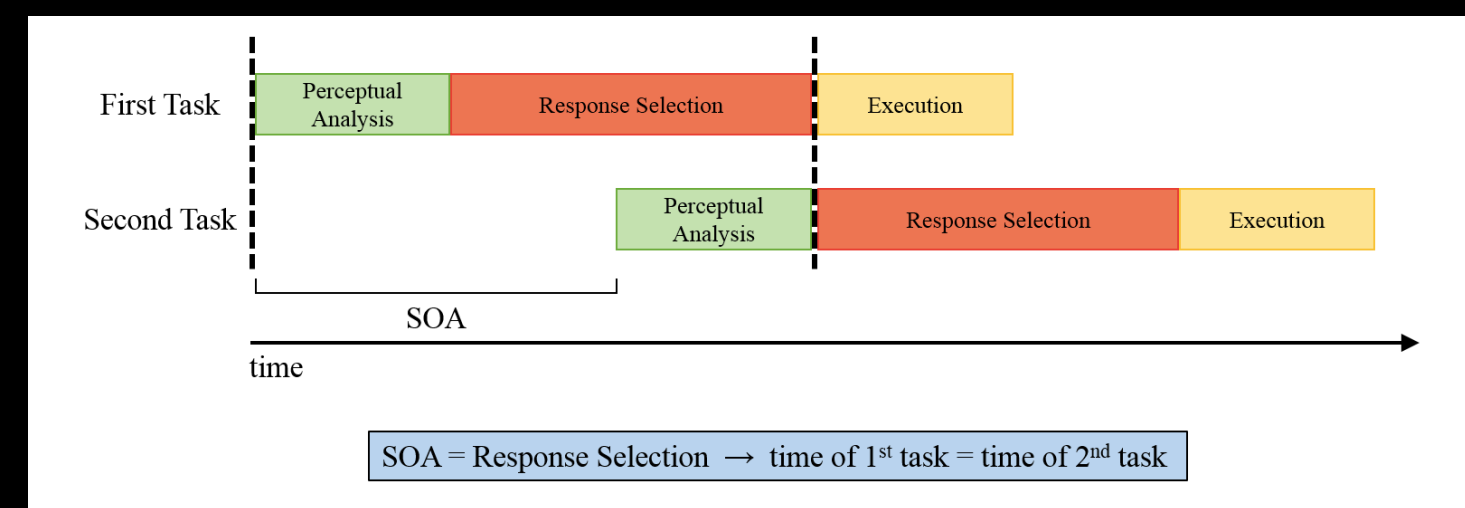
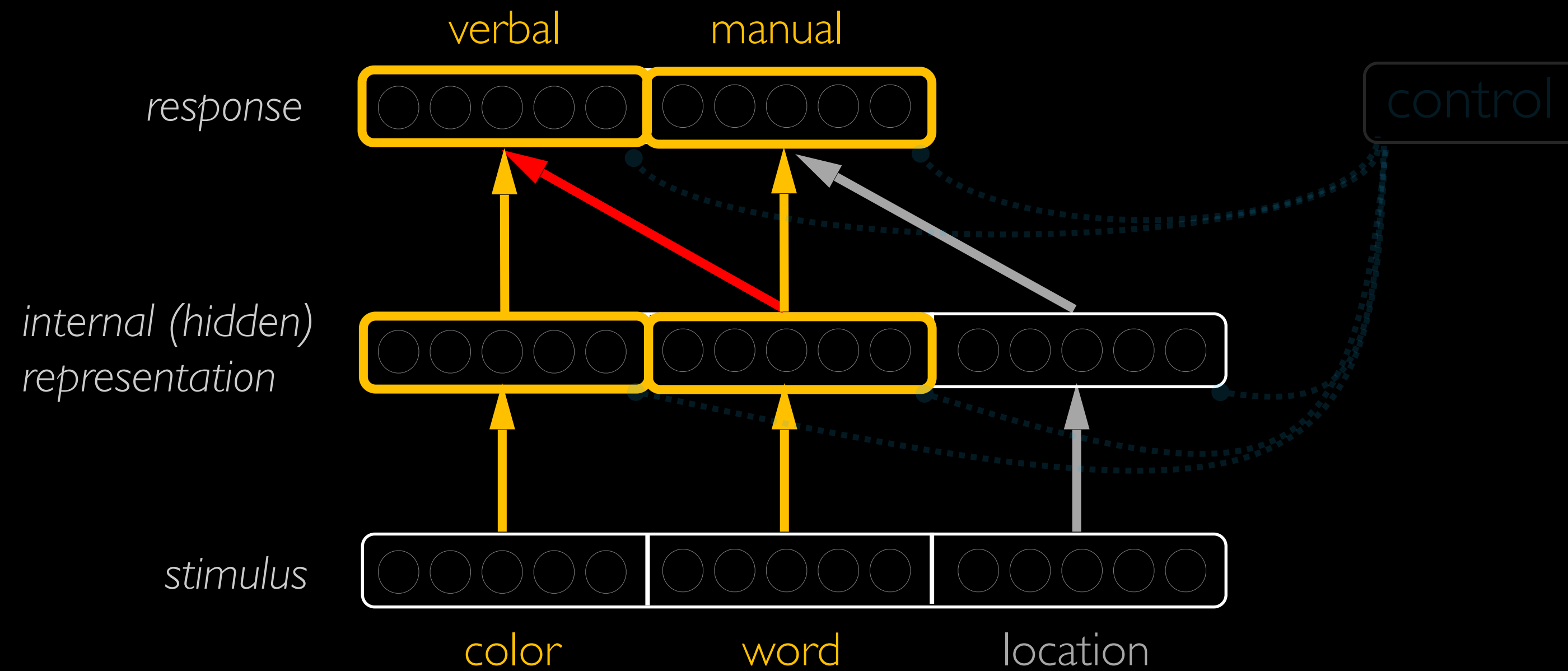
# Neural Network Model of Controlled Processing



## Multiple-Resource Theory

(Allport, 1972; Allport, 1980; Meyer & Kieras, 1997; Navon & Gopher, 1979; Wickens, 1984; Salvucci & Taatgen, 2008)

# Neural Network Model of Controlled Processing



Harold 1994

# Questions

# Questions

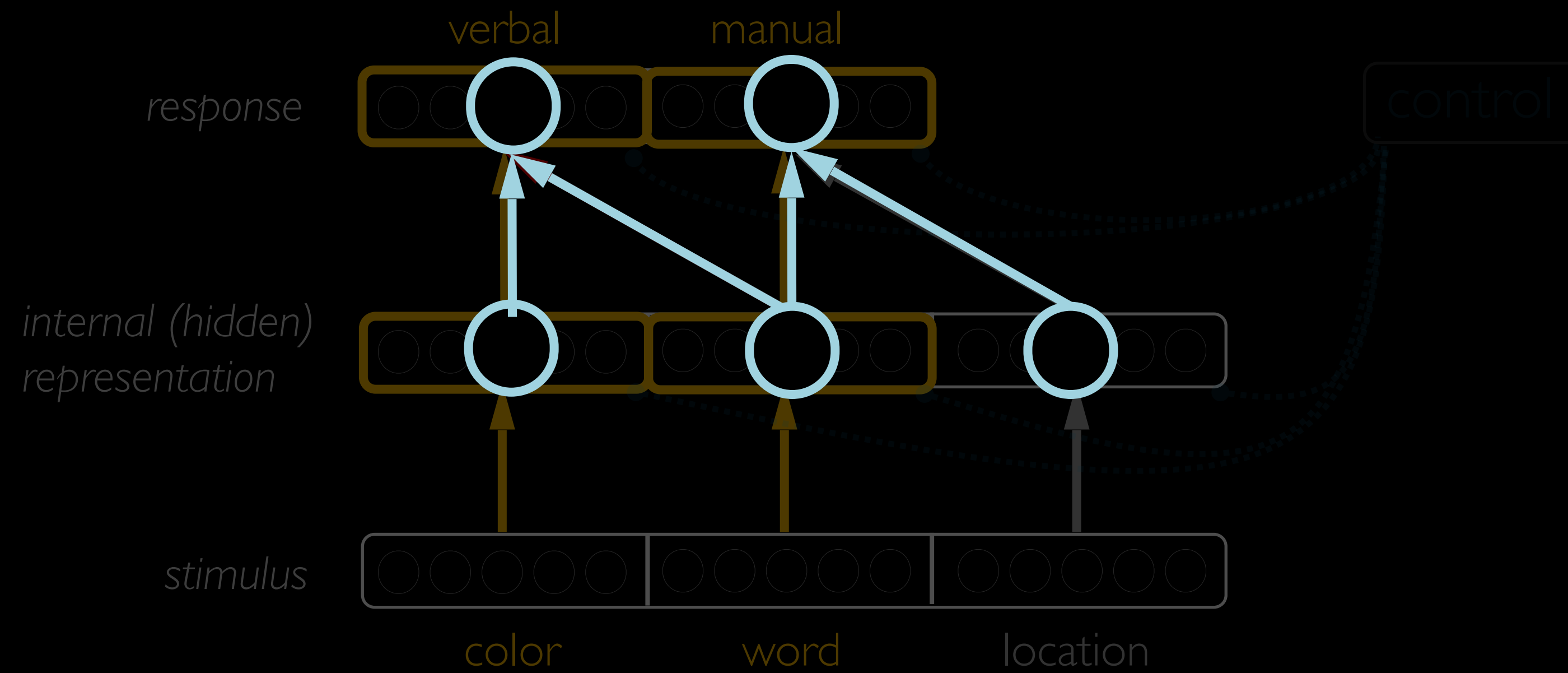
- 1. Do the limitations imposed by shared representations prevail in a system as large as the brain?**
- 2. Assuming that shared representation cause a lot of trouble, why do we use them in the first place? (Experiment is evidence we do)**

# Questions

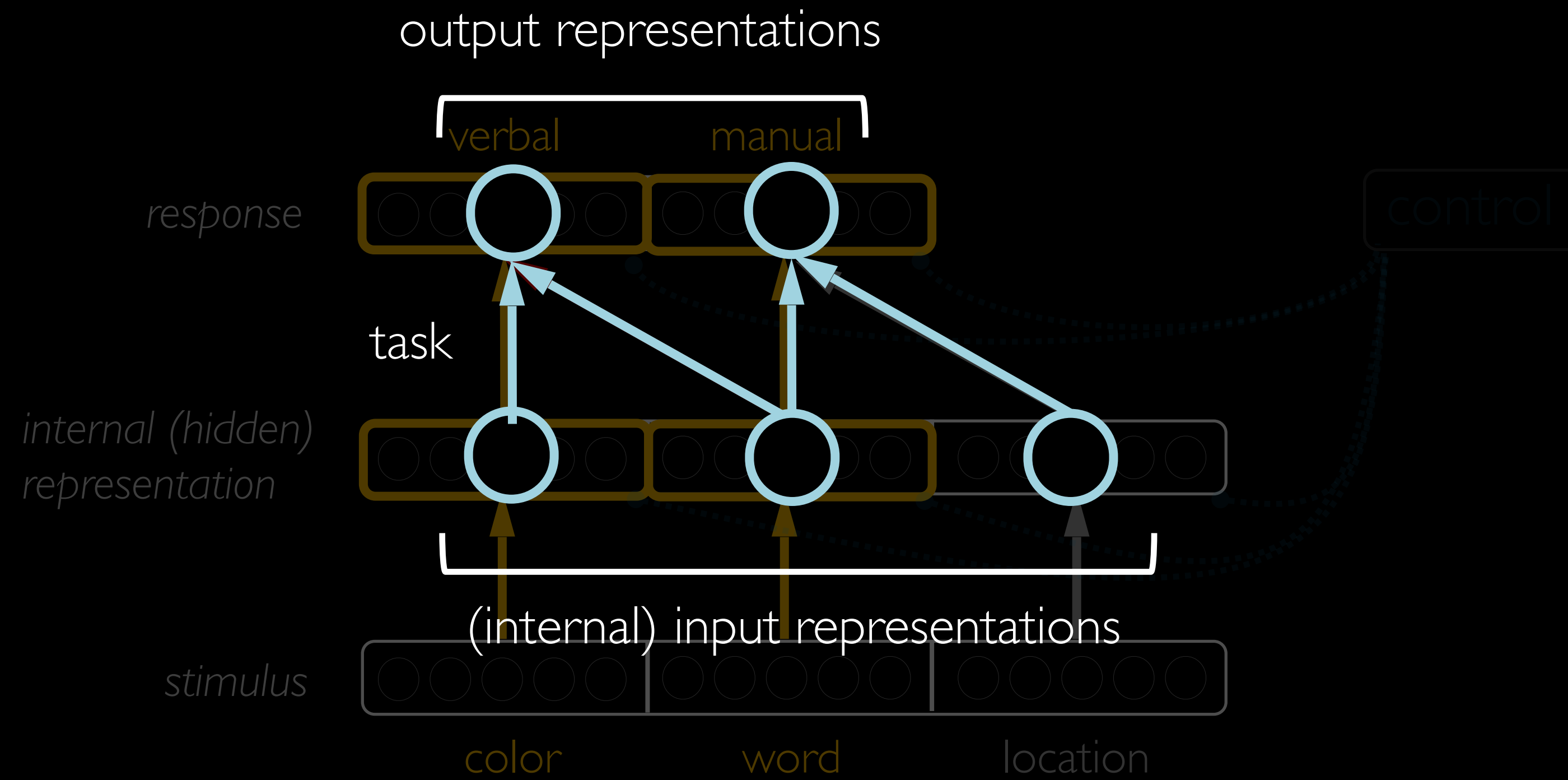
**1. Do the limitations imposed by shared representations prevail in a system as large as the brain?**

**2. Assuming that shared representation cause a lot of trouble, why do we use them in the first place? (Experiment is evidence we do)**

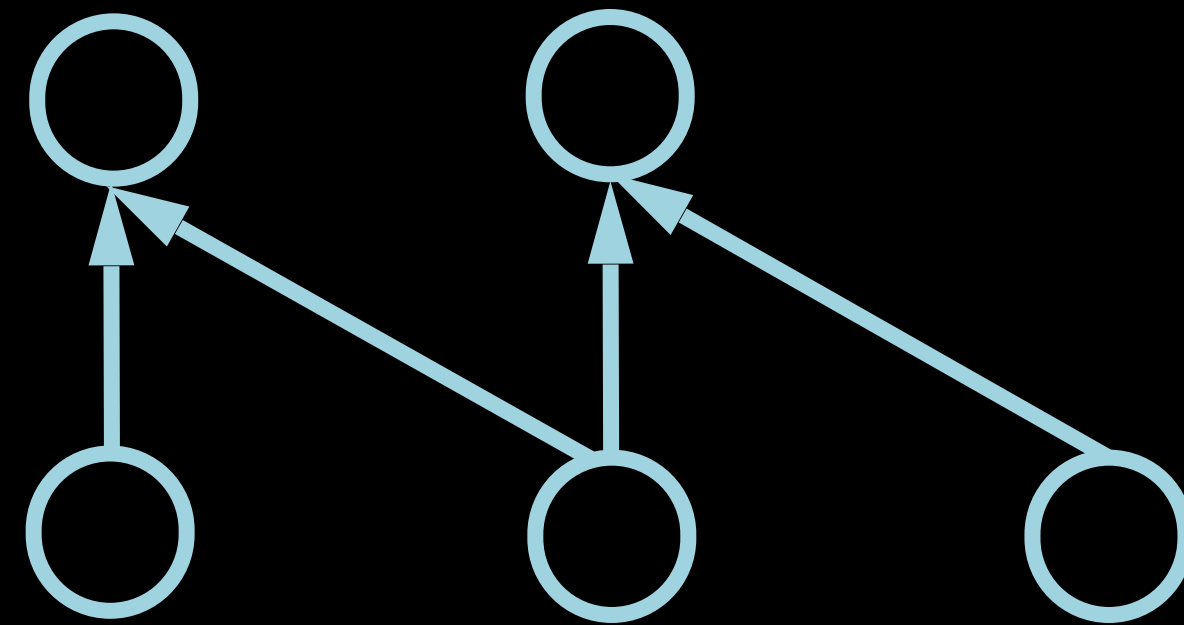
# Neural Network Model of Controlled Processing



# Neural Network Model of Controlled Processing

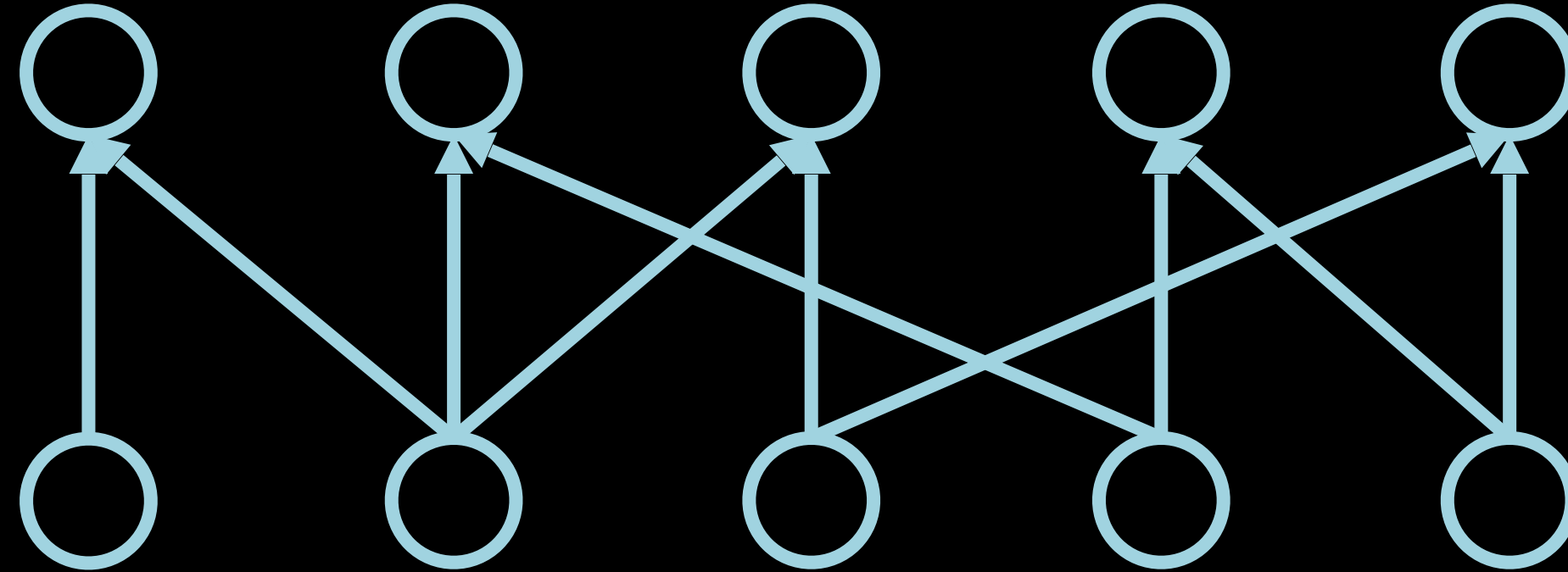


# Graph-Theoretic Analysis of Multitasking Capacity



What is the **maximum number of tasks** that the network can perform **in parallel without interference?**

# Graph-Theoretic Analysis of Multitasking Capacity

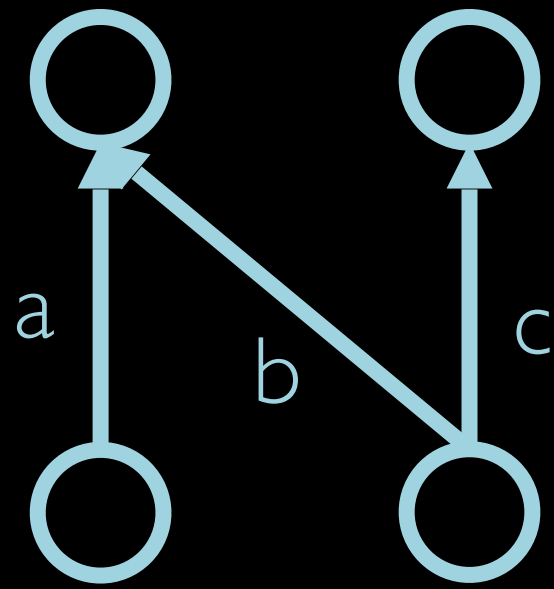


What is the **maximum number of tasks** that the network can perform **in parallel without interference**?

# Graph-Theoretic Approach

## Task Dependencies

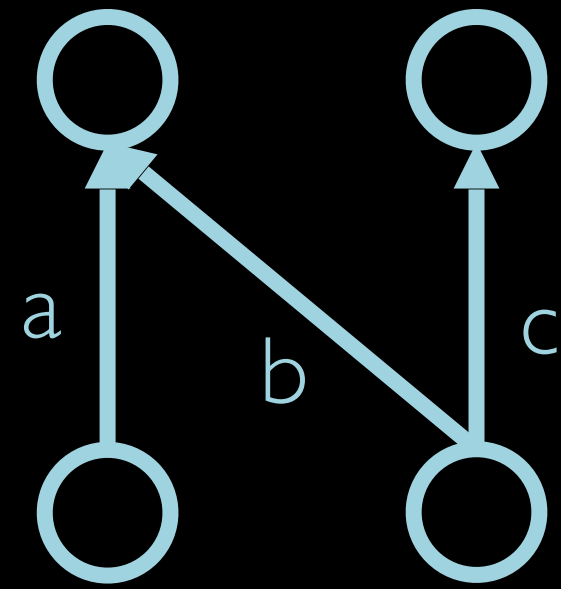
bipartite  
task graph



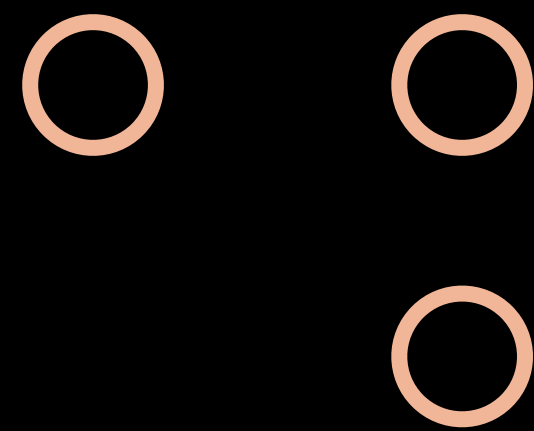
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph



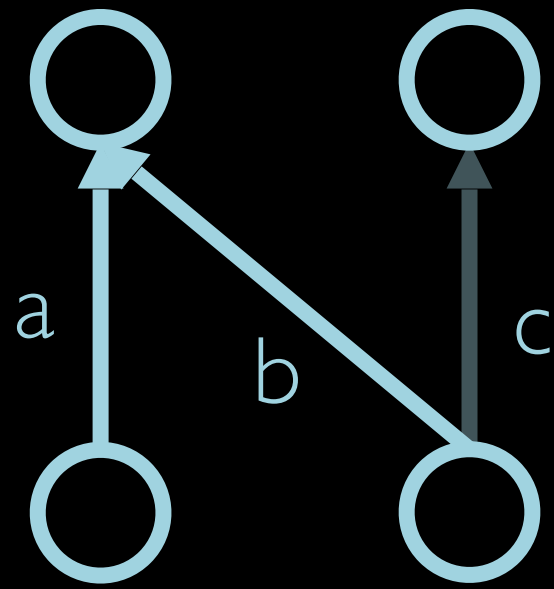
dependency  
graph



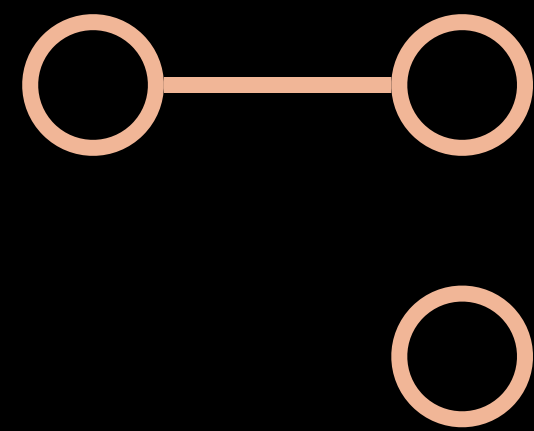
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph



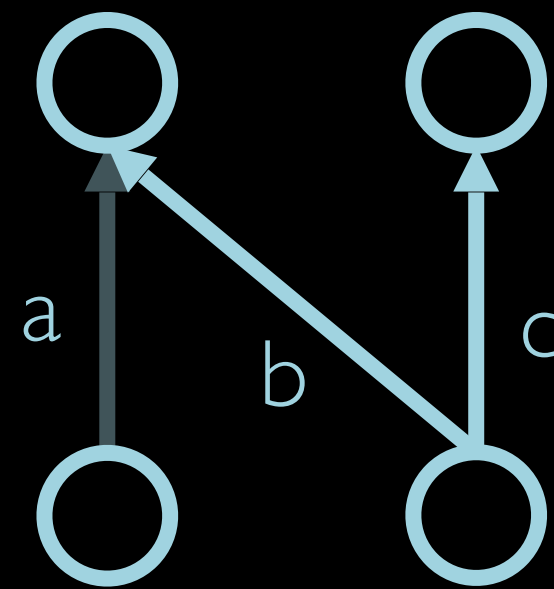
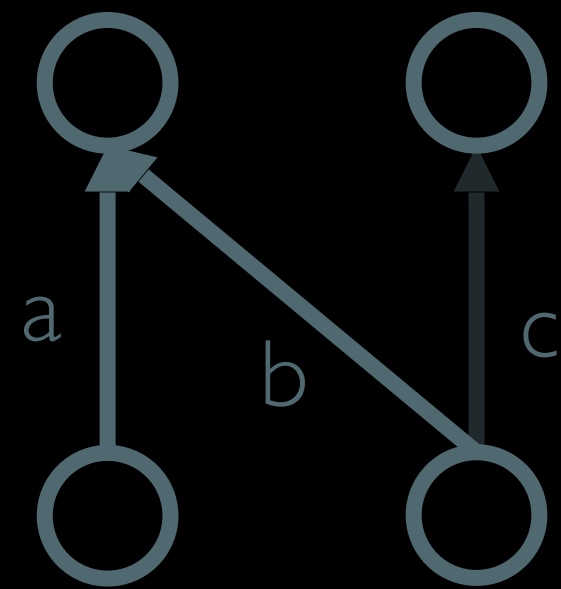
dependency  
graph



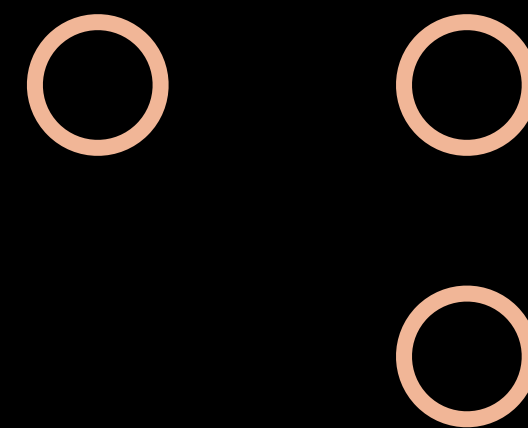
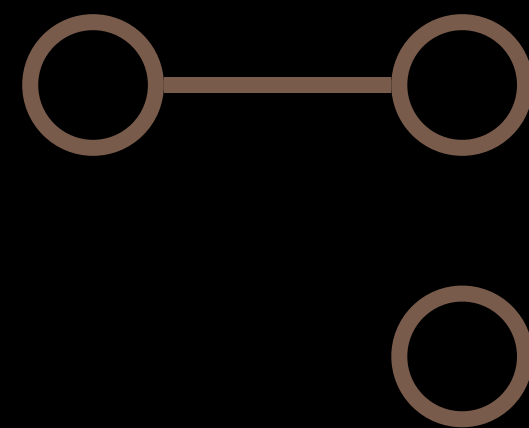
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph



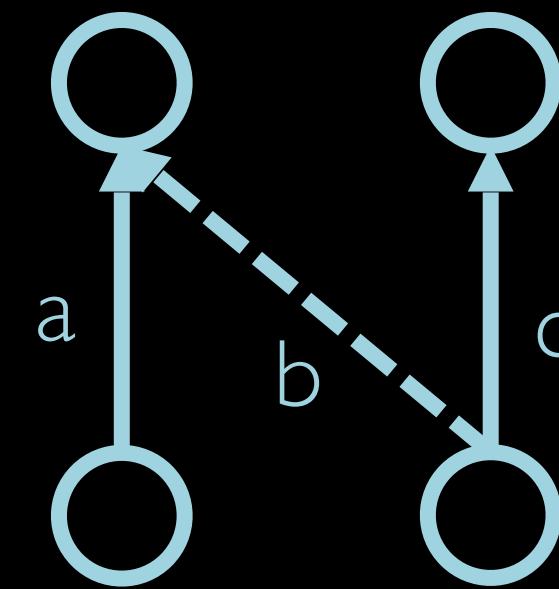
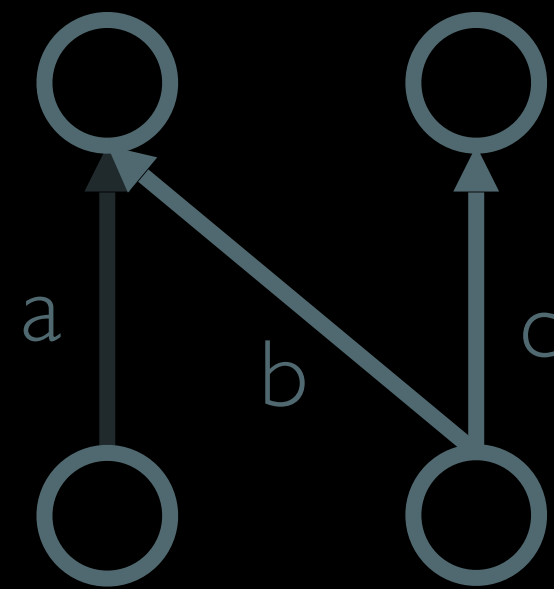
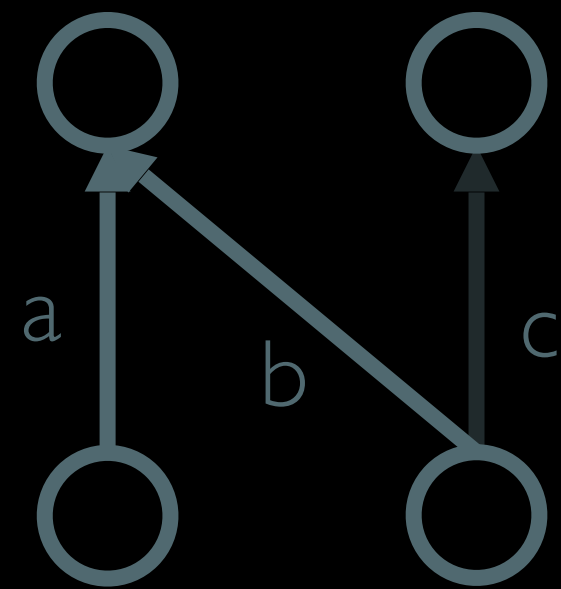
dependency  
graph



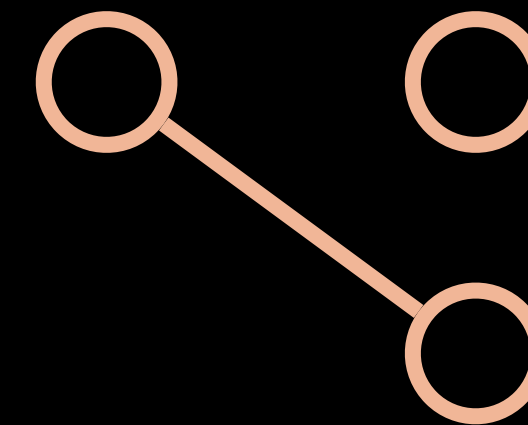
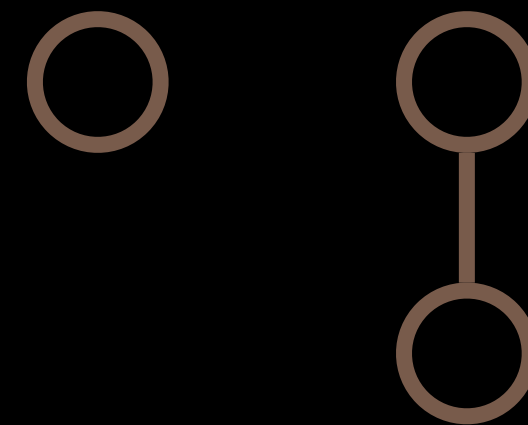
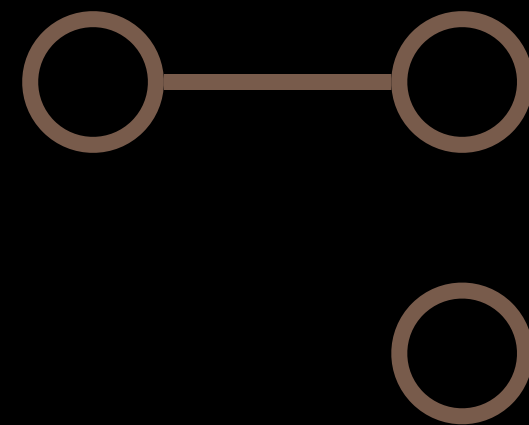
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph



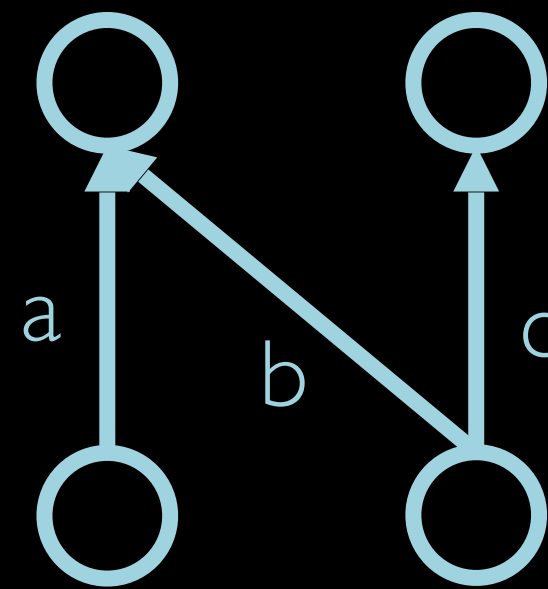
dependency  
graph



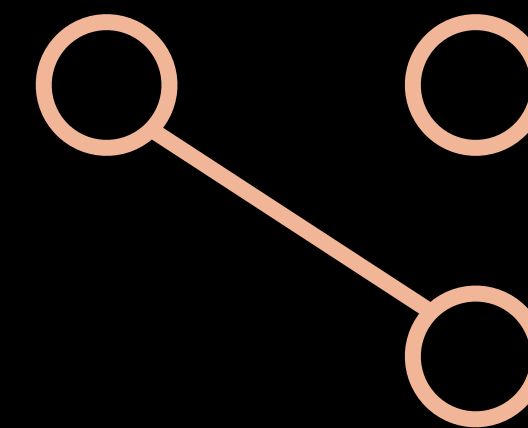
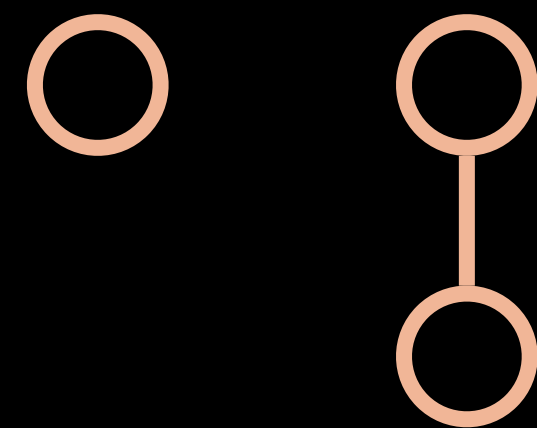
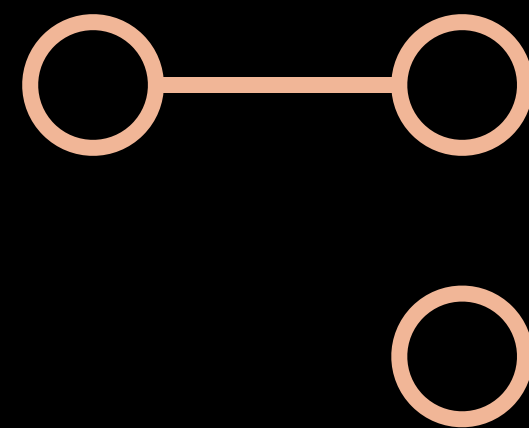
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph



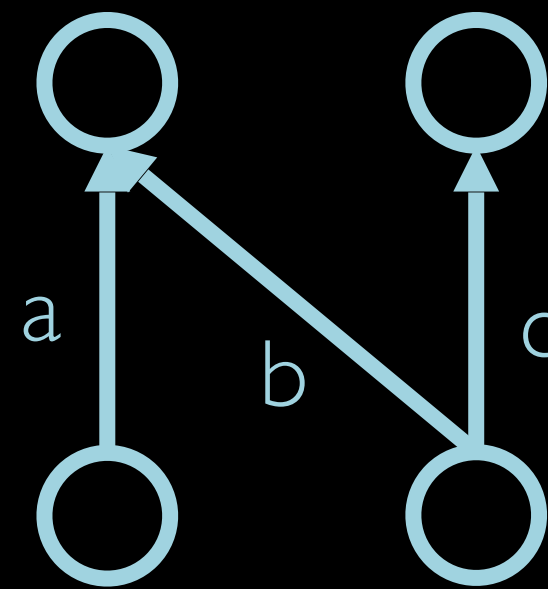
dependency  
graph



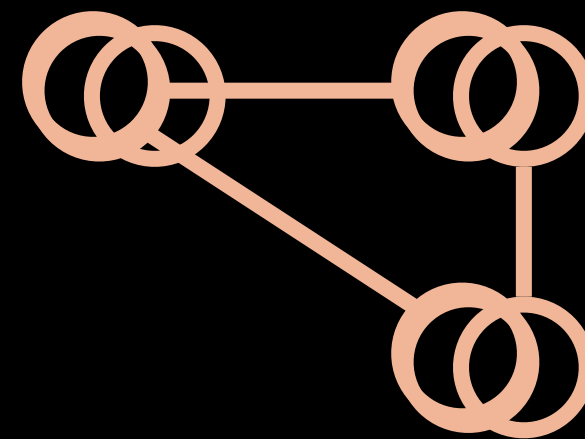
# Graph-Theoretic Approach

## Task Dependencies

bipartite  
task graph

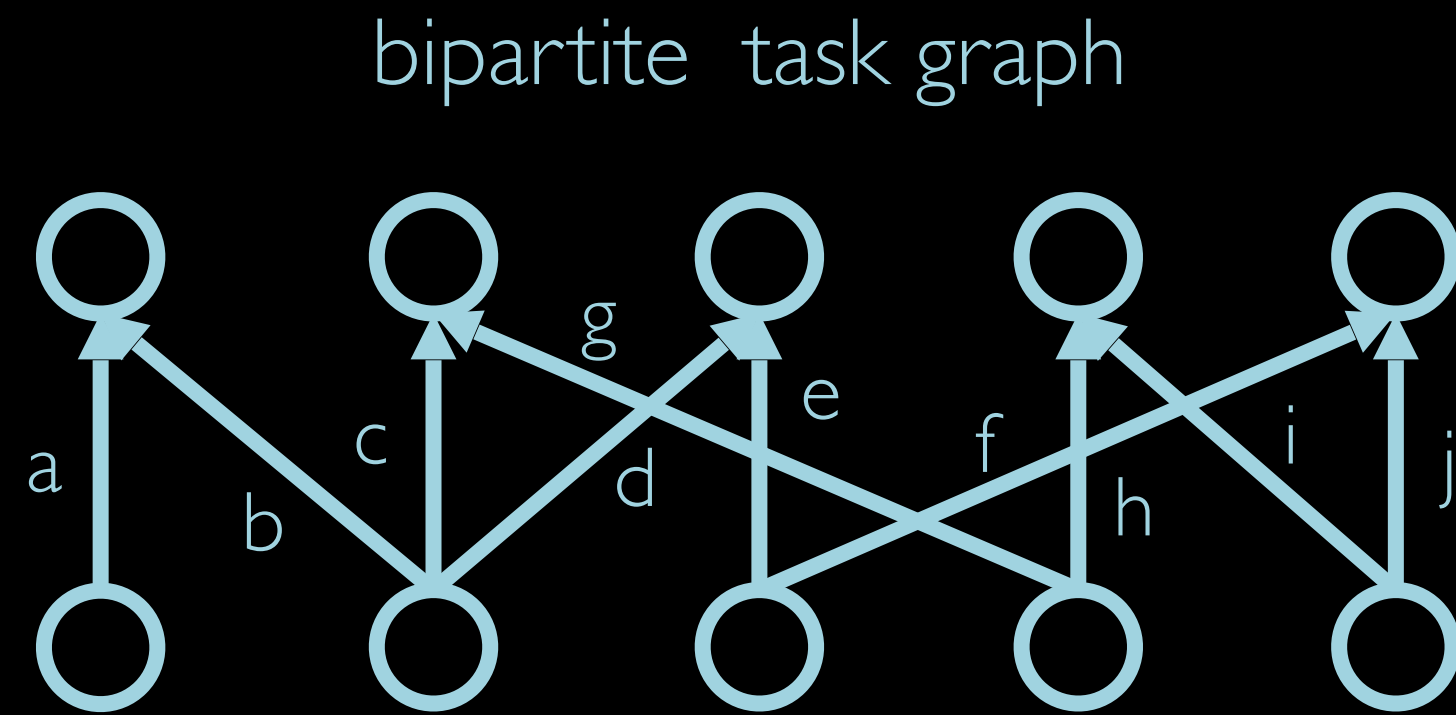


dependency  
graph

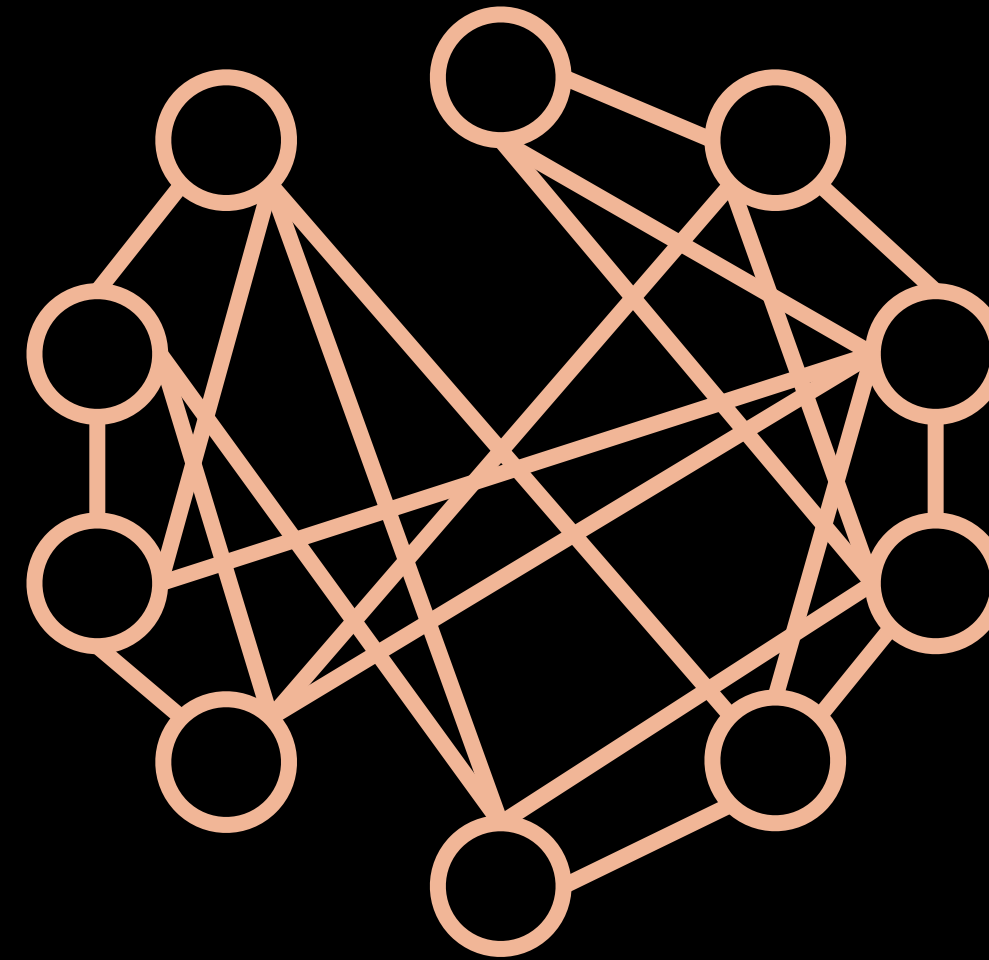


# Graph-Theoretic Approach

## Parallel Processing Capability

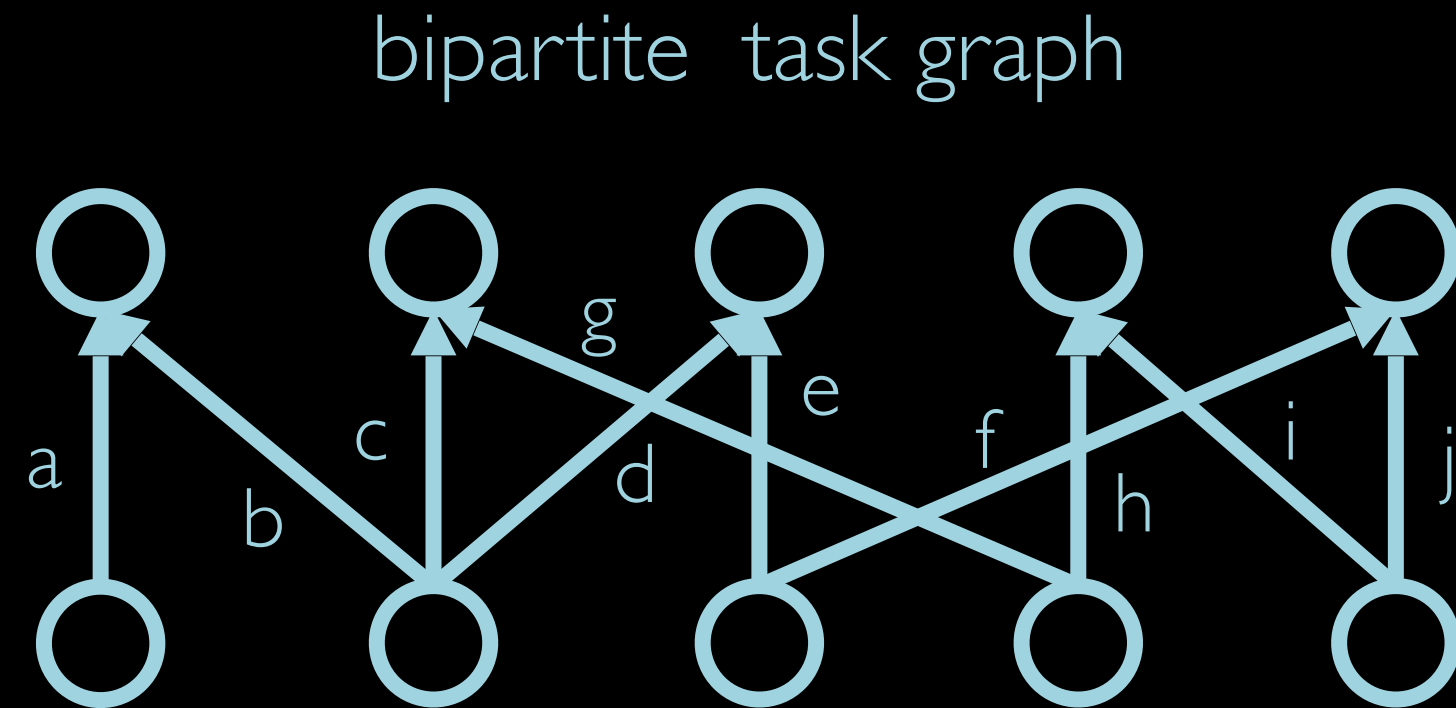


dependency graph

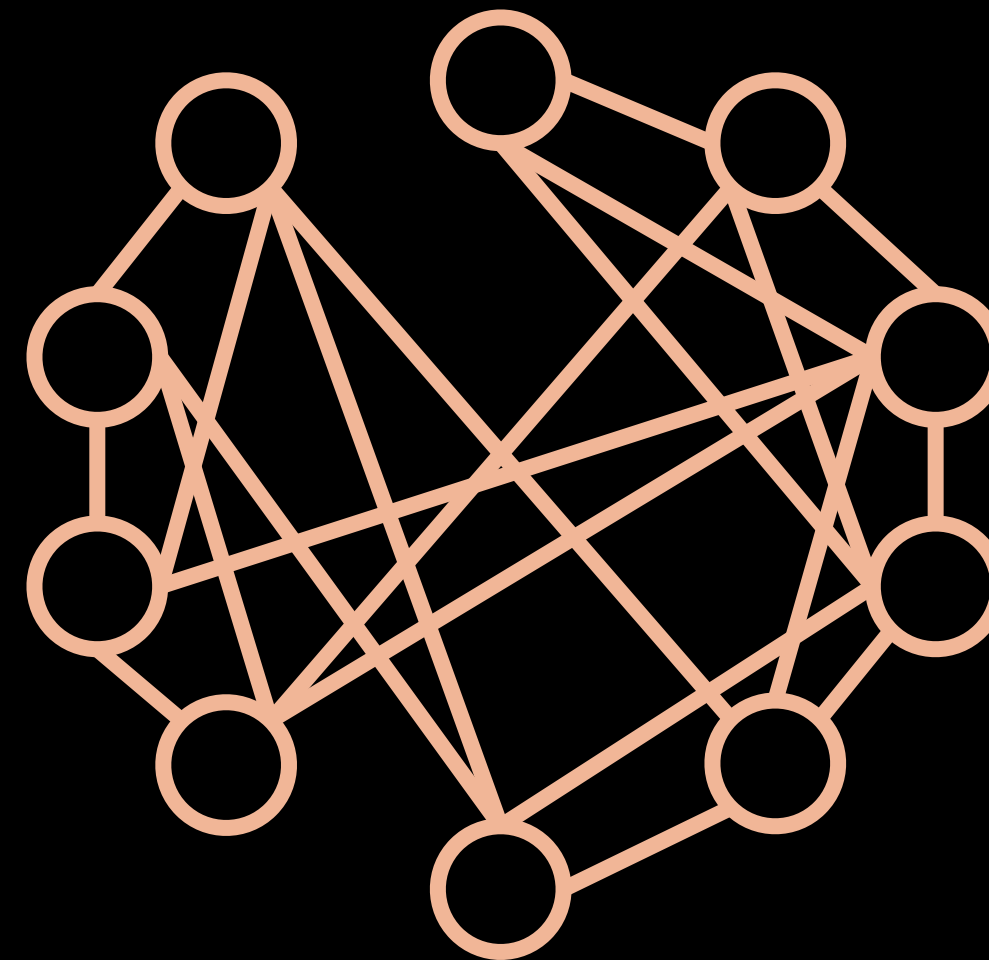


# Graph-Theoretic Approach

## Parallel Processing Capability



dependency graph

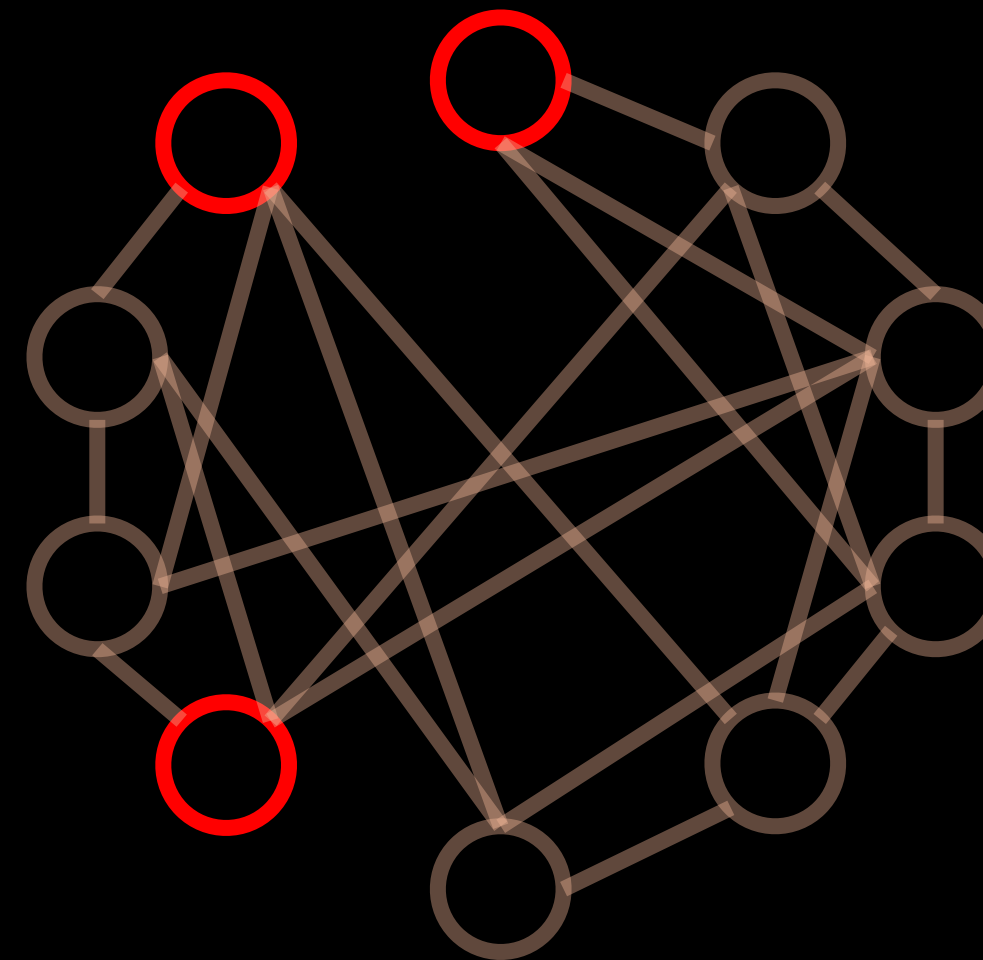


Maximum amount of tasks that can be performed in independently?

# Graph-Theoretic Approach

## Parallel Processing Capability

dependency graph



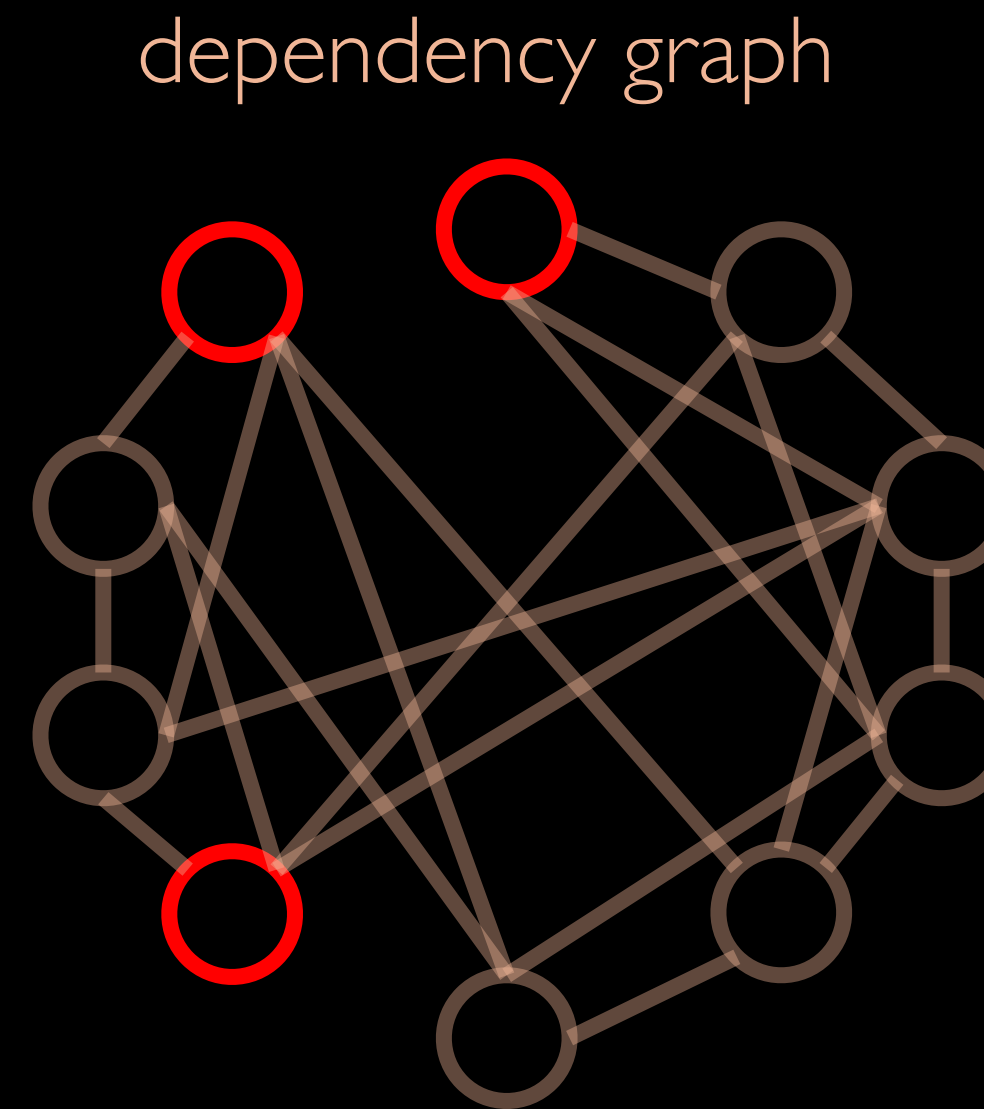
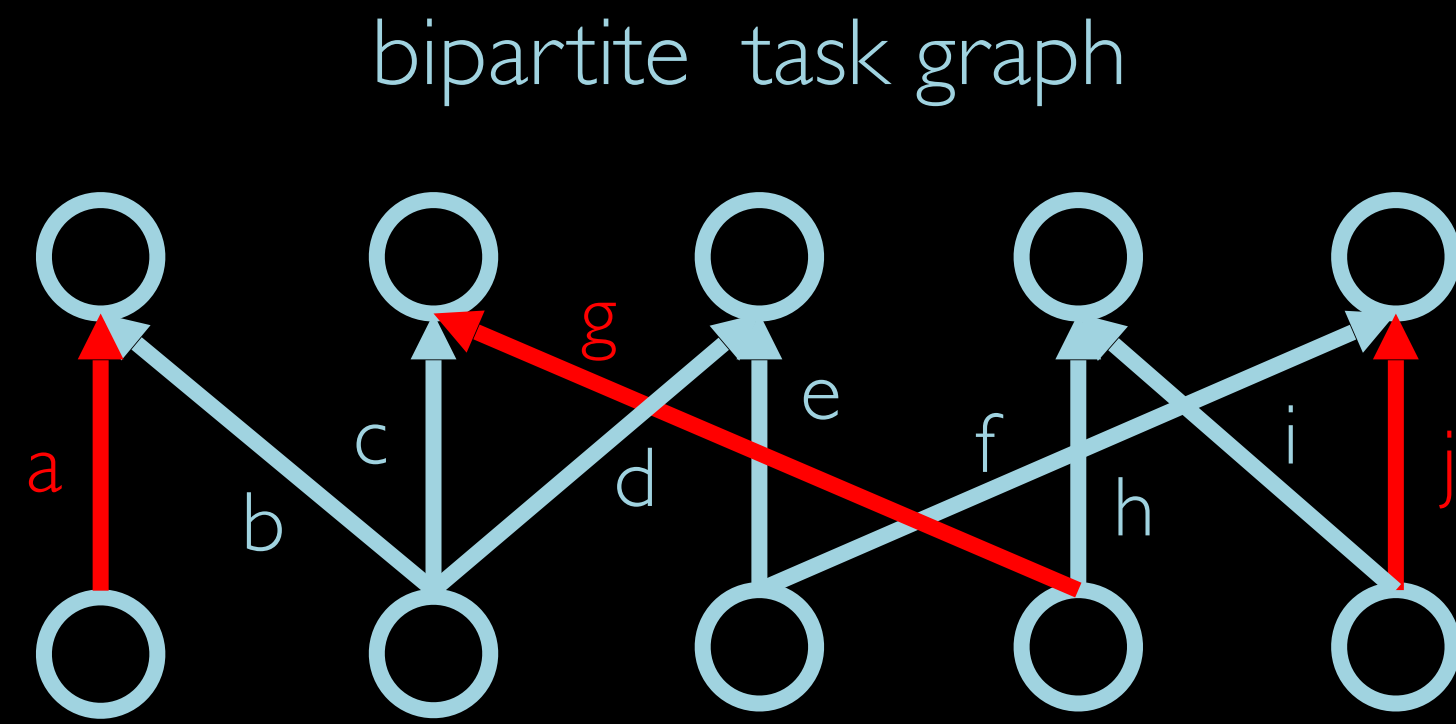
Maximum amount of tasks that can be performed independently?

An **independent vertex set** of a graph  $G$  is a subset of the vertices such that no two vertices in the subset are connected by an edge of  $G$ .

A **maximum independent vertex set** is an independent vertex set containing the largest possible number of vertices for a given graph.

# Graph-Theoretic Approach

## Parallel Processing Capability



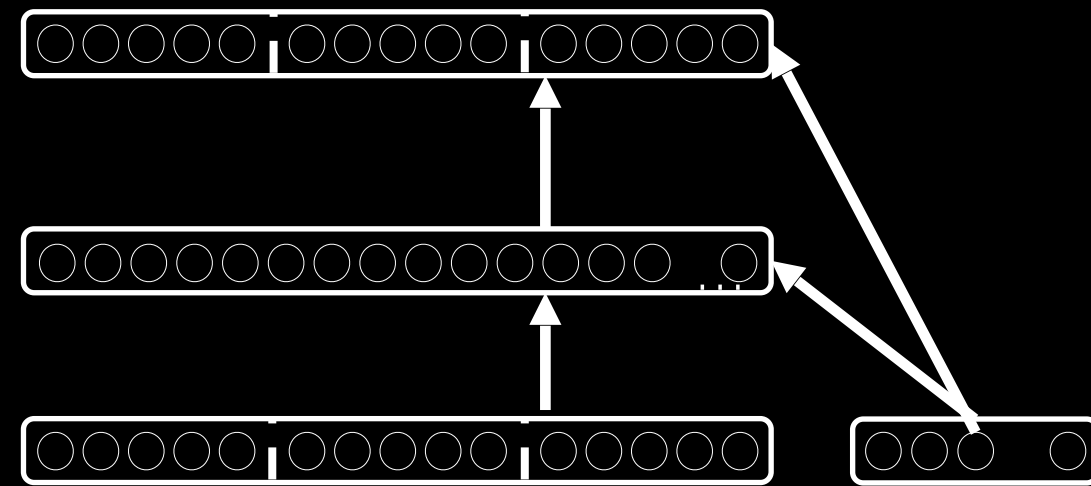
Maximum amount of tasks that can be performed in independently?

An **independent vertex set** of a graph  $G$  is a subset of the vertices such that no two vertices in the subset are connected by an edge of  $G$ .

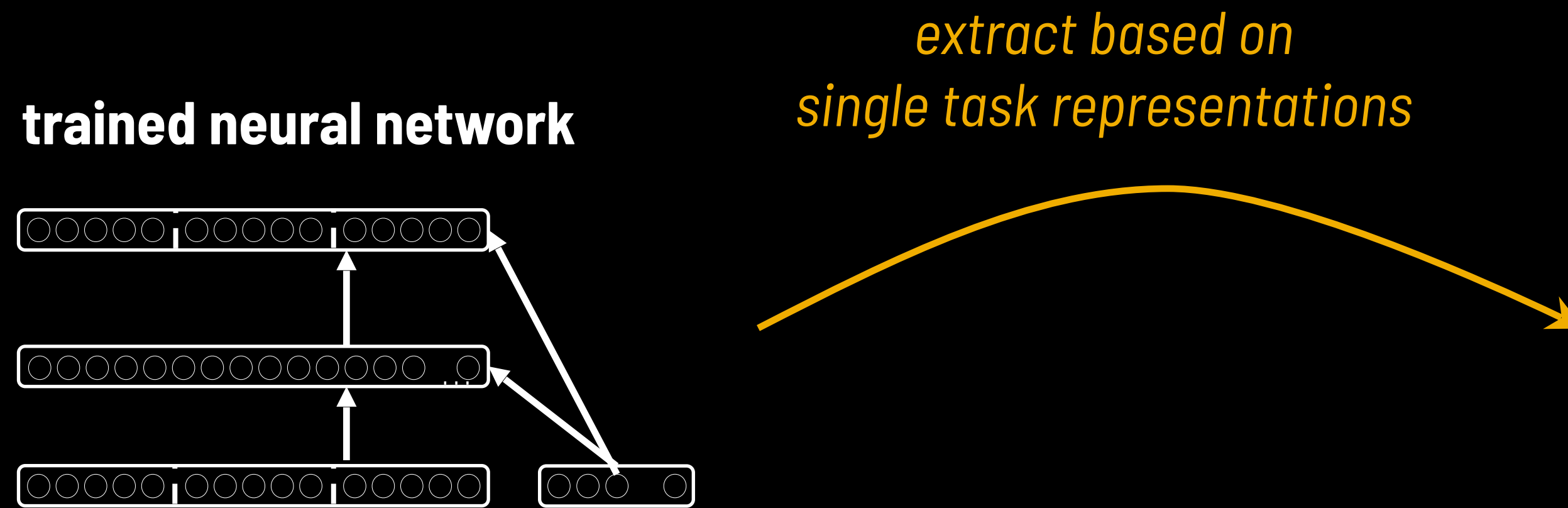
A **maximum independent vertex set** is an independent vertex set containing the largest possible number of vertices for a given graph.

# Application to Neural Systems

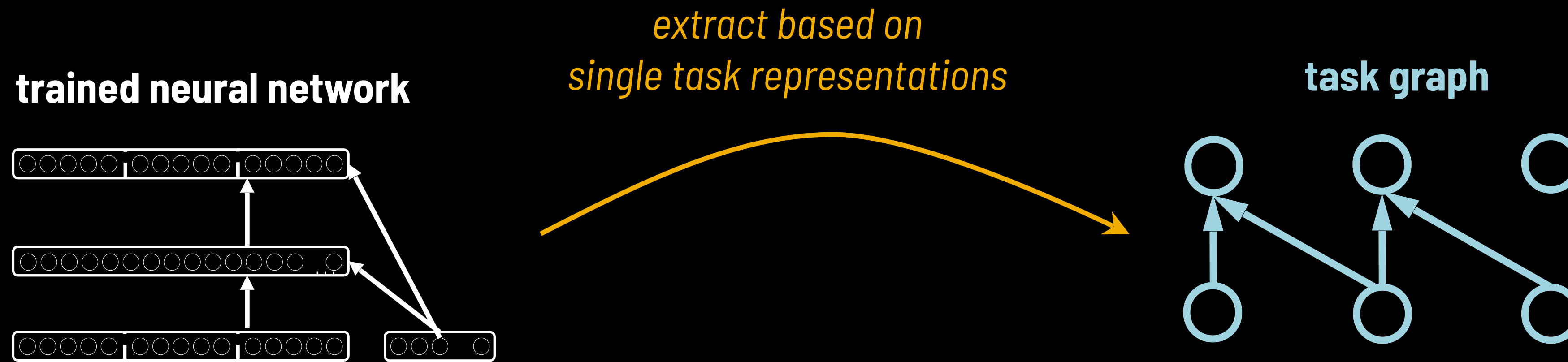
trained neural network



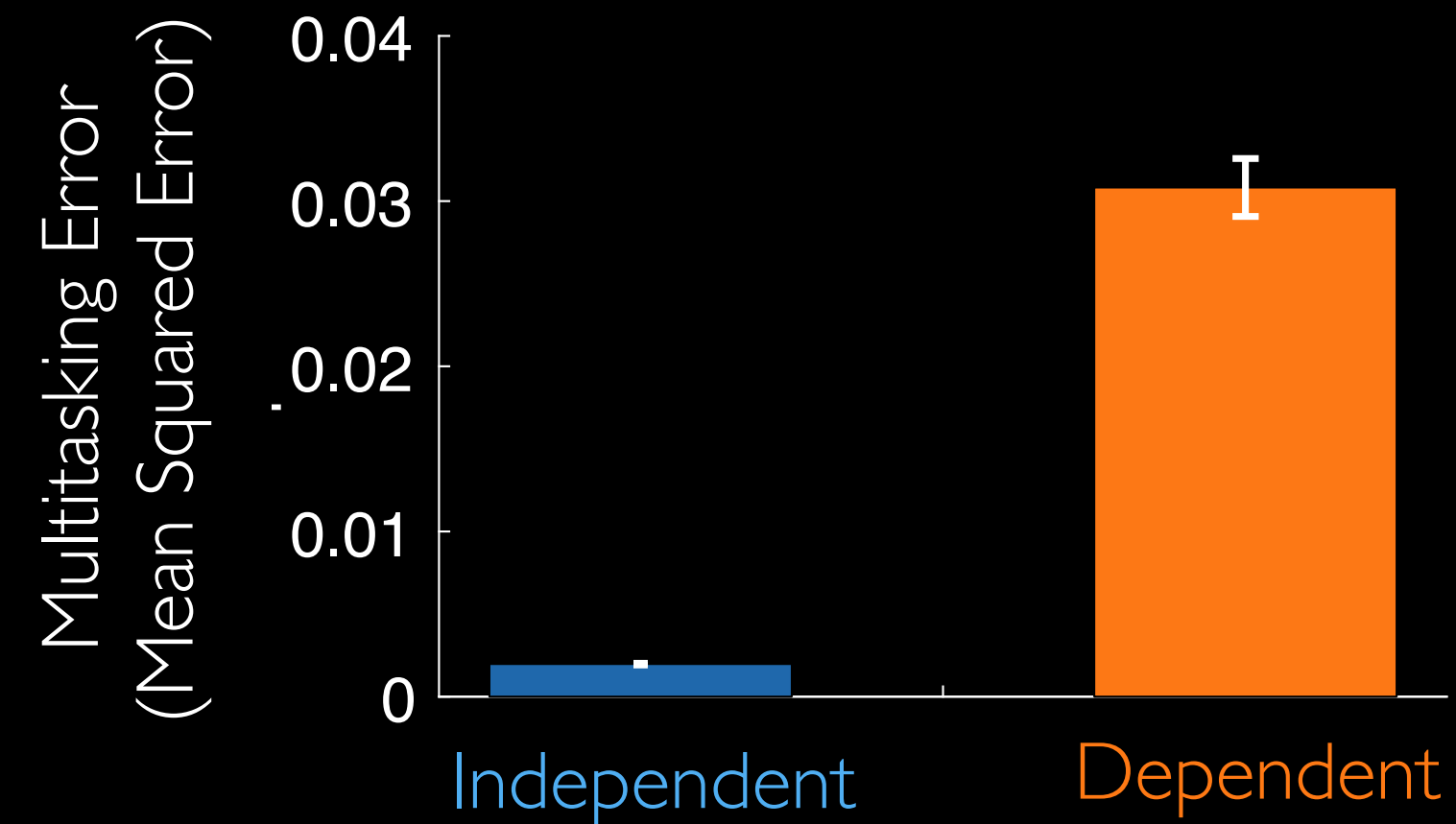
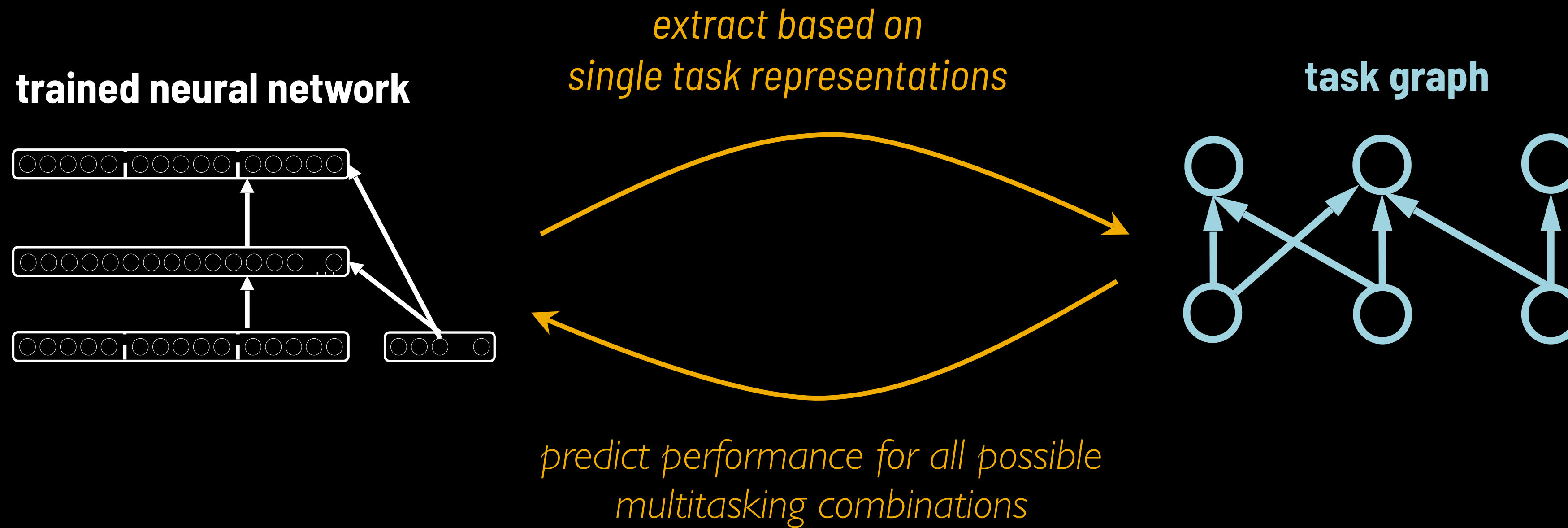
# Application to Neural Systems



# Application to Neural Systems

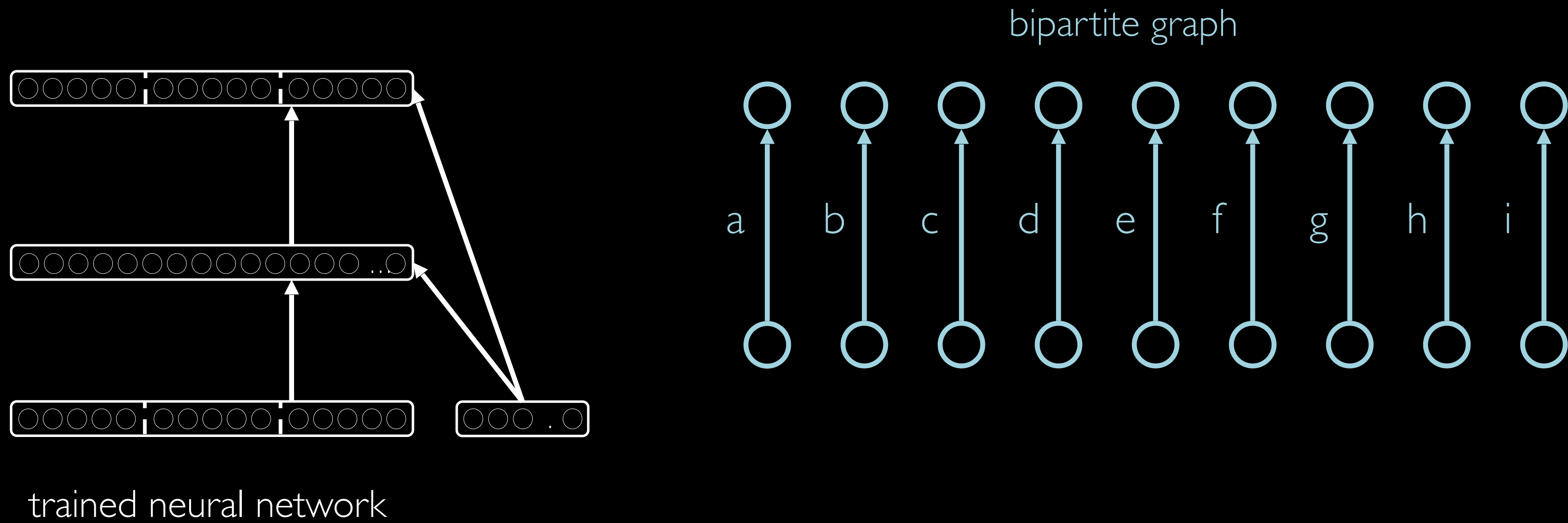


# Application to Neural Systems



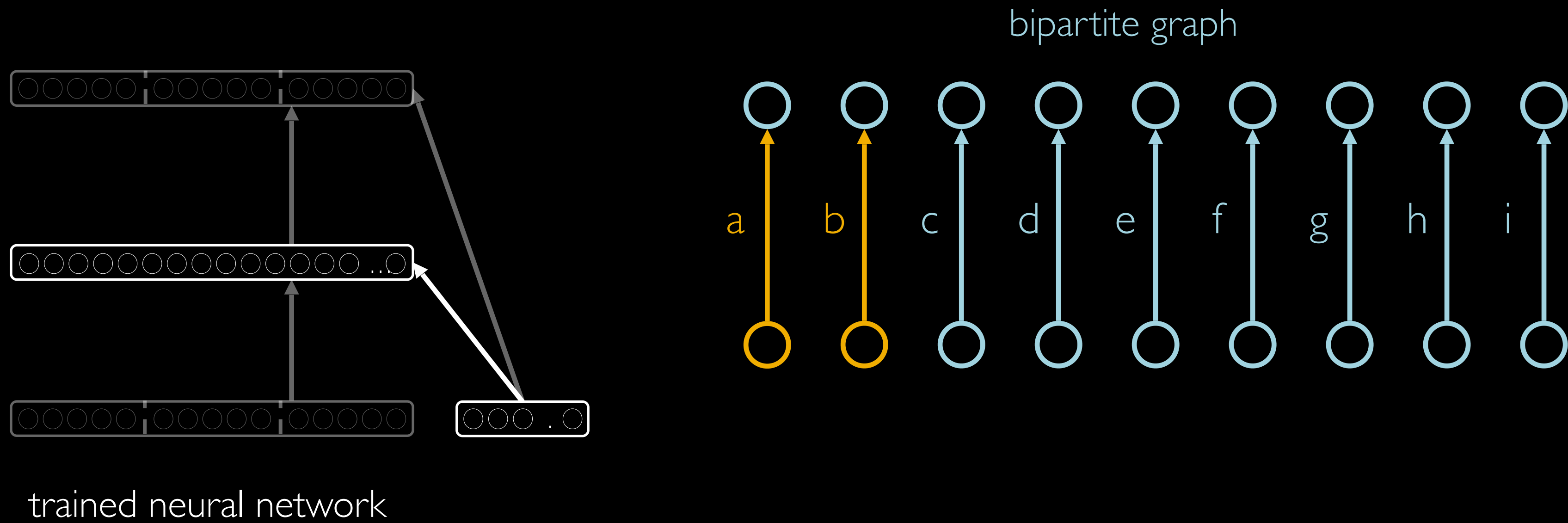
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



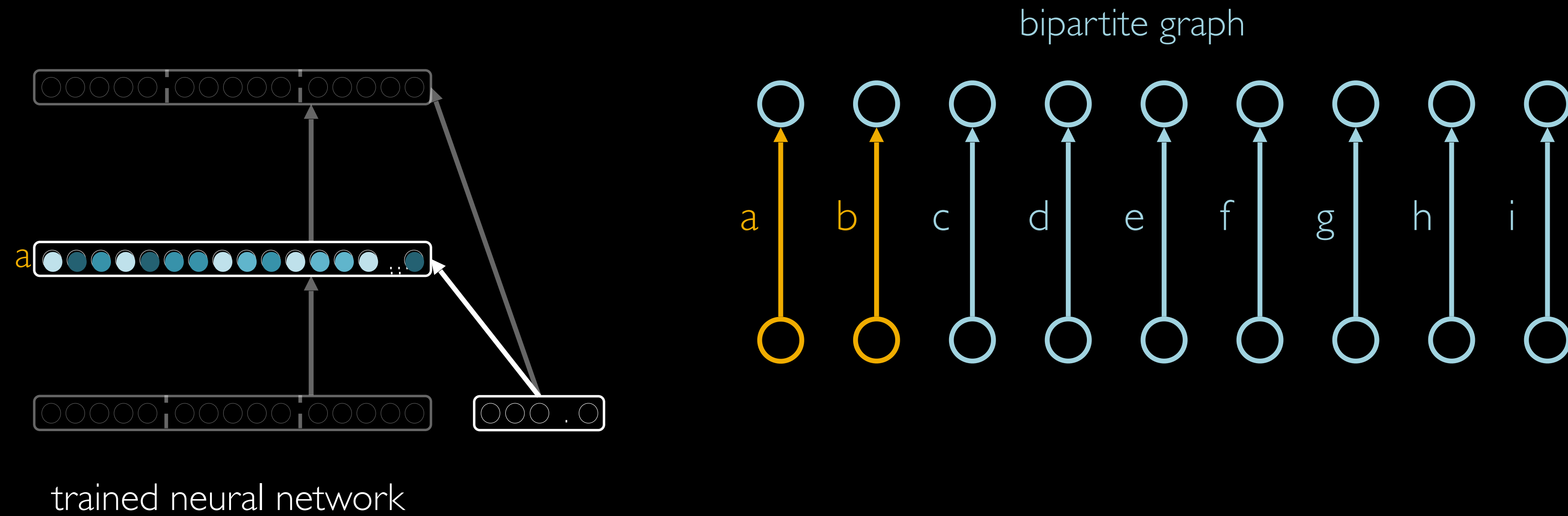
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



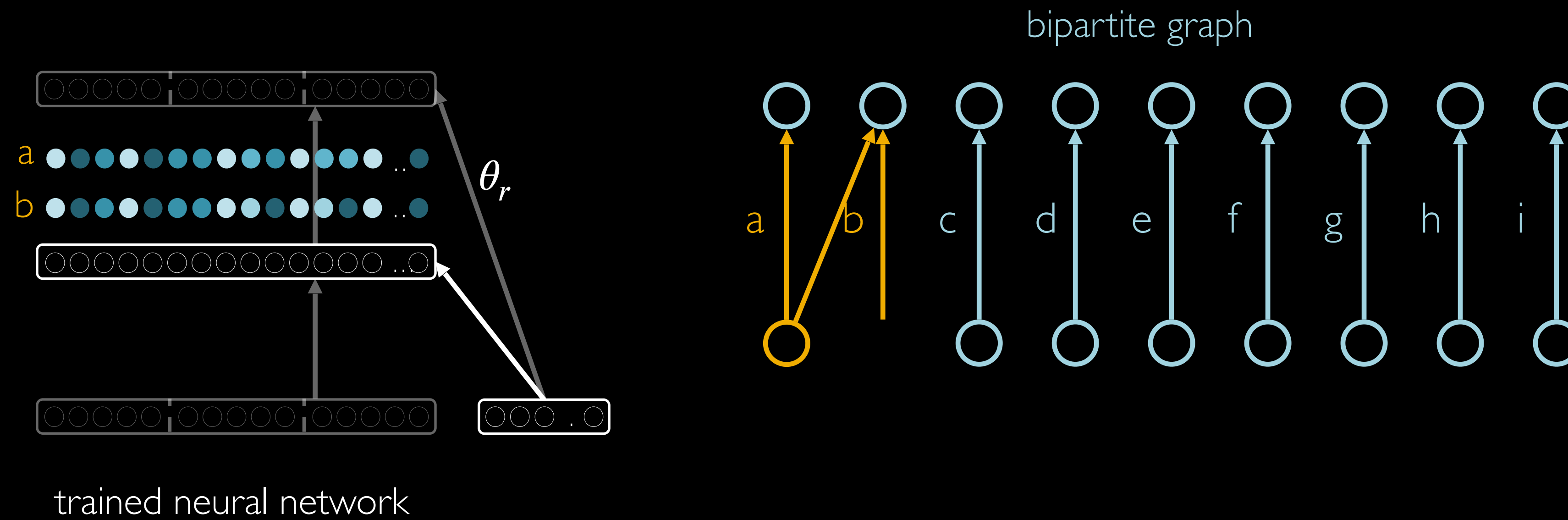
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



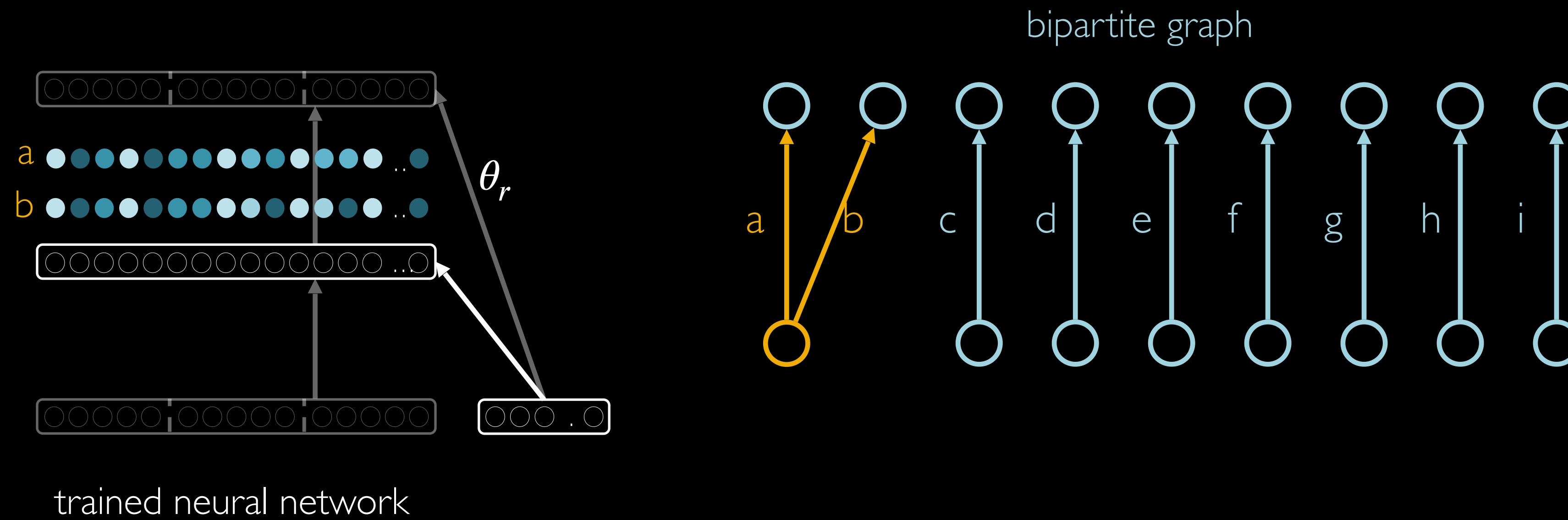
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



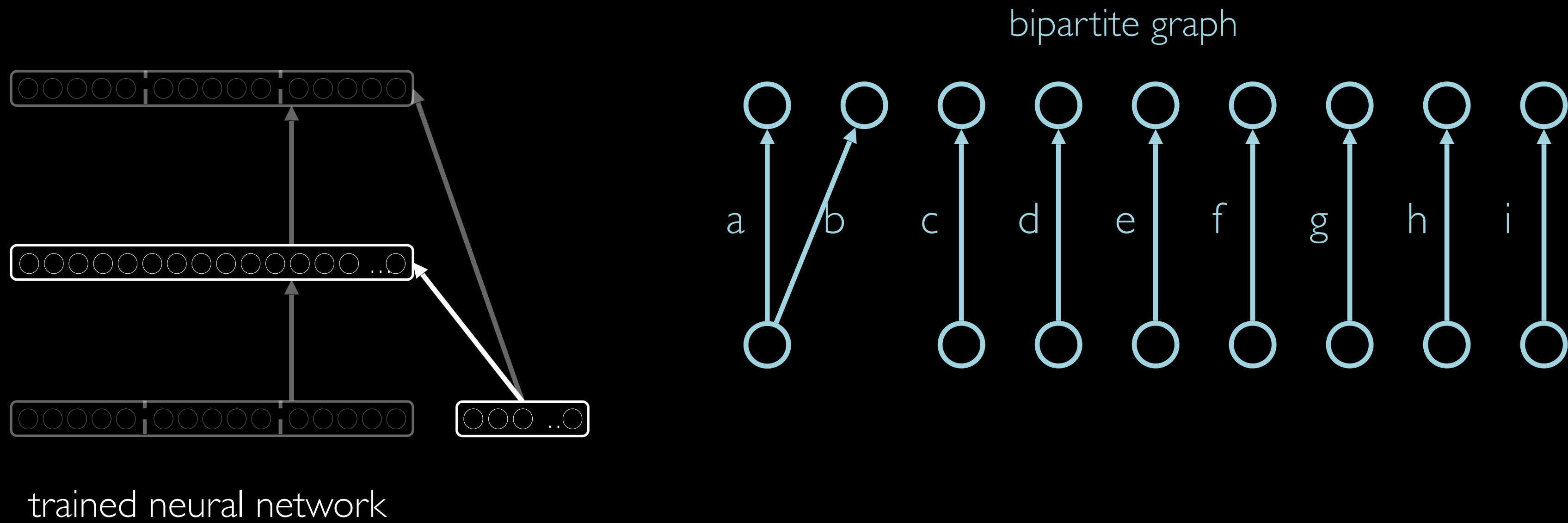
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



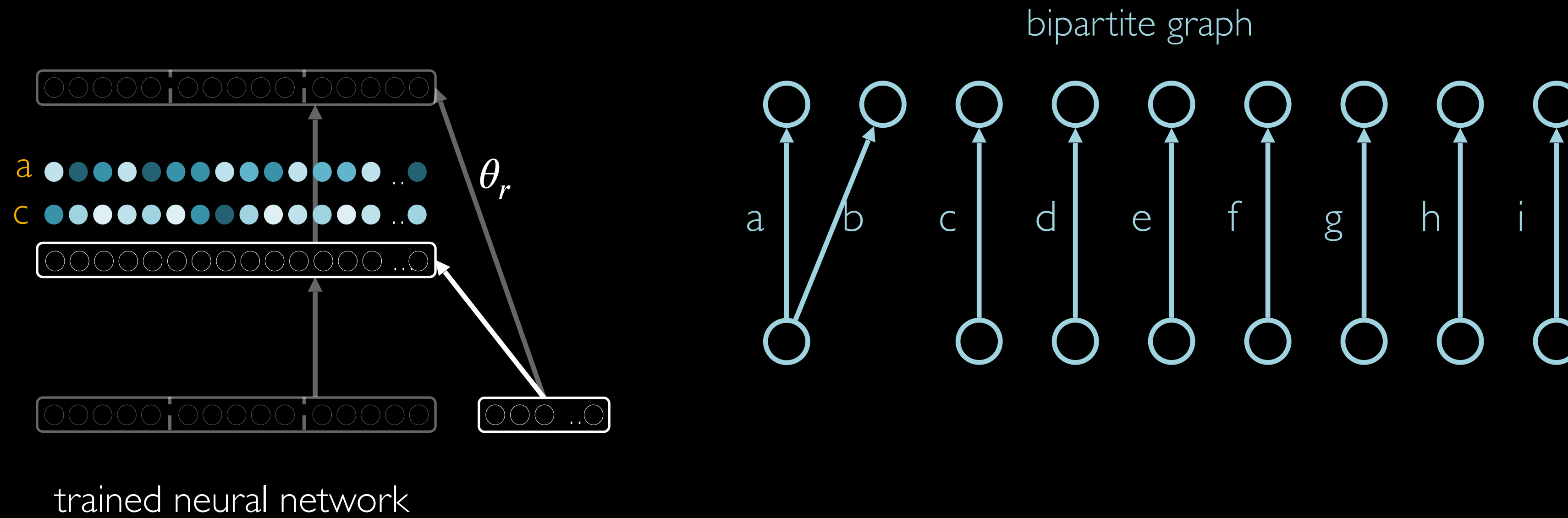
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



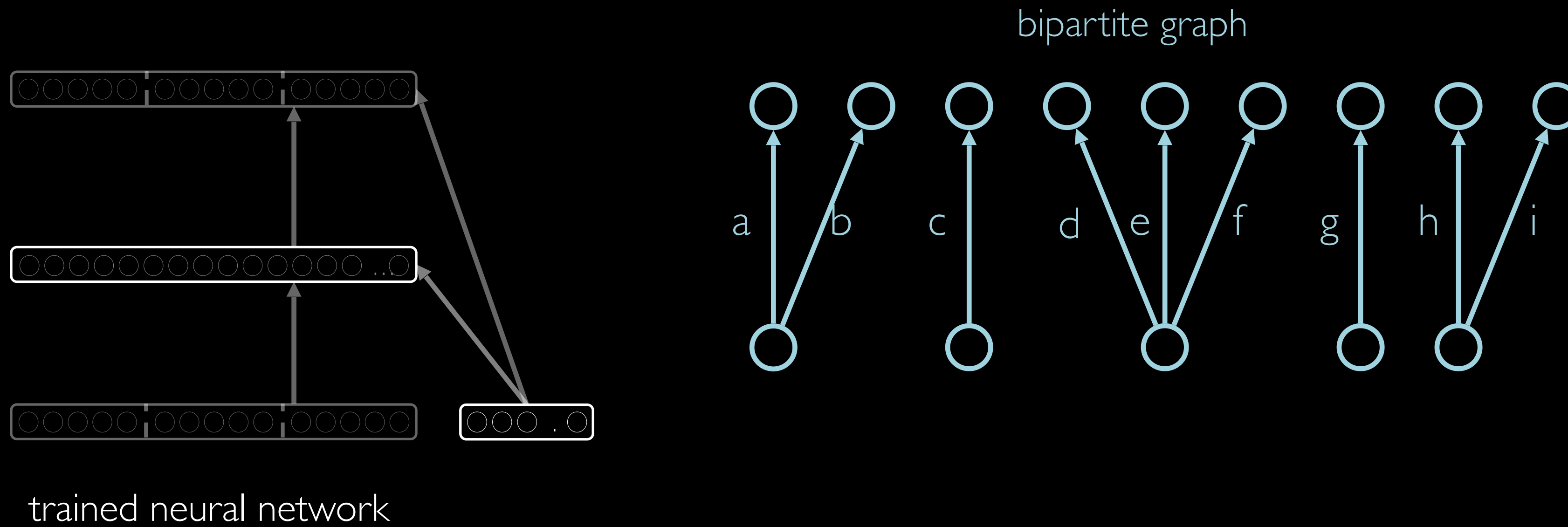
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



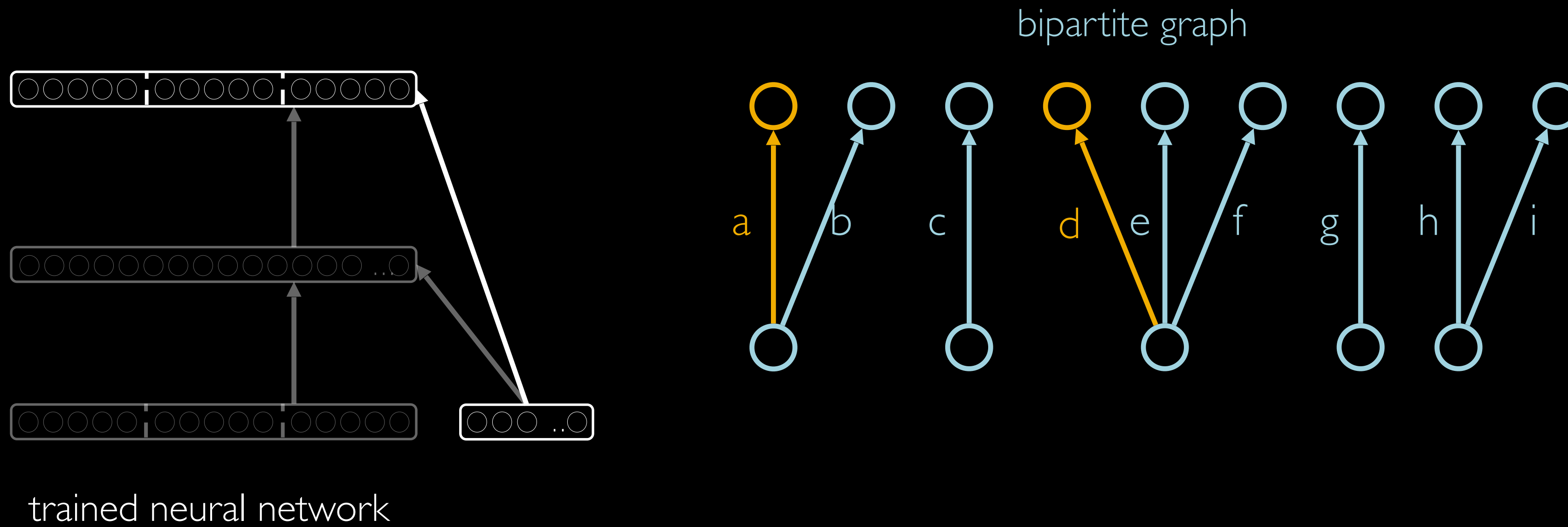
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



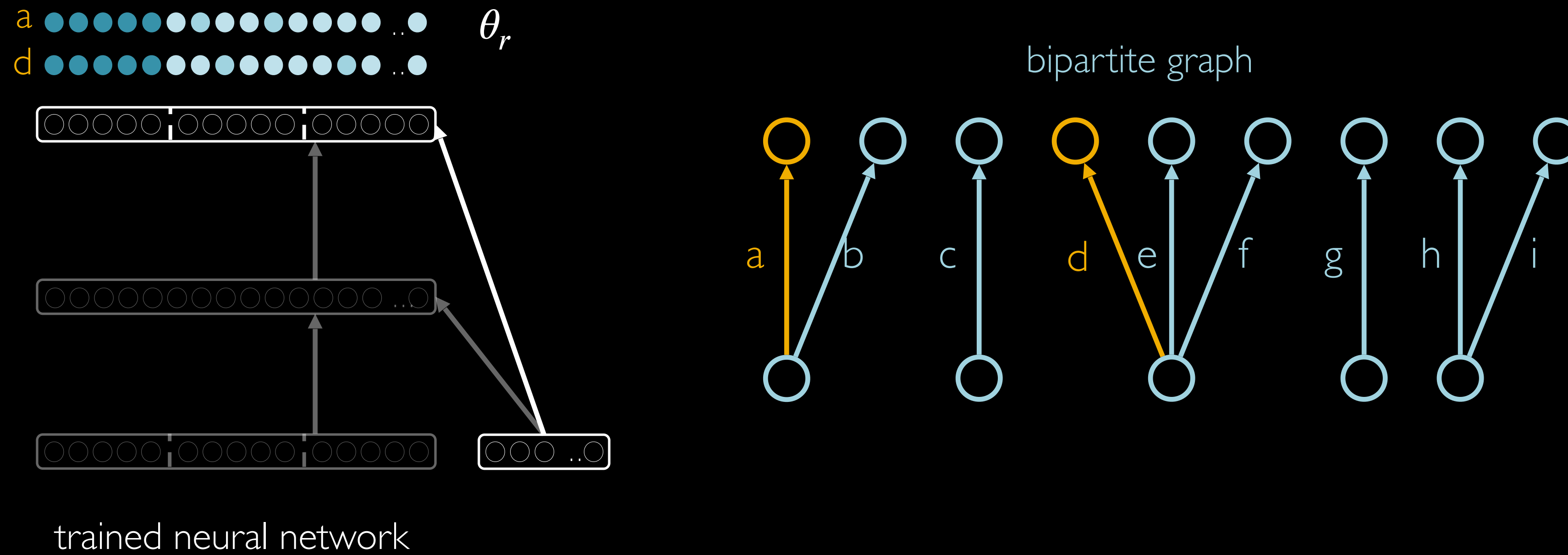
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



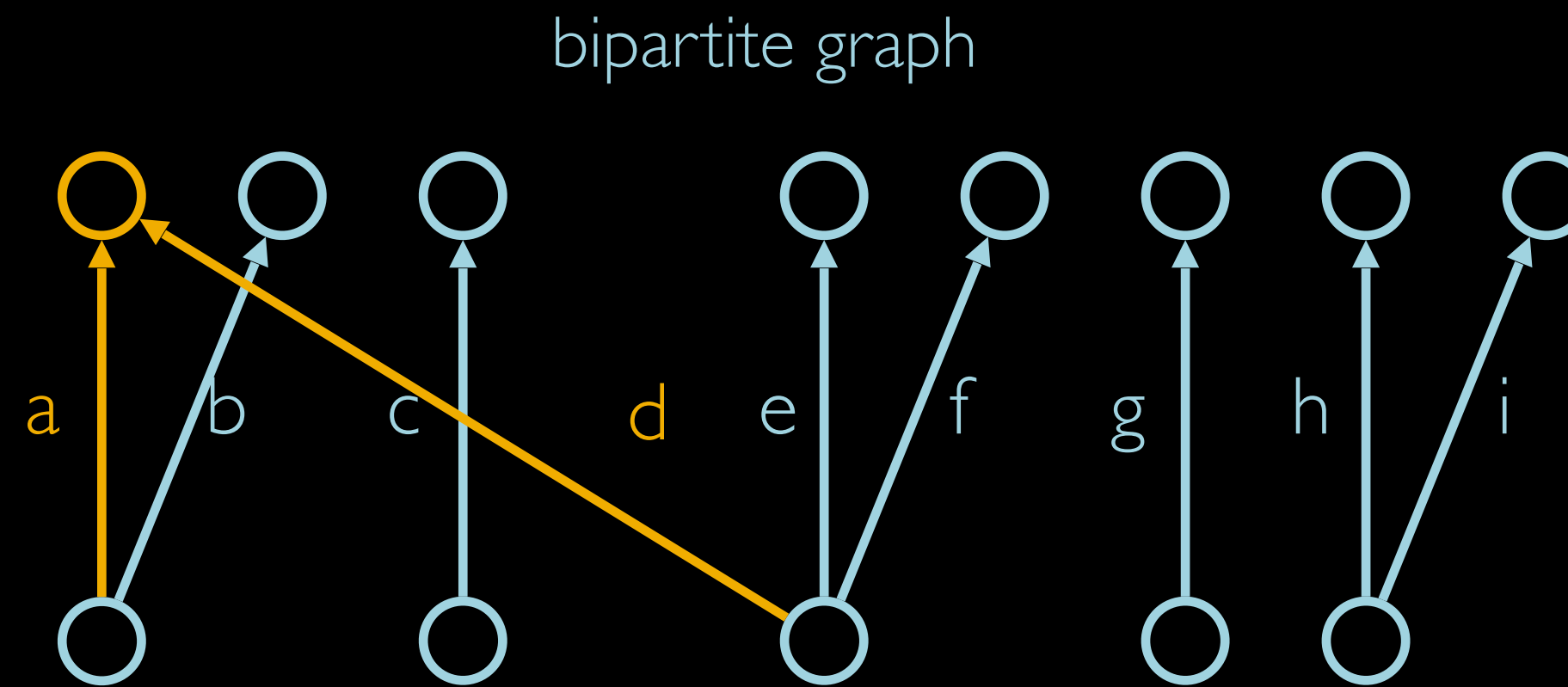
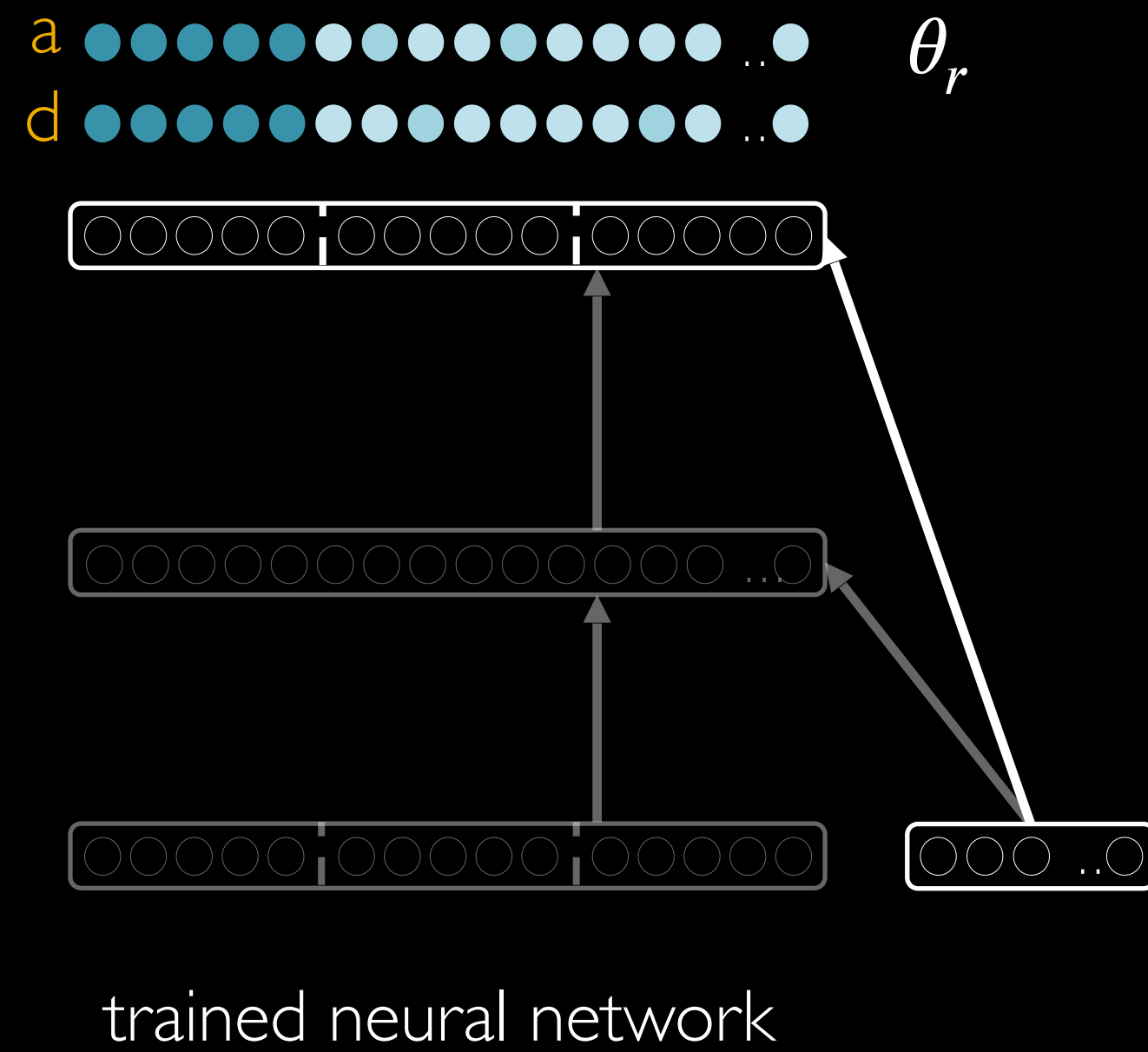
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



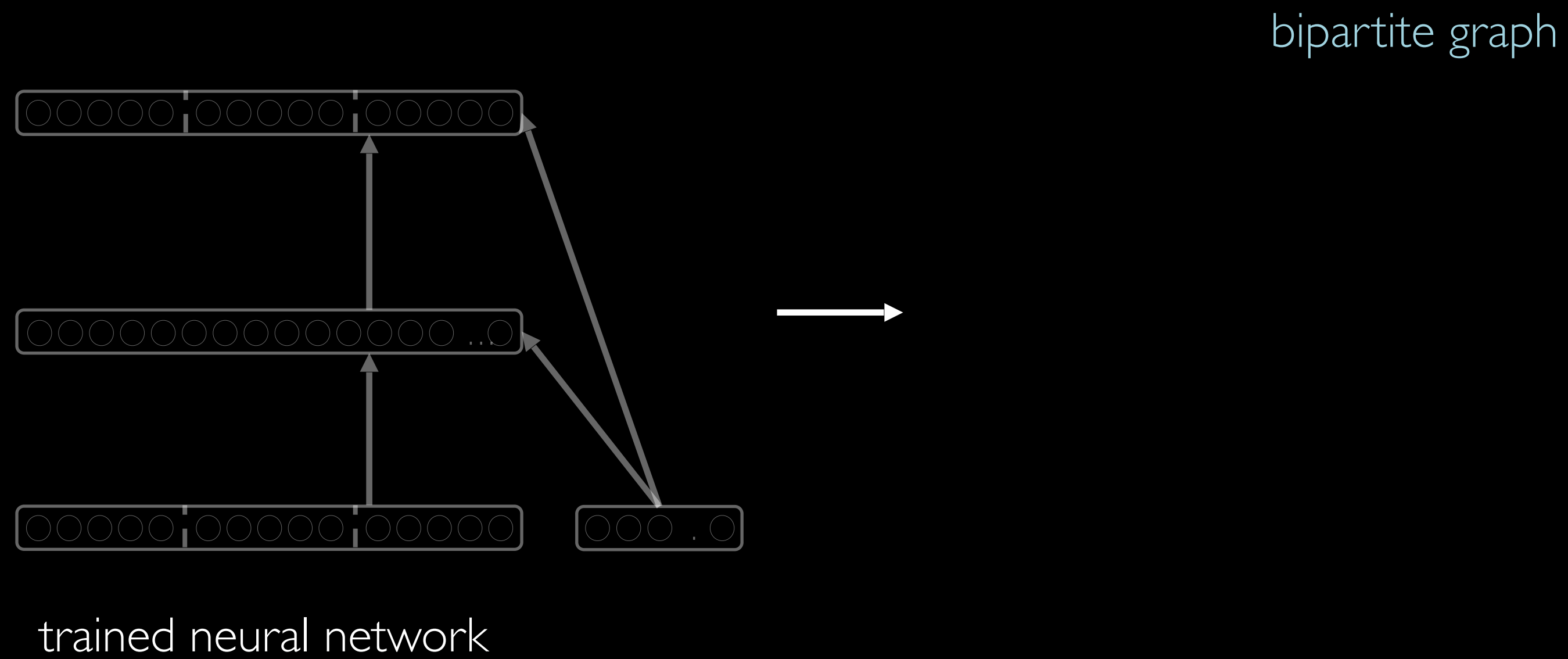
# Application to Neural Systems

## Extract Dependency Graph From Trained Neural Network



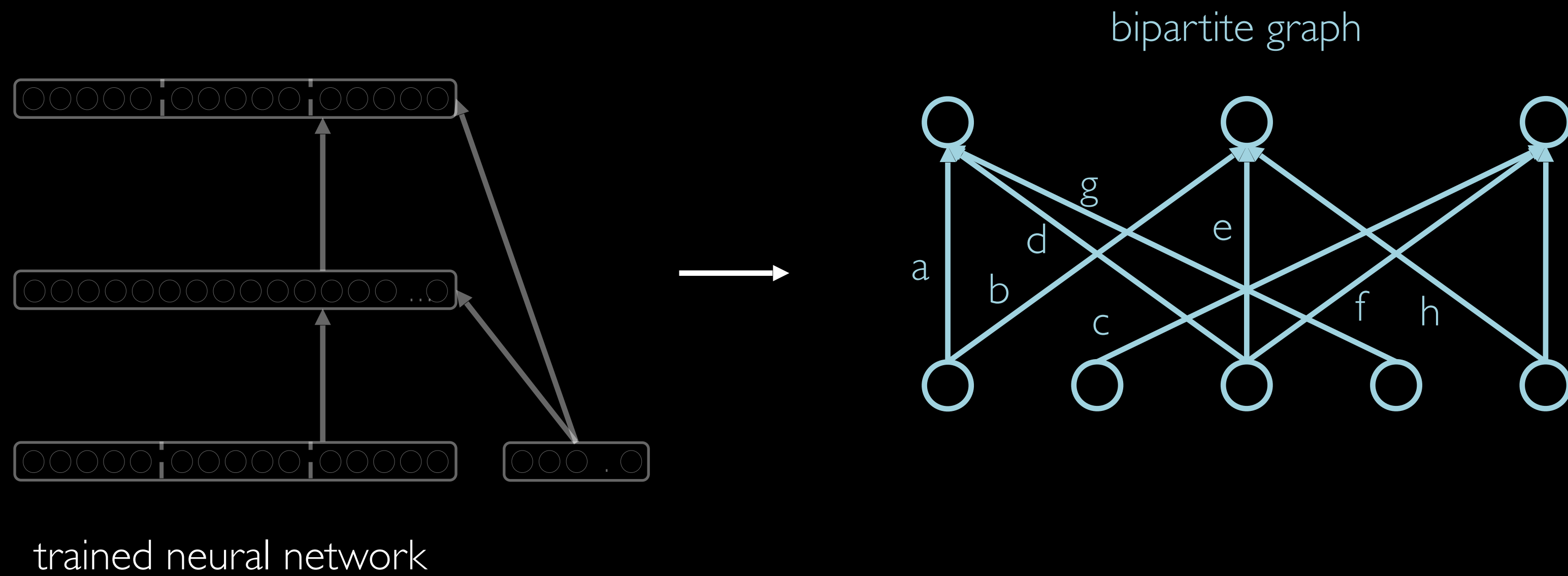
# Application to Neural Systems

## Predict Parallel Processing Performance Based On Dependency Graph



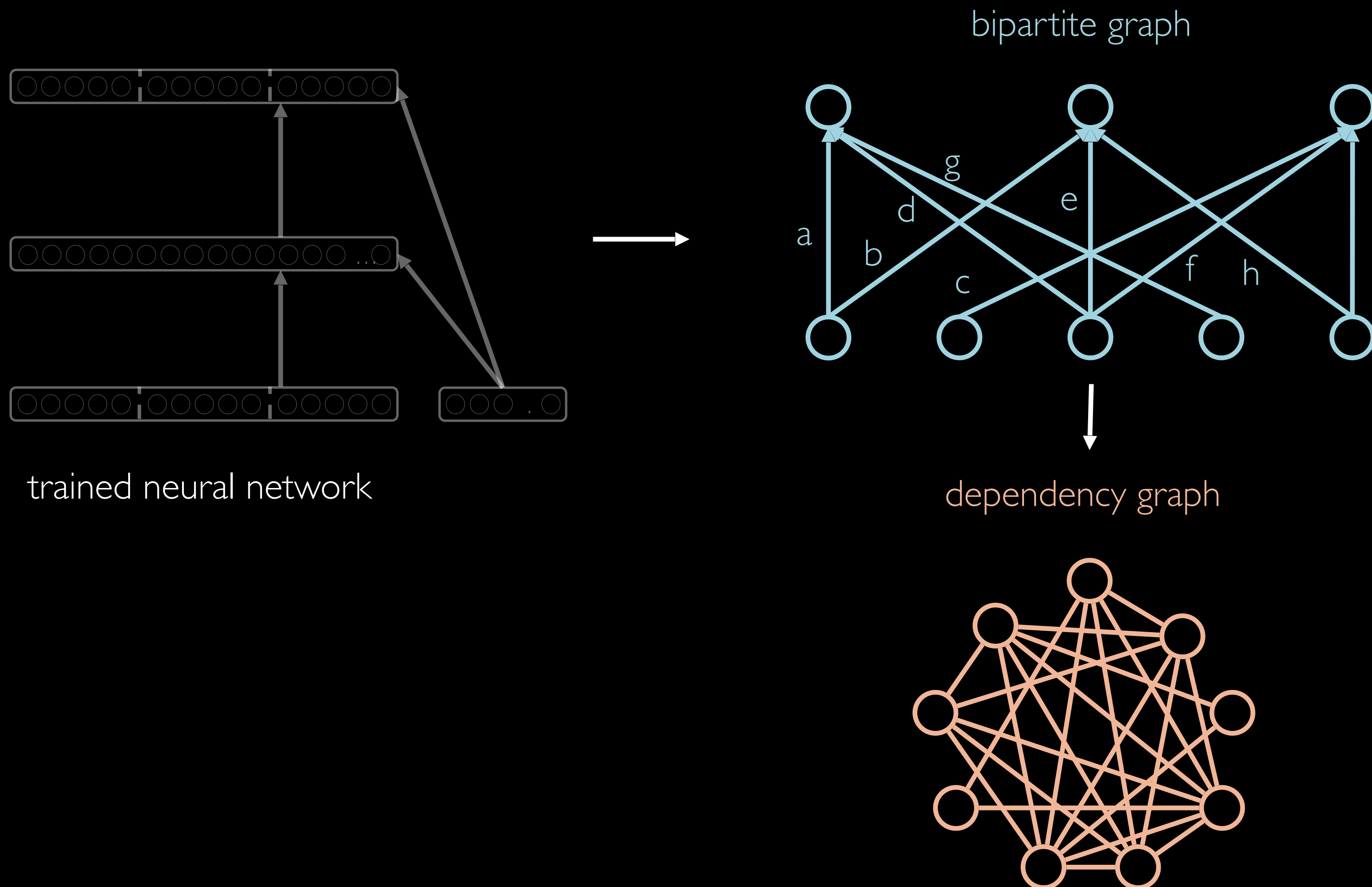
# Application to Neural Systems

## Predict Parallel Processing Performance Based On Dependency Graph



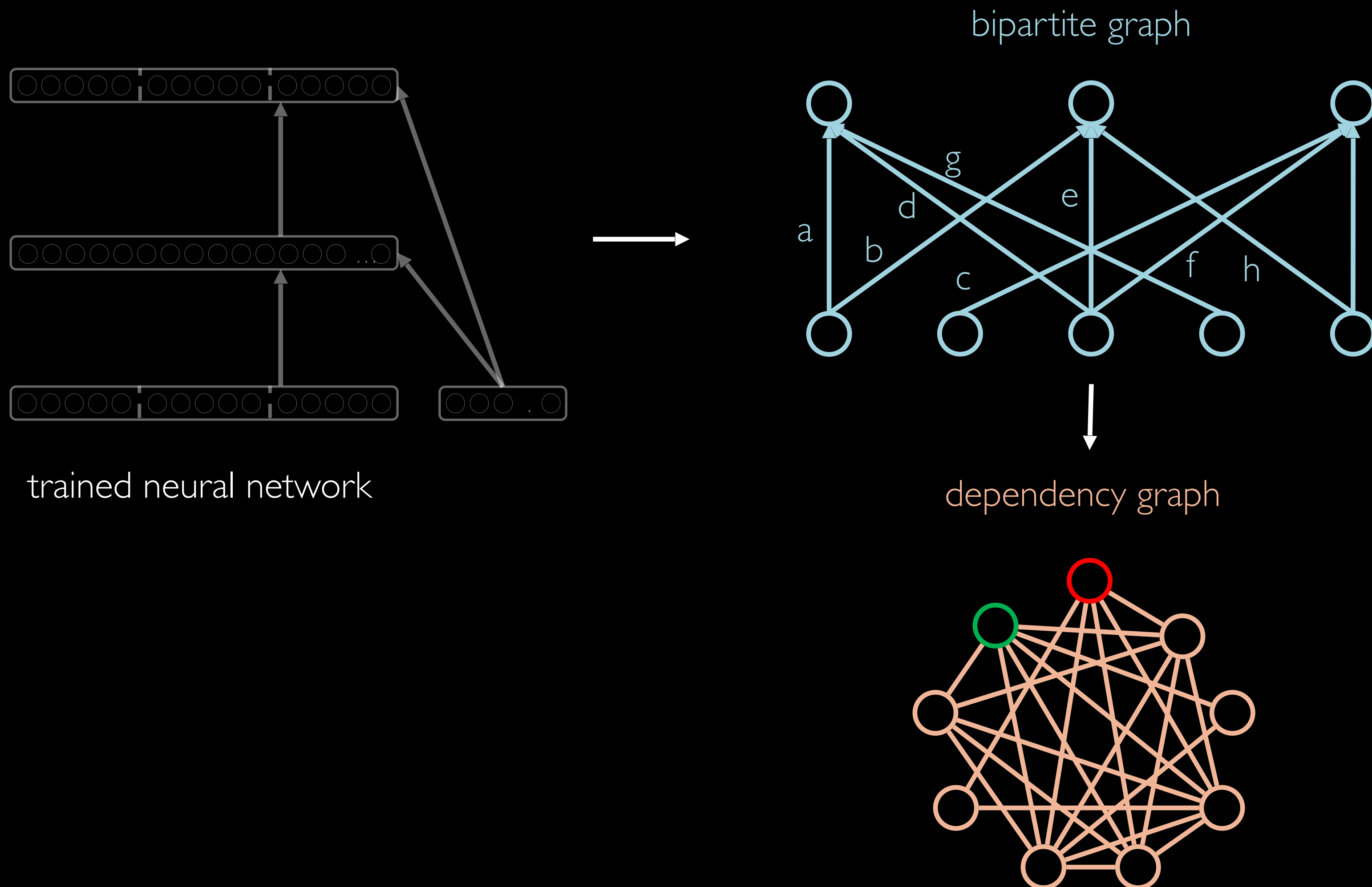
# Application to Neural Systems

## Predict Parallel Processing Performance Based On Dependency Graph



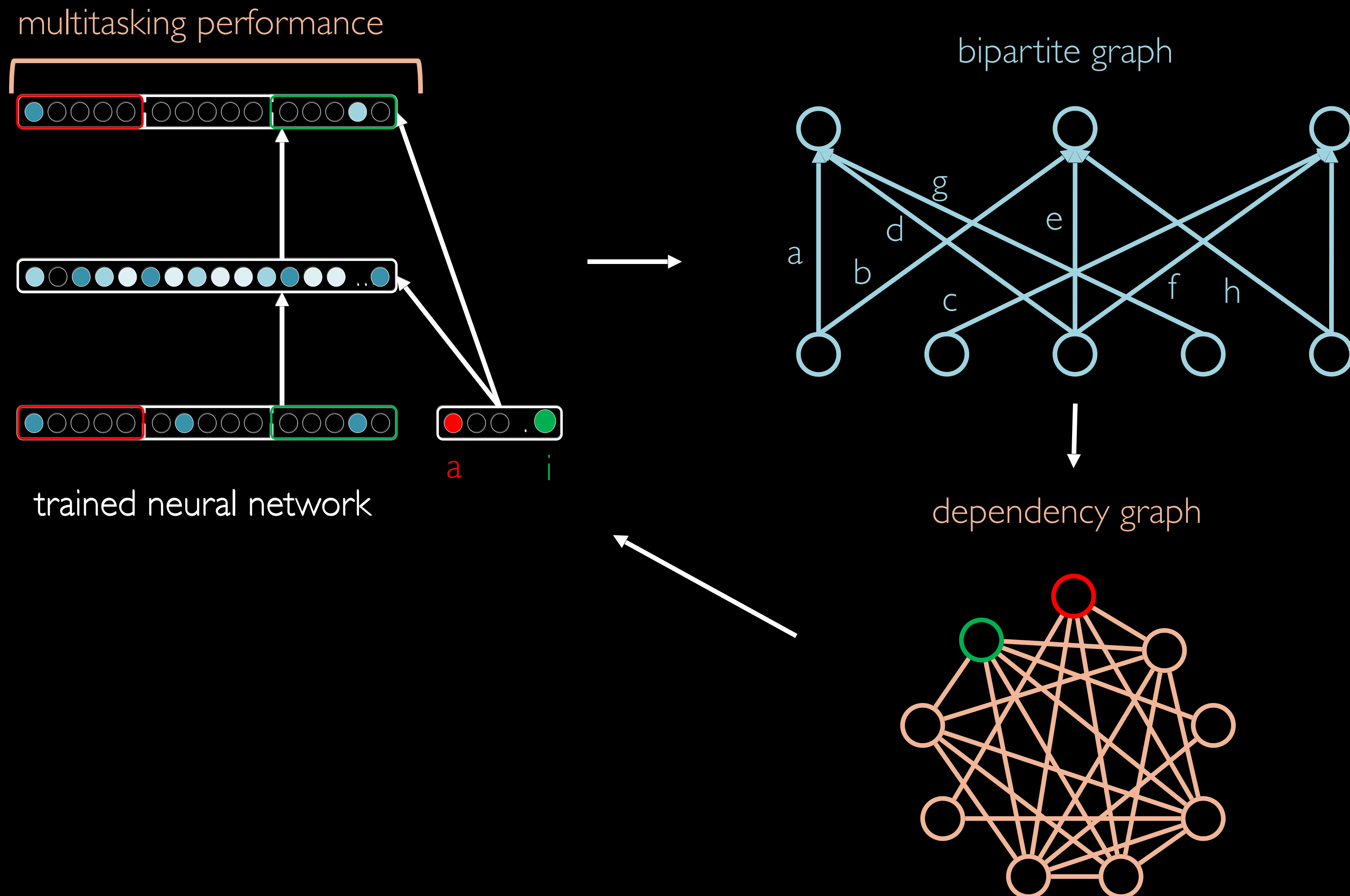
# Application to Neural Systems

## Predict Parallel Processing Performance Based On Dependency Graph



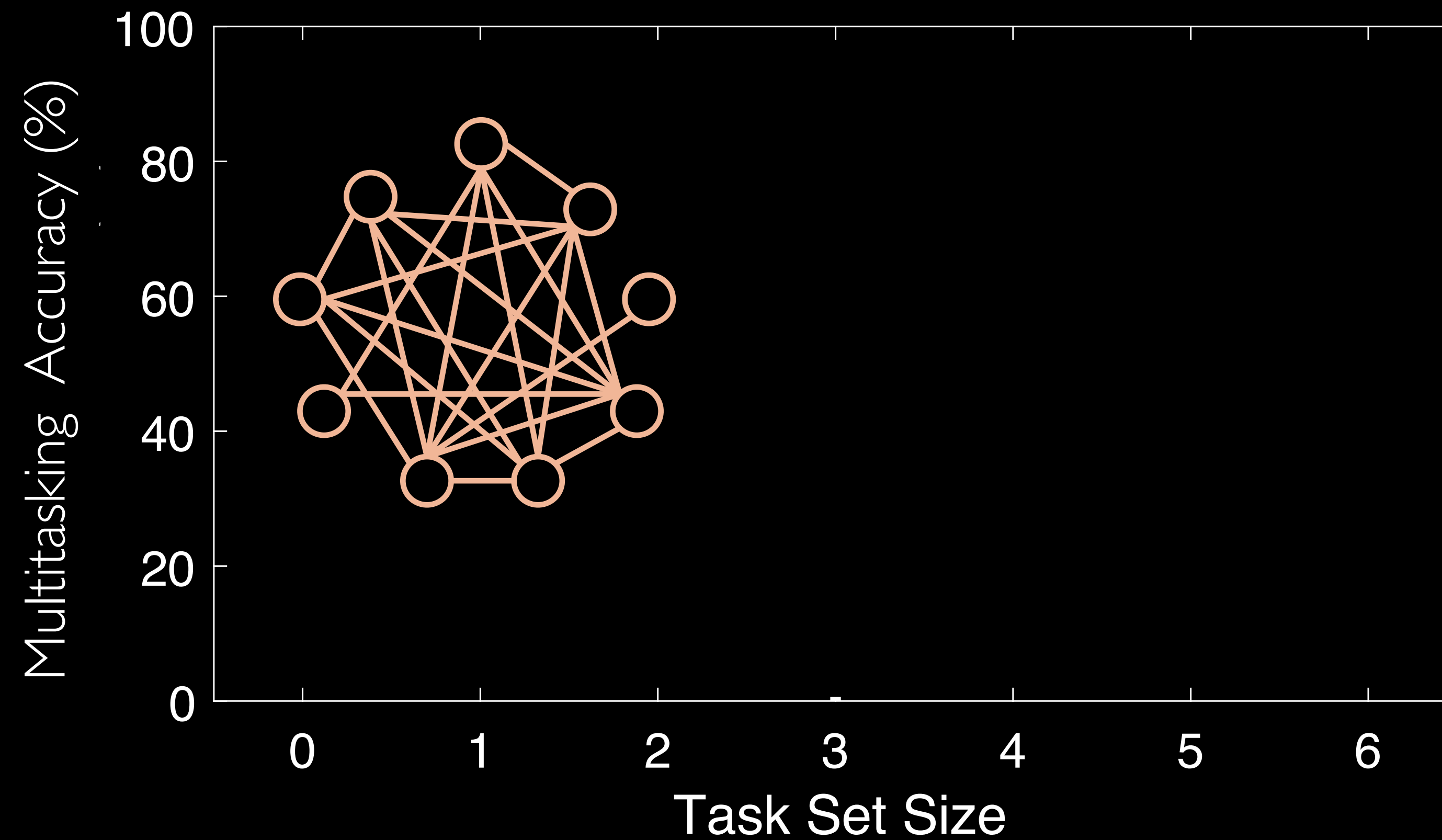
# Application to Neural Systems

## Predict Parallel Processing Performance Based On Dependency Graph



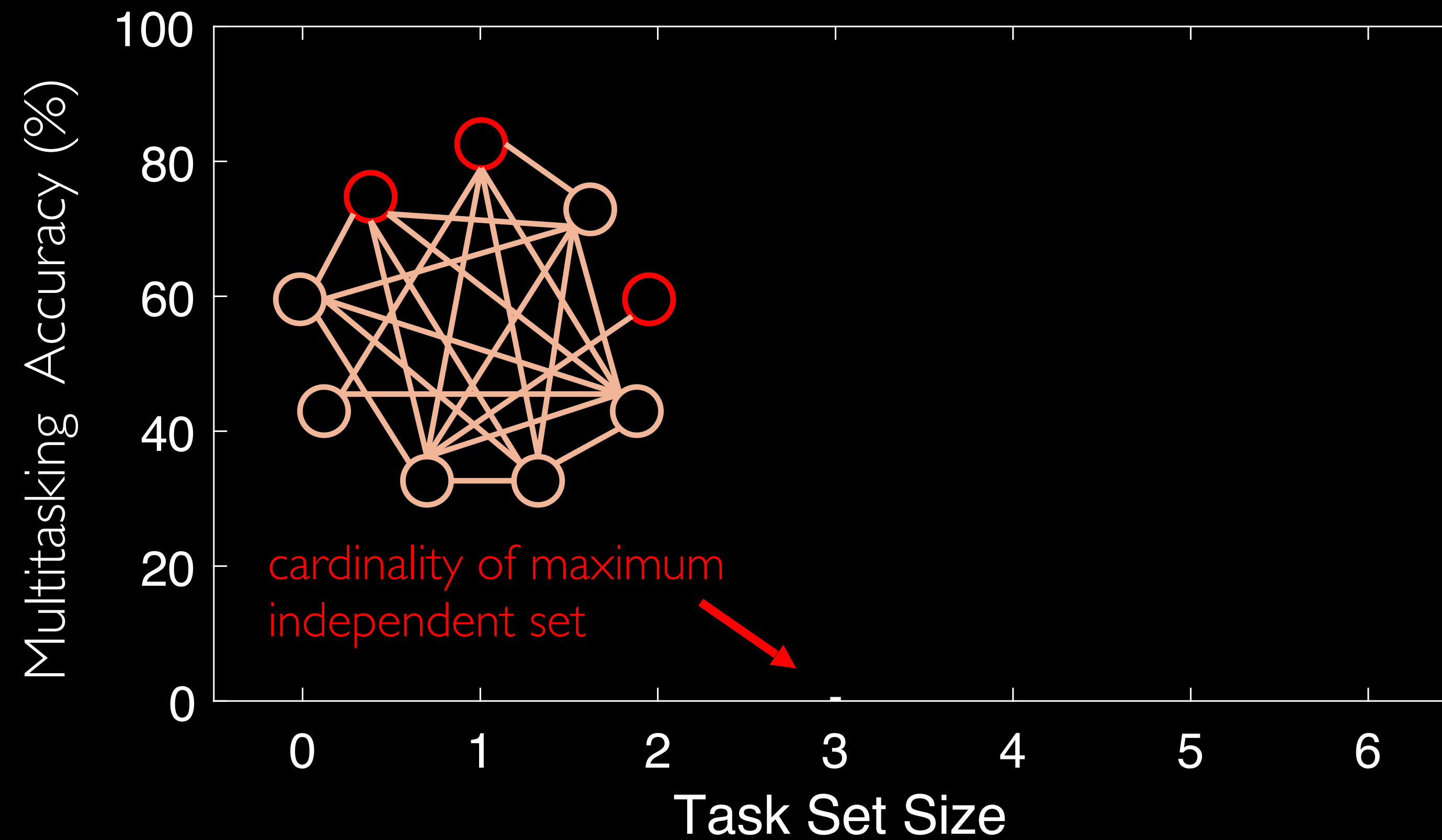
# Application to Neural Systems

Predict Parallel Processing Capacity Based on MIS



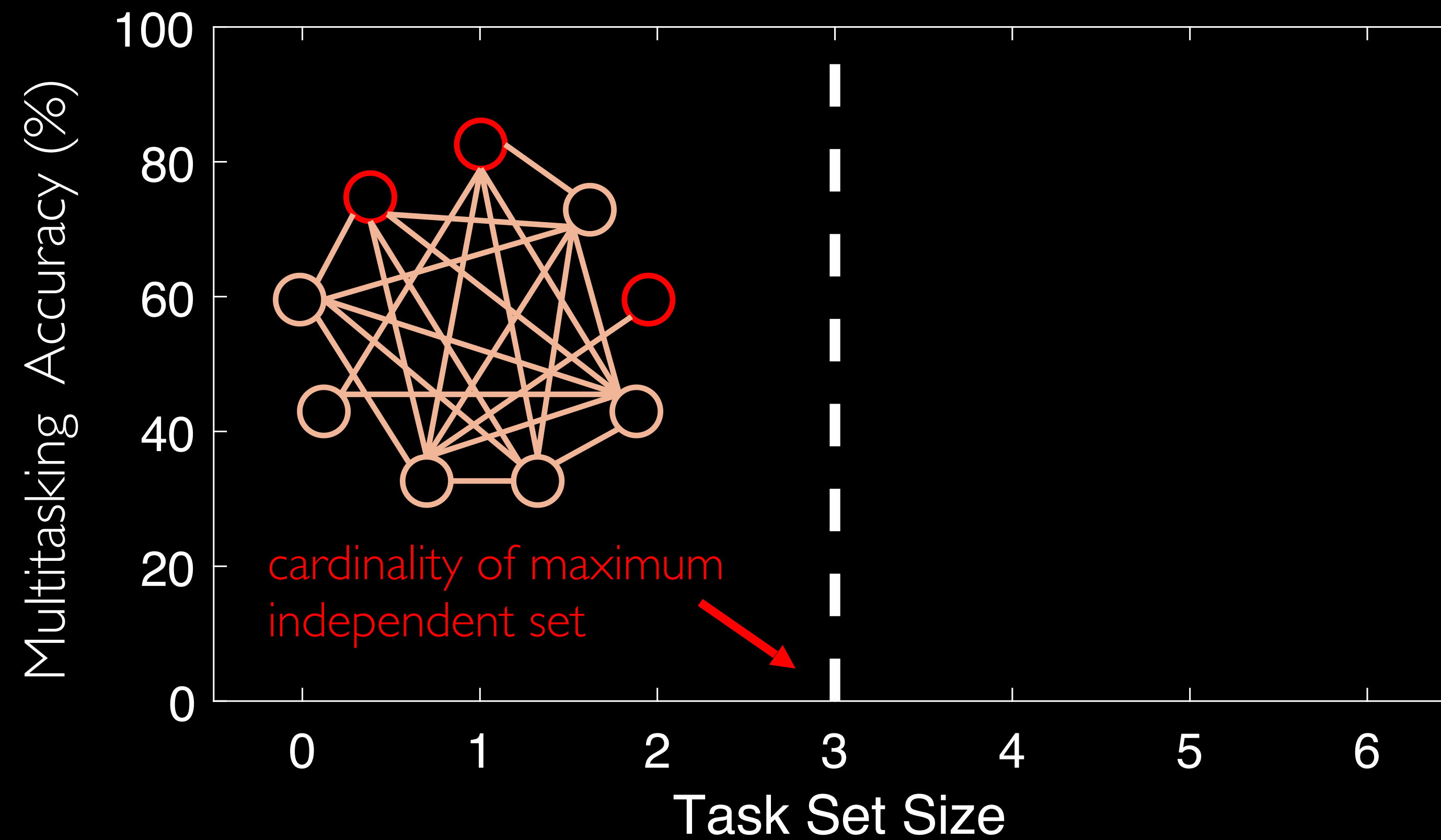
# Application to Neural Systems

## Predict Parallel Processing Capacity Based on MIS



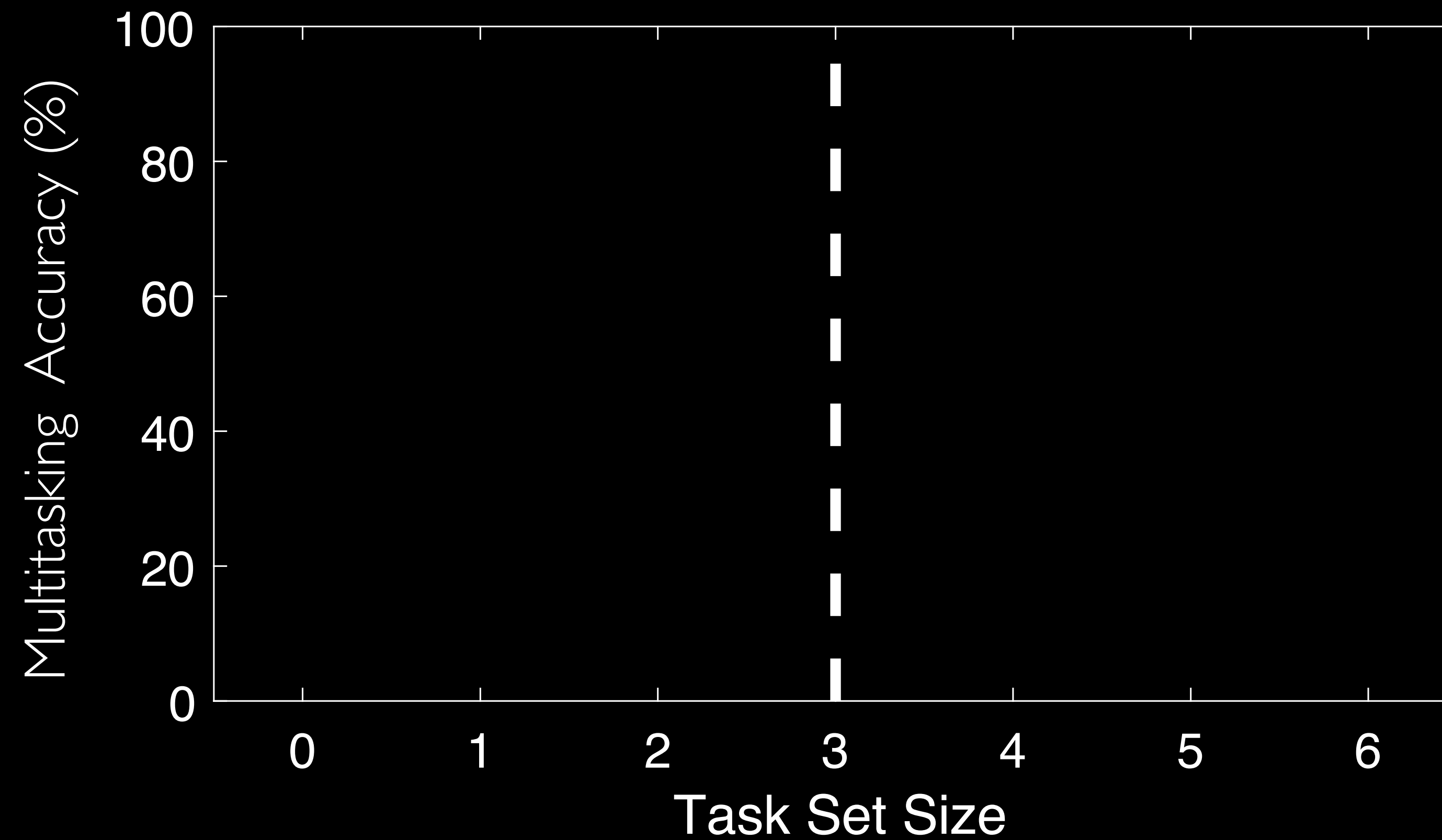
# Application to Neural Systems

## Predict Parallel Processing Capacity Based on MIS



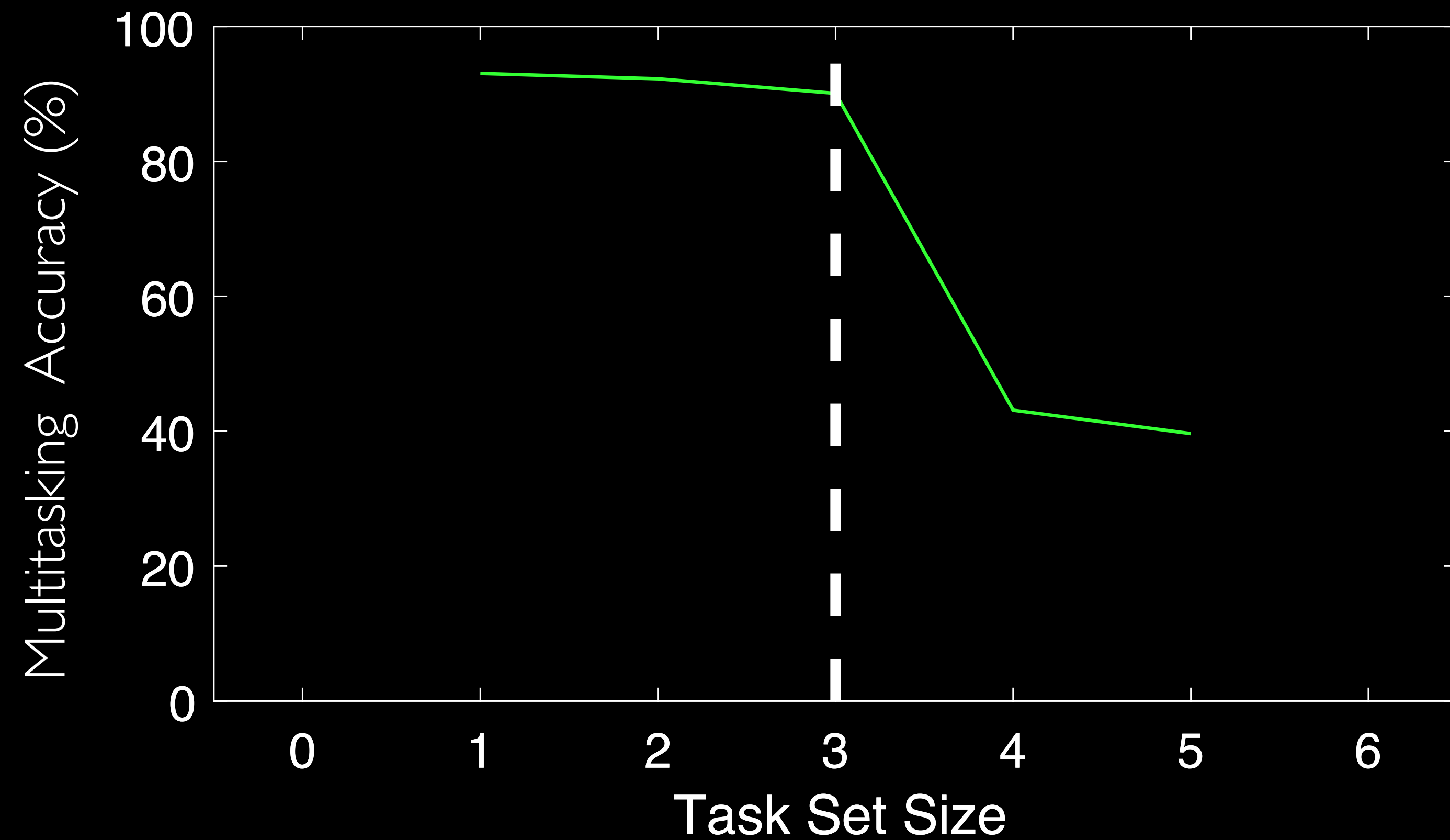
# Application to Neural Systems

Predict Parallel Processing Capacity Based on MIS



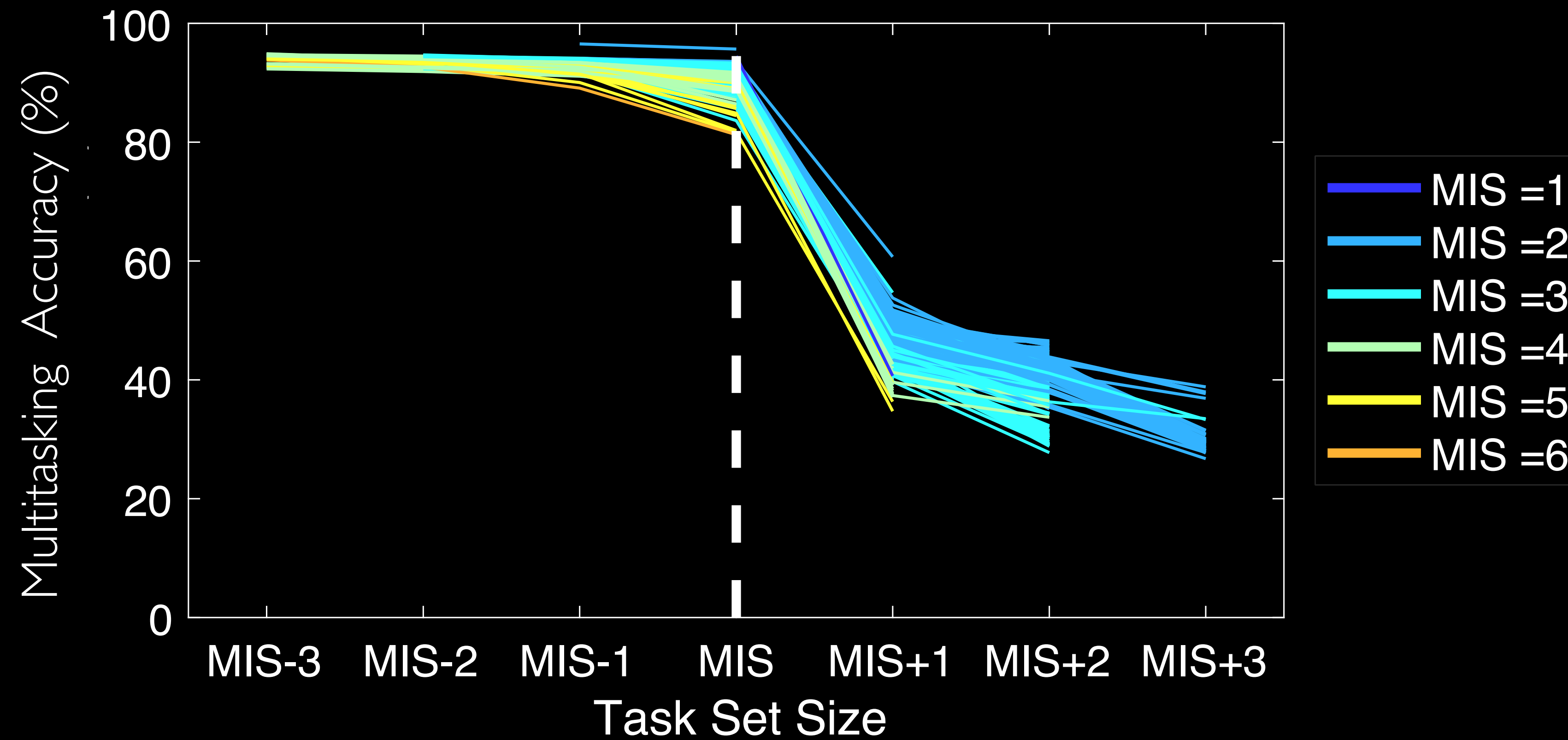
# Application to Neural Systems

Predict Parallel Processing Capacity Based on MIS



# Application to Neural Systems

Predict Parallel Processing Capacity Based on MIS



**Good.**

**Can we scale this up?**

# Stat. mech. approach to MIS estimation

## Max Set Packing in a nutshell

The solution comes out like this....

$$\rho_\alpha = \frac{\langle k \rangle}{\langle c \rangle} \left( 1 - p_*^{c/(c-1)} \right) + (M_k(t) - M'(t))$$

$$p_* = \mathbb{E}_{\tilde{c}} \left( 1 - \frac{1}{\langle k \rangle p_*} M'(t) \right)^{\tilde{c}}$$

# Stat. mech. approach to MIS estimation

## Max Set Packing in a nutshell

The solution comes out like this....

$$\rho_\alpha = \frac{\langle k \rangle}{\langle c \rangle} \left( 1 - p_*^{c/(c-1)} \right) + (M_k(t) - M'(t))$$

$$p_* = \mathbb{E}_{\tilde{c}} \left( 1 - \frac{1}{\langle k \rangle p_*} M'(t) \right)^{\tilde{c}}$$

Degree-generating  
function



# Stat. mech. approach to MIS estimation

## Max Set Packing in a nutshell

The solution comes out like this....

$$\rho_\alpha = \frac{\langle k \rangle}{\langle c \rangle} \left( 1 - p_*^{c/(c-1)} \right) + (M_k(t) - M'(t))$$

$$p_* = \mathbb{E}_{\tilde{c}} \left( 1 - \frac{1}{\langle k \rangle p_*} M'(t) \right)^{\tilde{c}}$$

Degree-generating function



Specifying to Gaussian degree distributions, it becomes...

$$M_k(t) = e^{\langle k \rangle t + \sigma^2 t^2 / 2}$$

# Stat. mech. approach to MIS estimation

## Max Set Packing in a nutshell

The solution comes out like this...

$$\rho_\alpha = \frac{\langle k \rangle}{\langle c \rangle} \left( 1 - p_*^{c/(c-1)} \right) + (M_k(t) - M'(t))$$

$$p_* = \mathbb{E}_{\tilde{c}} \left( 1 - \frac{1}{\langle k \rangle p_*} M'(t) \right)^{\tilde{c}}$$

Degree-generating function



Specifying to Gaussian degree distributions, it becomes...

$$M_k(t) = e^{\langle k \rangle t + \sigma^2 t^2 / 2}$$

$$\rho_\alpha = \frac{\langle k \rangle}{\langle c \rangle} \left( 1 - p_*^{c/(c-1)} \right) + M_k(\ln p_*) (1 - \langle k \rangle - \sigma^2 \ln p_*)$$

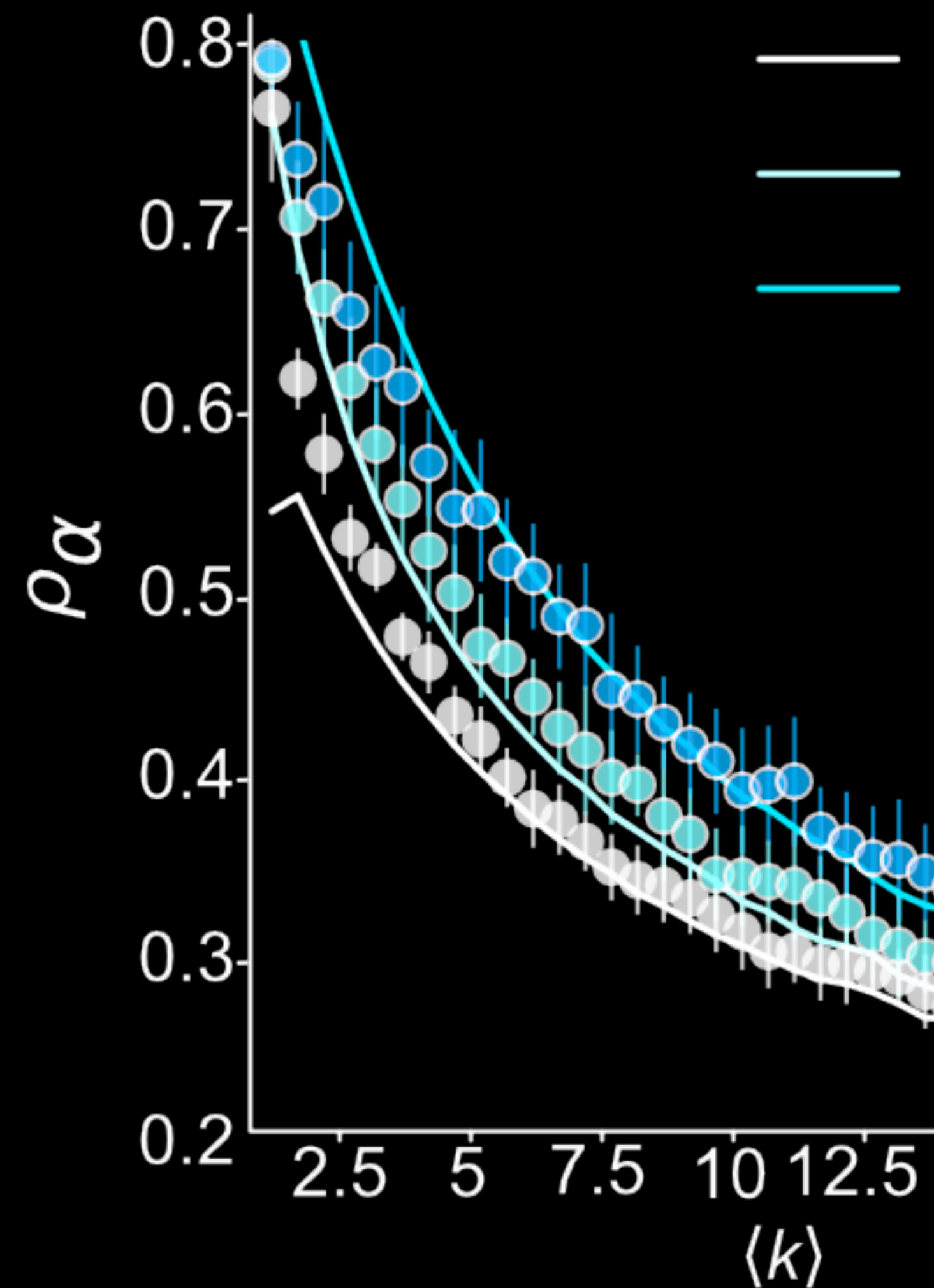
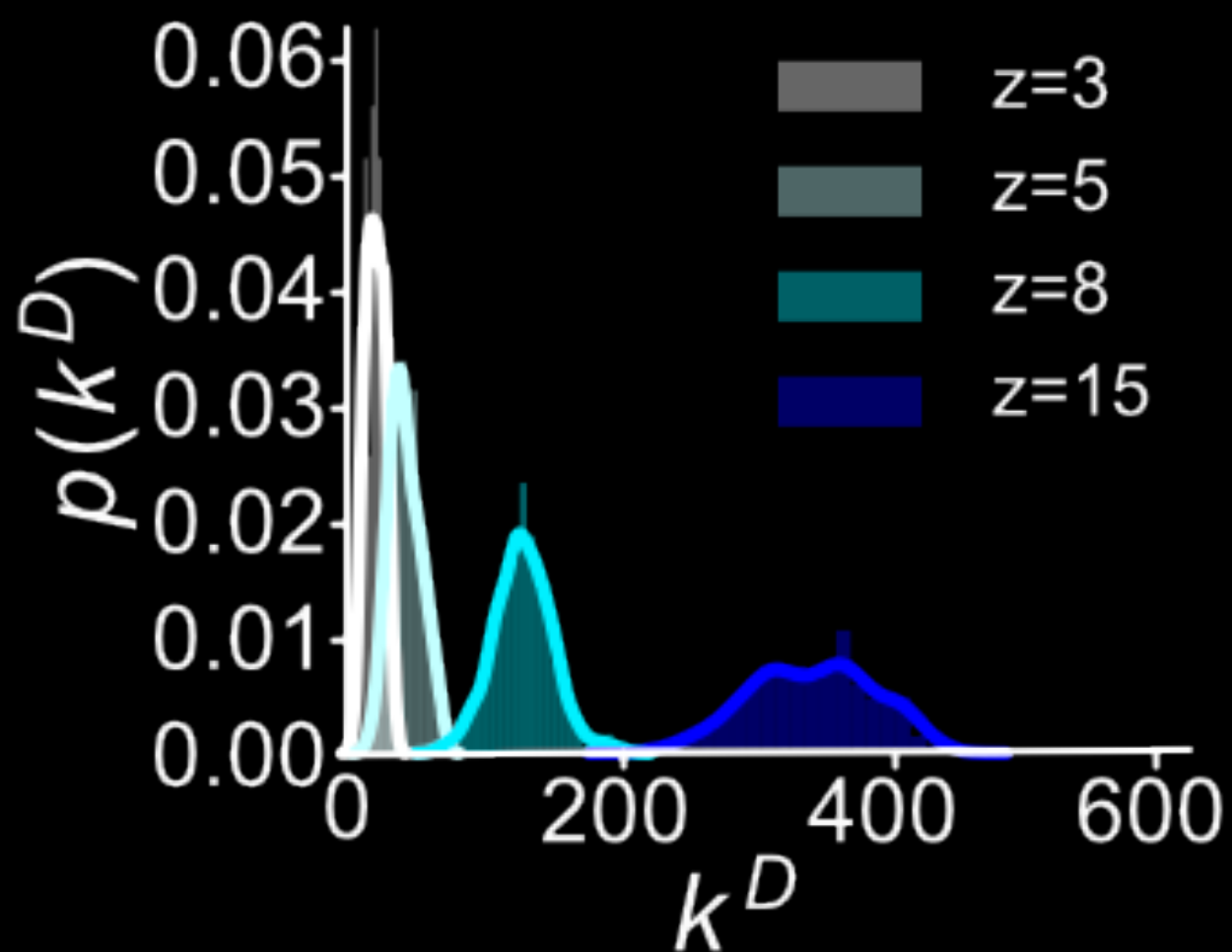
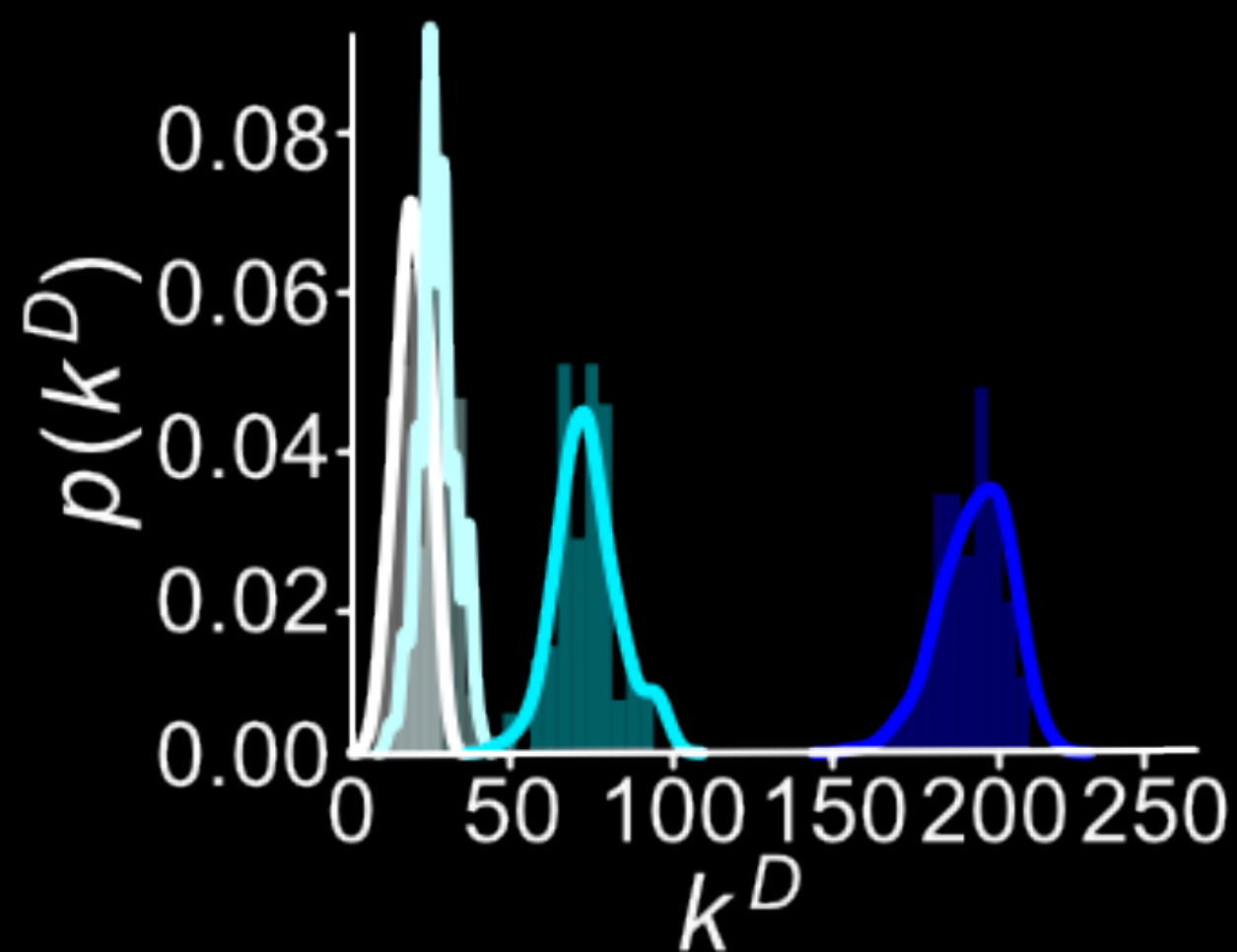
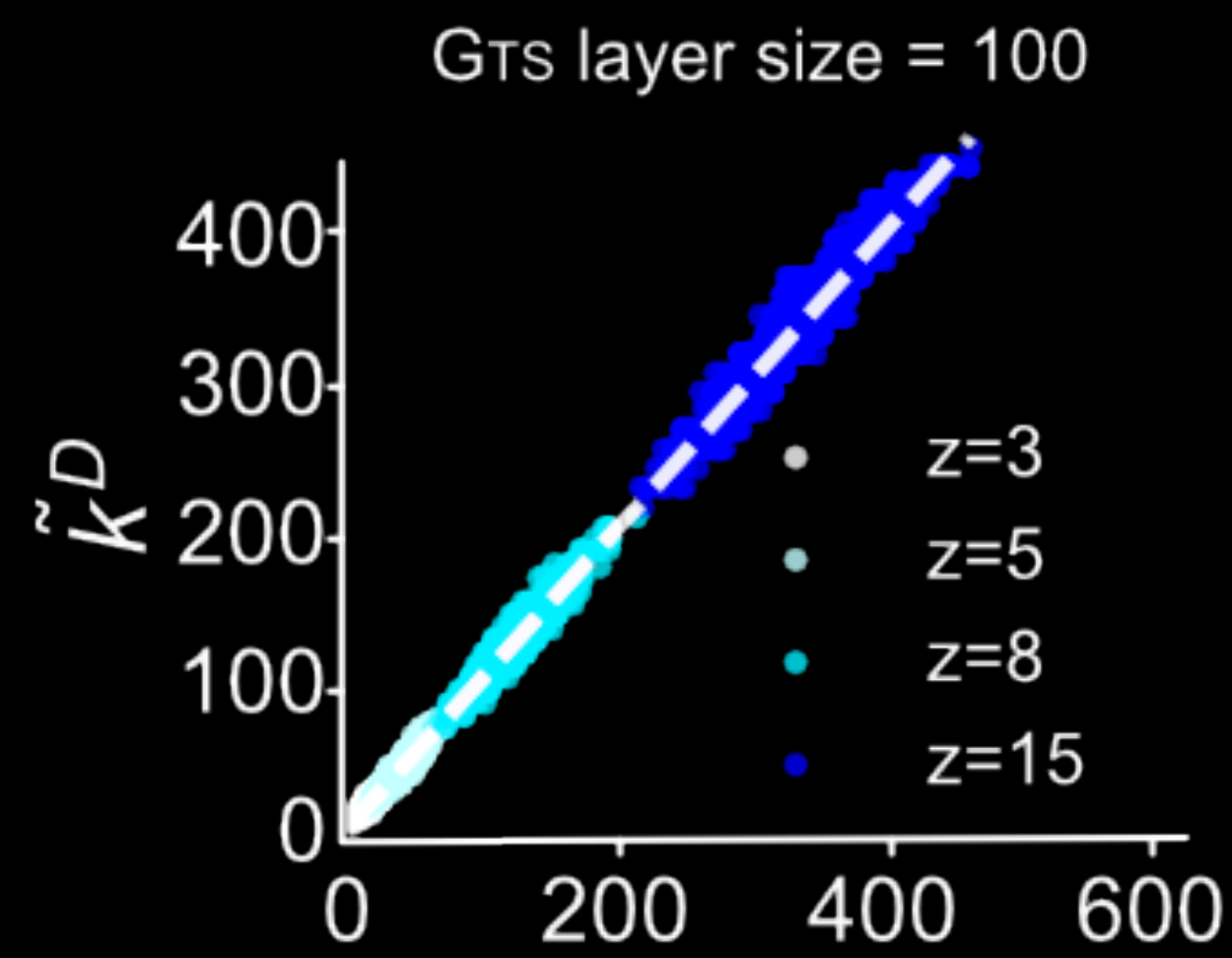
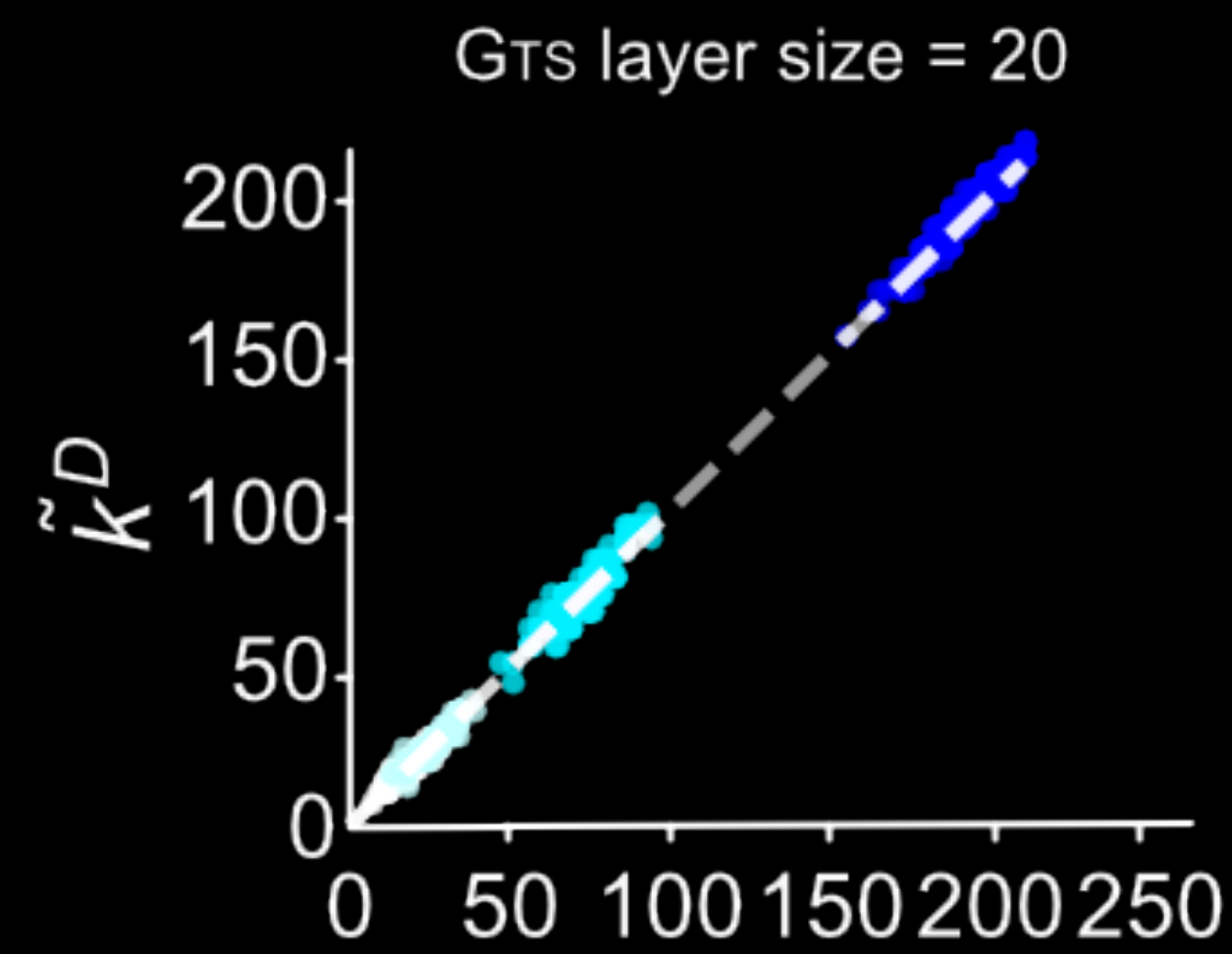
# Stat. mech. approach to MIS estimation

MIS results

# Stat. mech. approach to MIS estimation

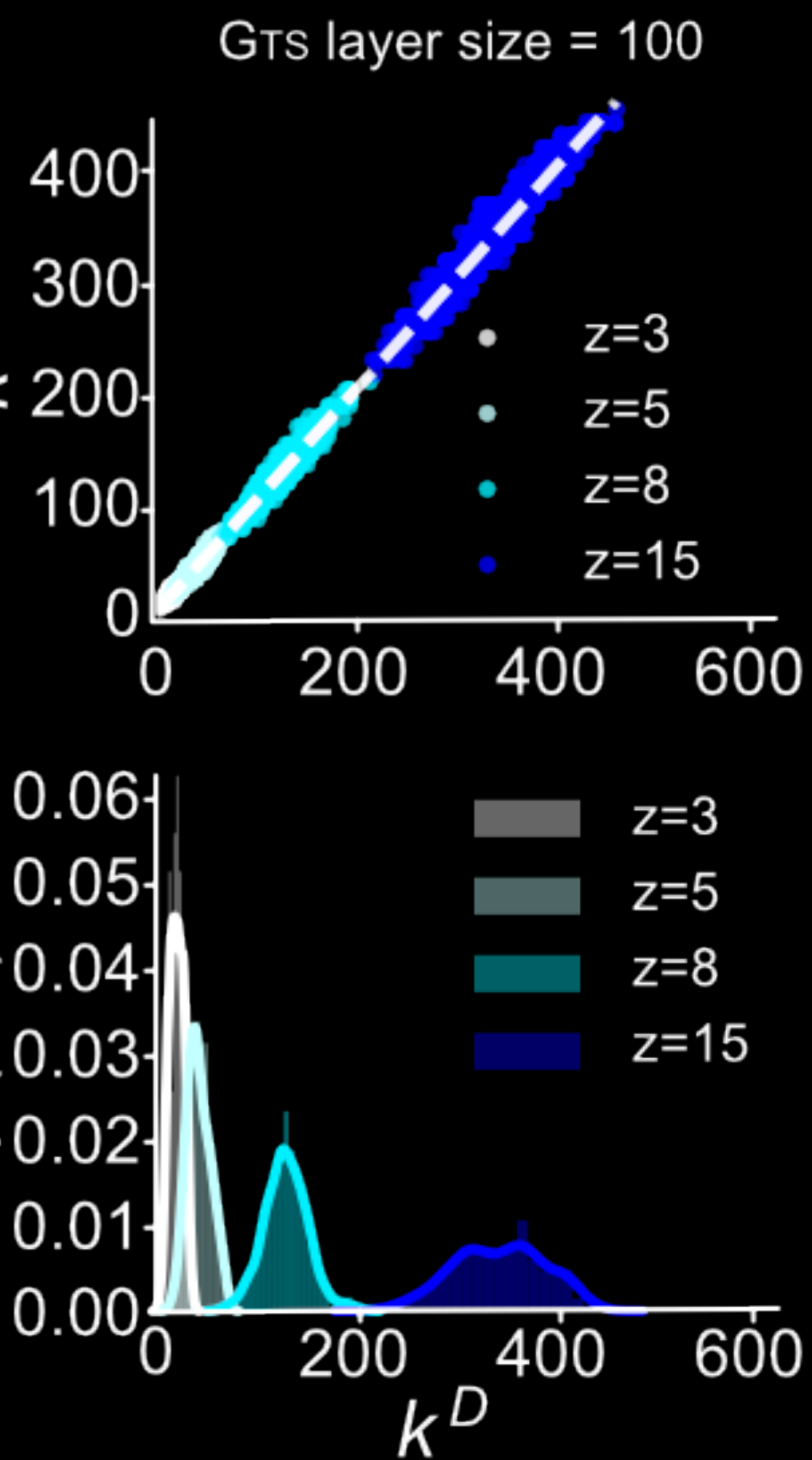
## MIS results

Dependency graph degree estimation

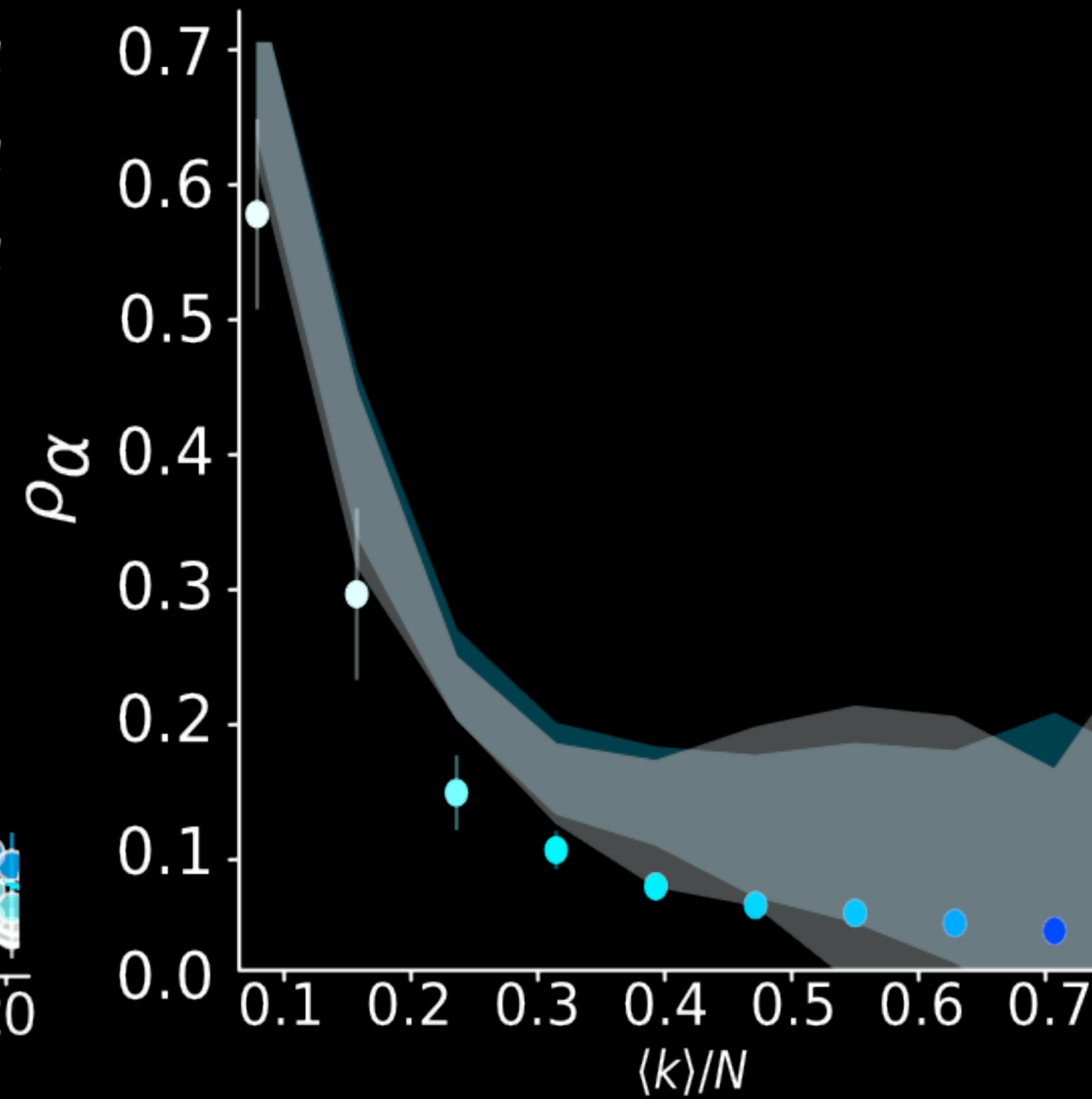
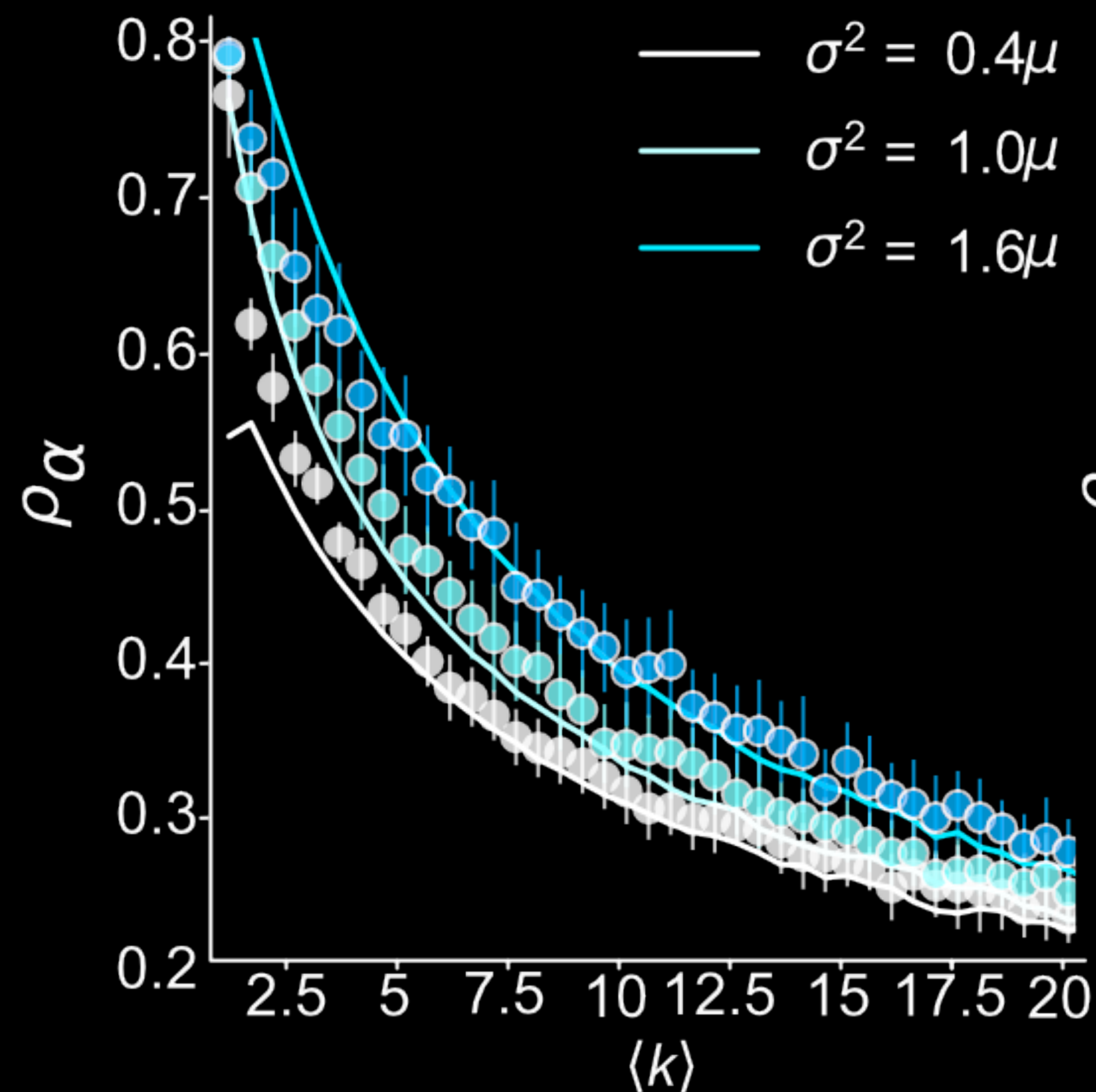


# Stat. mech. approach to MIS estimation

## MIS results

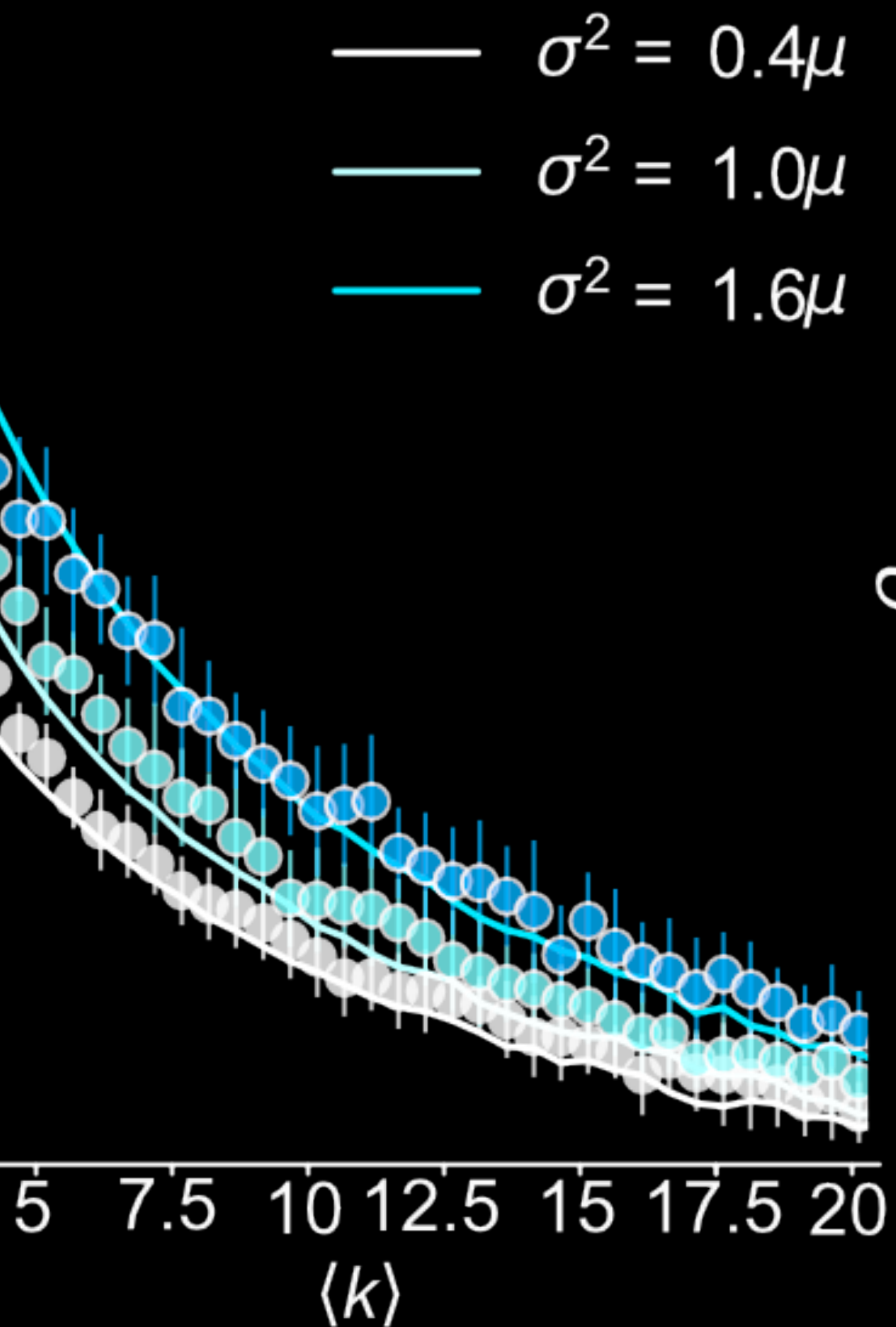


### MIS prediction as function of density and heterogeneity

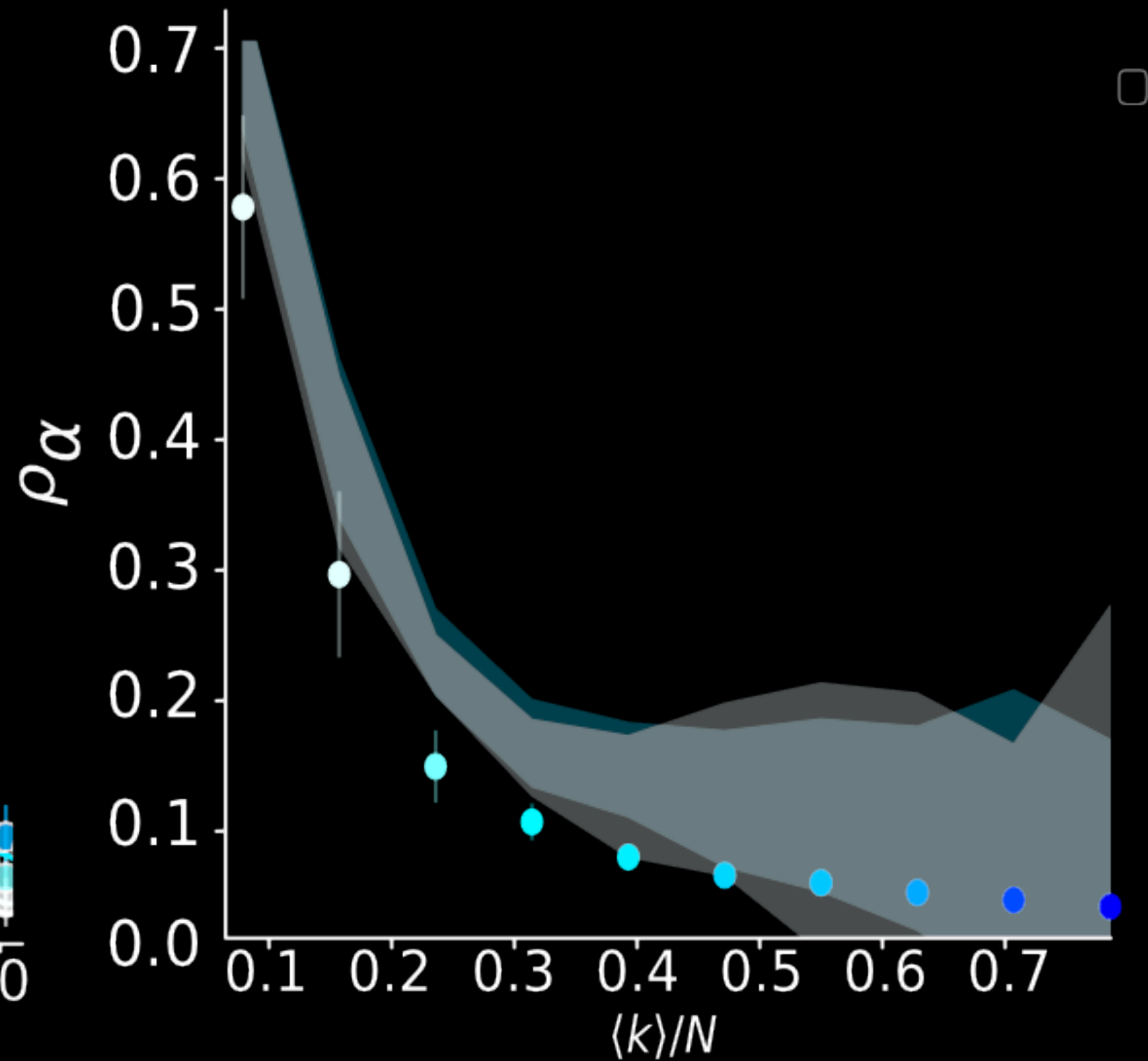


# Stat. mech. approach to MIS estimation

MIS results



All together!



**This holds for the MIS.  
That is, a specific set.**

**This holds for the MIS.  
That is, a specific set.**

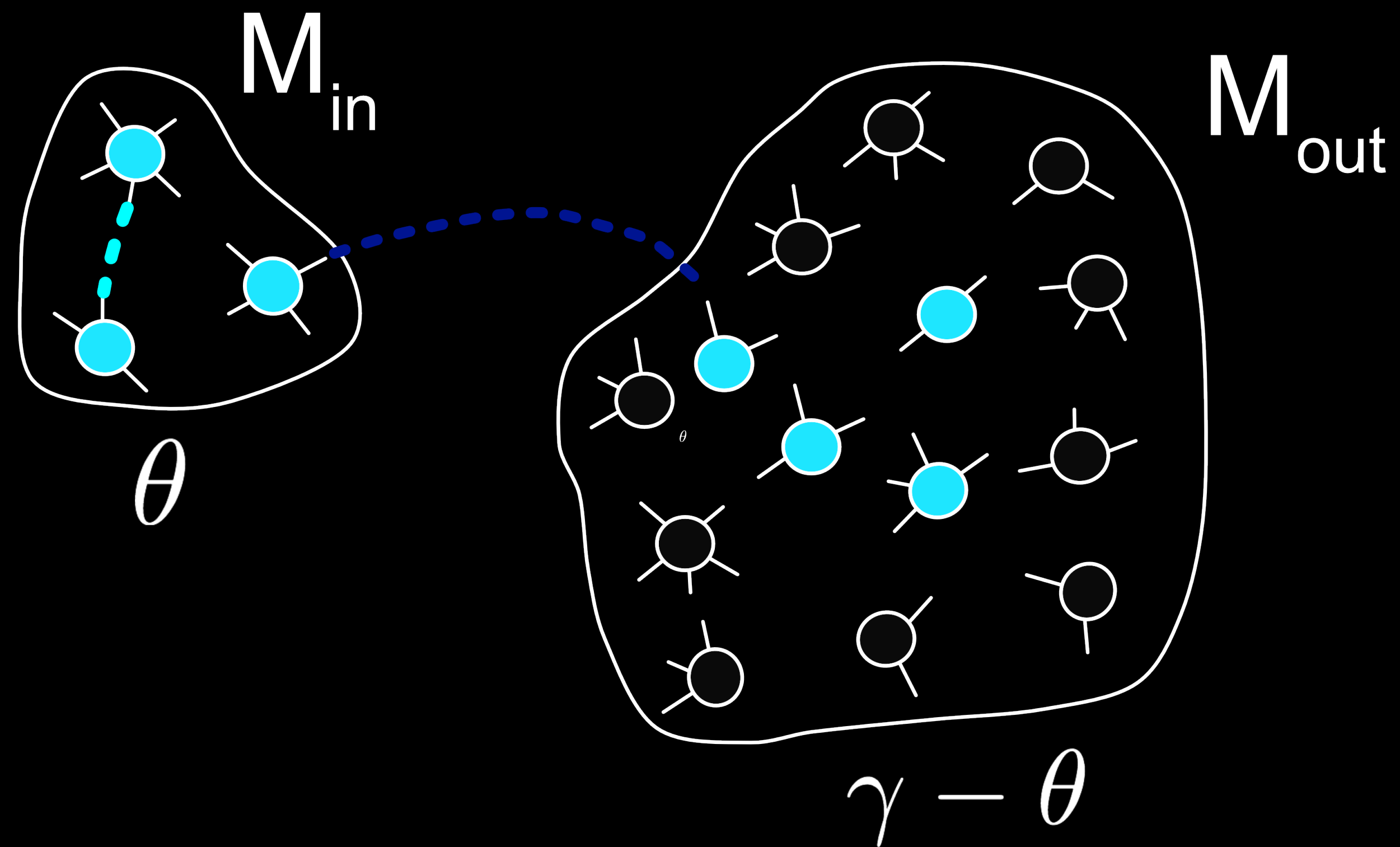
**What happens if I give you  
 $\gamma$  randomly chosen tasks?**

**This holds for the MIS.  
That is, a specific set.**

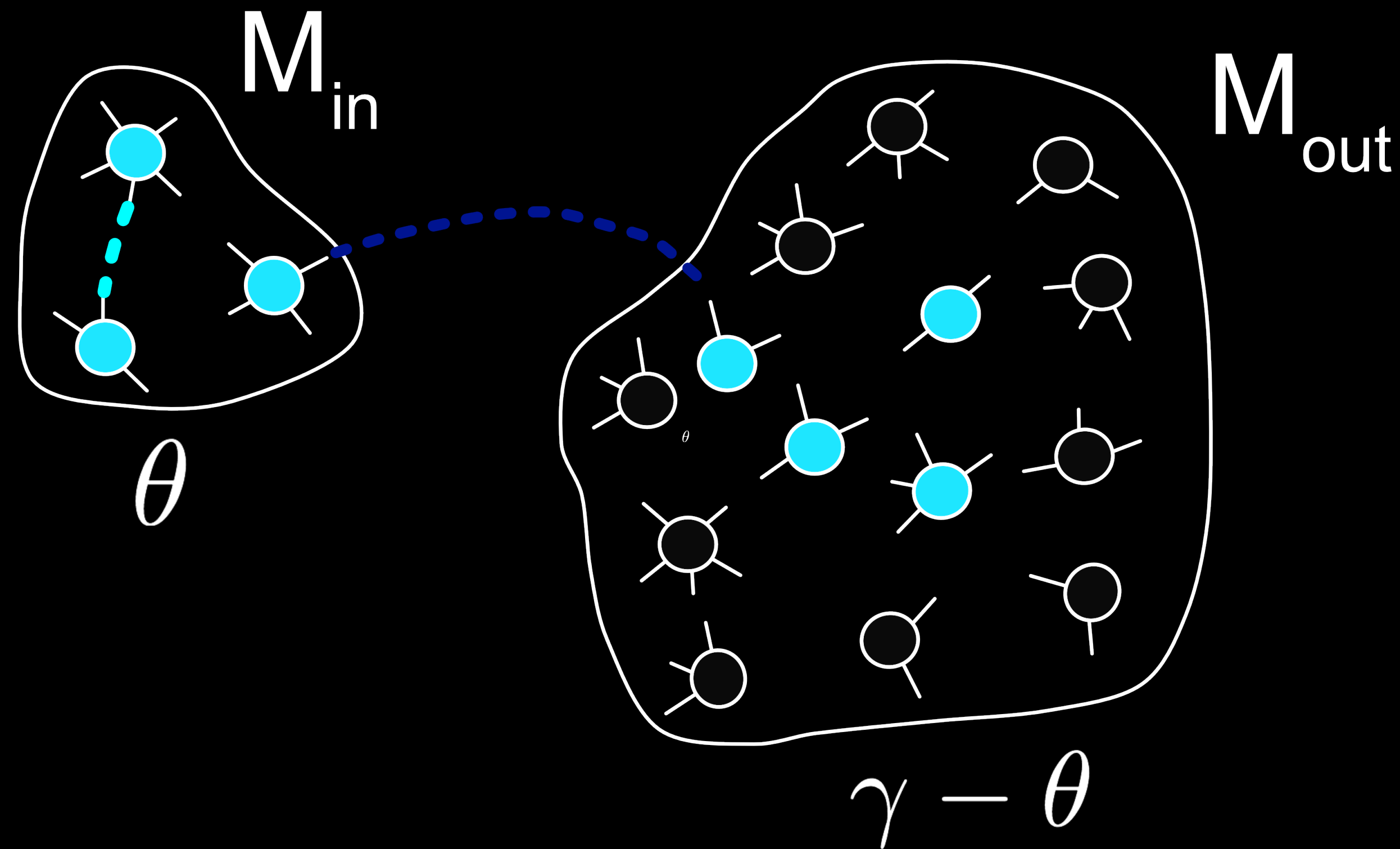
**What happens if I give you  
 $\gamma$  randomly chosen tasks?**

**expected capacity vs  
maximum capacity**

# Estimation of expected parallel capacity



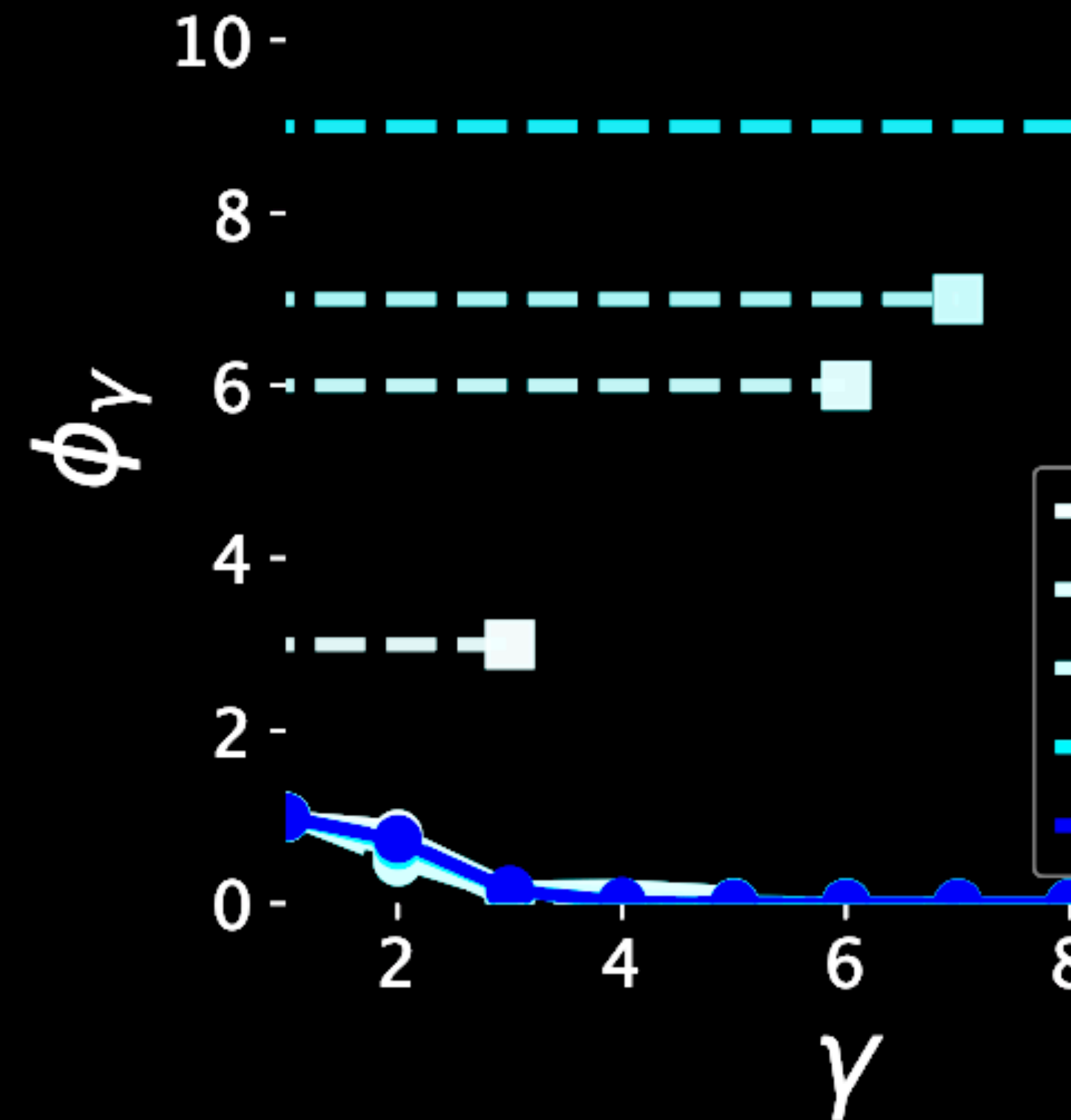
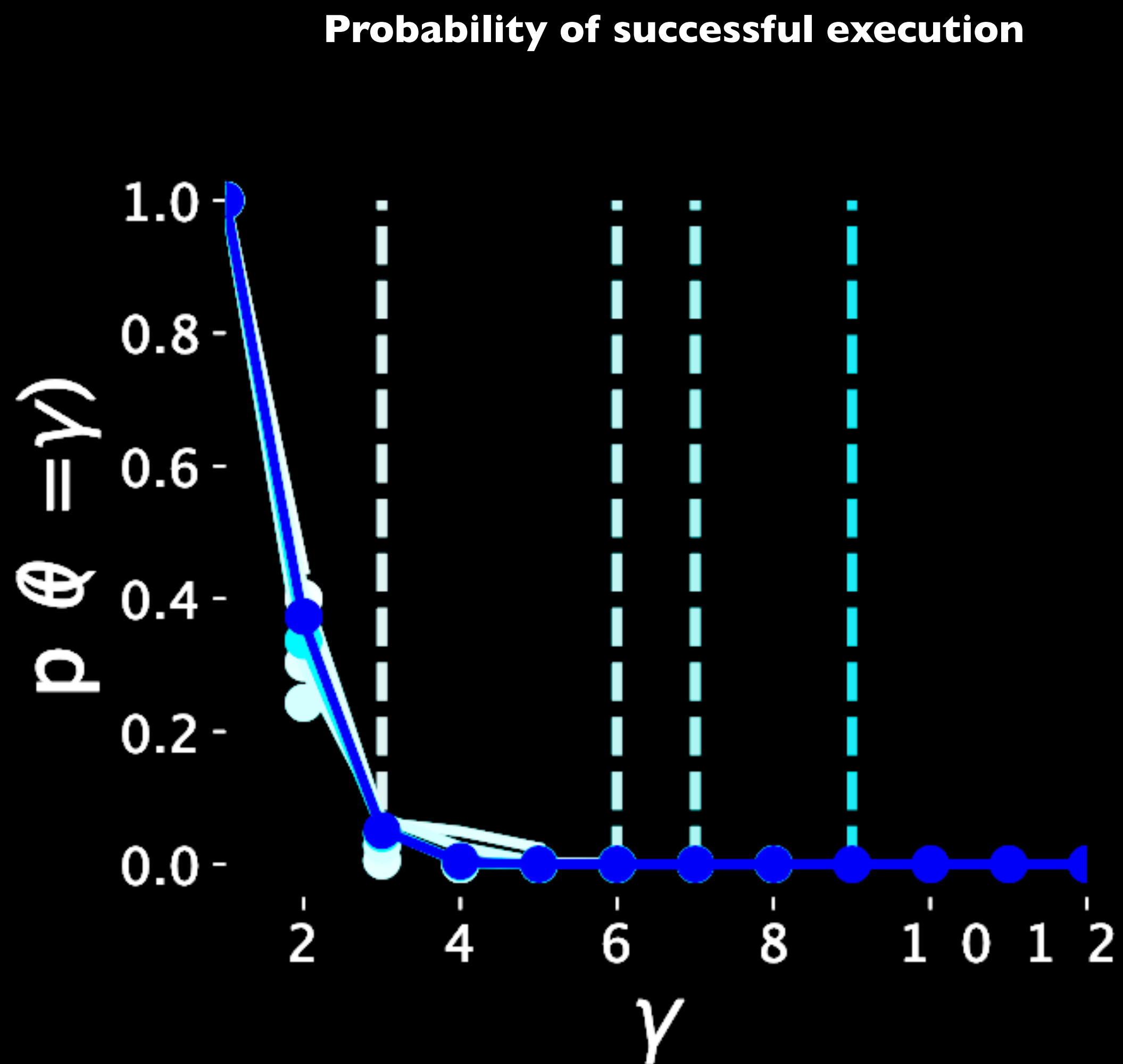
# Estimation of expected parallel capacity



$$P(\theta; \gamma, \mathcal{G}_D) \simeq \left(1 - \frac{\langle k^2 \rangle}{2M_D}\right)^{\binom{\theta}{2}} \left(\frac{\theta \langle k \rangle^2}{2M_D}\right)^{\gamma - \theta}$$

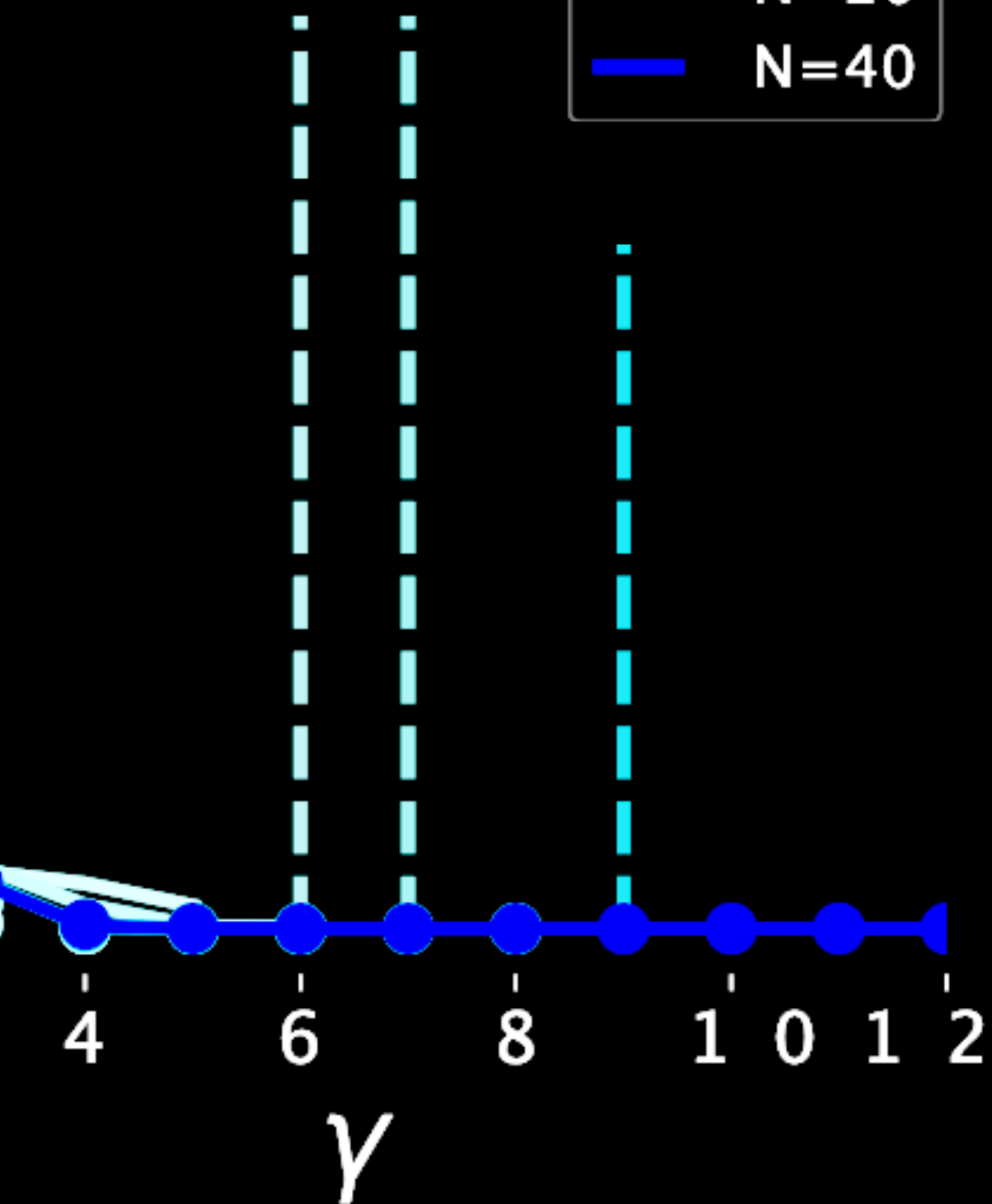
# Estimation of expected parallel capacity

Number of representations  
in task graph layers

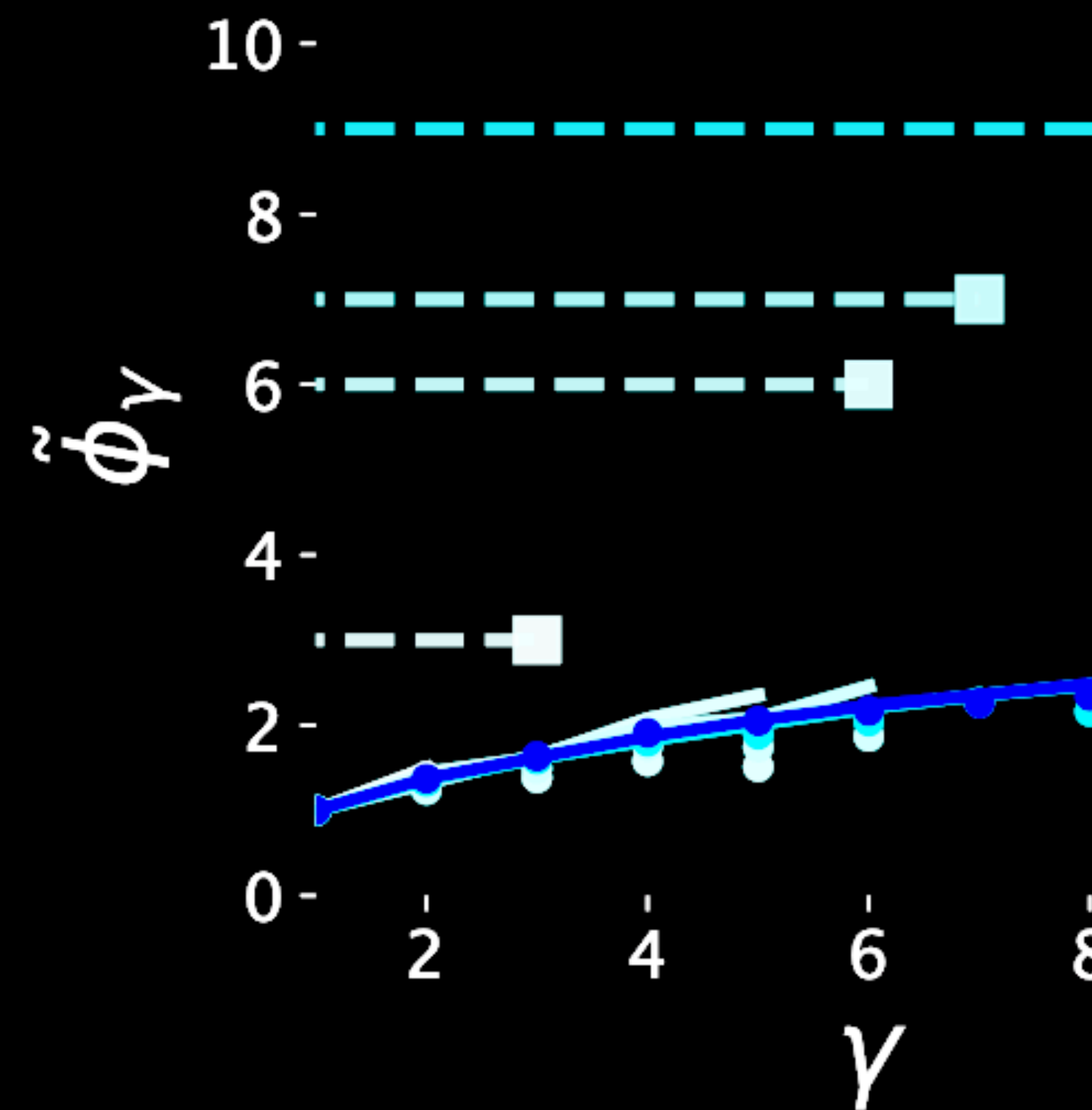
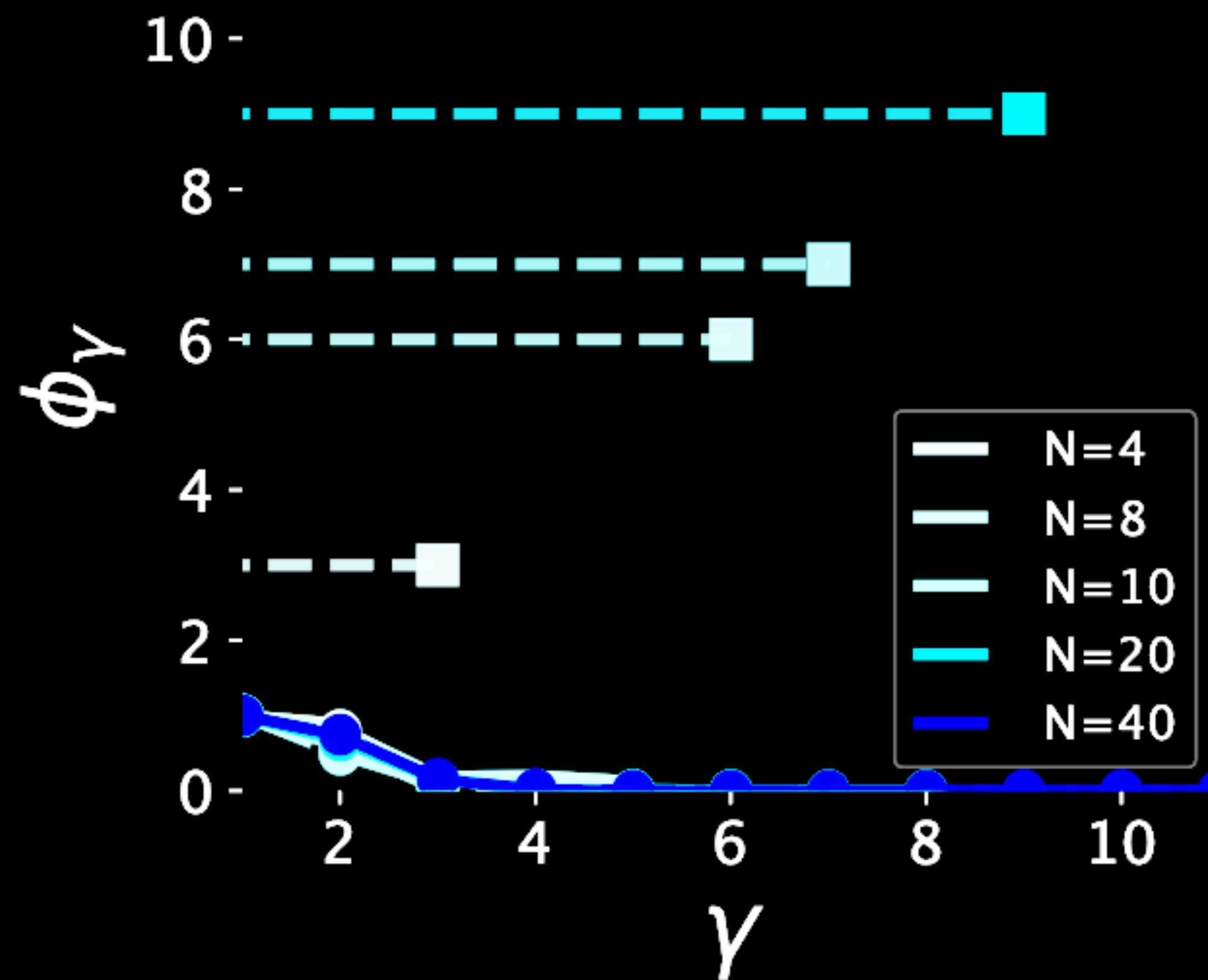


# Estimation of expected parallel capacity

Number of representations  
in task graph layers

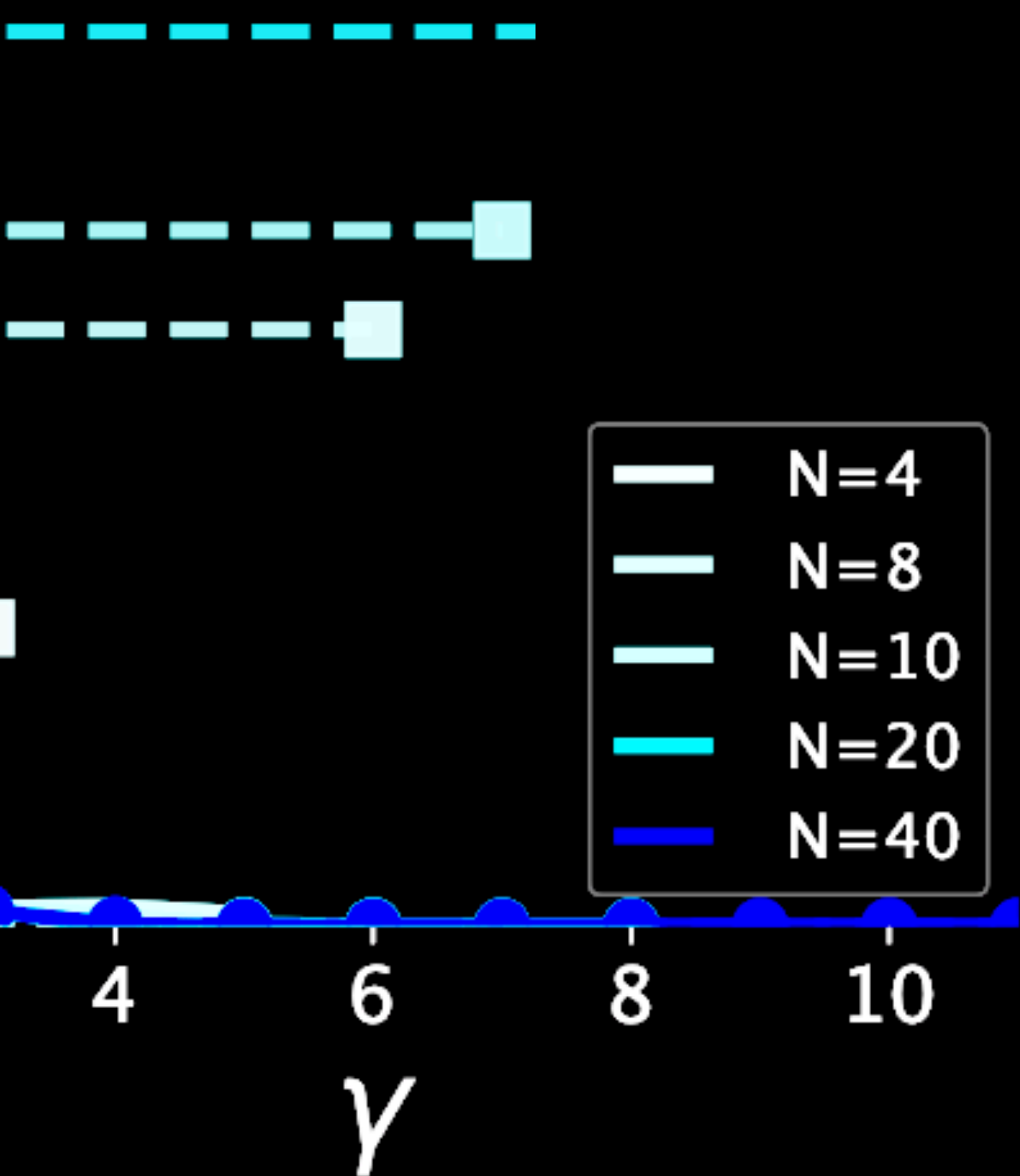


## Hard regime (all or nothing)

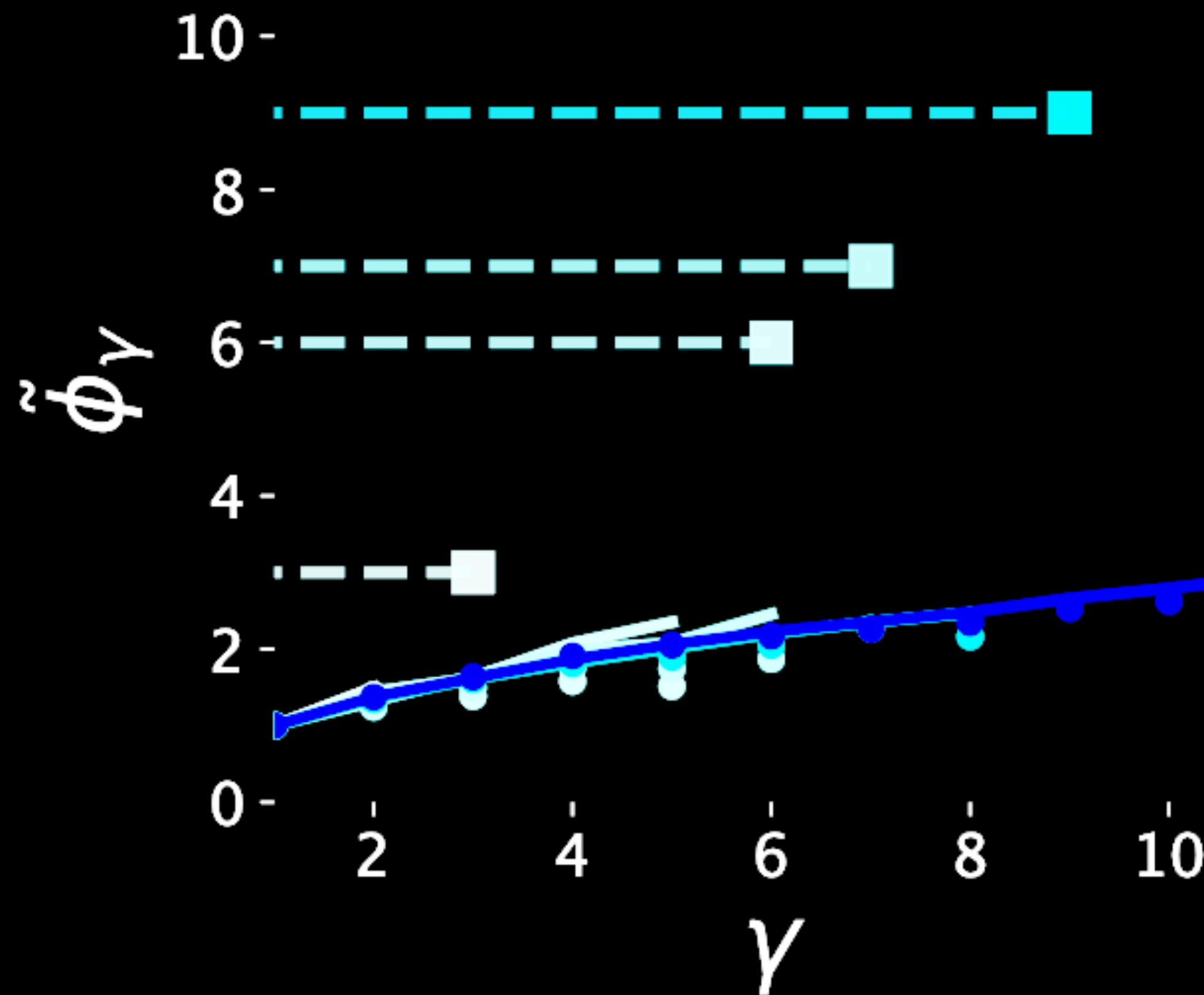


# Estimation of expected parallel capacity

Number of representations  
in task graph layers



Soft regime (one is enough)



# Take-home #1

# Take-home #1

## Results:

- we can extract dependency patterns from data
- severe limits to max. capacity under modest repr. sharing
- expected capacity is even worse

# Take-home #1

## Results:

- we can extract dependency patterns from data
- severe limits to max. capacity under modest repr. sharing
- expected capacity is even worse

A little too rigid...

No weights, simple networks,  
no control!!

# Can we extend this?

 eLife |  Home Browse Magazine Community About

[Neuroscience](#)

## **An Information-Theoretic Approach to Reward Rate Optimization in the Tradeoff Between Controlled and Automatic Processing in Neural Network Architectures**

Giovanni Petri , Sebastian Musslick, Jonathan D. Cohen

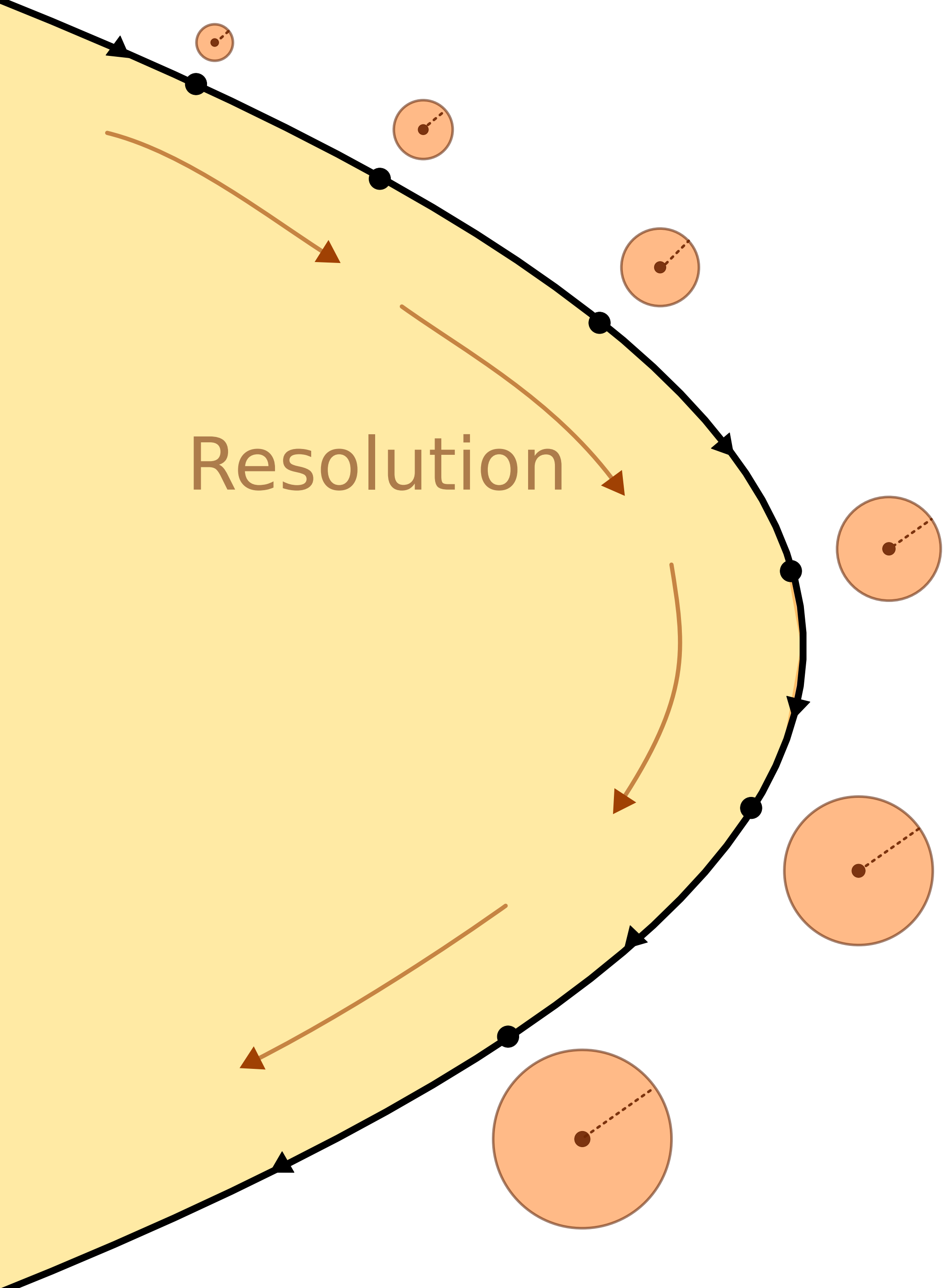
# Questions

- 1. Do the limitations imposed by shared representations prevail in a system as large as the brain?**
- 2. Assuming that shared representation cause a lot of trouble, why do we use them in the first place? (Experiment is evidence we do)**

# Questions

**1. Do the limitations imposed by shared representations prevail in a system as large as the brain?**

**2. Assuming that shared representation cause a lot of trouble, why do we use them in the first place? (Experiment is evidence we do)**



# Miller's Law



**Generalization**



**Parallel Capacity**

---

# Toward a Universal Law of Generalization for Psychological Science

ROGER N. SHEPARD

Generalization



Parallel Capacity

---

# Toward a Universal Law of Generalization for Psychological Science

ROGER N. SHEPARD

Generalization



Parallel Capacity

## THE PSYCHOLOGICAL REVIEW

---

THE MAGICAL NUMBER SEVEN, PLUS OR MINUS TWO:  
SOME LIMITS ON OUR CAPACITY FOR  
PROCESSING INFORMATION <sup>1</sup>

GEORGE A. MILLER

*Harvard University*

---

# **Toward a Universal Law of Generalization for Psychological Science**

ROGER N. SHEPARD

---

# Toward a Universal Law of Generalization for Psychological Science

ROGER N. SHEPARD

$$\mathcal{G}_{x\hat{x}} \triangleq \left( \frac{p_{x\hat{x}} \cdot p_{\hat{x}x}}{p_{\hat{x}\hat{x}} \cdot p_{xx}} \right)^{\frac{1}{2}}$$

---

# Toward a Universal Law of Generalization for Psychological Science

ROGER N. SHEPARD

$$\mathcal{G}_{x\hat{x}} \triangleq \left( \frac{p_{x\hat{x}} \cdot p_{\hat{x}x}}{p_{\hat{x}\hat{x}} \cdot p_{xx}} \right)^{\frac{1}{2}}$$

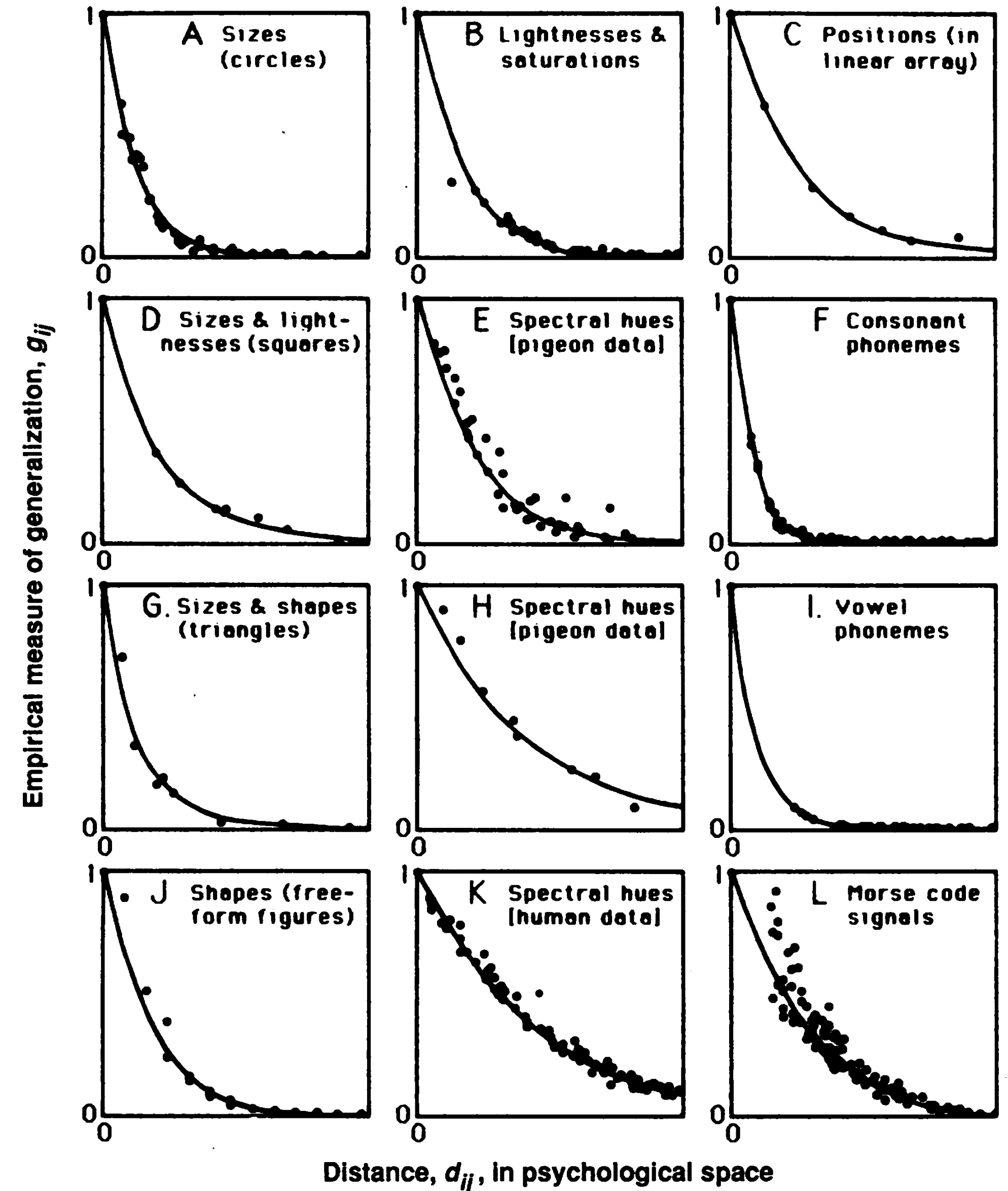
$$D_{ik}^R \cong -\log \left\{ (1 + C^R) \left( \frac{P_{ik}^R \cdot P_{ki}^R}{P_{ii}^R \cdot P_{kk}^R} \right)^{1/2} - C^R \right\}.$$

# Toward a Universal Law of Generalization for Psychological Science

ROGER N. SHEPARD

$$G_{x\hat{x}} \triangleq \left( \frac{P_{x\hat{x}} \cdot P_{\hat{x}x}}{P_{\hat{x}\hat{x}} \cdot P_{xx}} \right)^{\frac{1}{2}}$$

$$D_{ik}^R \cong -\log \left\{ (1 + C^R) \left( \frac{P_{ik}^R \cdot P_{ki}^R}{P_{ii}^R \cdot P_{kk}^R} \right)^{1/2} - C^R \right\}$$



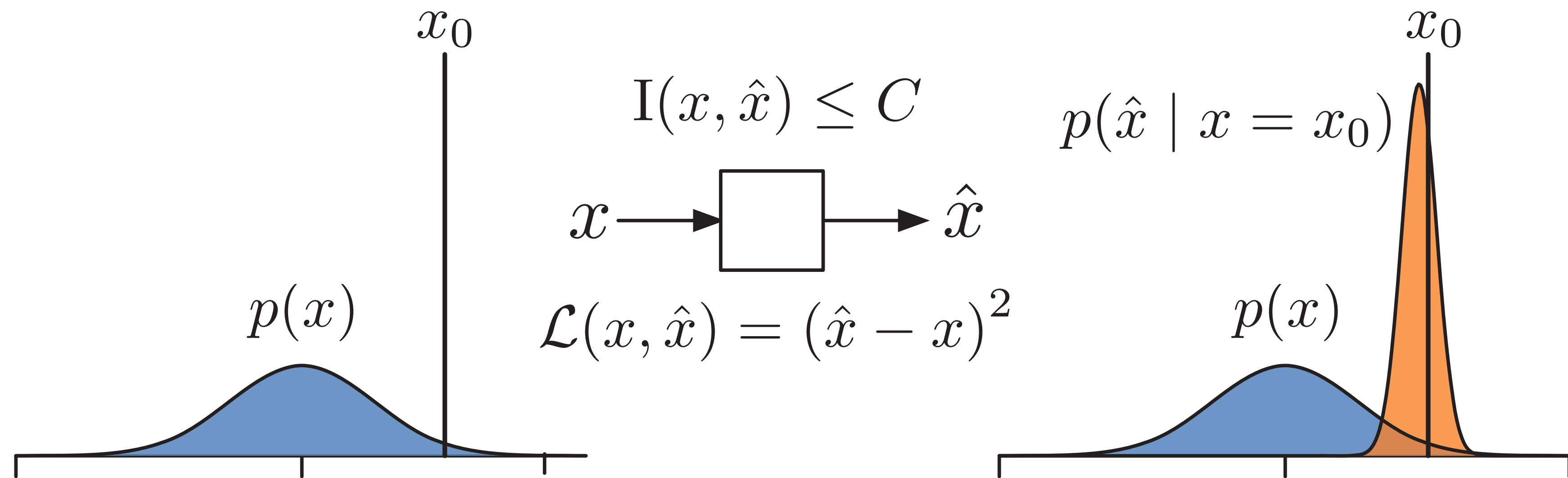
**COGNITIVE PSYCHOLOGY**

# **Efficient coding explains the universal law of generalization in human perception**

**Chris R. Sims\***

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*



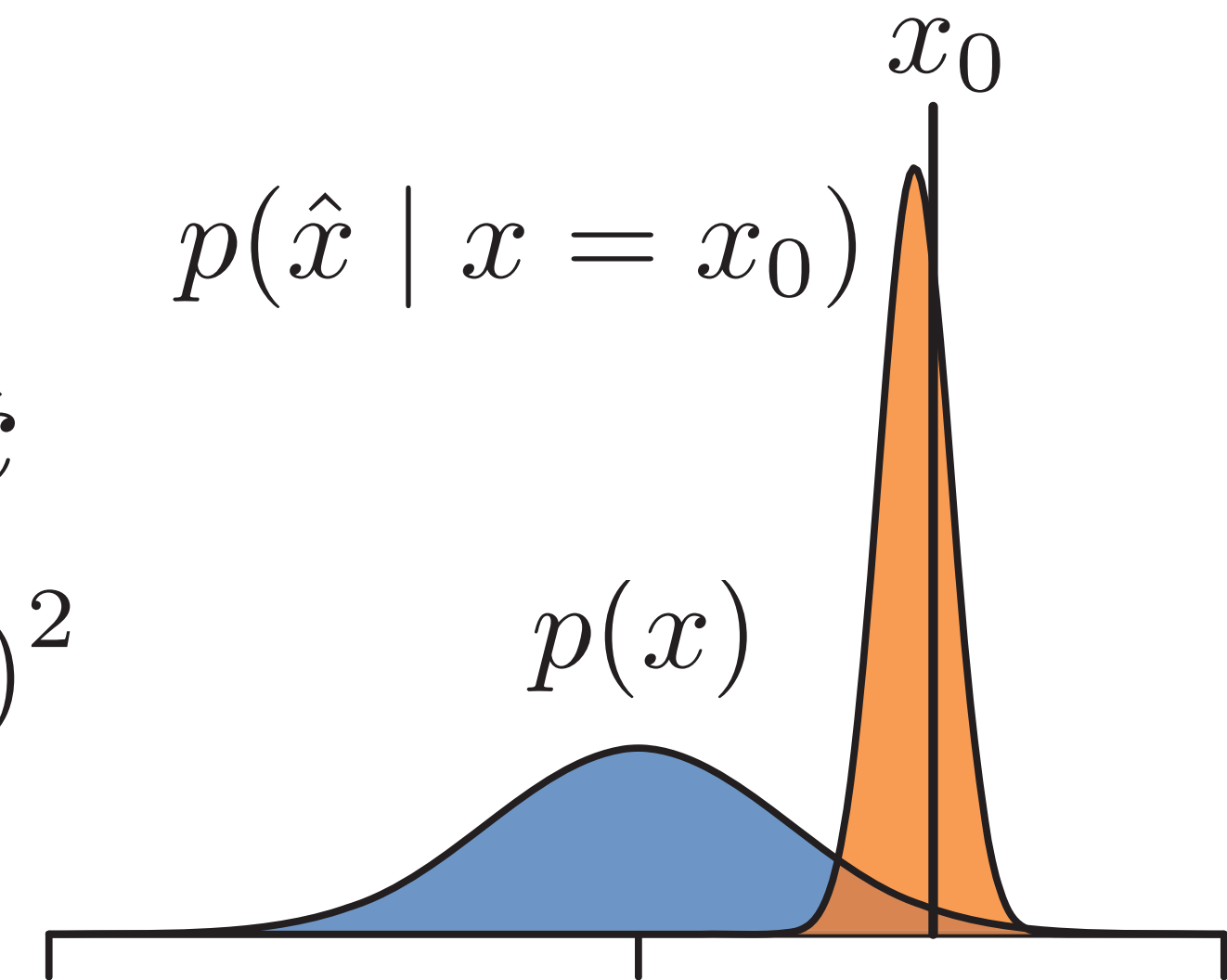
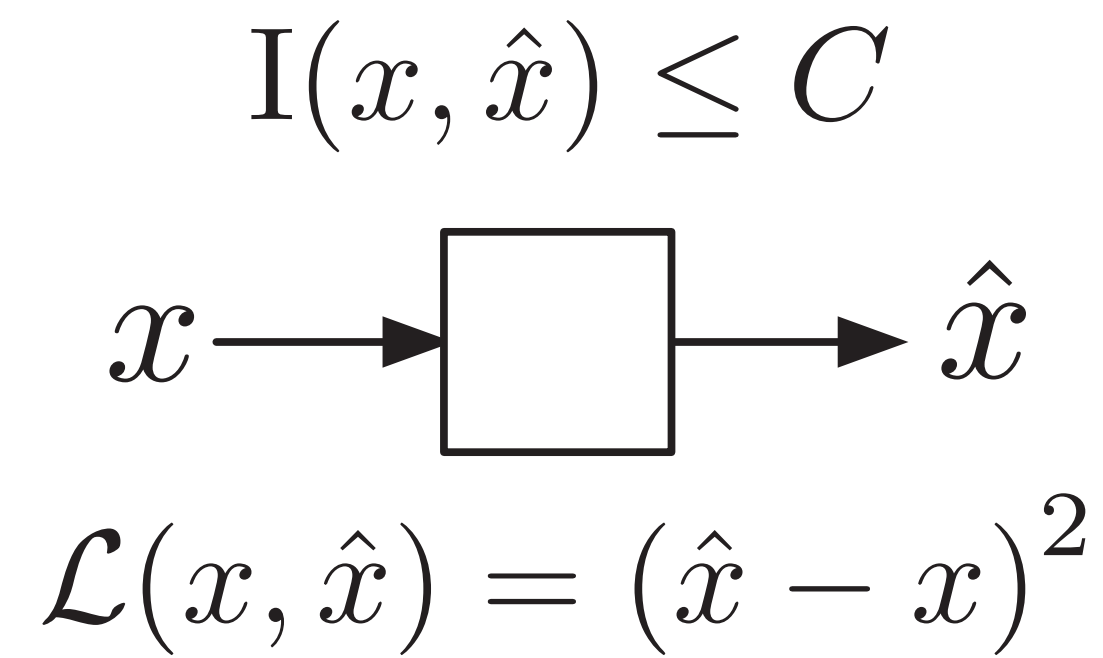
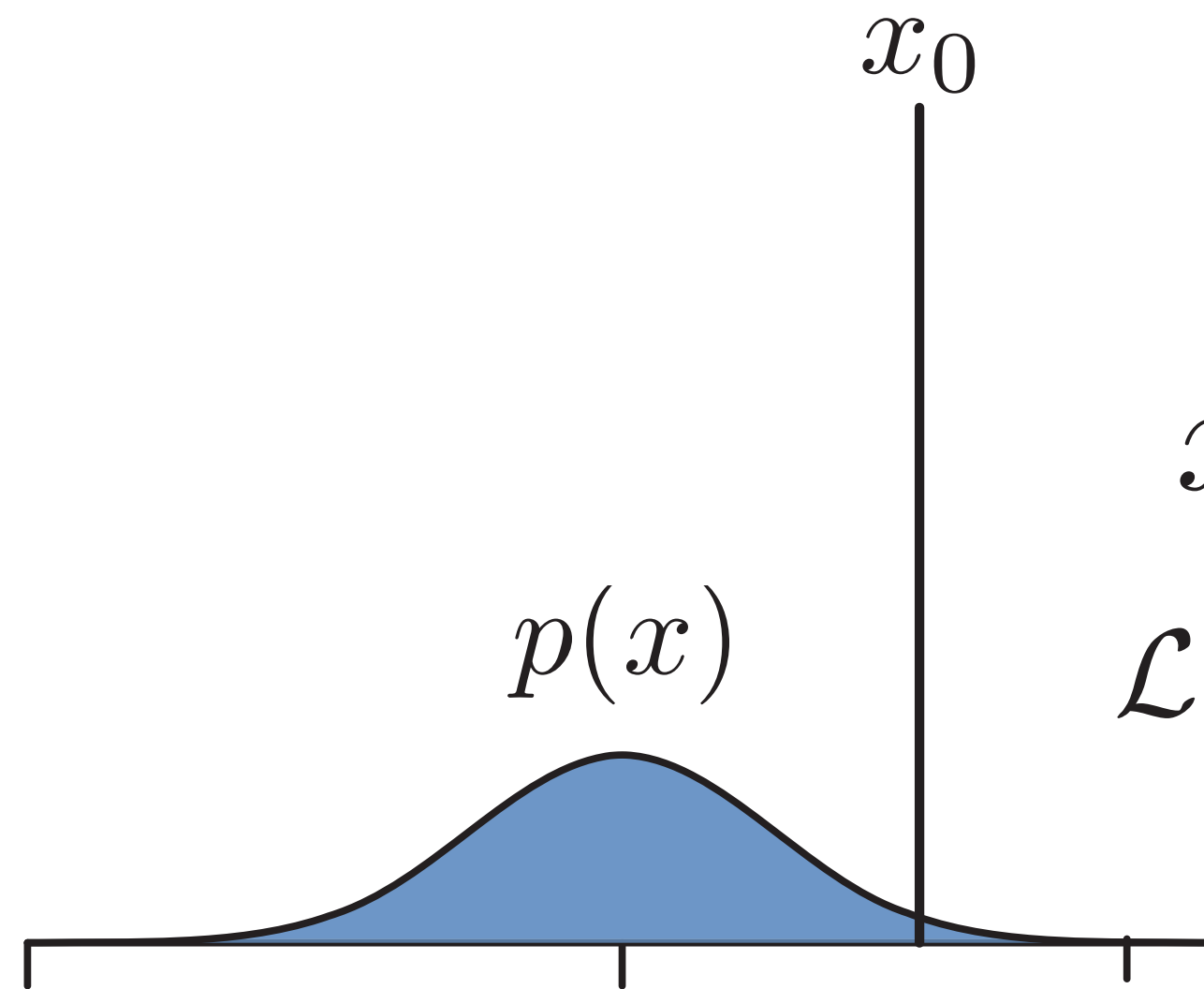
# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}$$

$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$



COGNITIVE PSYCHOLOGY

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}$$

$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}$$

$$\mathcal{G}_{x\hat{x}} \triangleq \left( \frac{\mathcal{P}_{x\hat{x}} \cdot \mathcal{P}_{\hat{x}x}}{\mathcal{P}_{\hat{x}\hat{x}} \cdot \mathcal{P}_{xx}} \right)^{\frac{1}{2}}$$

$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}, \quad \mathcal{G}_{x\hat{x}} \triangleq \left( \frac{\mathcal{P}_{x\hat{x}} \cdot \mathcal{P}_{\hat{x}x}}{\mathcal{P}_{\hat{x}\hat{x}} \cdot \mathcal{P}_{xx}} \right)^{\frac{1}{2}}$$

$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$

$$\mathcal{G}_{x\hat{x}} = \exp \left[ s \frac{1}{2} \left( \mathcal{L}(x, \hat{x}) + \mathcal{L}(\hat{x}, x) - \mathcal{L}(x, x) - \mathcal{L}(\hat{x}, \hat{x}) \right) \right]$$

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}, \quad \mathcal{G}_{x\hat{x}} \triangleq \left( \frac{p_{x\hat{x}} \cdot p_{\hat{x}x}}{p_{\hat{x}\hat{x}} \cdot p_{xx}} \right)^{\frac{1}{2}}$$

$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$

$$\mathcal{G}_{x\hat{x}} = \exp \left[ s \frac{1}{2} \left( \mathcal{L}(x, \hat{x}) + \mathcal{L}(\hat{x}, x) - \mathcal{L}(x, x) - \mathcal{L}(\hat{x}, \hat{x}) \right) \right] \quad \mathcal{G}_{x\hat{x}} = \exp[s \mathcal{L}(x, \hat{x})]$$

COGNITIVE PSYCHOLOGY

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}$$

$$\mathcal{G}_{x\hat{x}} \triangleq \left( \frac{p_{x\hat{x}} \cdot p_{\hat{x}x}}{p_{\hat{x}\hat{x}} \cdot p_{xx}} \right)^{\frac{1}{2}}$$

$$p(\hat{x}) = \int p(x)p(\hat{x} | x)dx.$$

$$\mathcal{G}_{x\hat{x}} = \exp \left[ s \frac{1}{2} \left( \mathcal{L}(x, \hat{x}) + \mathcal{L}(\hat{x}, x) - \mathcal{L}(x, x) - \mathcal{L}(\hat{x}, \hat{x}) \right) \right]$$

$$\mathcal{G}_{x\hat{x}} = \exp[s \mathcal{L}(x, \hat{x})]$$



this can be computed iteratively  
(Blahut-Arimoto)

COGNITIVE PSYCHOLOGY

# Efficient coding explains the universal law of generalization in human perception

Chris R. Sims\*

RDT

$$p(\hat{x} | x) = \frac{p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x}))}{\int p(\hat{x}) \exp(s \mathcal{L}(x, \hat{x})) d\hat{x}}$$

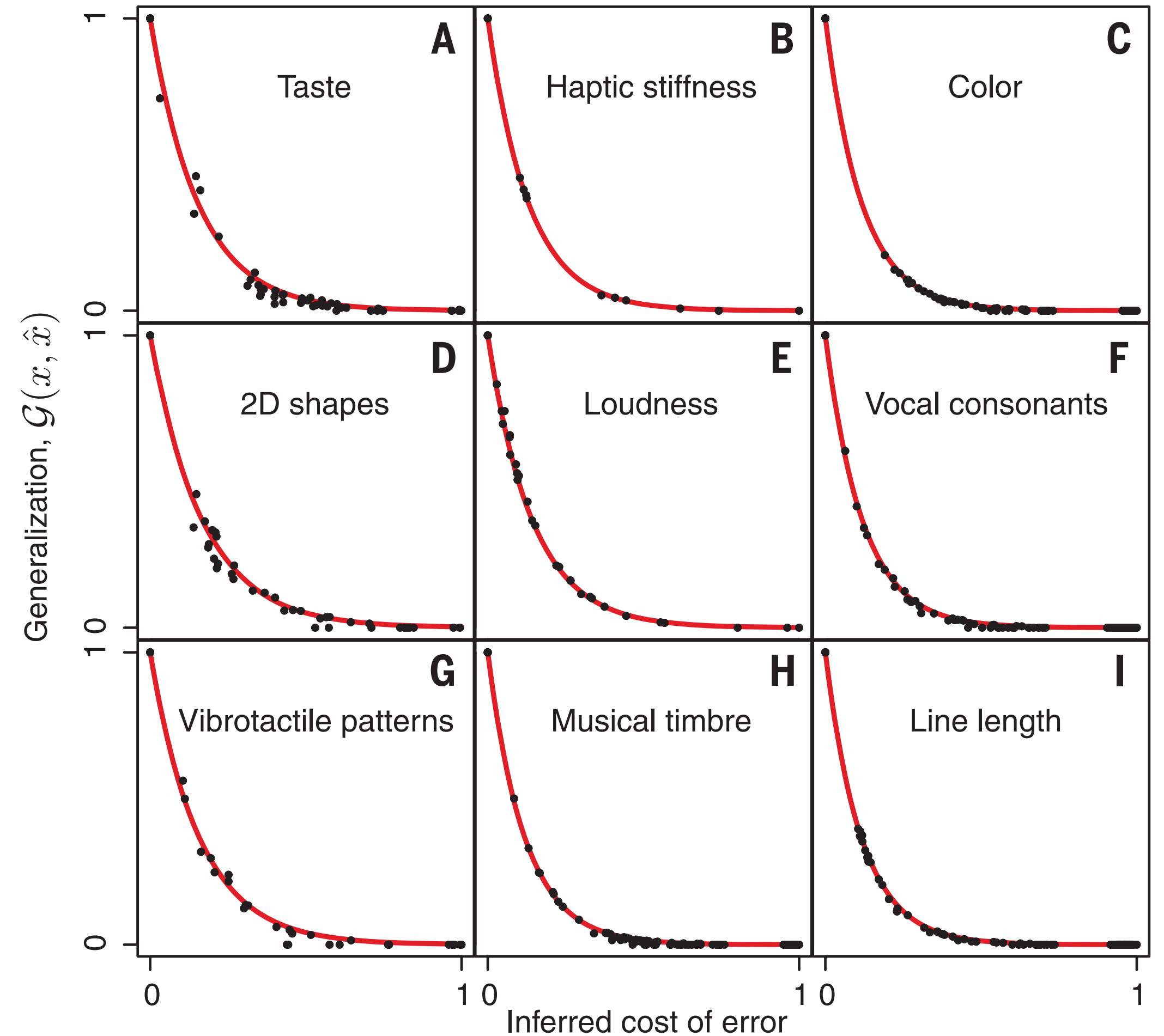
$$p(\hat{x}) = \int p(x) p(\hat{x} | x) dx.$$

$$\mathcal{G}_{x\hat{x}} \triangleq \left( \frac{p_{x\hat{x}} \cdot p_{\hat{x}x}}{p_{\hat{x}\hat{x}} \cdot p_{xx}} \right)^{\frac{1}{2}}$$

$$\mathcal{G}_{x\hat{x}} = \exp \left[ s \frac{1}{2} \left( \mathcal{L}(x, \hat{x}) + \mathcal{L}(\hat{x}, x) - \mathcal{L}(x, x) - \mathcal{L}(\hat{x}, \hat{x}) \right) \right]$$

$$\mathcal{G}_{x\hat{x}} = \exp[s \mathcal{L}(x, \hat{x})]$$

this can be computed iteratively (Blahut-Arimoto)



# THE PSYCHOLOGICAL REVIEW

---

THE MAGICAL NUMBER SEVEN, PLUS OR MINUS TWO:  
SOME LIMITS ON OUR CAPACITY FOR  
PROCESSING INFORMATION <sup>1</sup>

GEORGE A. MILLER

*Harvard University*

# THE PSYCHOLOGICAL REVIEW

---

**THE MAGICAL NUMBER SEVEN, PLUS OR MINUS TWO:  
SOME LIMITS ON OUR CAPACITY FOR  
PROCESSING INFORMATION <sup>1</sup>**

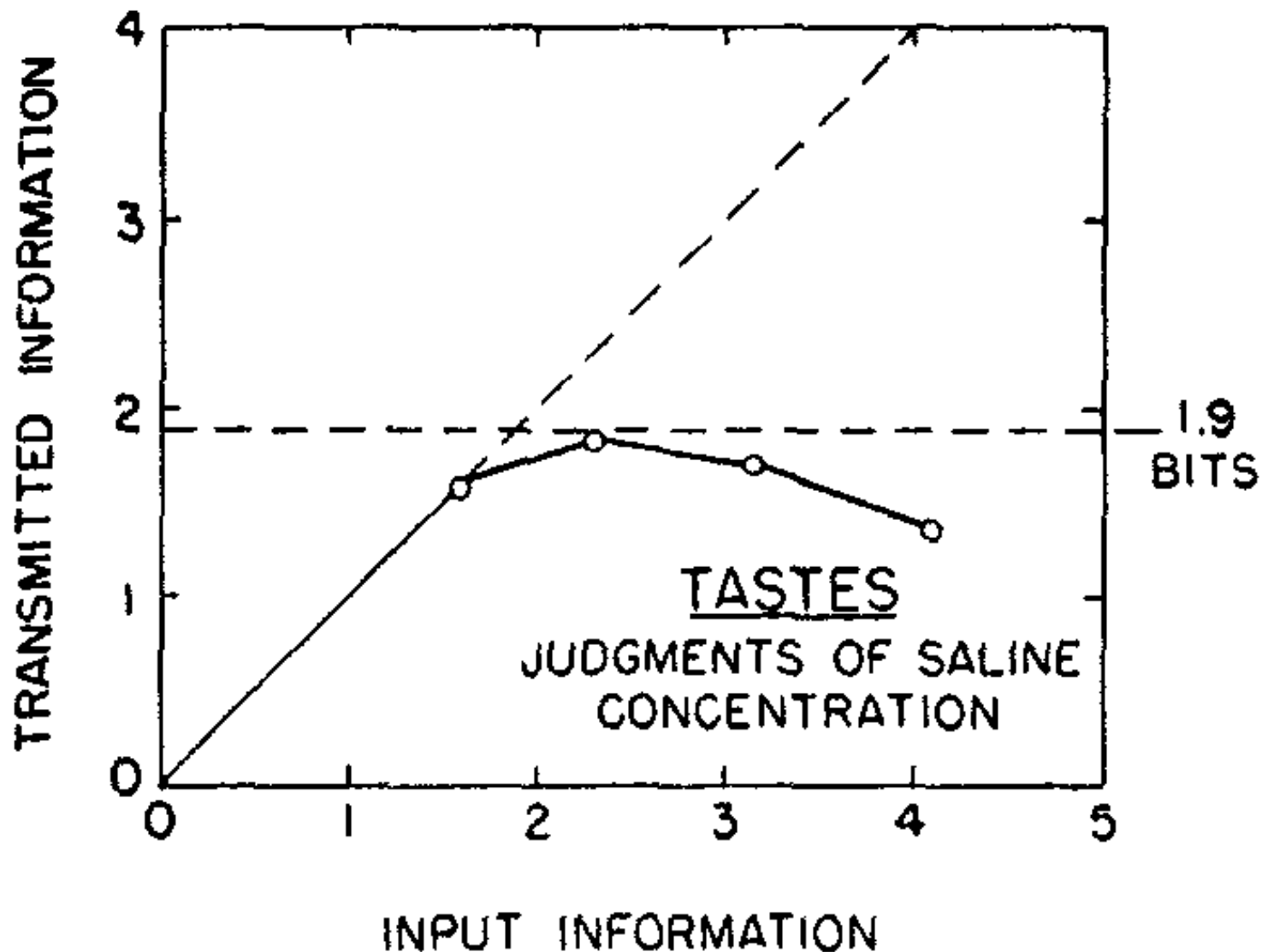
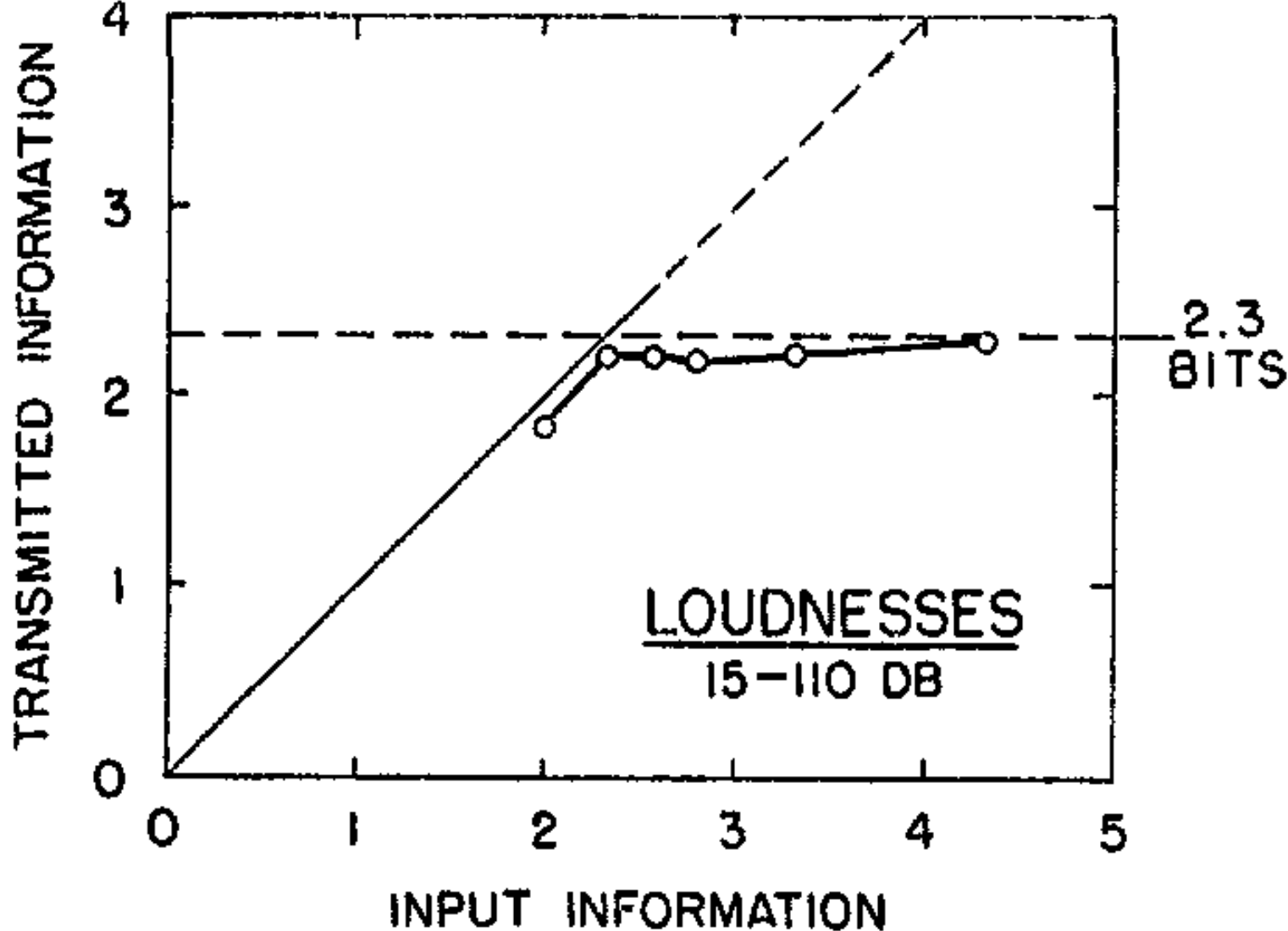
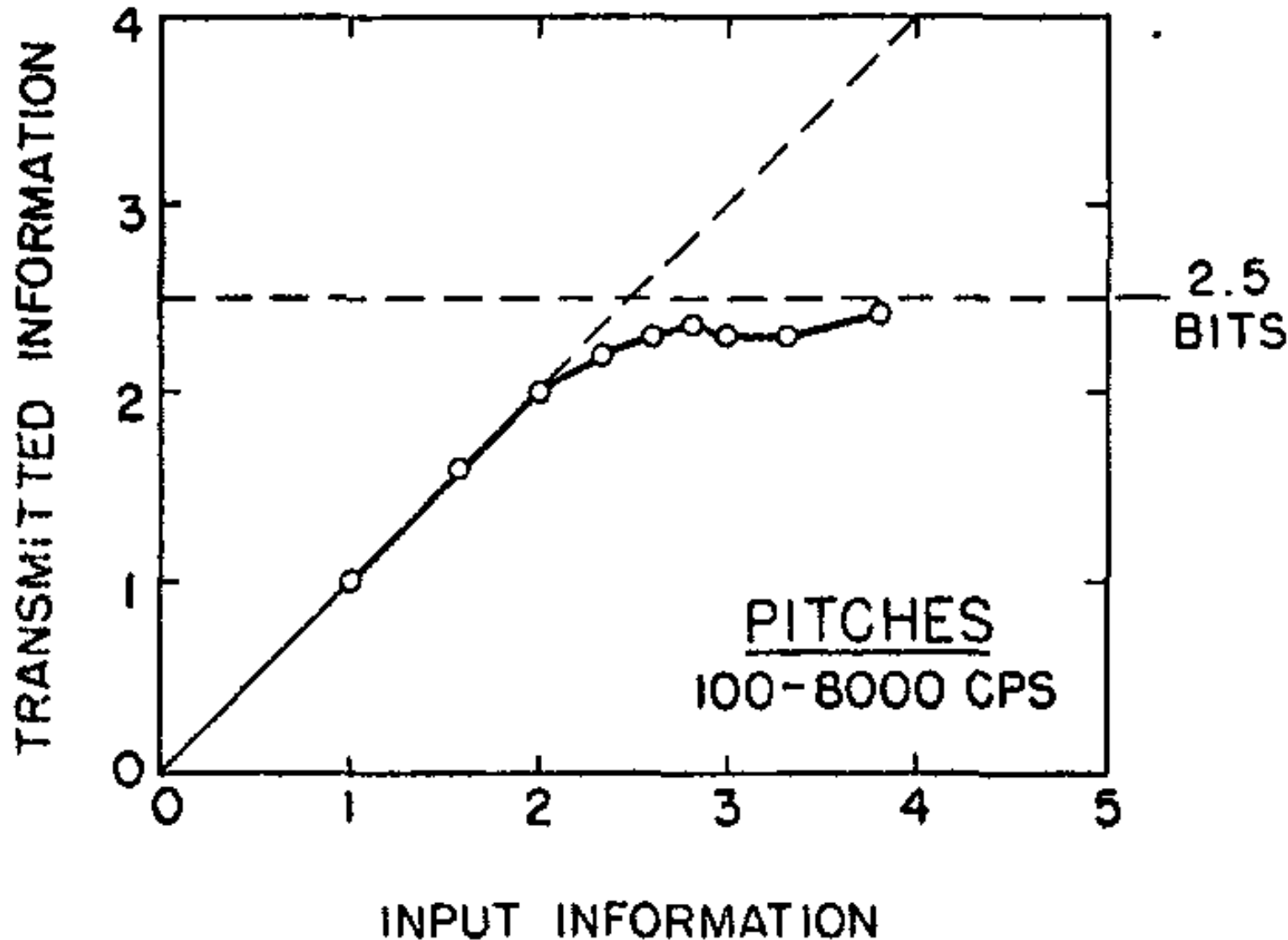
**GEORGE A. MILLER**

*Harvard University*

# THE PSYCHOLOGICAL REVIEW

## THE MAGICAL NUMBER SEVEN, PLUS OR MINUS TWO: SOME LIMITS ON OUR CAPACITY FOR PROCESSING INFORMATION <sup>1</sup>

GEORGE A. MILLER  
*Harvard University*

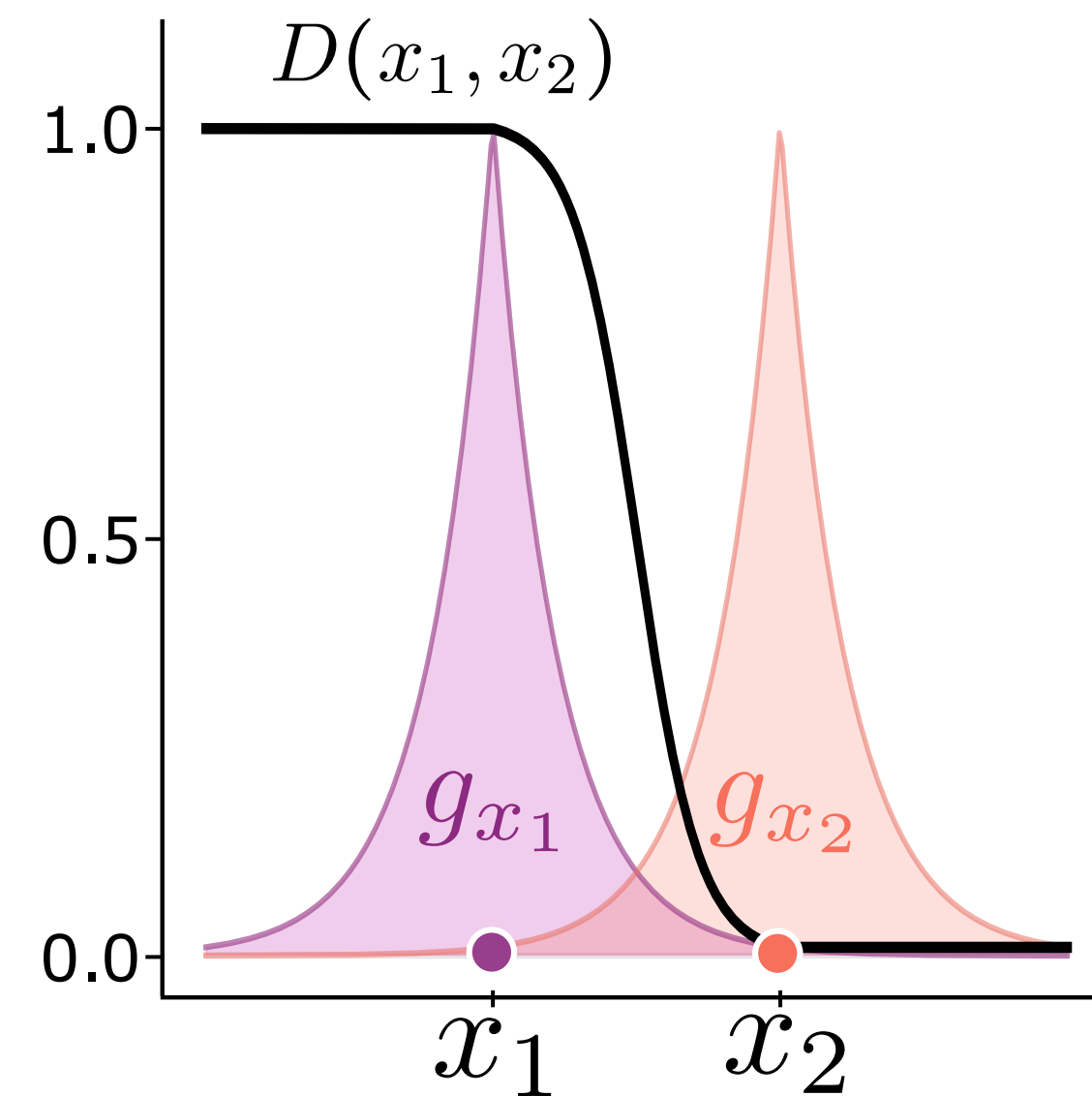


# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

**a** No resolution

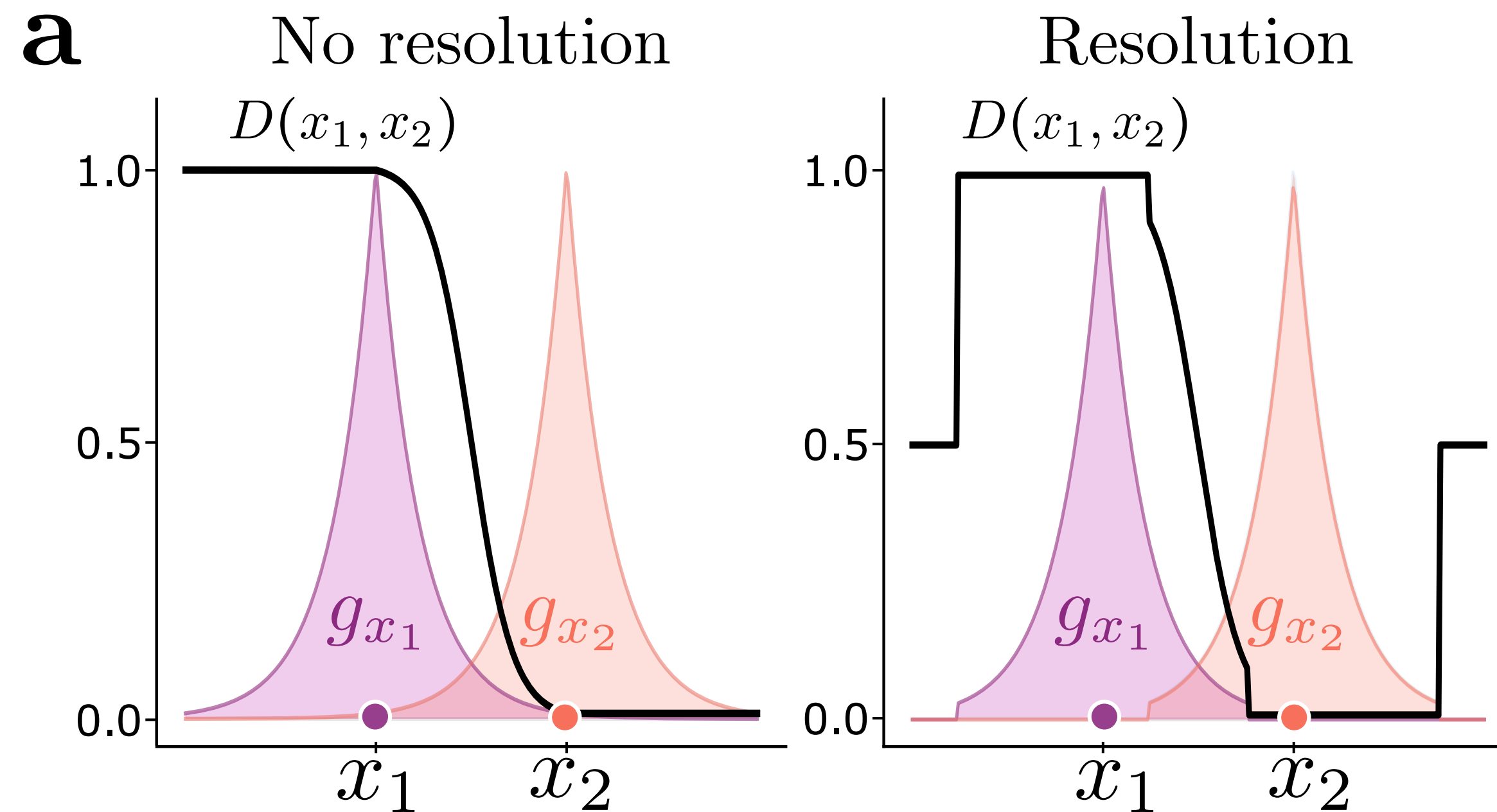


$$D(Y = x_1 | p) = \frac{G(x_1, p)}{G(x_1, p) + G(x_2, p)}$$

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

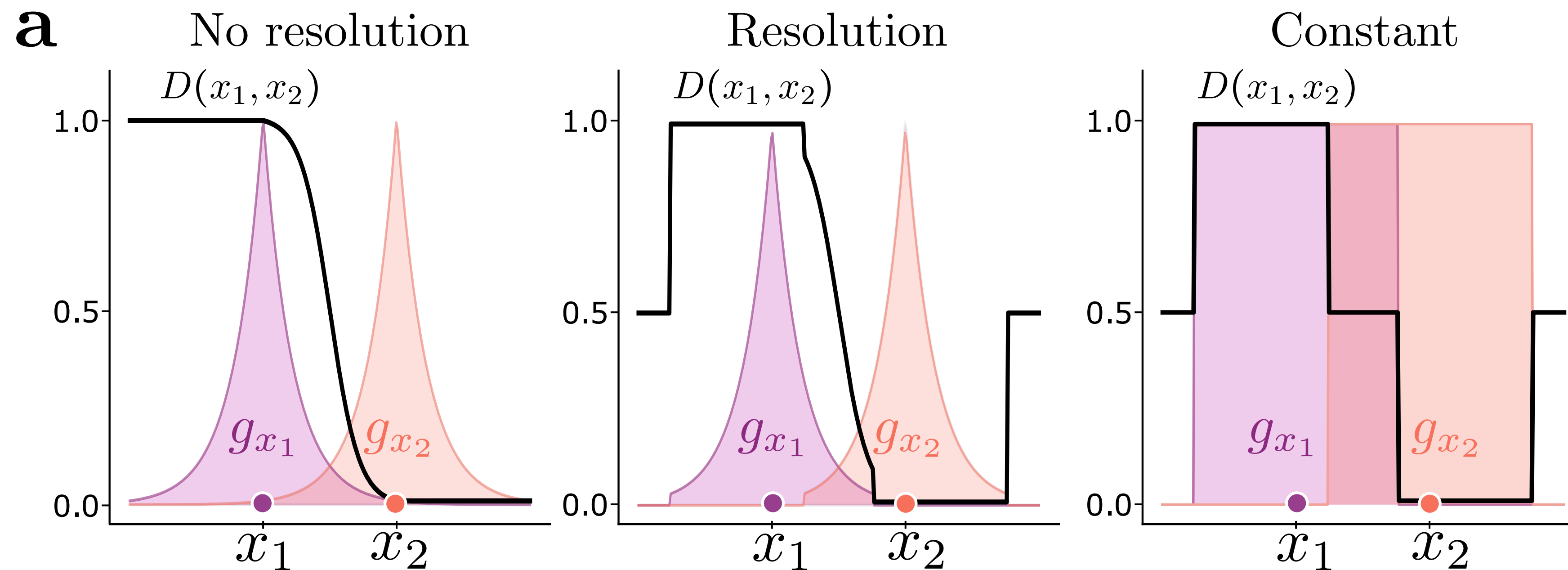


$$D(Y = x_1 | p) = \frac{G(x_1, p)}{G(x_1, p) + G(x_2, p)}$$

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

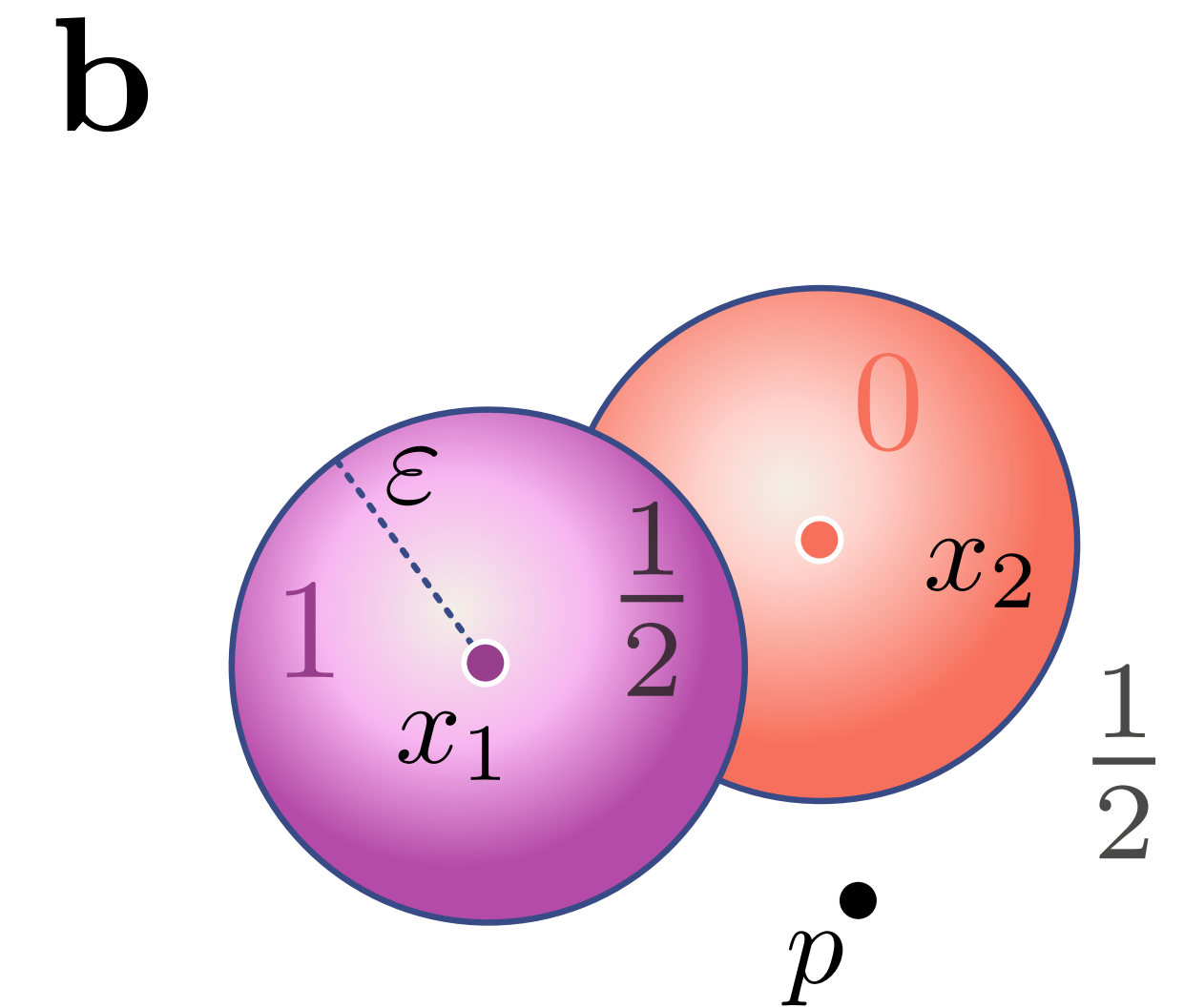
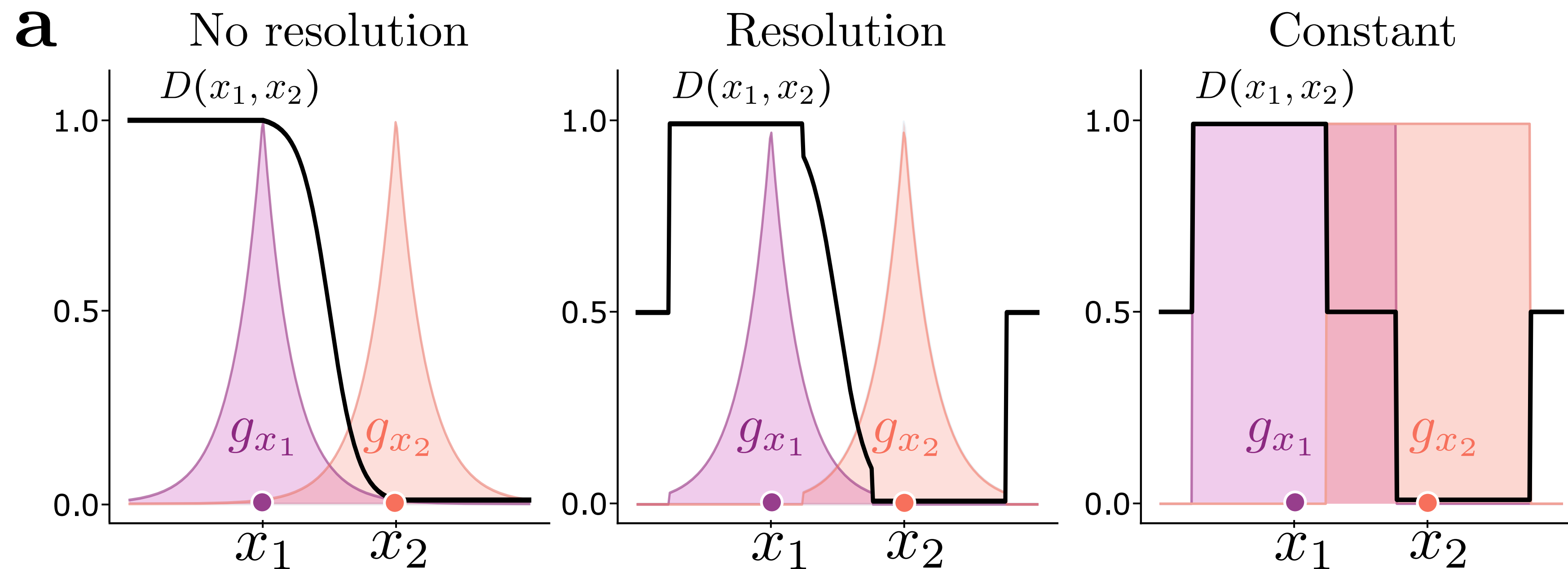


$$D(Y = x_1 | p) = \frac{G(x_1, p)}{G(x_1, p) + G(x_2, p)}$$

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



$$D(Y = x_1 | p) = \frac{G(x_1, p)}{G(x_1, p) + G(x_2, p)}$$

# Generalization vs Processing

---

**Bound by semanticity: universal laws governing the  
generalization-identification tradeoff**

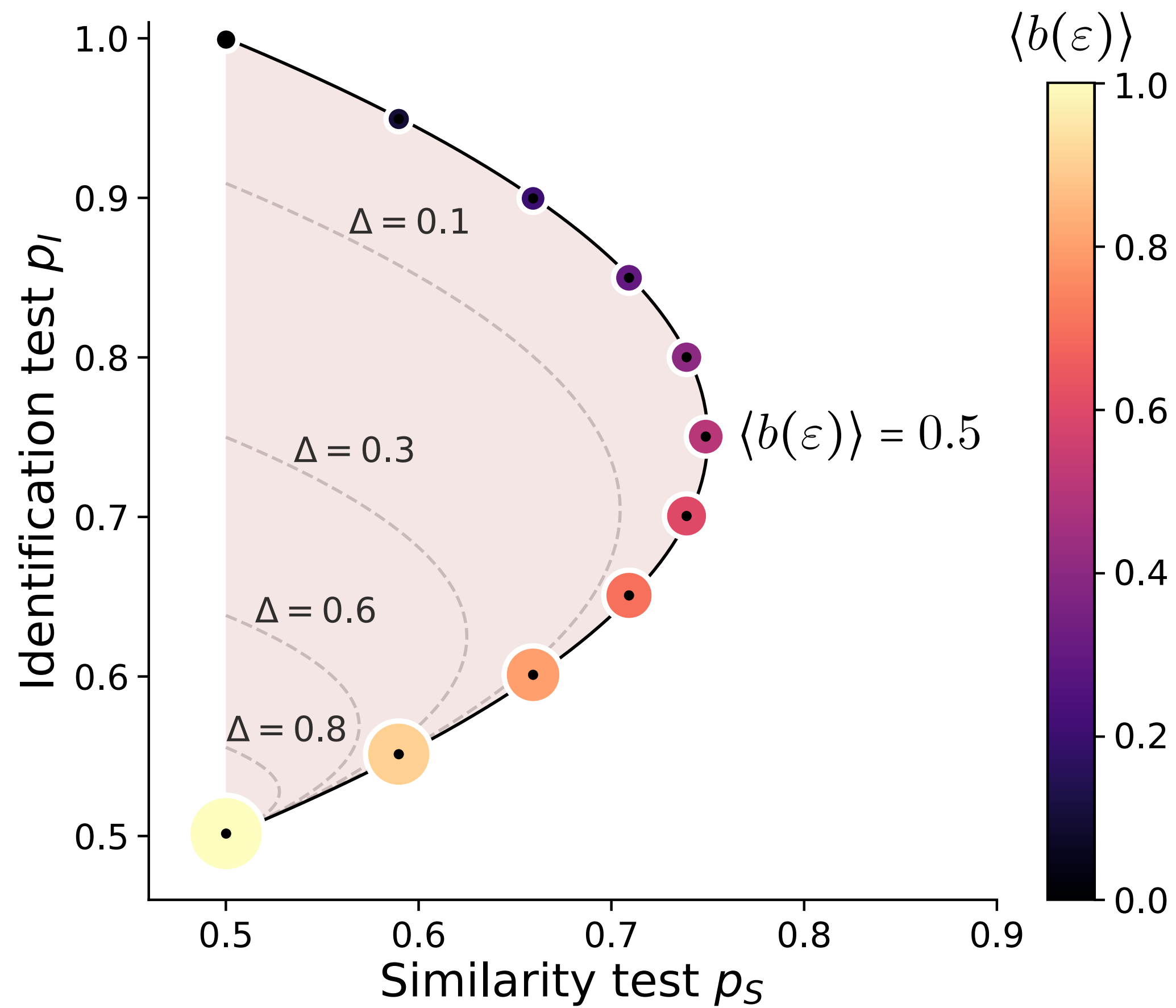
---

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

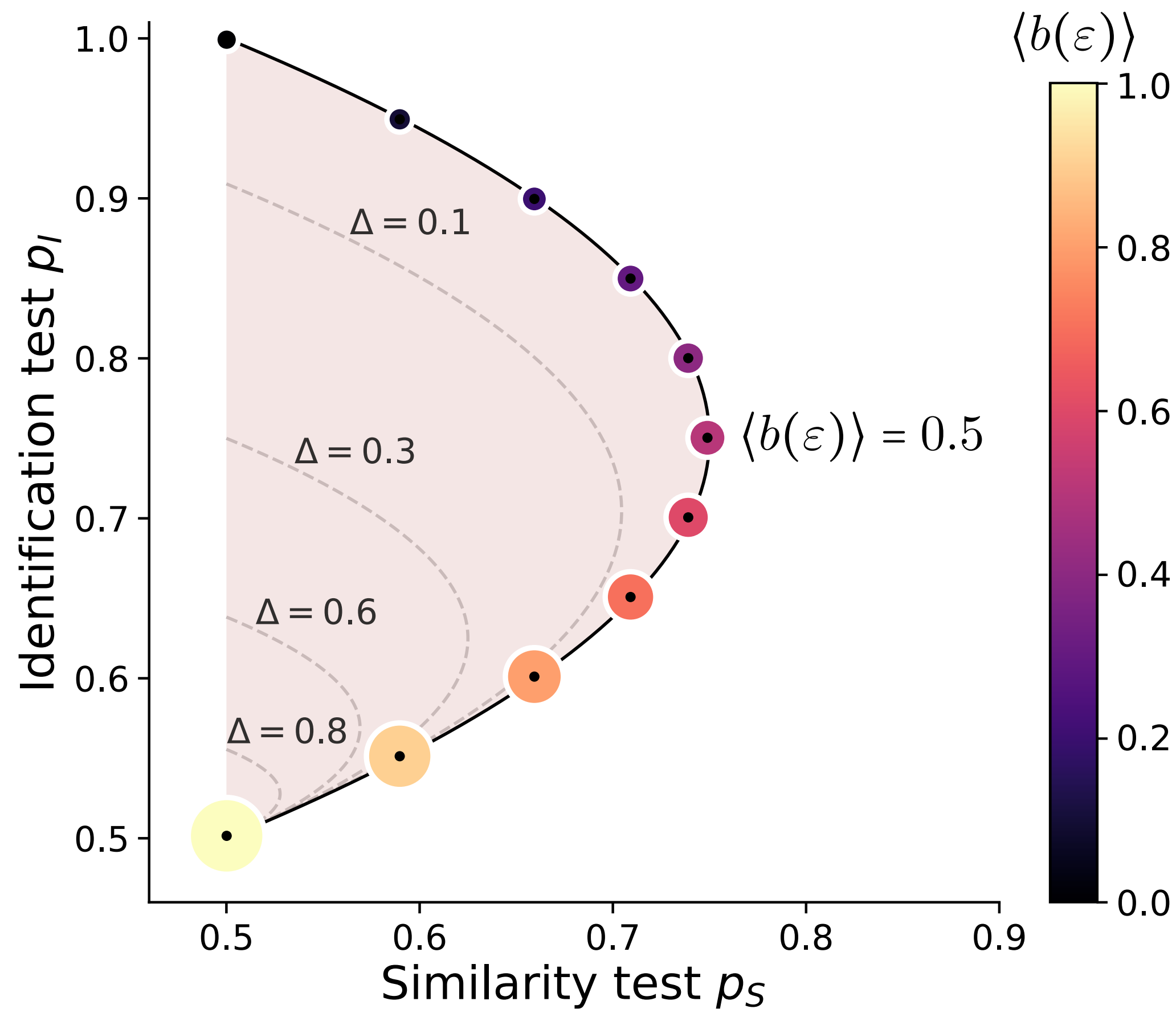
Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



**Theorem 1** (2-item tests). *Let  $(M, d, \Sigma, \nu)$  be a separable metric probability space. If, for every  $p \in M$ ,  $b_p$  is absolutely continuous on every closed sub-interval of  $[0, \infty)$ , then, for the noise-free constant similarity function  $g = g_{\varepsilon;0}$  it holds that*

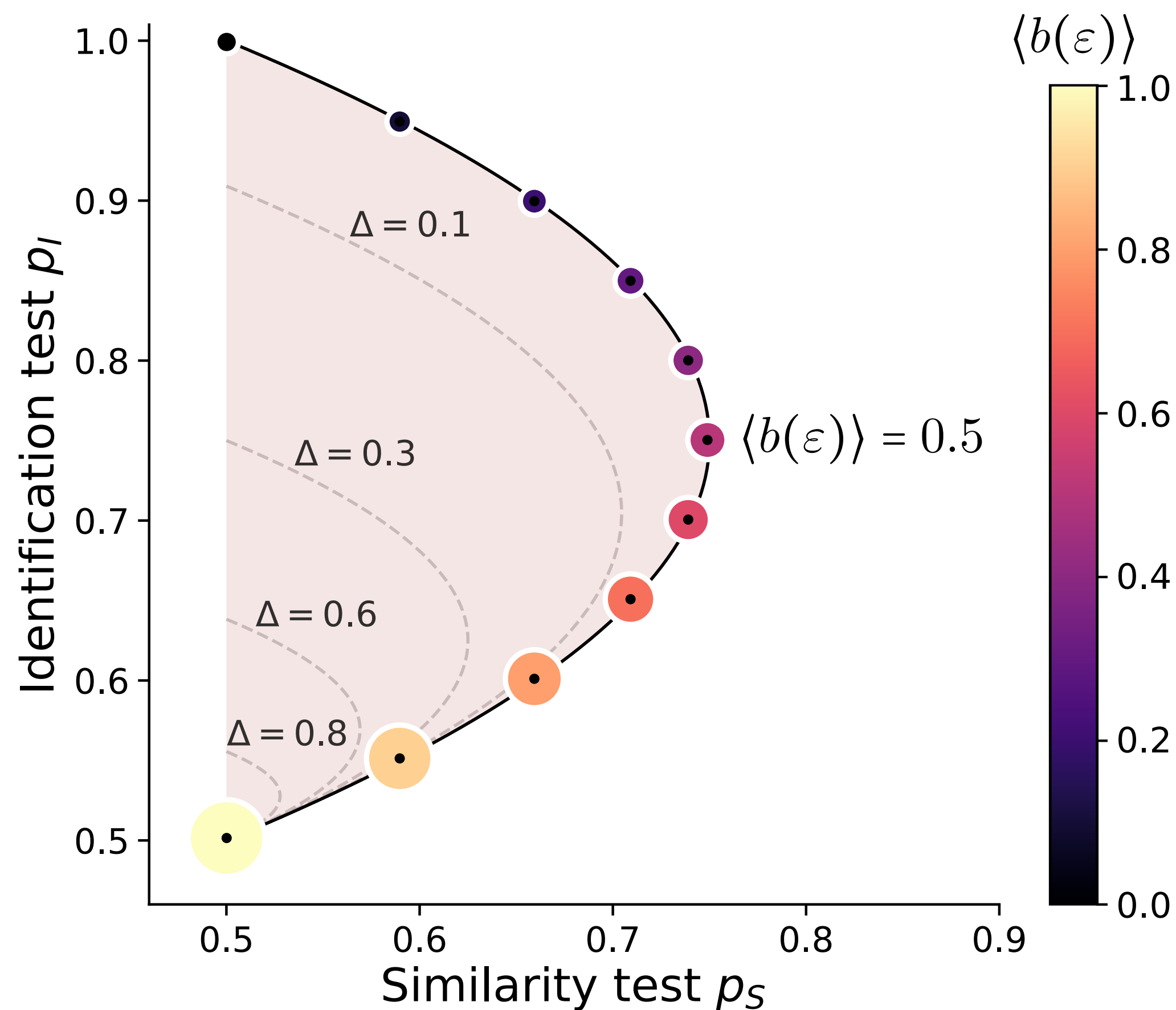
$$p_S(\varepsilon) = \frac{1}{2} + \langle b(\varepsilon) \rangle - \langle b(\varepsilon) \rangle^2 - \text{Var}(b(\varepsilon)), \quad (3)$$

$$p_I(\varepsilon) = 1 - \frac{1}{2} \langle b(\varepsilon) \rangle. \quad (4)$$

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

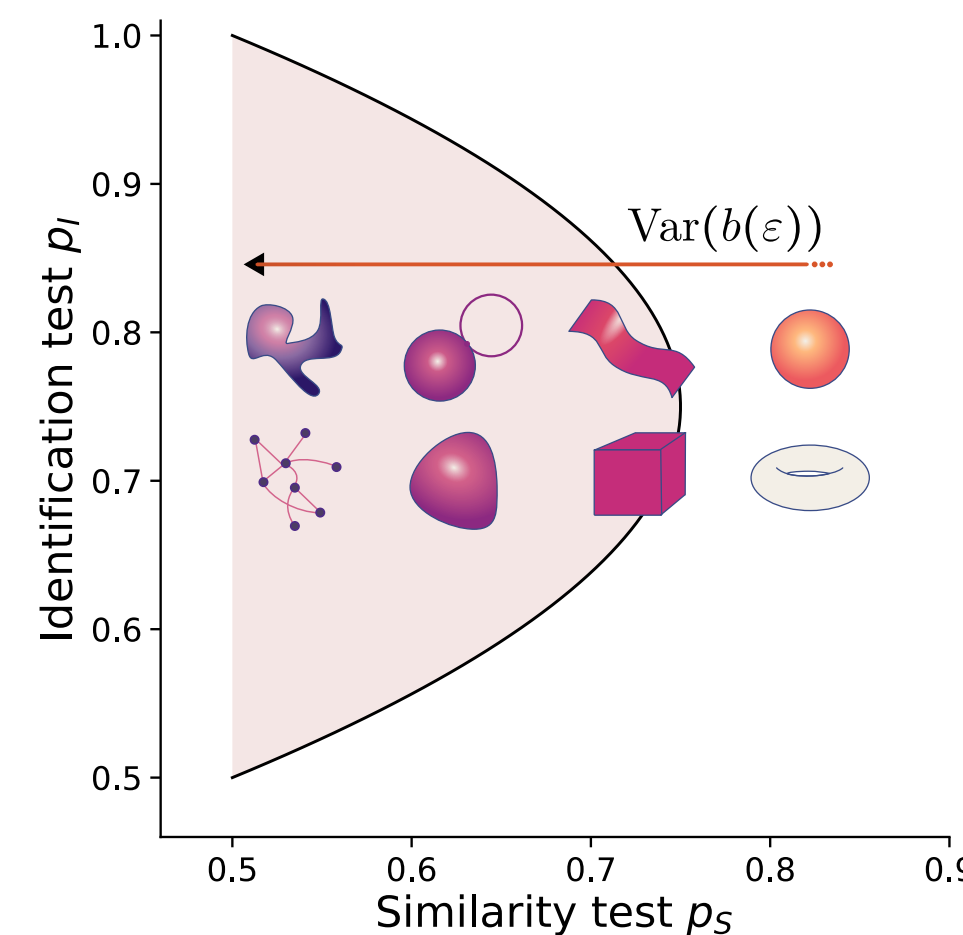
Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



**Theorem 1** (2-item tests). *Let  $(M, d, \Sigma, \nu)$  be a separable metric probability space. If, for every  $p \in M$ ,  $b_p$  is absolutely continuous on every closed sub-interval of  $[0, \infty)$ , then, for the noise-free constant similarity function  $g = g_{\epsilon;0}$  it holds that*

$$p_S(\epsilon) = \frac{1}{2} + \langle b(\epsilon) \rangle - \langle b(\epsilon) \rangle^2 - \text{Var}(b(\epsilon)), \quad (3)$$

$$p_I(\epsilon) = 1 - \frac{1}{2} \langle b(\epsilon) \rangle. \quad (4)$$



**Theorem 2** (Noise). *Under the same assumptions of Theorem 1, for the two-item similarity and identification tests with constant similarity functions  $g = g_{\epsilon;\Delta}$  with noise level  $\Delta \geq 0$  it holds that*

$$p_S(\epsilon, \Delta) = \frac{1}{2} + \frac{1 - \Delta}{1 + \Delta} (\langle b(\epsilon) \rangle - \langle b(\epsilon)^2 \rangle), \quad (5)$$

$$p_I(\epsilon, \Delta) = \frac{2 - (1 - \Delta) \langle b(\epsilon) \rangle}{2 + 2\Delta}. \quad (6)$$

# Generalization vs Processing

Multi-item identification decays!

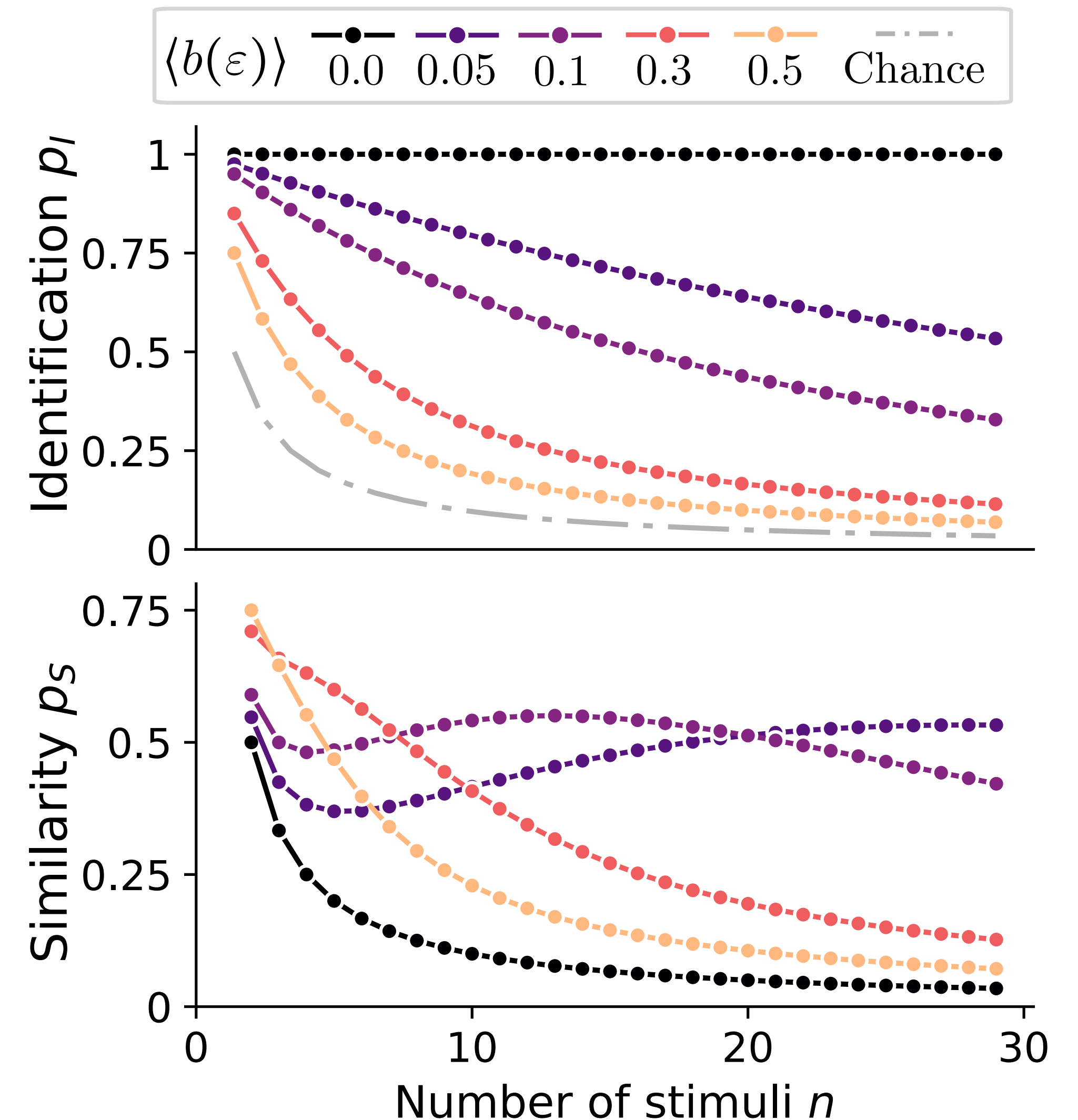
Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

**Theorem 3** ( $n$ -item tests). Under the same assumptions of Theorem 1, for the constant noise-free ( $\Delta = 0$ ) similarity function  $g = g_{\varepsilon;0}$  we have that

$$p_S^n(\varepsilon) = \mathbb{E}_{p \sim \nu} \left[ \frac{1}{n} + \sum_{k=1}^{n-1} \frac{(1 - b_p(\varepsilon))^{n-k} - (1 - b_p(\varepsilon))^n}{k} \right], \quad (7)$$

$$p_I^n(\varepsilon) = \mathbb{E}_{p \sim \nu} \left[ \frac{1 - (1 - b_p(\varepsilon))^n}{n b_p(\varepsilon)} \right]. \quad (8)$$

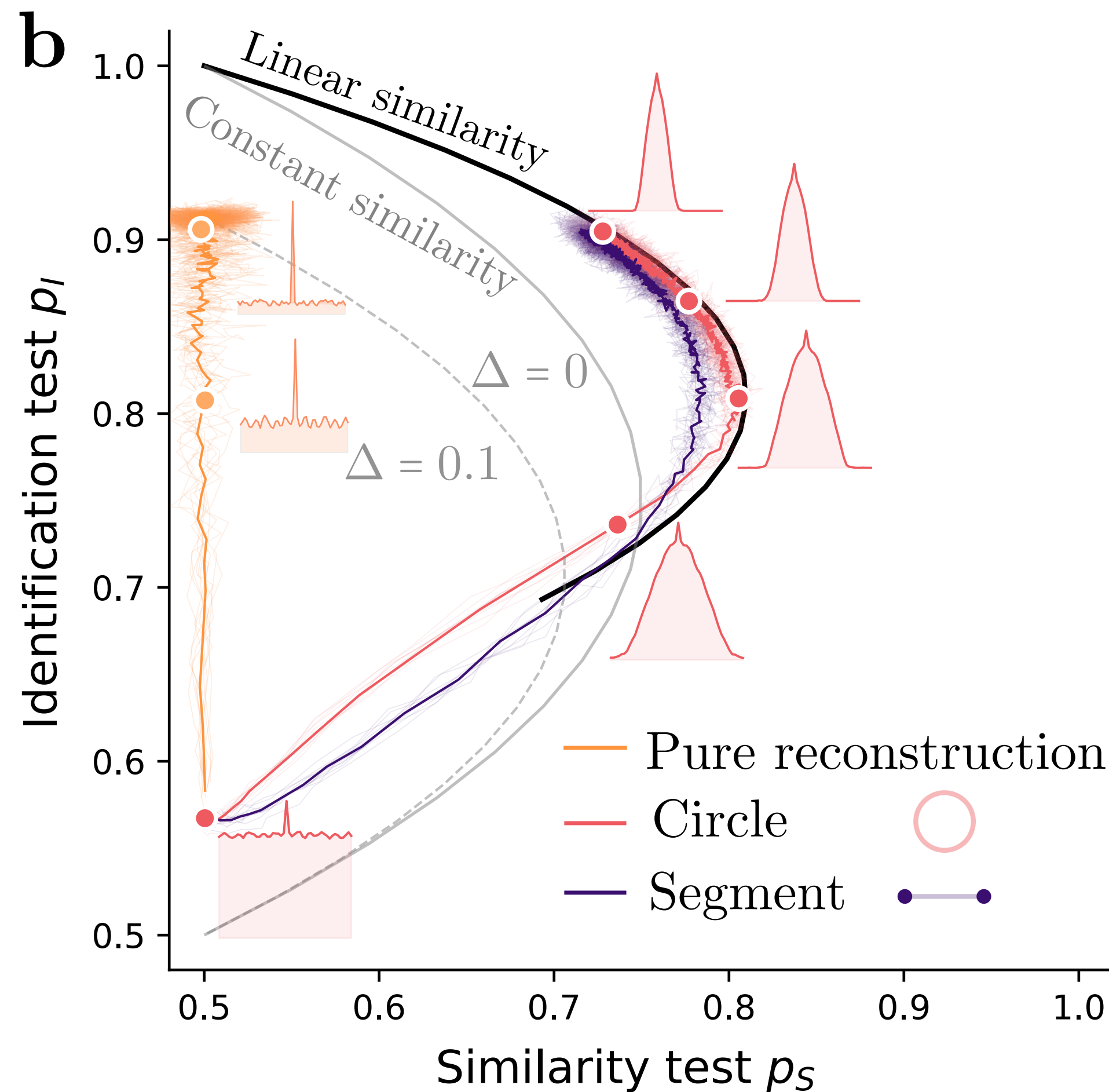
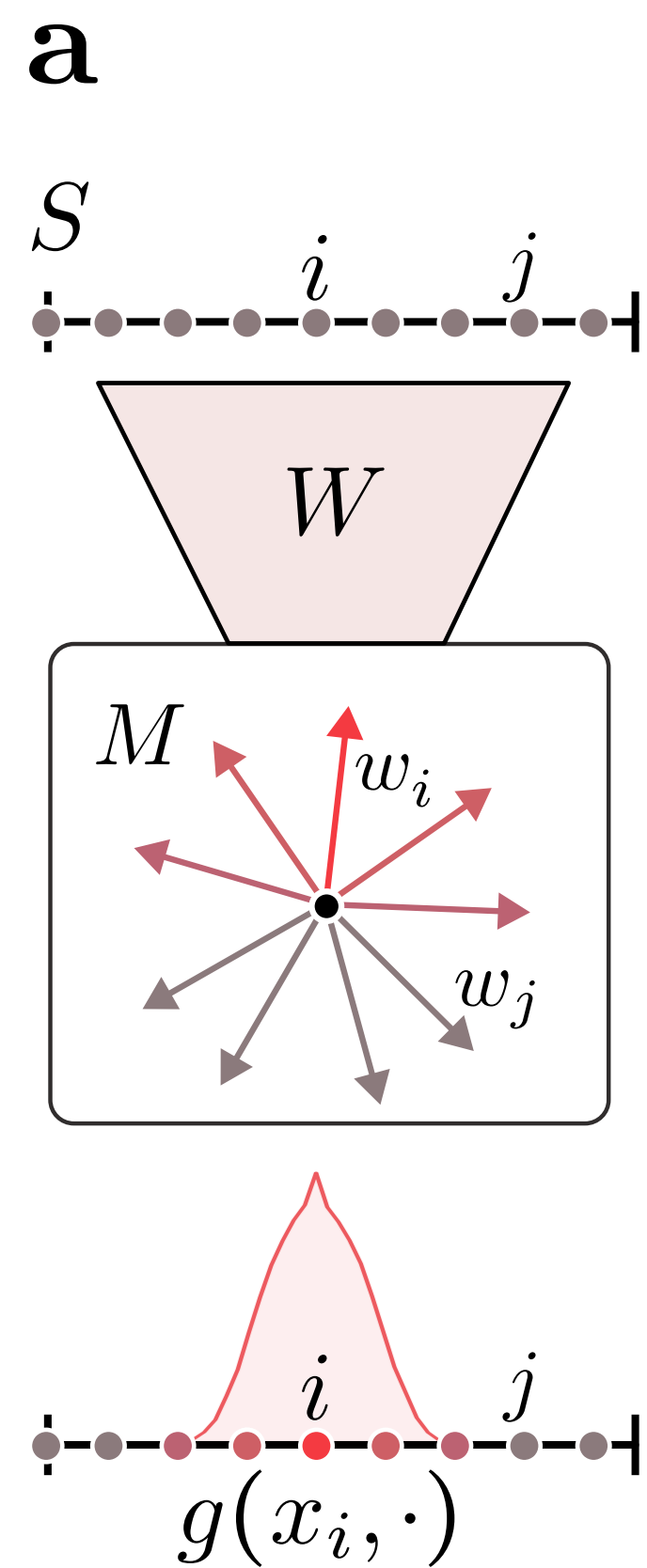


Is there an emergent resolution?

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



$$f(x) = \sigma(W^\top W x),$$

$$L_{\text{rec}} = \sum_{i=1}^l \|e_i - \sigma(W^\top W e_i)\|^2 = \sum_{i=1}^l \|e_i - \sigma(W^\top w_i)\|^2.$$

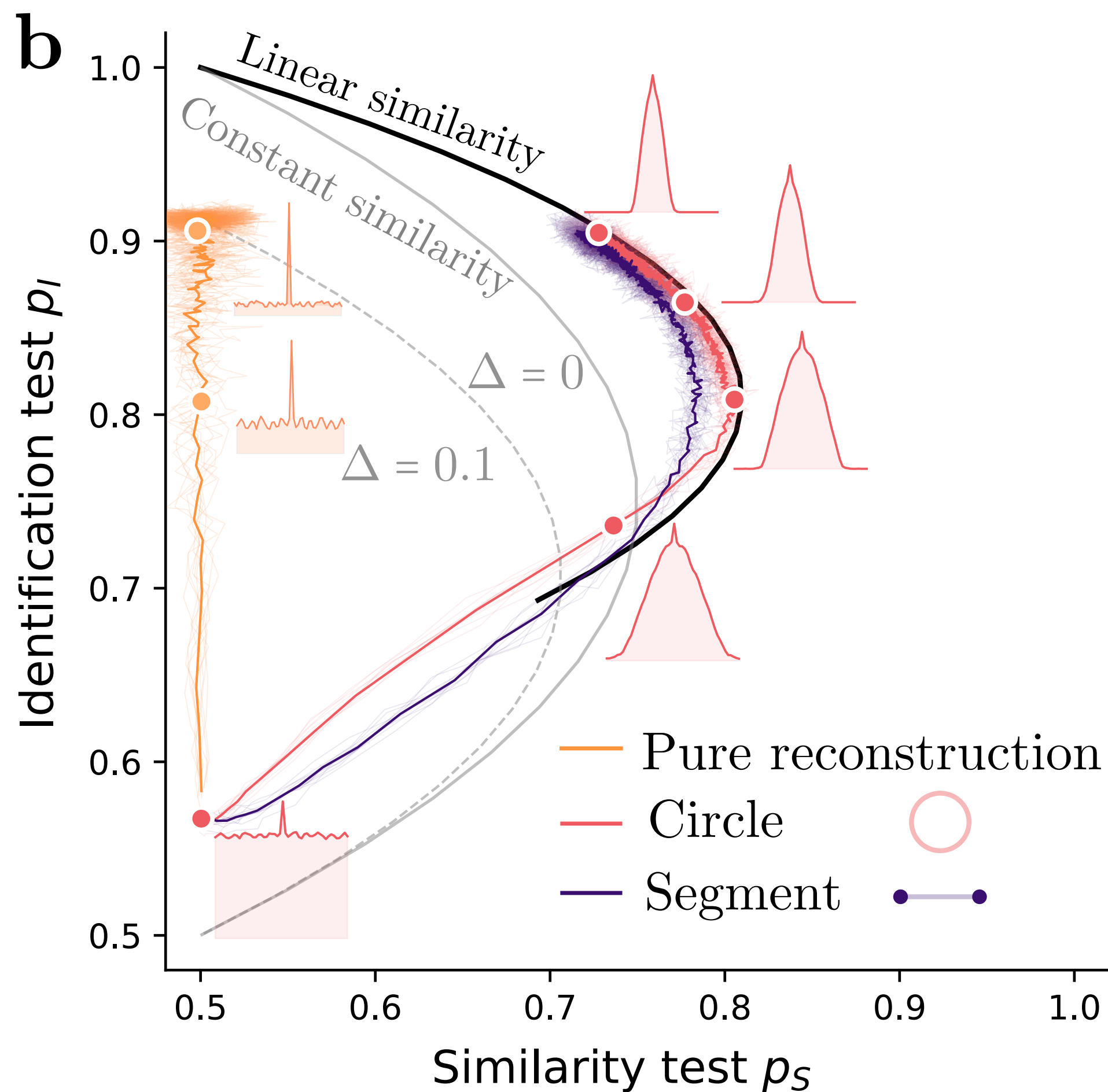
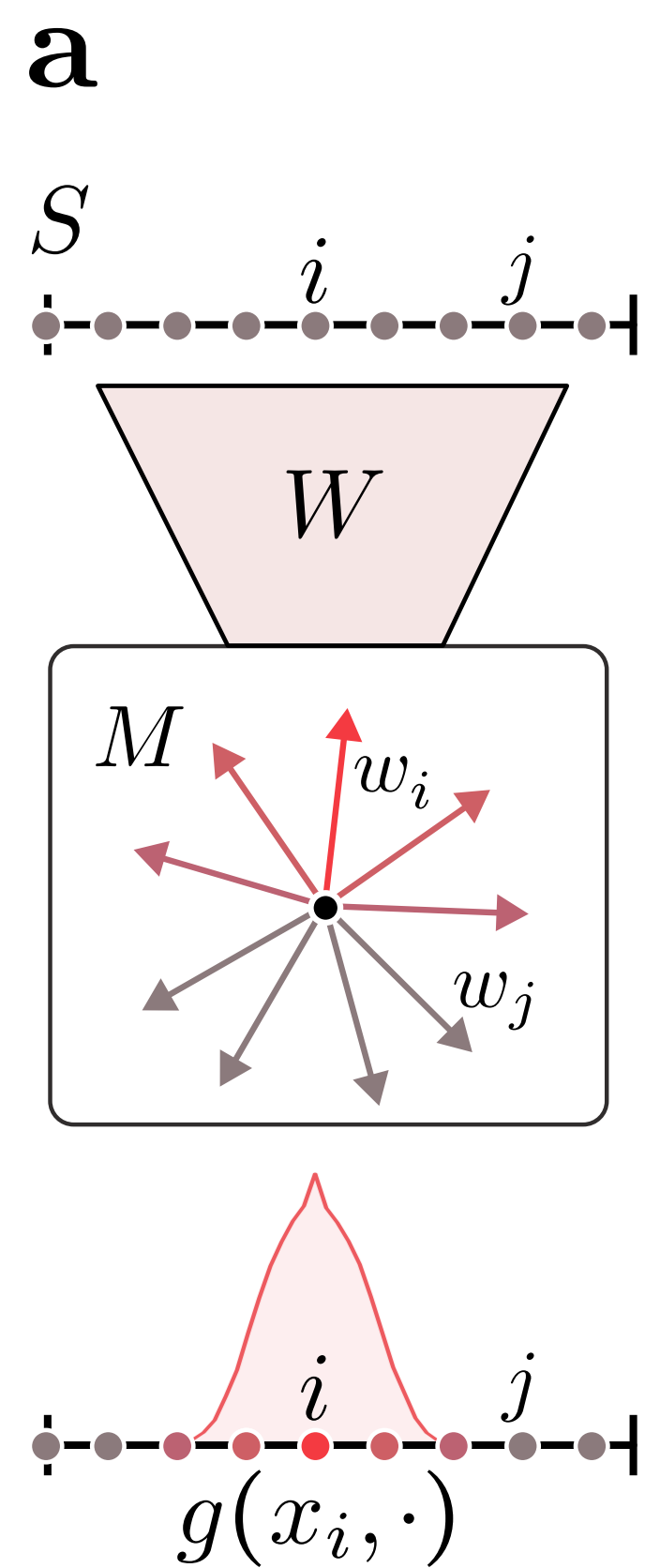
$$D_i = \frac{\sigma(w_i^\top w_k)}{\sigma(w_i^\top w_k) + \sigma(w_j^\top w_k)}, \quad D_j = \frac{\sigma(w_j^\top w_k)}{\sigma(w_i^\top w_k) + \sigma(w_j^\top w_k)}.$$

$$L_{\text{sim}} = -\frac{1}{2} D_{\hat{i}}.$$

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

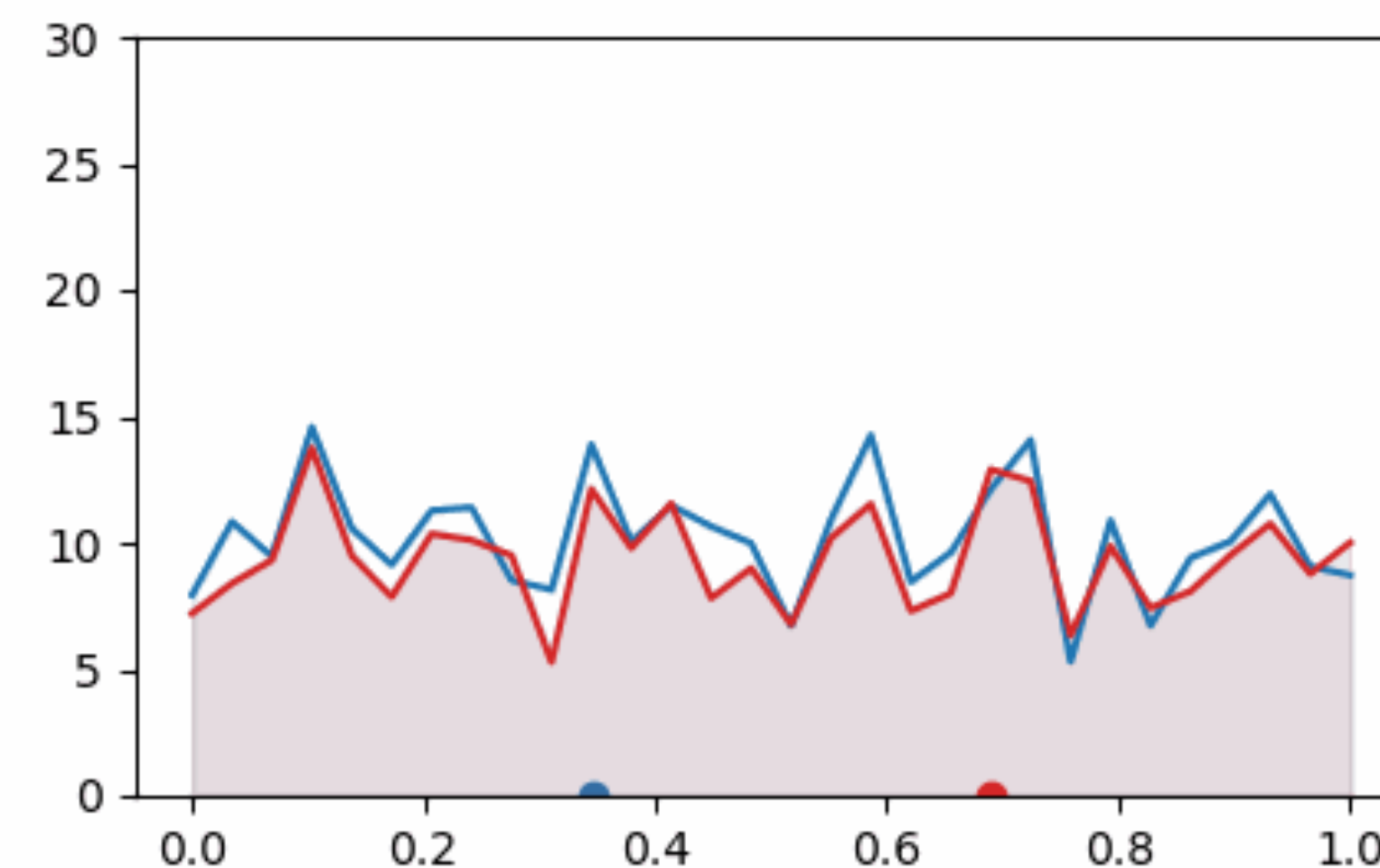


$$f(x) = \sigma(W^\top Wx),$$

$$L_{\text{rec}} = \sum_{i=1}^l \|e_i - \sigma(W^\top W e_i)\|^2 = \sum_{i=1}^l \|e_i - \sigma(W^\top w_i)\|^2.$$

$$D_i = \frac{\sigma(w_i^\top w_k)}{\sigma(w_i^\top w_k) + \sigma(w_j^\top w_k)}, \quad D_j = \frac{\sigma(w_j^\top w_k)}{\sigma(w_i^\top w_k) + \sigma(w_j^\top w_k)}.$$

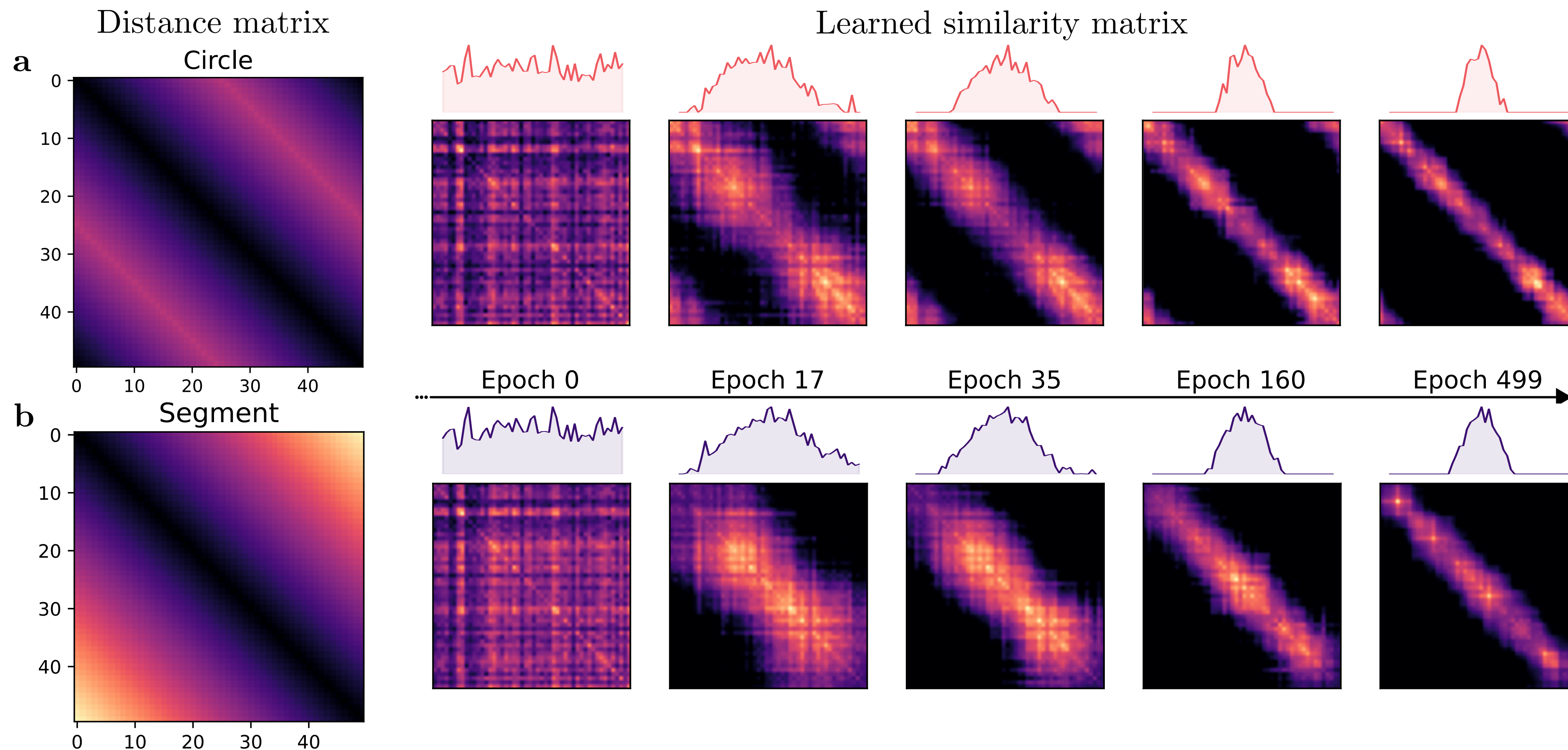
$$L_{\text{sim}} = -\frac{1}{2} D_{\hat{i}}.$$



# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>



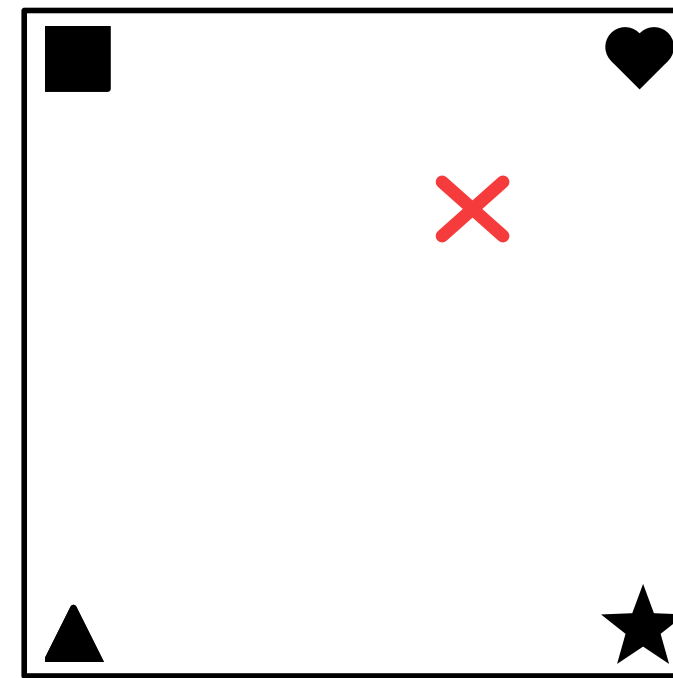
Link to percolation on random lattice models?

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

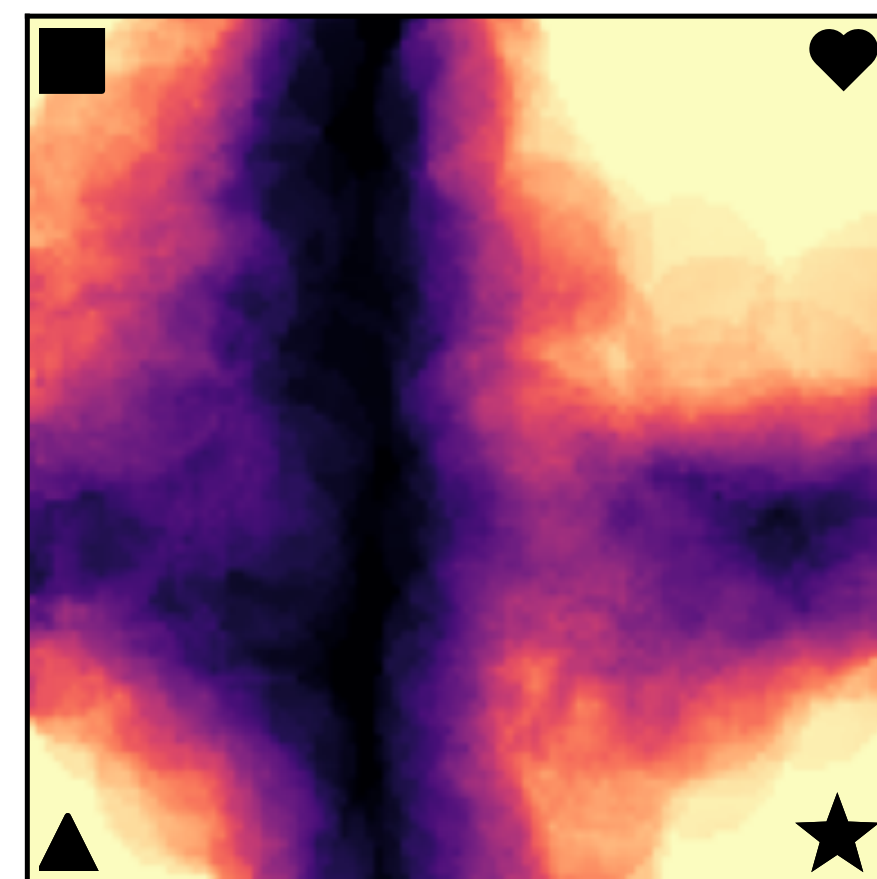
Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

## VLM spatial similarity task

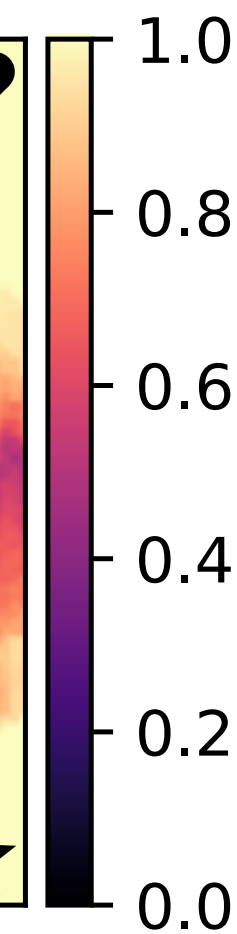
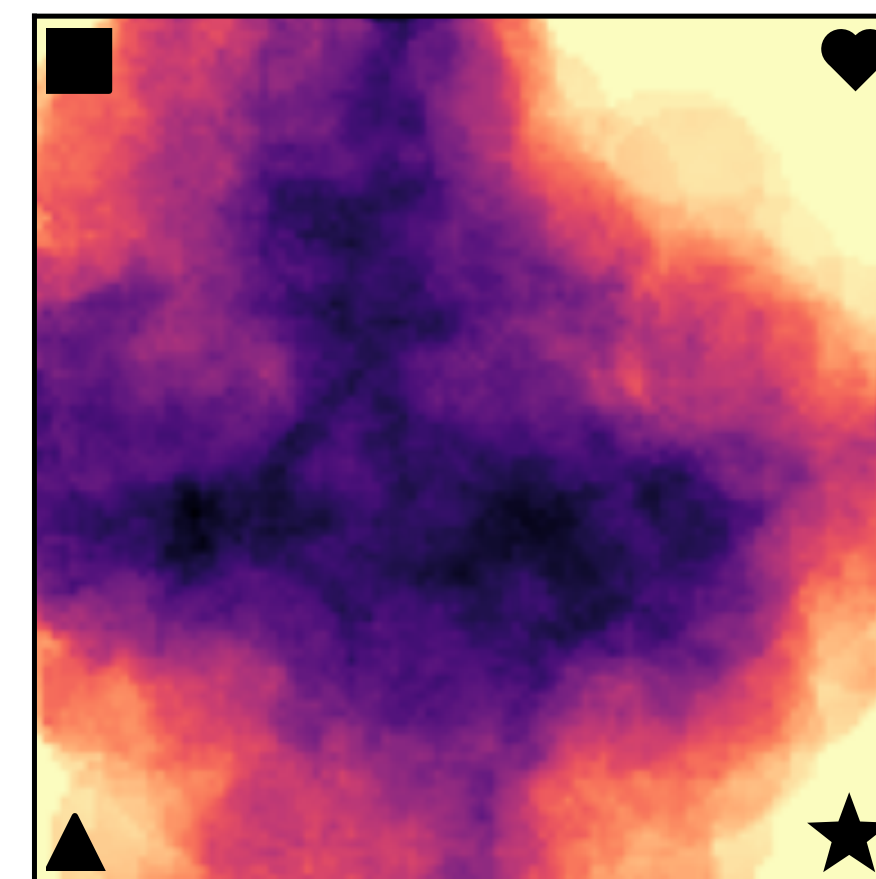


Which **black shape** is closest to the **red X**?

gemma-3-12b-it



Qwen2.5-VL-7B-Instruct  $\mathbb{P}$  correct



# Generalization vs Processing

---

**Bound by semanticity: universal laws governing the  
generalization-identification tradeoff**

---

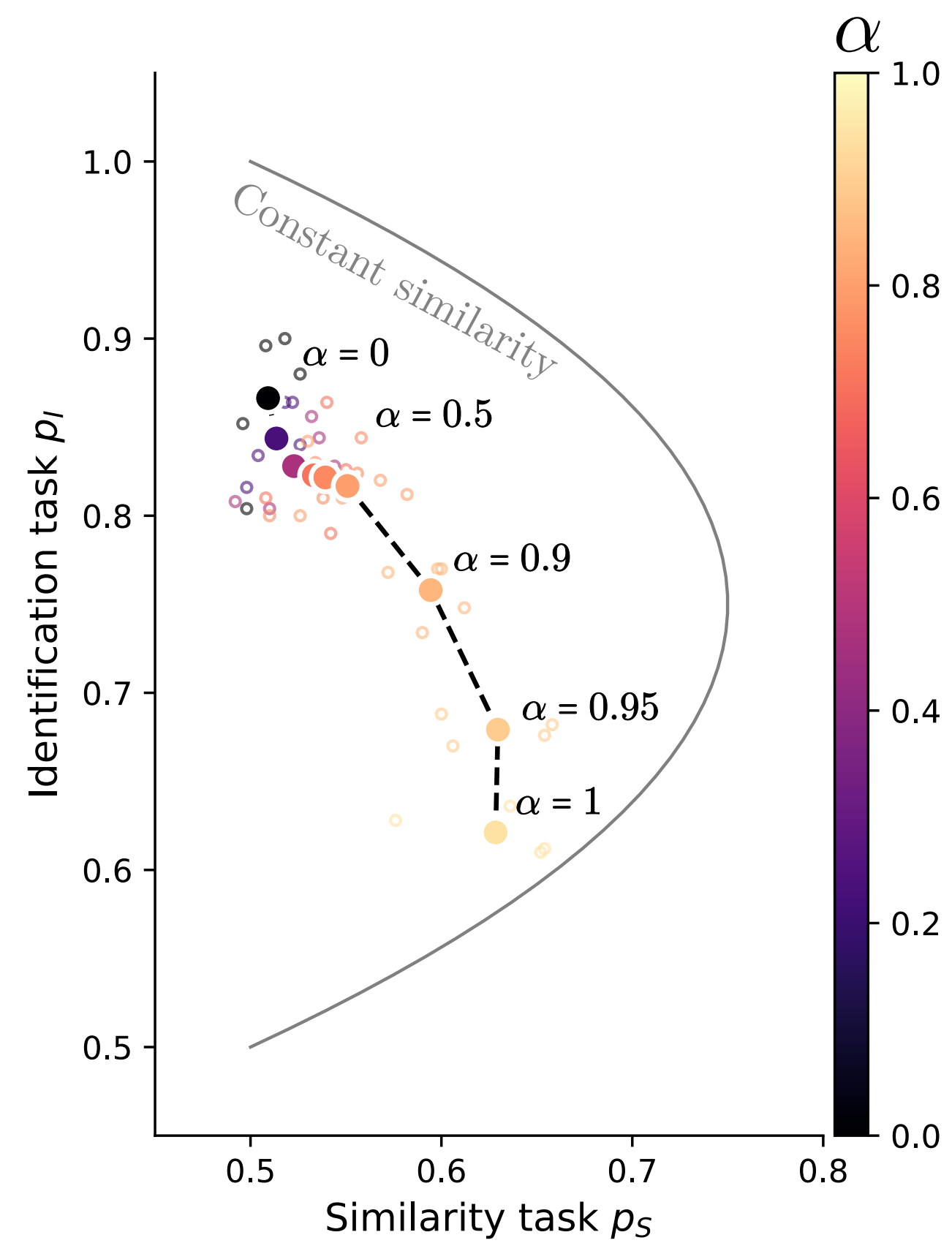
Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

# Generalization vs Processing

Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

## CNN finetuning



# Generalization vs Processing

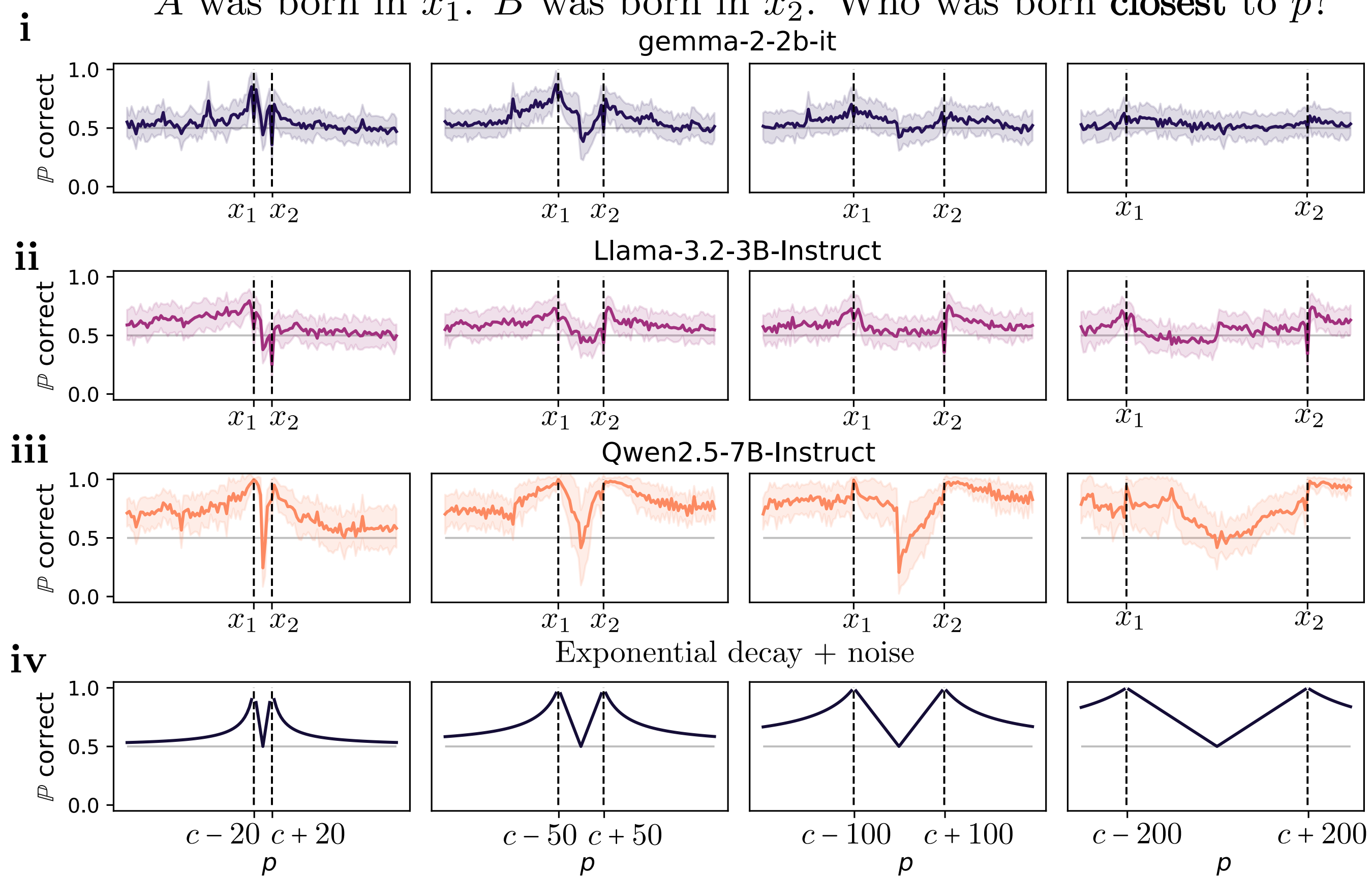
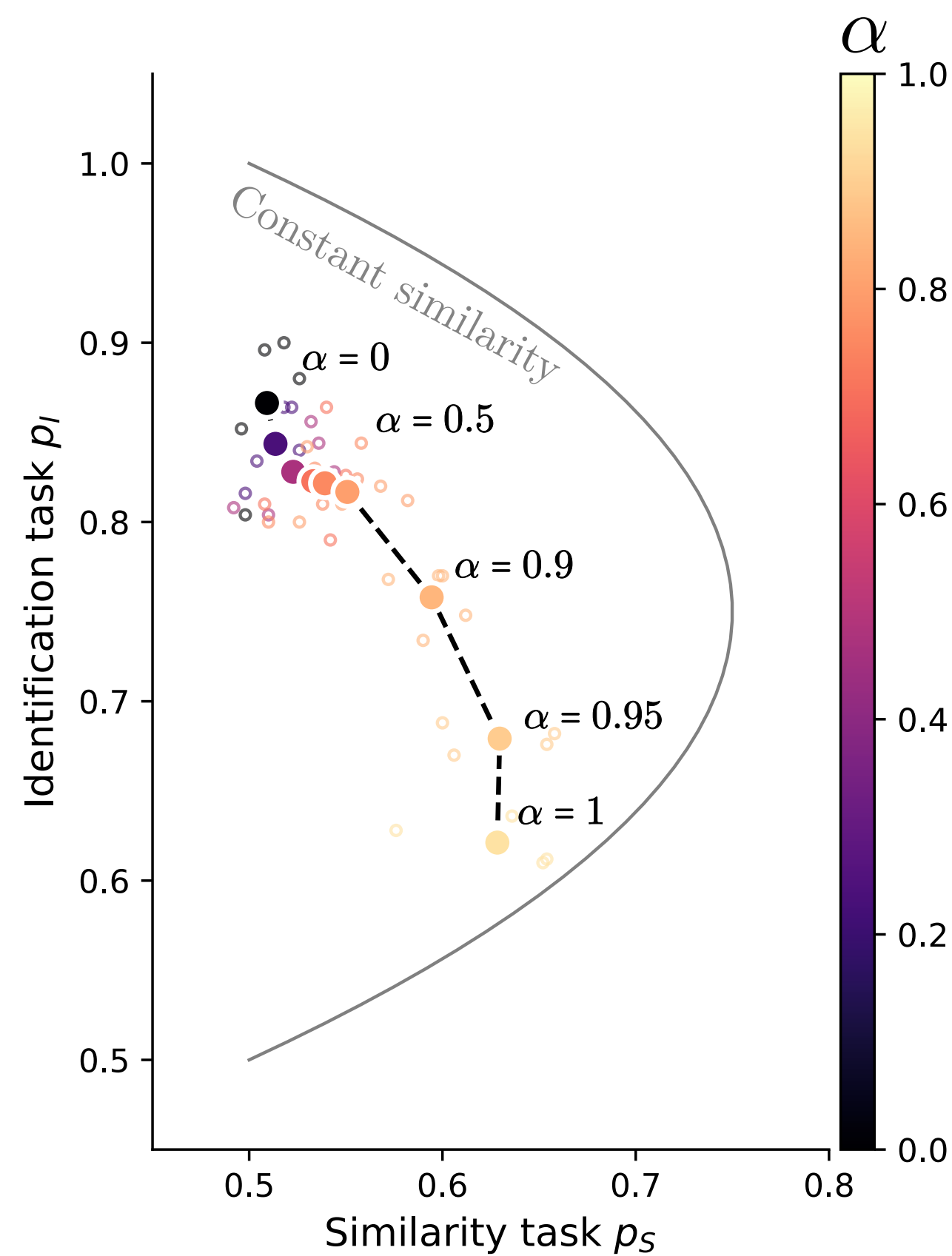
Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurisso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>,  
Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>,  
Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

## LLM year similarity task

$A$  was born in  $x_1$ .  $B$  was born in  $x_2$ . Who was born closest to  $p$ ?

## CNN finetuning



# Generalization vs Processing

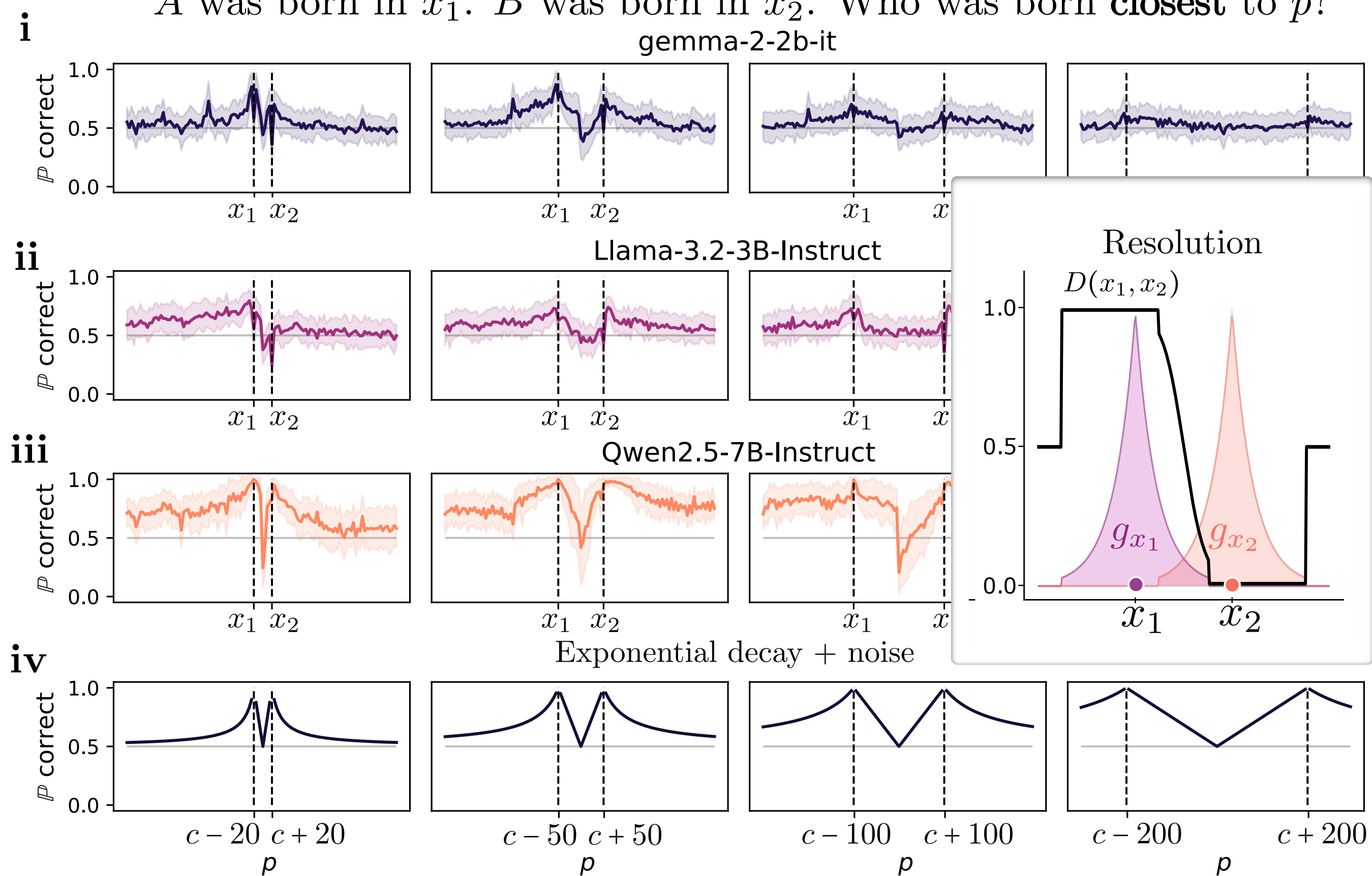
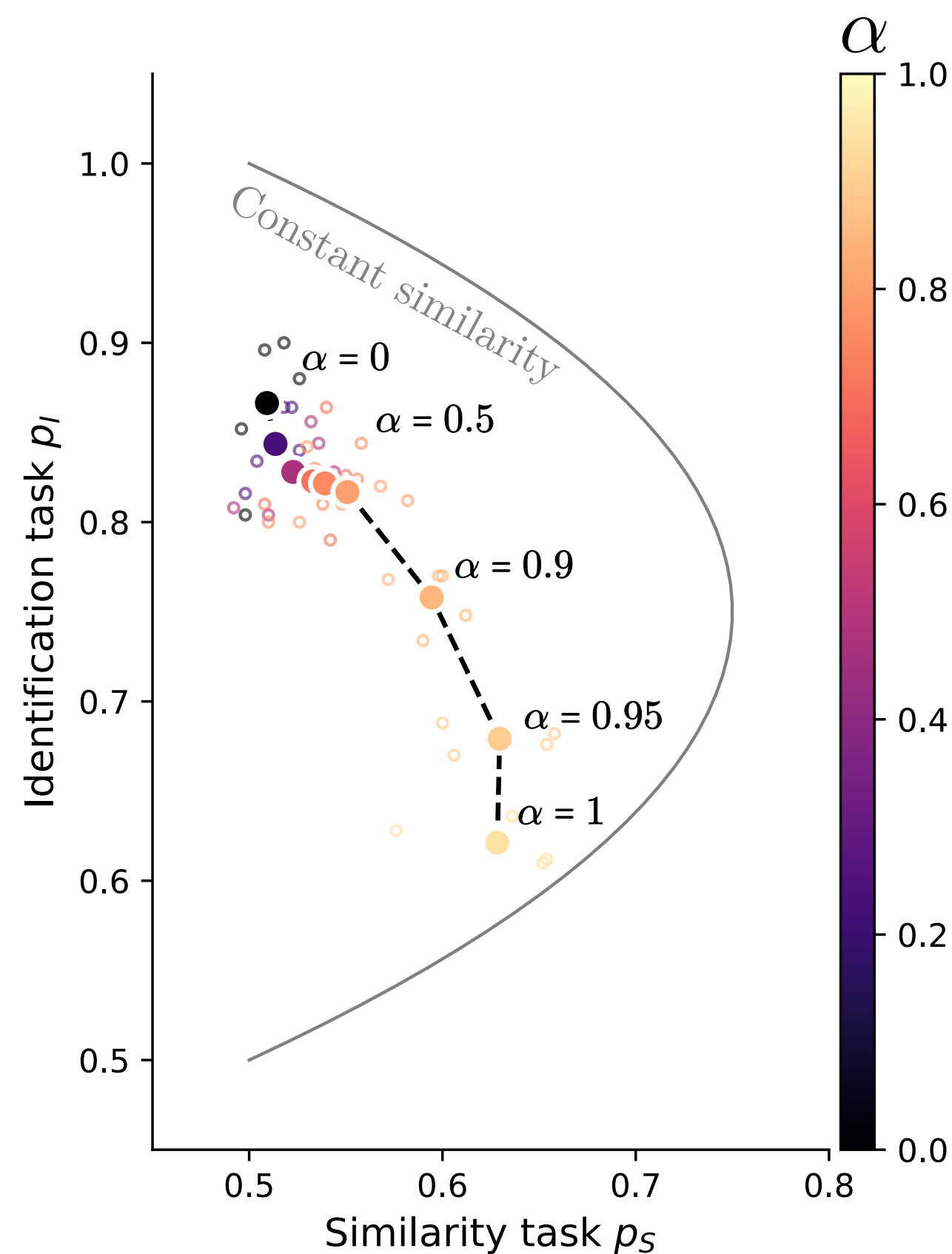
Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

## LLM year similarity task

A was born in  $x_1$ . B was born in  $x_2$ . Who was born closest to  $p$ ?

## CNN finetuning



# Take-home

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

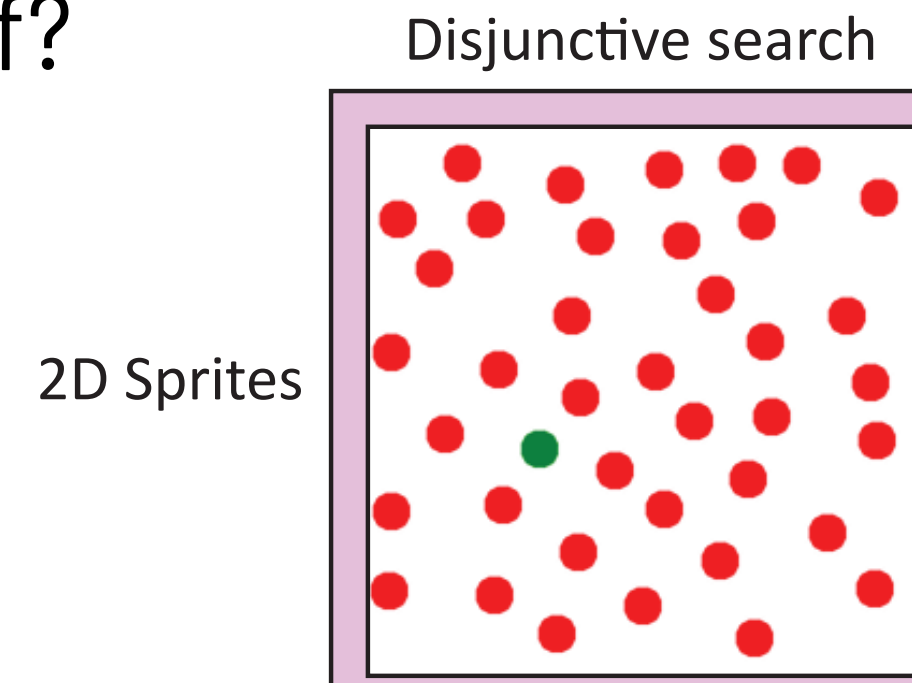
# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?



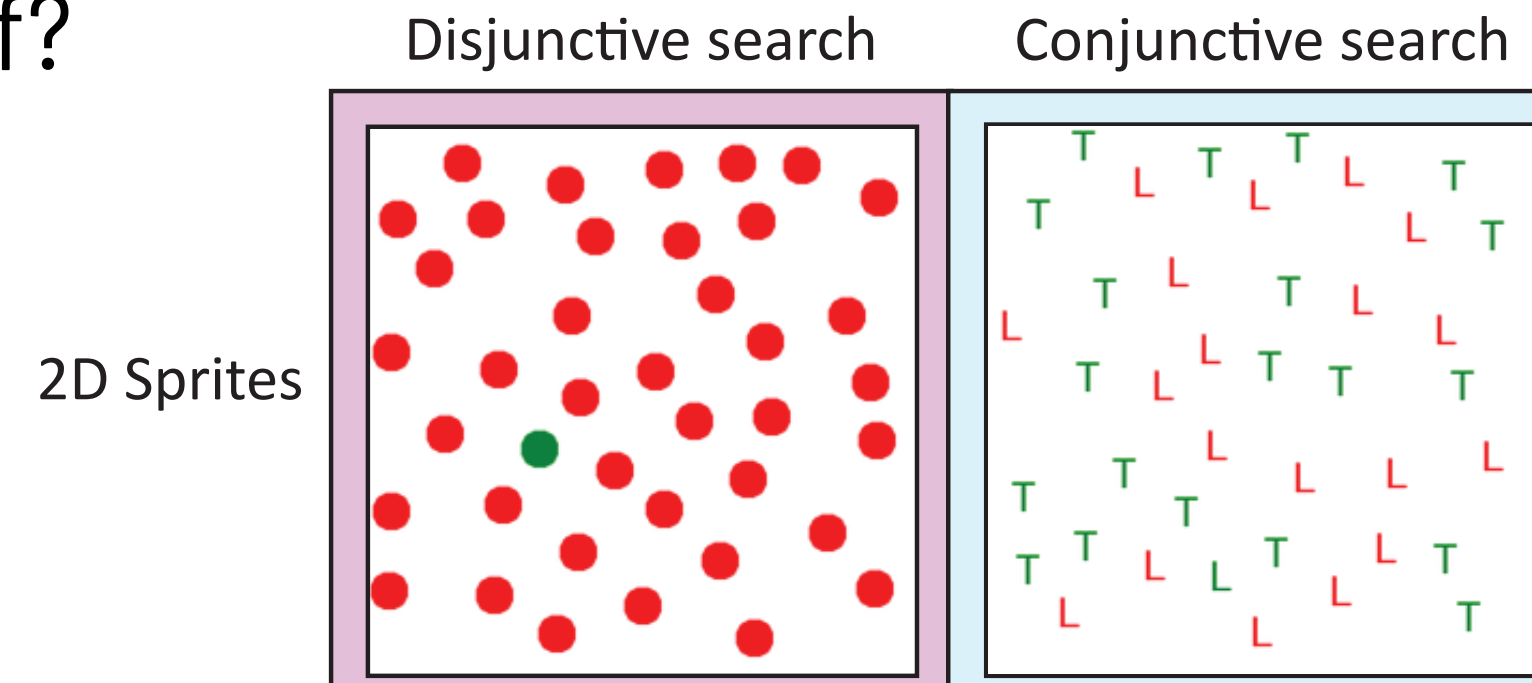
# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?



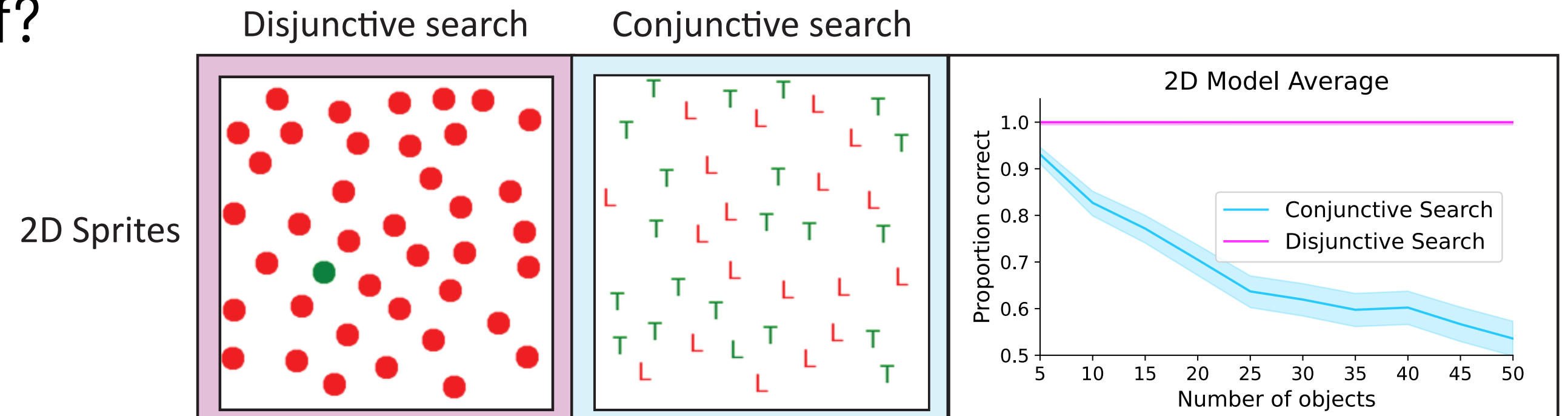
# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?



# Take-home

## Results:

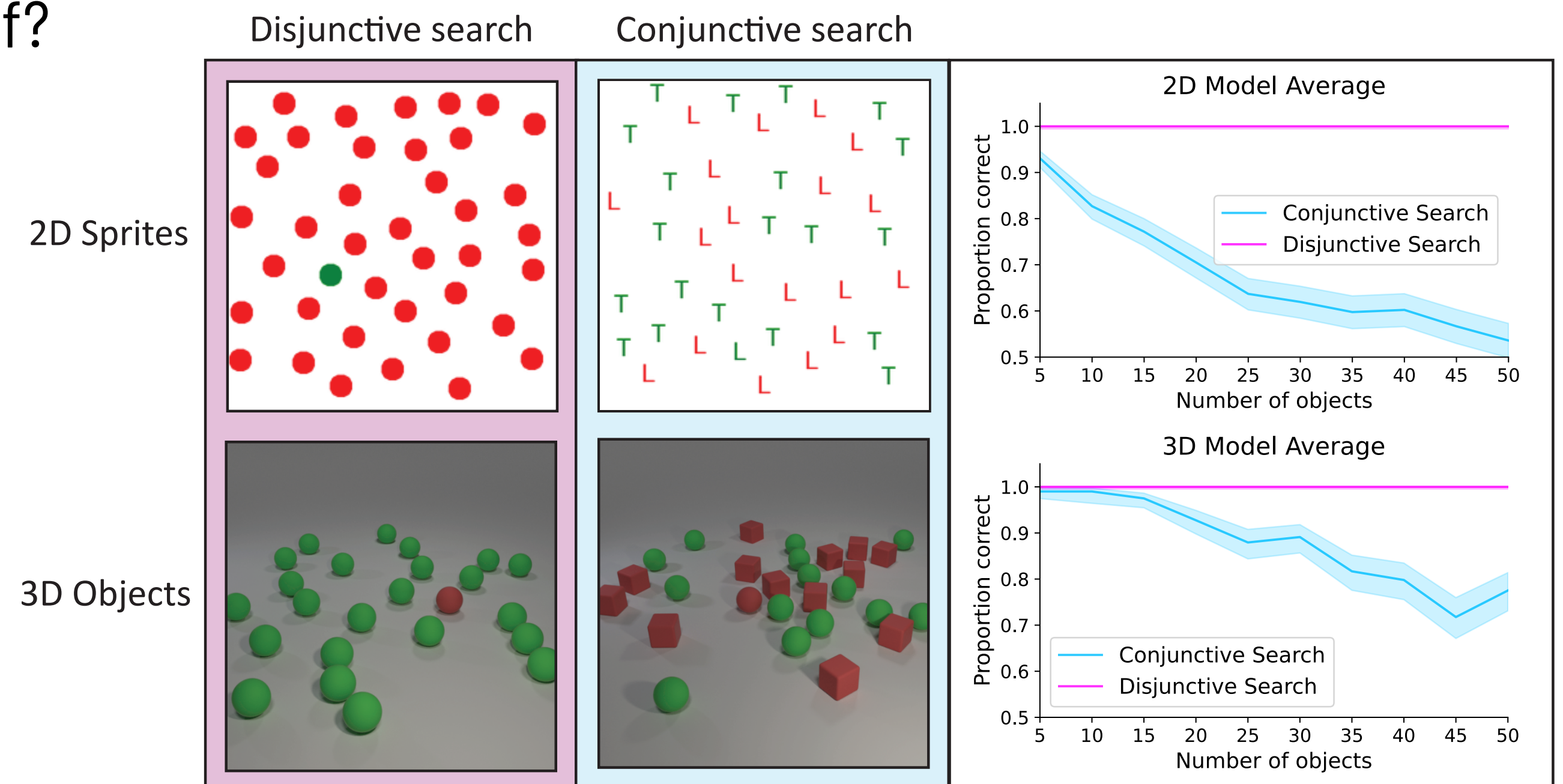
- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

Understanding the Limits of Vision Language Models  
Through the Lens of the Binding Problem

Declan Campbell<sup>1</sup>, Sunayana Rane<sup>2</sup>, Tyler Giallanza<sup>2</sup>, Nicolò De Sabbata<sup>3</sup>, Kia Ghods<sup>1</sup>,  
Amogh Joshi<sup>1</sup>, Alexander Ku<sup>2</sup>, Steven M. Frankland<sup>4</sup>, Thomas L. Griffiths<sup>2,5\*</sup>,  
Jonathan D. Cohen<sup>1,2\*</sup>, and Taylor W. Webb<sup>6\*</sup>



# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

filler  
(red, square,  
etc..)

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

filler  
(red, square, etc..)

Role  
(abstract object slot)

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

filler  
(red, square, etc..)

Role  
(abstract object slot)

$$\begin{aligned}\hat{f} &= \Phi r = FR^{\top} r^{(k)} \\ &= f^{(k)} + \sum_{j \neq k} f^{(j)} (r^{(j)})^{\top} r^{(k)} \\ &= f^{(k)} + \sum_{j \neq k} G_{jk}^R f^{(j)},\end{aligned}$$

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

filler  
(red, square, etc..)

Role  
(abstract object slot)

$$\begin{aligned}\hat{f} &= \Phi r = FR^{\top} r^{(k)} \\ &= f^{(k)} + \sum_{j \neq k} f^{(j)} (r^{(j)})^{\top} r^{(k)} \\ &= \underset{\text{feature}}{f^{(k)}} + \sum_{j \neq k} G_{jk}^R f^{(j)},\end{aligned}$$

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

TPR formulation of Miller-Shepard's law

$$\Phi = \sum_{i=1}^n f^{(i)} \otimes r^{(i)} = FR^{\top},$$

filler  
(red, square, etc..)  
Role  
(abstract object slot)

$$\begin{aligned}\hat{f} &= \Phi r = FR^{\top} r^{(k)} \\ &= f^{(k)} + \sum_{j \neq k} f^{(j)} (r^{(j)})^{\top} r^{(k)} \\ &= \underbrace{f^{(k)}}_{\text{feature}} + \sum_{j \neq k} G_{jk}^R \underbrace{f^{(j)}}_{\text{interference}},\end{aligned}$$

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

# Take-home

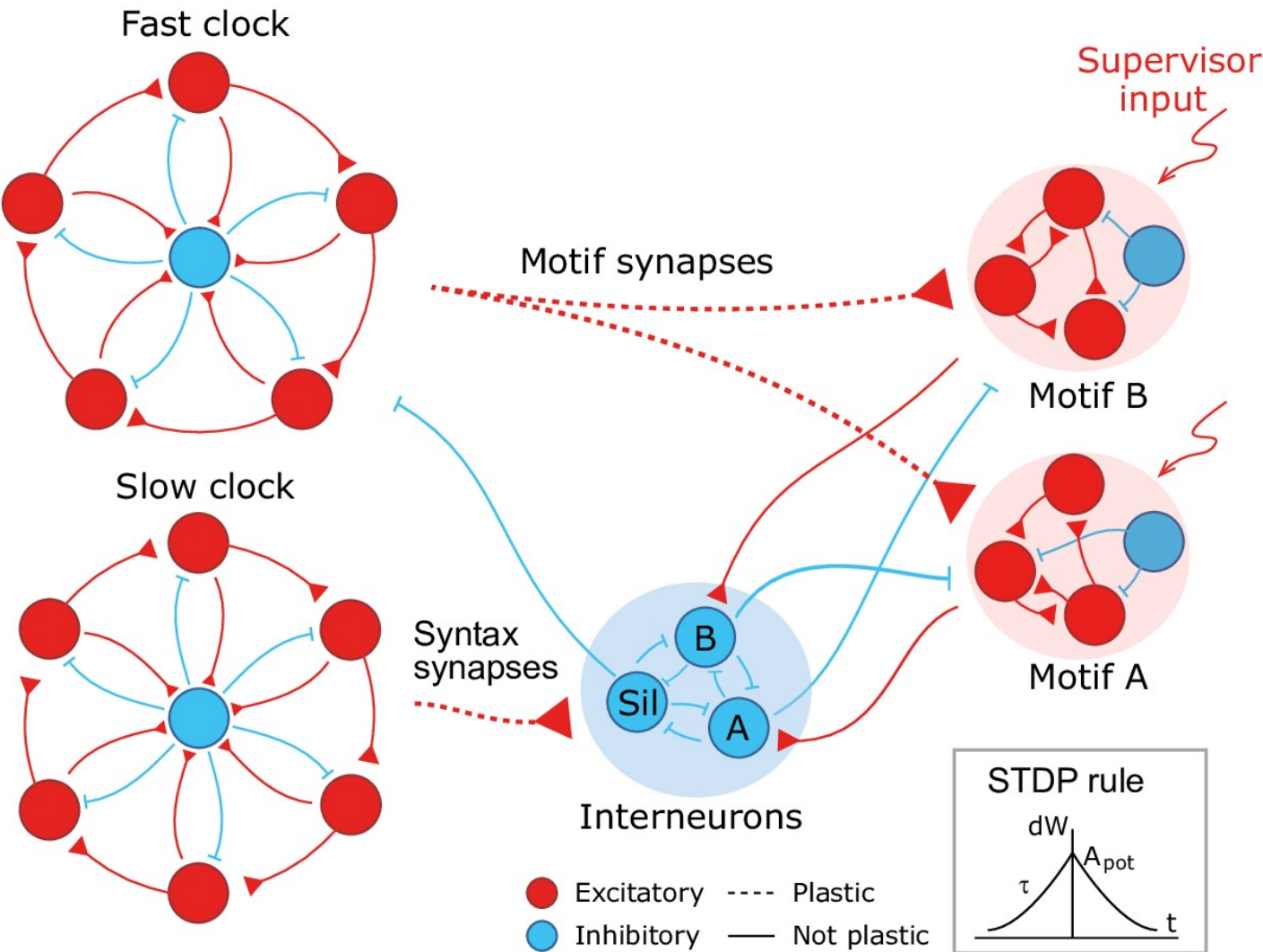
## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

RESEARCH ARTICLE  
Learning compositional sequences with multiple time scales through a hierarchical network of spiking neurons  
Amadeus Maes<sup>1</sup>, Mauricio Barahona<sup>2</sup>, Claudia Clopath<sup>1\*</sup>



# Take-home

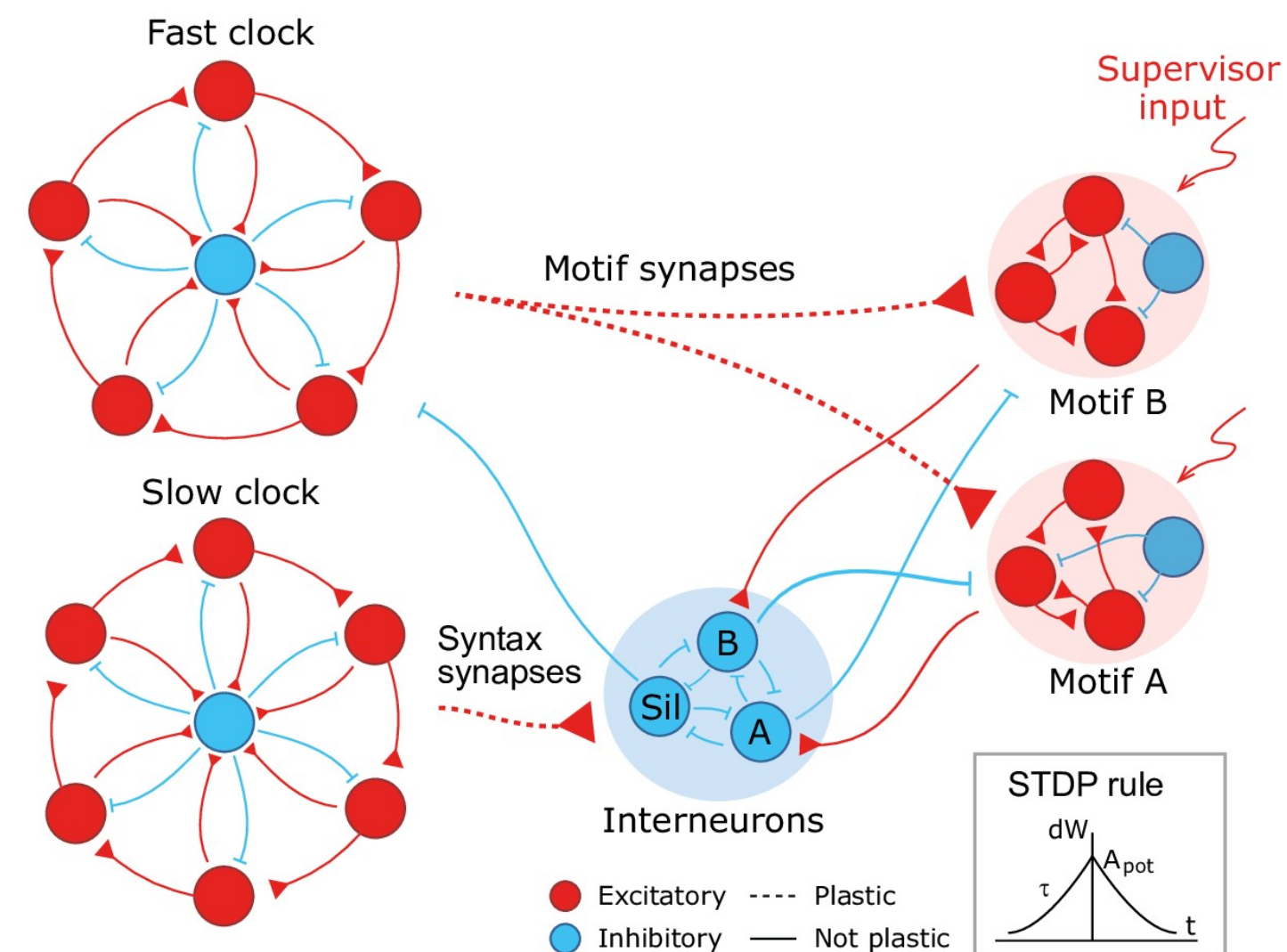
## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

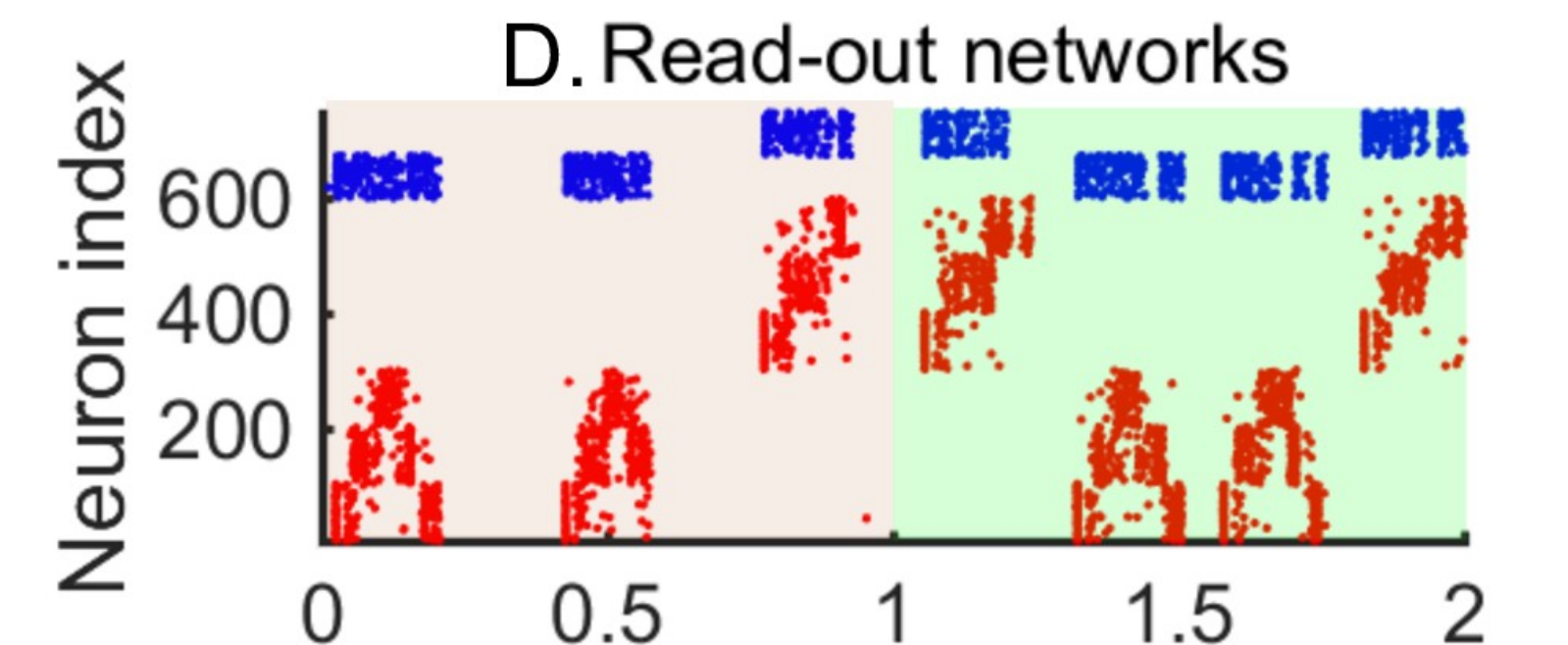
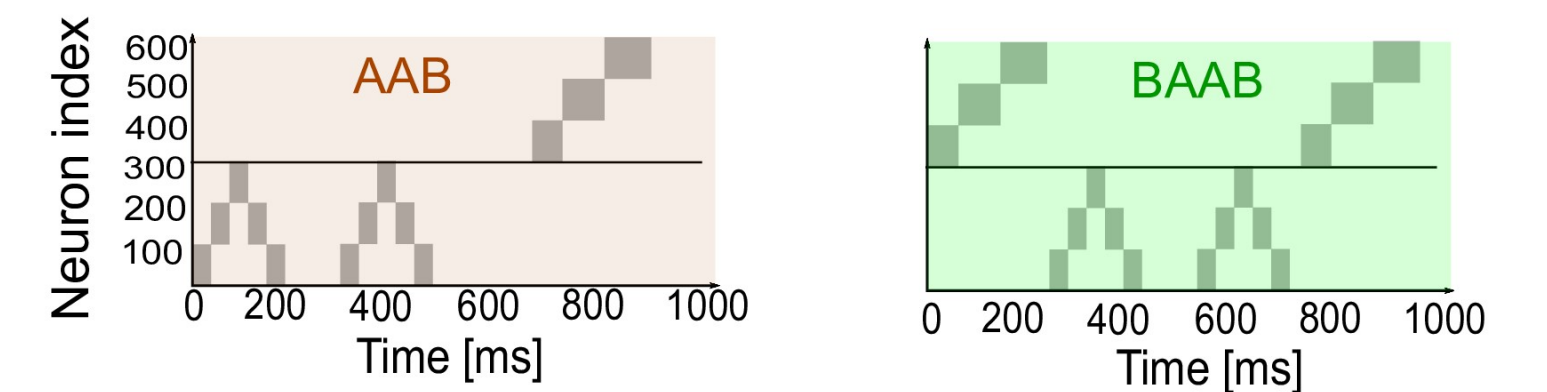
## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

RESEARCH ARTICLE  
Learning compositional sequences with multiple time scales through a hierarchical network of spiking neurons  
Amadeus Maes<sup>1</sup>, Mauricio Barahona<sup>2</sup>, Claudia Clopath<sup>1\*</sup>



A. Target sequences



# Take-home

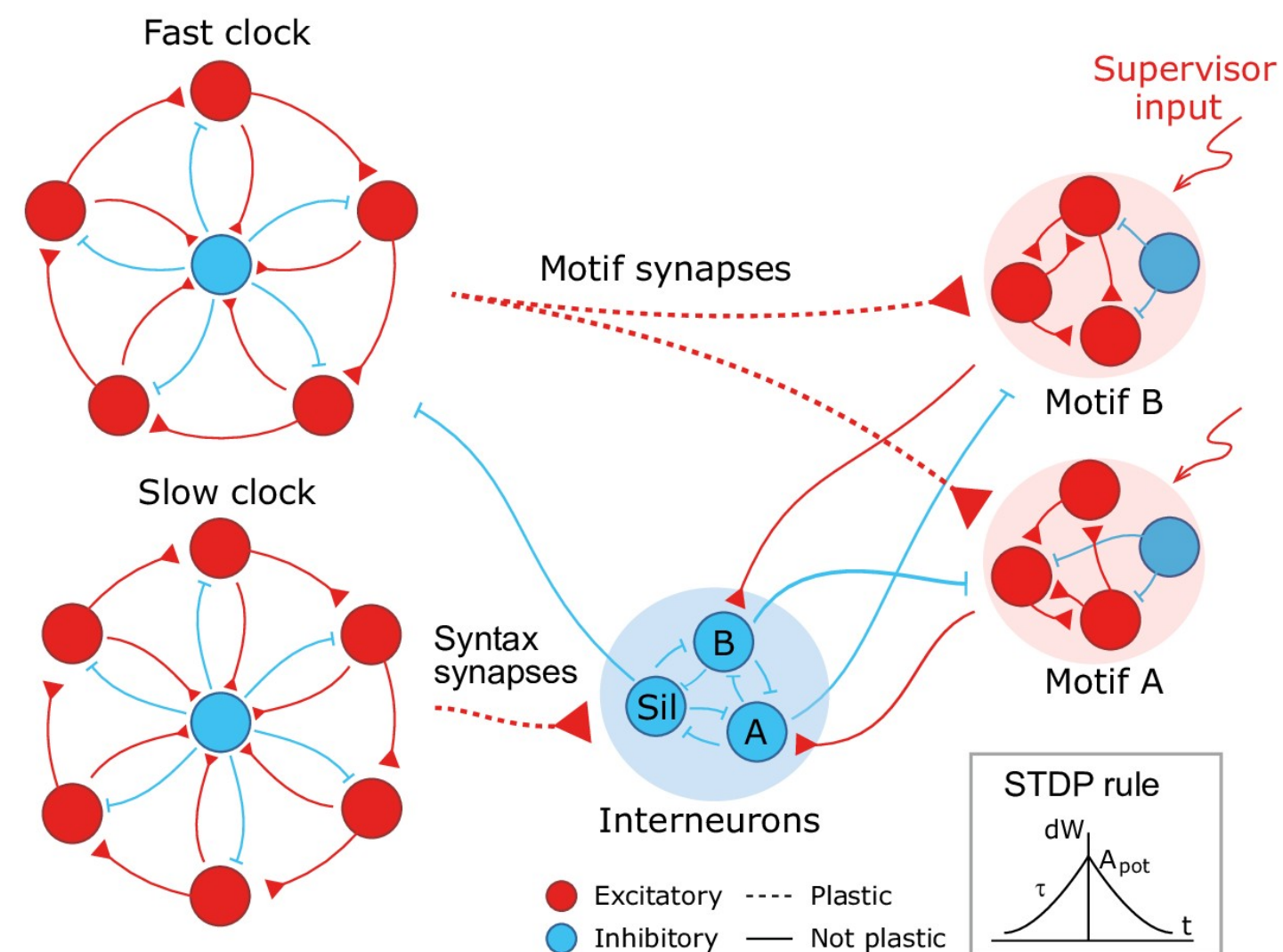
## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

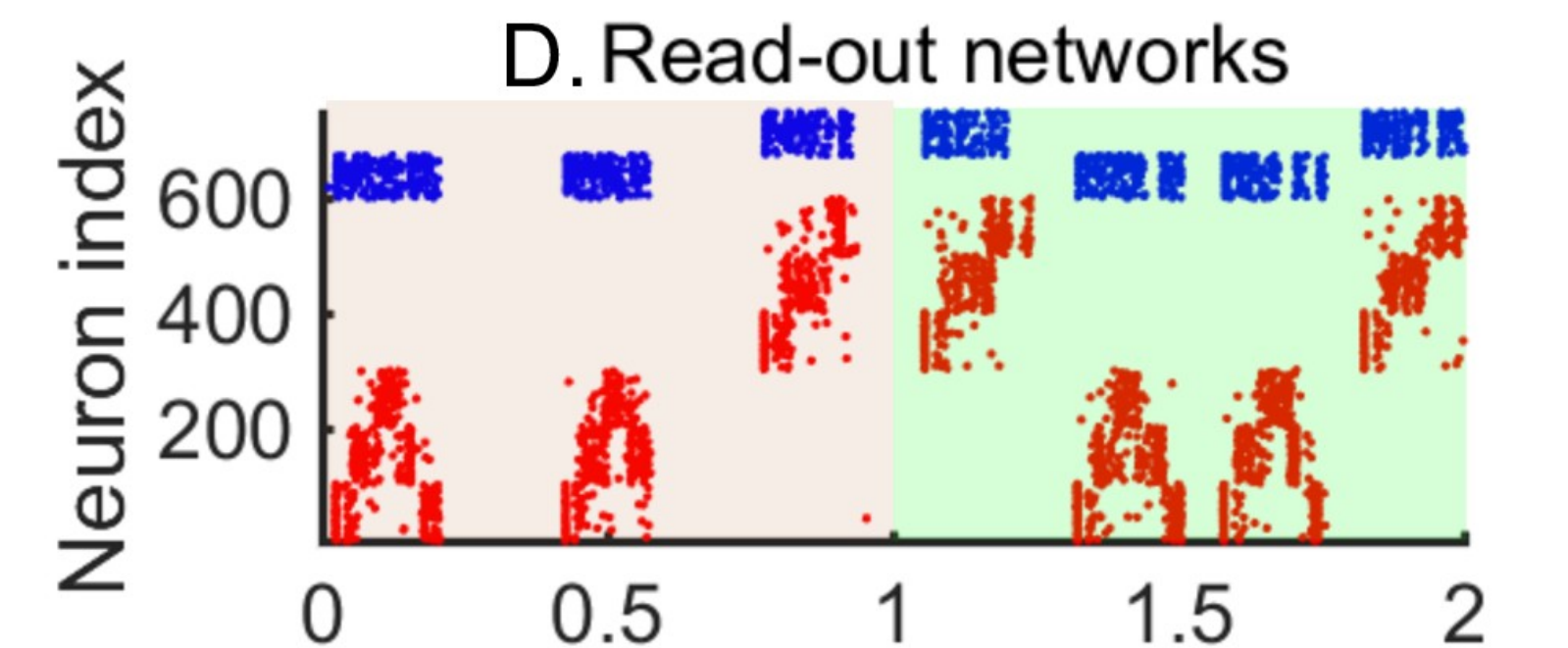
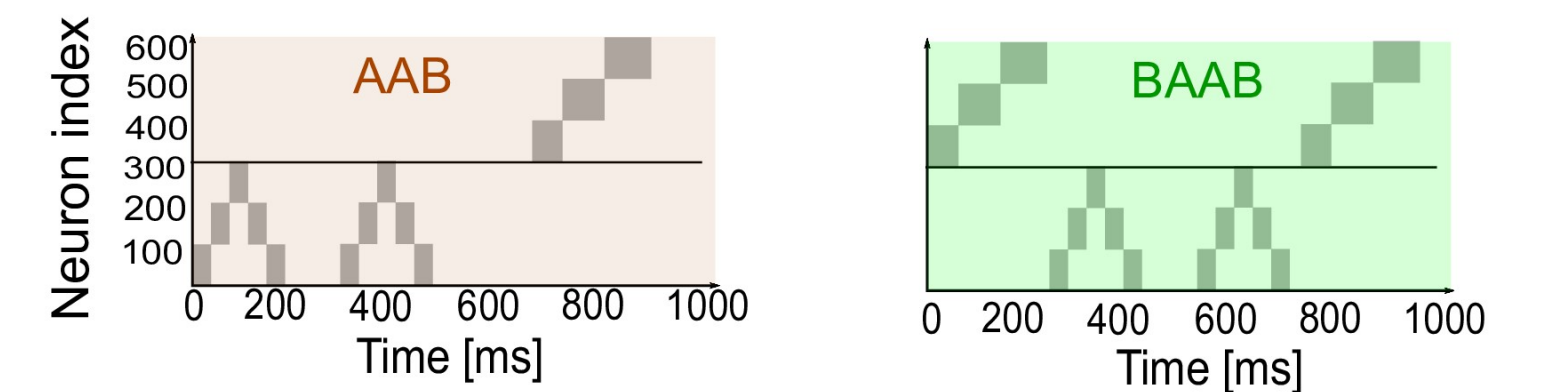
- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?

RESEARCH ARTICLE  
Learning compositional sequences with multiple time scales through a hierarchical network of spiking neurons  
Amadeus Maes<sup>1</sup>, Mauricio Barahona<sup>2</sup>, Claudia Clopath<sup>1\*</sup>



Jesseba Fernando

A. Target sequences



# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

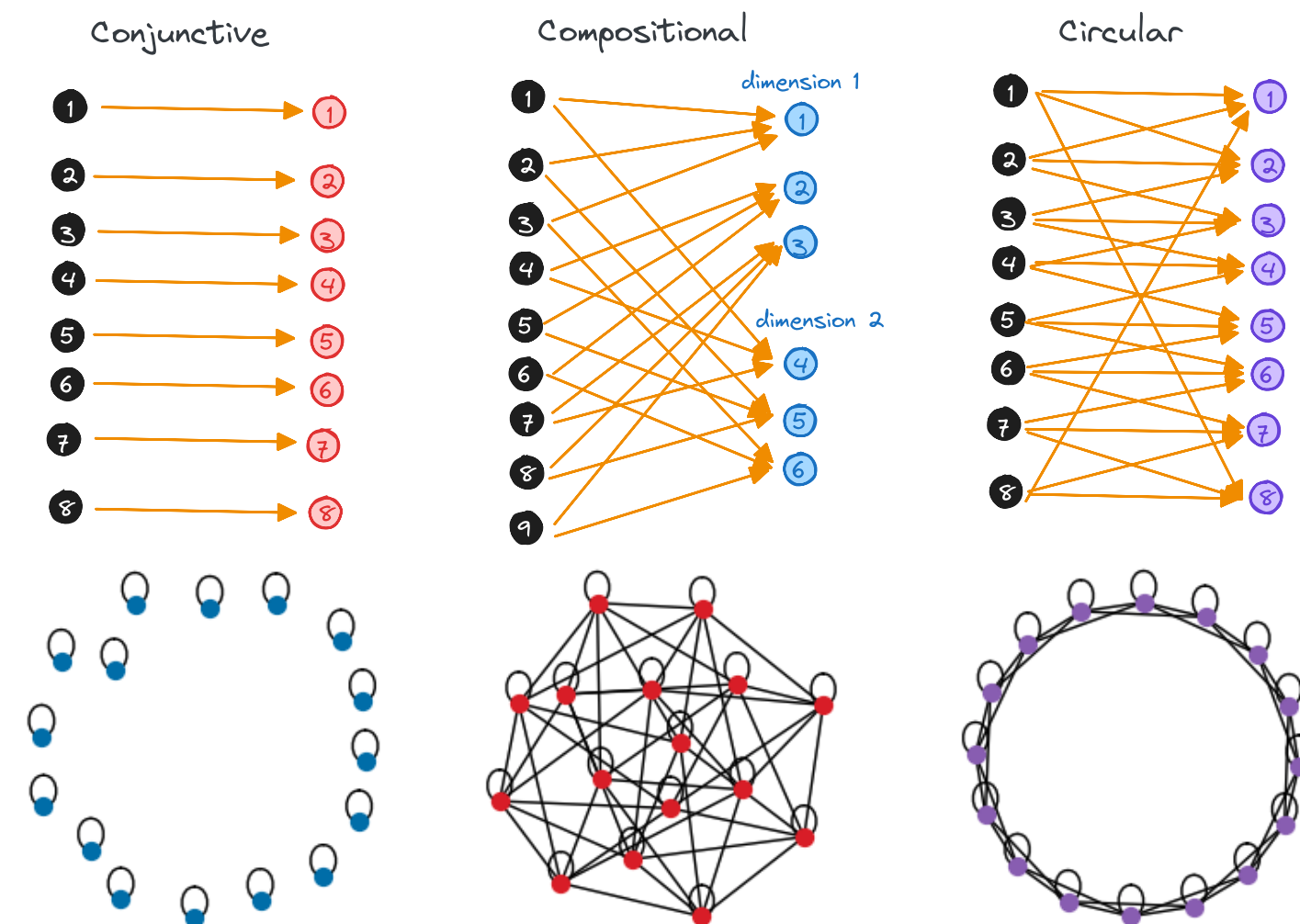
# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)



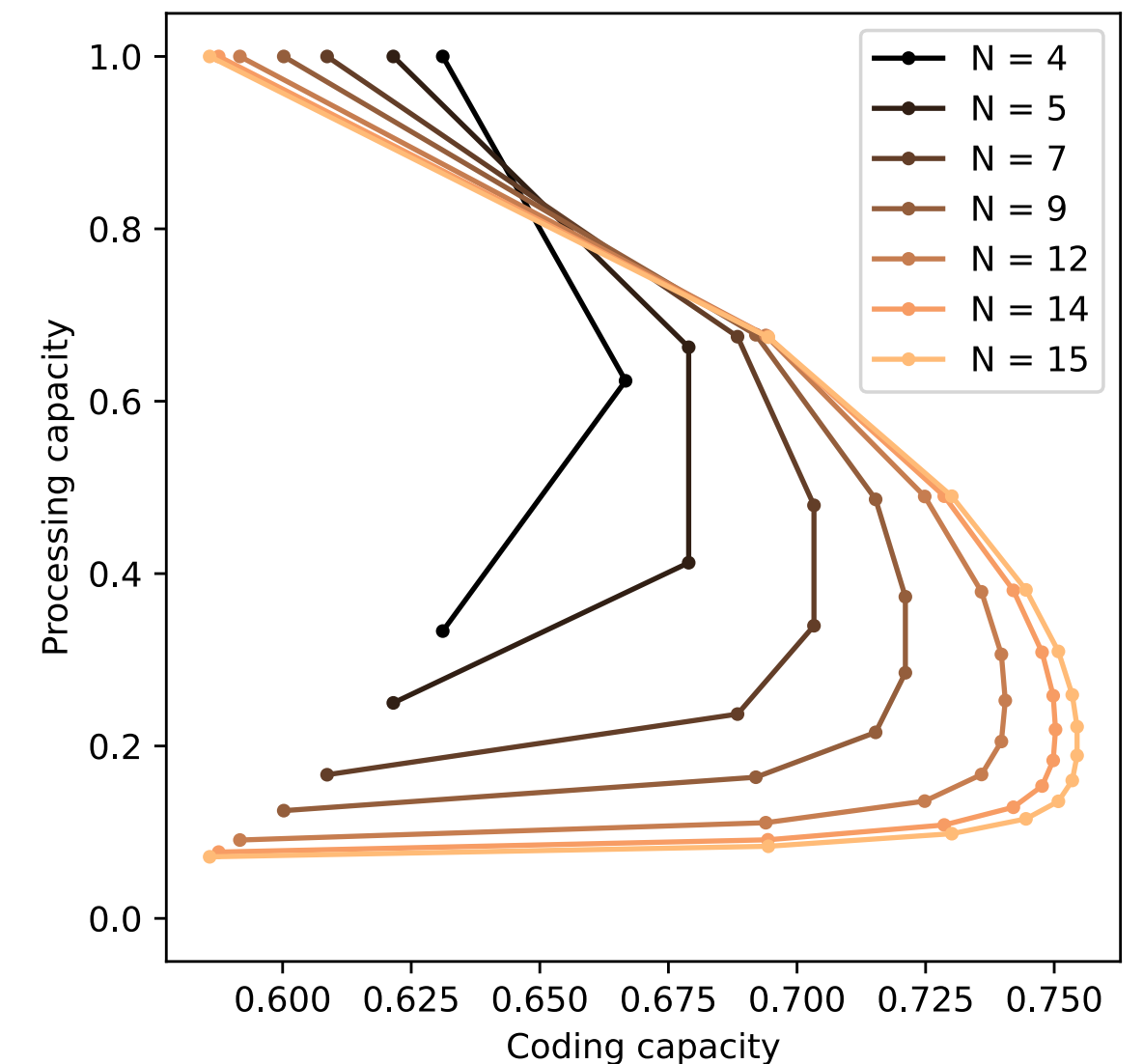
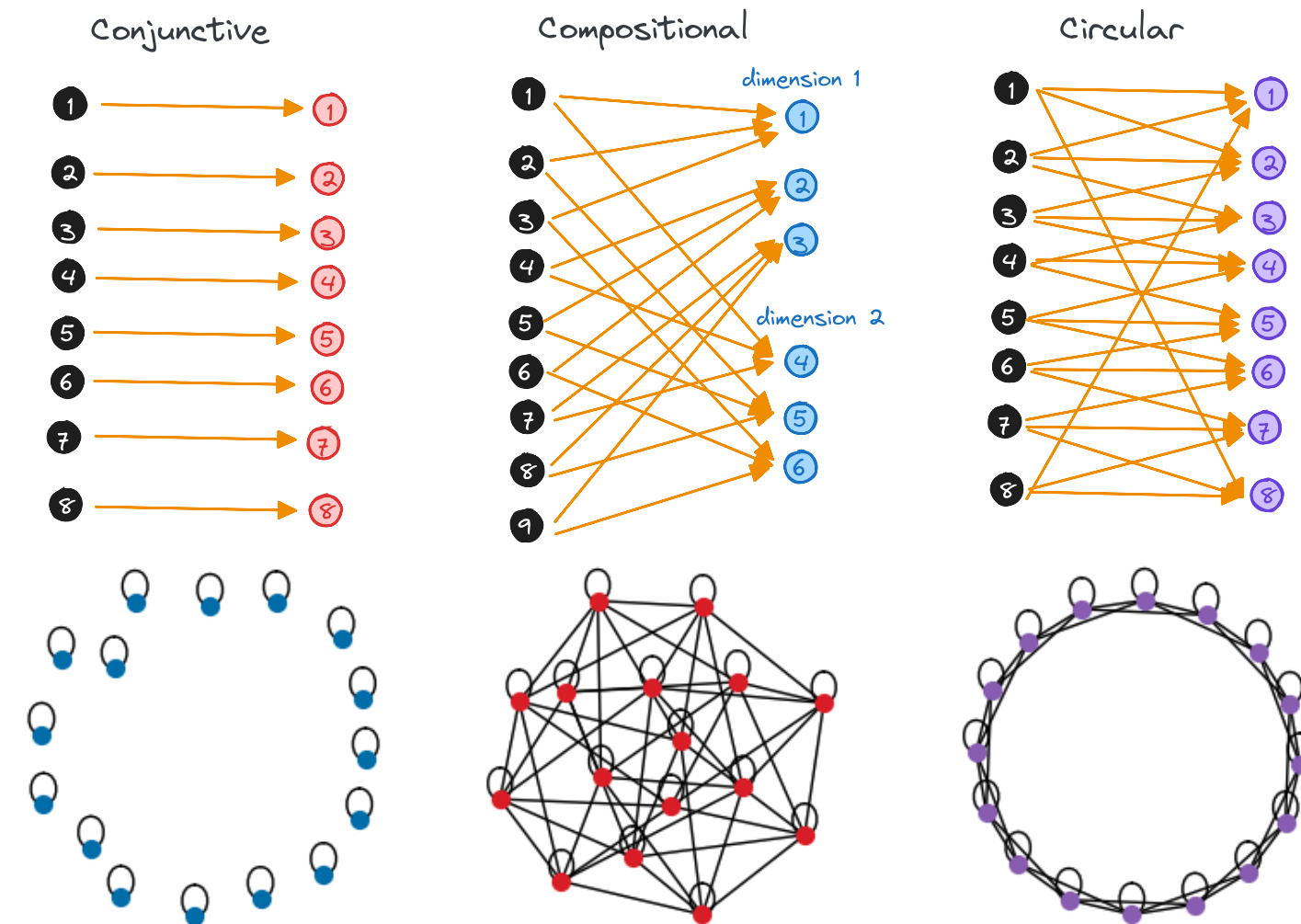
# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)



# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?

# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

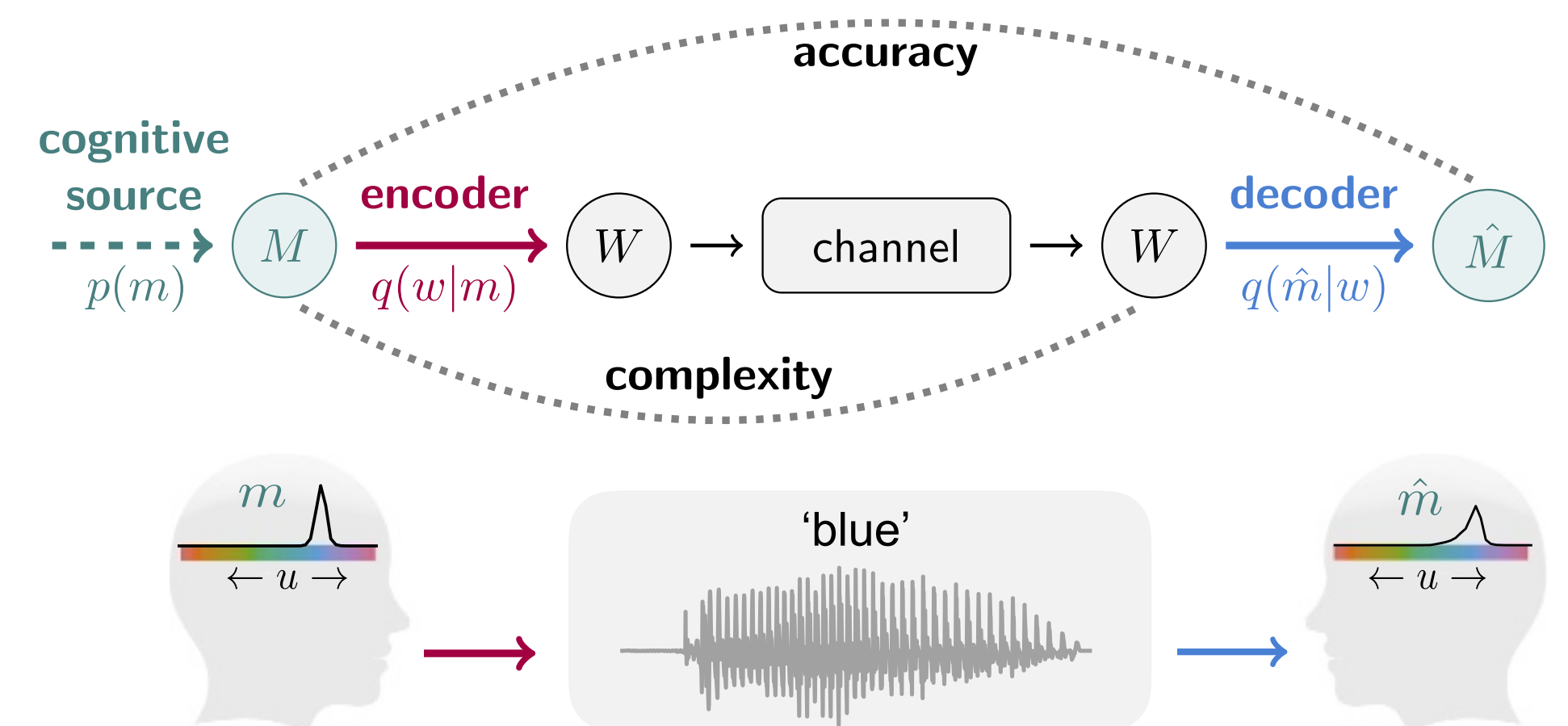
- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?

## Efficient compression in color naming and its evolution

Noga Zaslavsky<sup>a,b,1</sup>, Charles Kemp<sup>c,2</sup>, Terry Regier<sup>b,d</sup>, and Naftali Tishby<sup>a,e</sup>



# Take-home

## Results:

- Any resolution limit  $\rightarrow$  tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

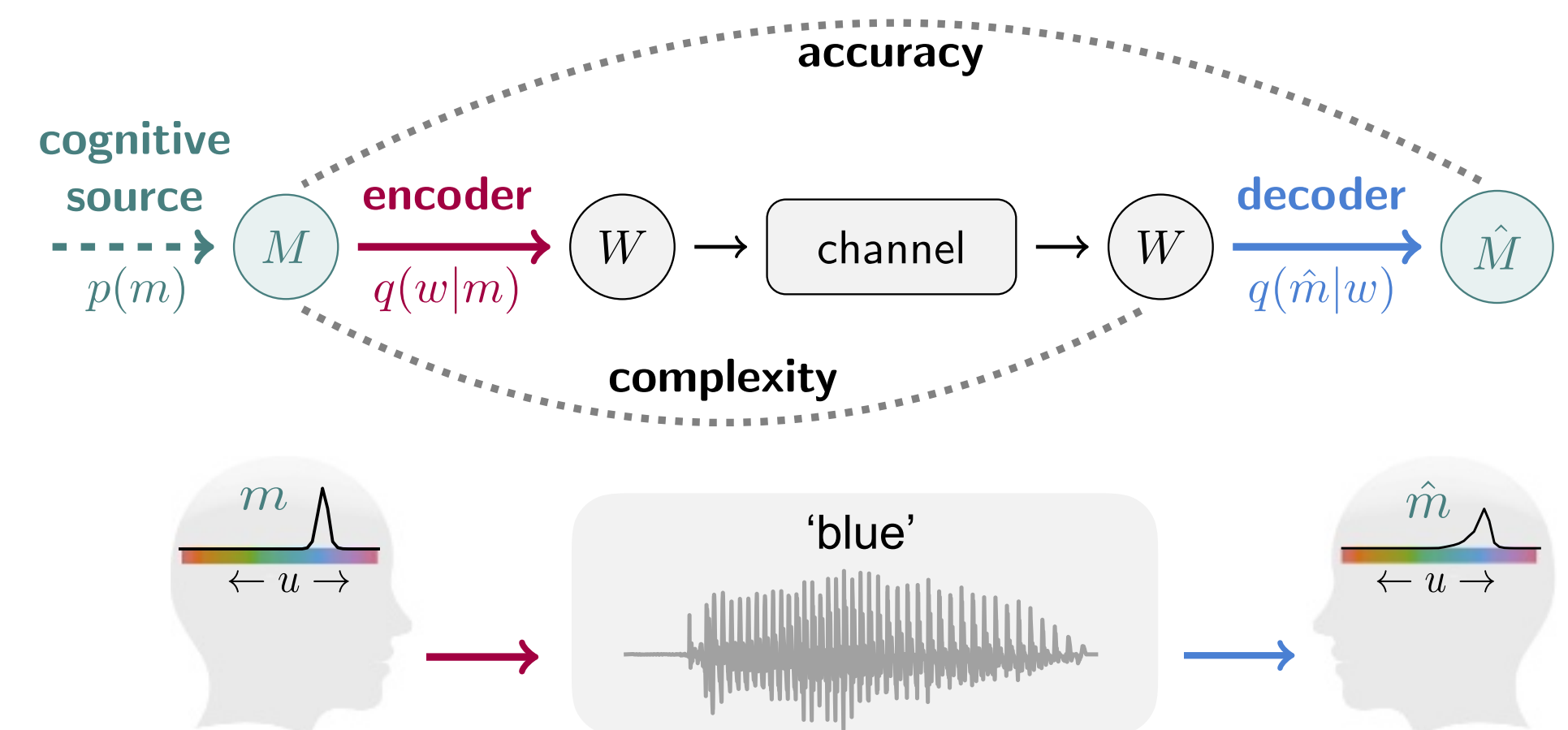
- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?

### Perspective

## Language is primarily a tool for communication rather than thought

Efficient compression in color naming and its evolution

Noga Zaslavsky<sup>a,b,1</sup>, Charles Kemp<sup>c,2</sup>, Terry Regier<sup>b,d</sup>, and Naftali Tishby<sup>a,e</sup>



# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?

# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?



Biswadip Dey



H. Kayan Oczimder



Nesreen Ahmed



Ted Willke



Topological limits to the parallel processing capability of network architectures



eLife | Home Magazine Community About

Neuroscience

An Information-Theoretic Approach to Reward Rate Optimization in the Tradeoff Between Controlled and Automatic Processing in Neural Network Architectures



Giovanni Petri, Sebastian Musslick, Jonathan D. Cohen



Sebastian Musslick



Jonathan Cohen

---

**Bound by semanticity: universal laws governing the generalization-identification tradeoff**

---



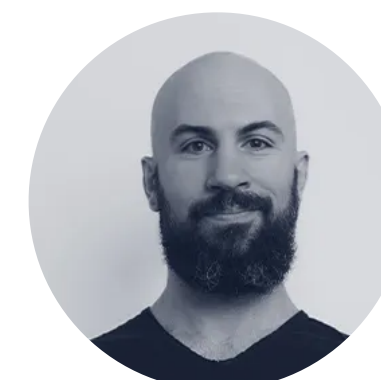
Marco Nurisso



Raj Deshpande



Jesseba Fernando



Alan Perotti

Raja Marjeh, Steven M. Frankland, Richard L. Lewis, Taylor W. Webb, Declan Campbell, Francesco Vaccarino

---

**The Geometry of Representational Failures in Vision Language Models**

---



# Take-home

## Results:

- Any resolution limit → tradeoff generalization vs processing
- Resolution emerges independently!

## Next:

- We can steer representations to reduce/increase performances (Savietto, ICML 2026)
- How does compositionality affects the tradeoff?
- Does this depend on scale? (i.e. What if we renormalize representations?)

## Hot takes:

- All living/cog systems (immune system, cells, your football team, etc..) face this tradeoff! What does it imply?

# Thank you!



Biswadip Dey



H. Kayan Oczimder



Nesreen Ahmed



Ted Willke



Topological limits to the parallel processing capability of network architectures



eLife | Home Magazine Community About

Neuroscience

An Information-Theoretic Approach to Reward Rate Optimization in the Tradeoff Between Controlled and Automatic Processing in Neural Network Architectures



Giovanni Petri, Sebastian Musslick, Jonathan D. Cohen



Sebastian Musslick



Jonathan Cohen

---

**Bound by semanticity: universal laws governing the generalization-identification tradeoff**

---



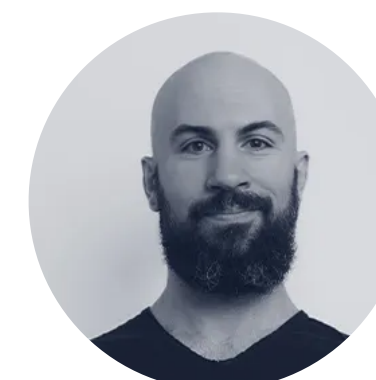
Marco Nurisso



Raj Deshpande



Jesseba Fernando



Alan Perotti

Raja Marjeh, Steven M. Frankland, Richard L. Lewis, Taylor W. Webb, Declan Campbell, Francesco Vaccarino

---

**The Geometry of Representational Failures in Vision Language Models**

---



# Backup

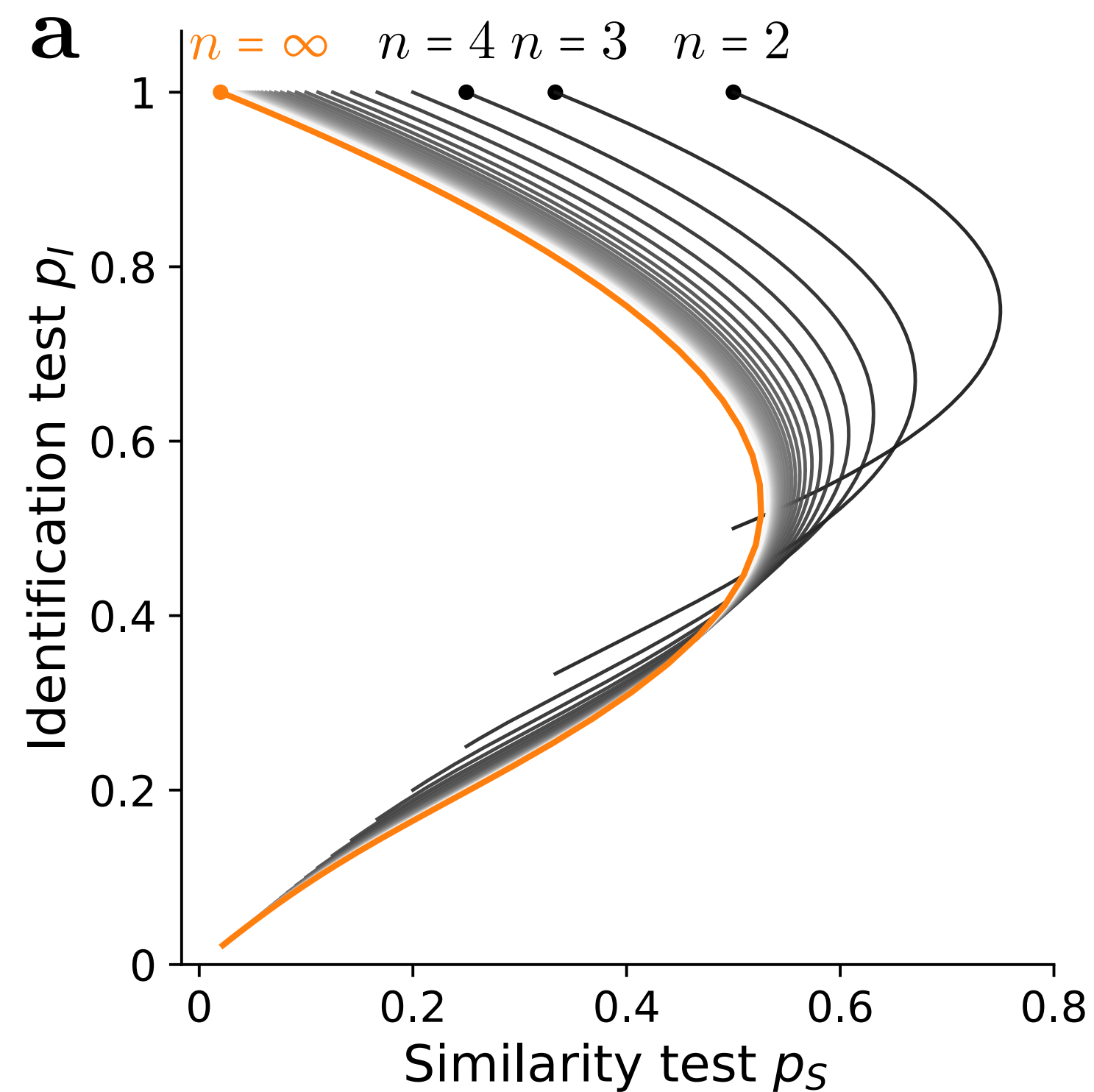
# Generalization vs Processing

## Multiple objects

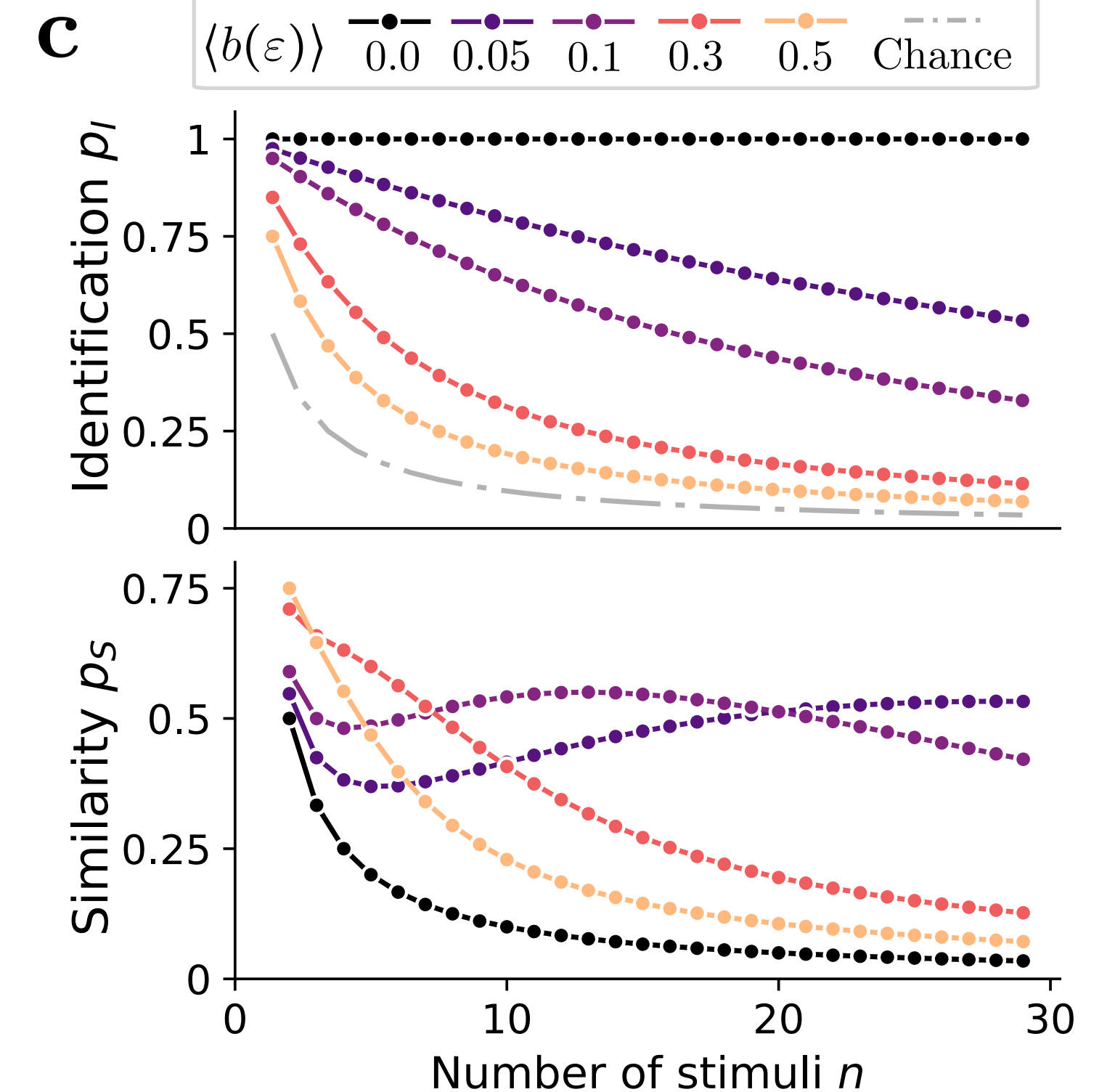
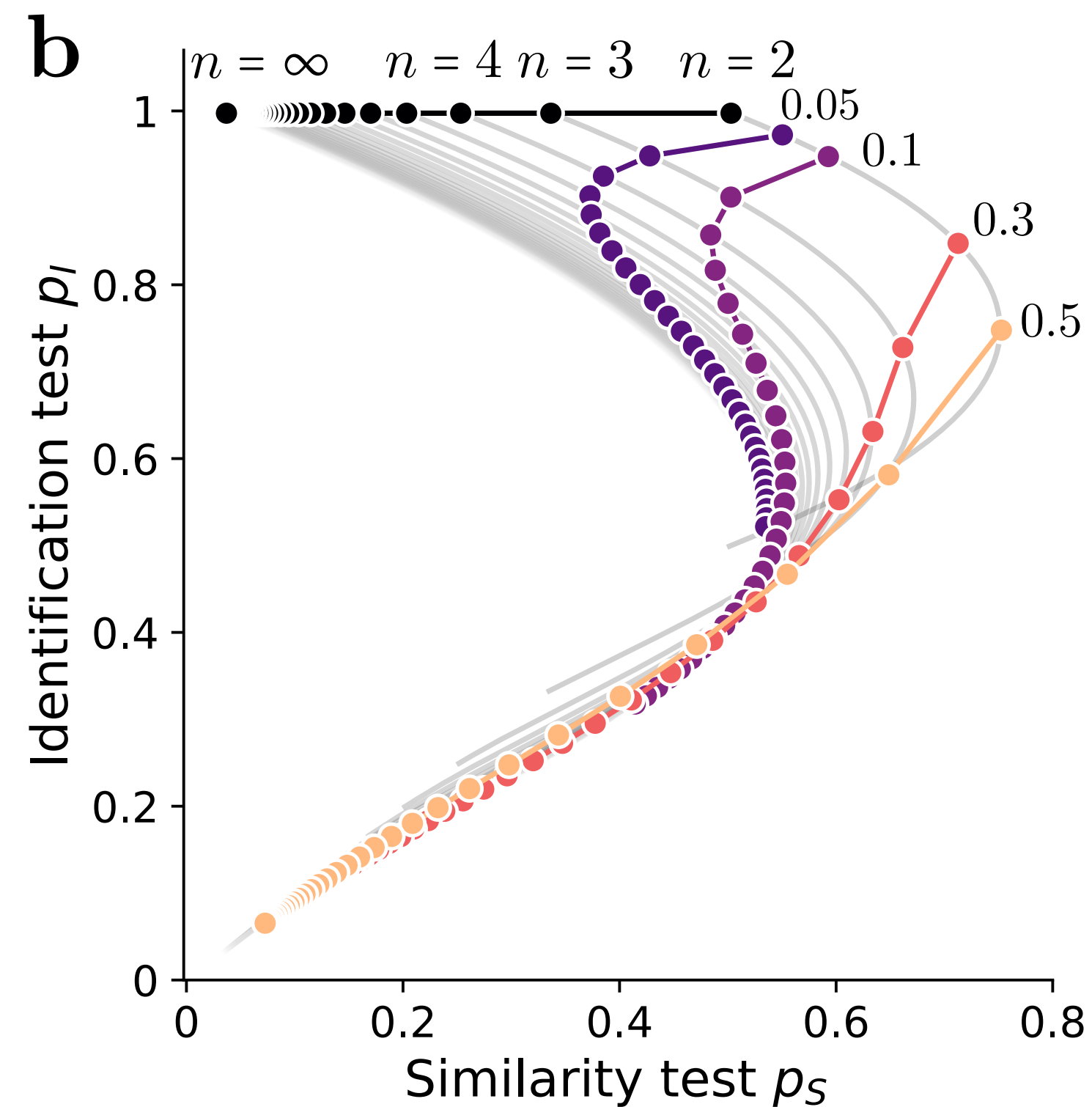
Bound by semanticity: universal laws governing the generalization-identification tradeoff

Marco Nurişso<sup>1,2</sup>, Jesseba Fernando<sup>3,4</sup>, Raj Deshpande<sup>5</sup>, Alan Perotti<sup>2</sup>, Raja Marjich<sup>6</sup>, Steven M. Frankland<sup>7</sup>, Richard L. Lewis<sup>8</sup>, Taylor W. Webb<sup>9</sup>, Declan Campbell<sup>10</sup>, Francesco Vaccarino<sup>1</sup>, Jonathan D. Cohen<sup>6,10</sup>, Giovanni Petri<sup>2,5,11</sup>

Multiple  $n$  - All resolutions



Multiple  $n$  - Fixed resolutions



**Theorem 3** ( $n$ -item tests). Under the same assumptions of Theorem 1, for the constant noise-free ( $\Delta = 0$ ) similarity function  $g = g_{\epsilon;0}$  we have that

$$p_S^n(\epsilon) = \mathbb{E}_{p \sim \nu} \left[ \frac{1}{n} + \sum_{k=1}^{n-1} \frac{(1 - b_p(\epsilon))^{n-k} - (1 - b_p(\epsilon))^n}{k} \right], \quad (7)$$

$$p_I^n(\epsilon) = \mathbb{E}_{p \sim \nu} \left[ \frac{1 - (1 - b_p(\epsilon))^n}{n b_p(\epsilon)} \right]. \quad (8)$$

# Binary code networks

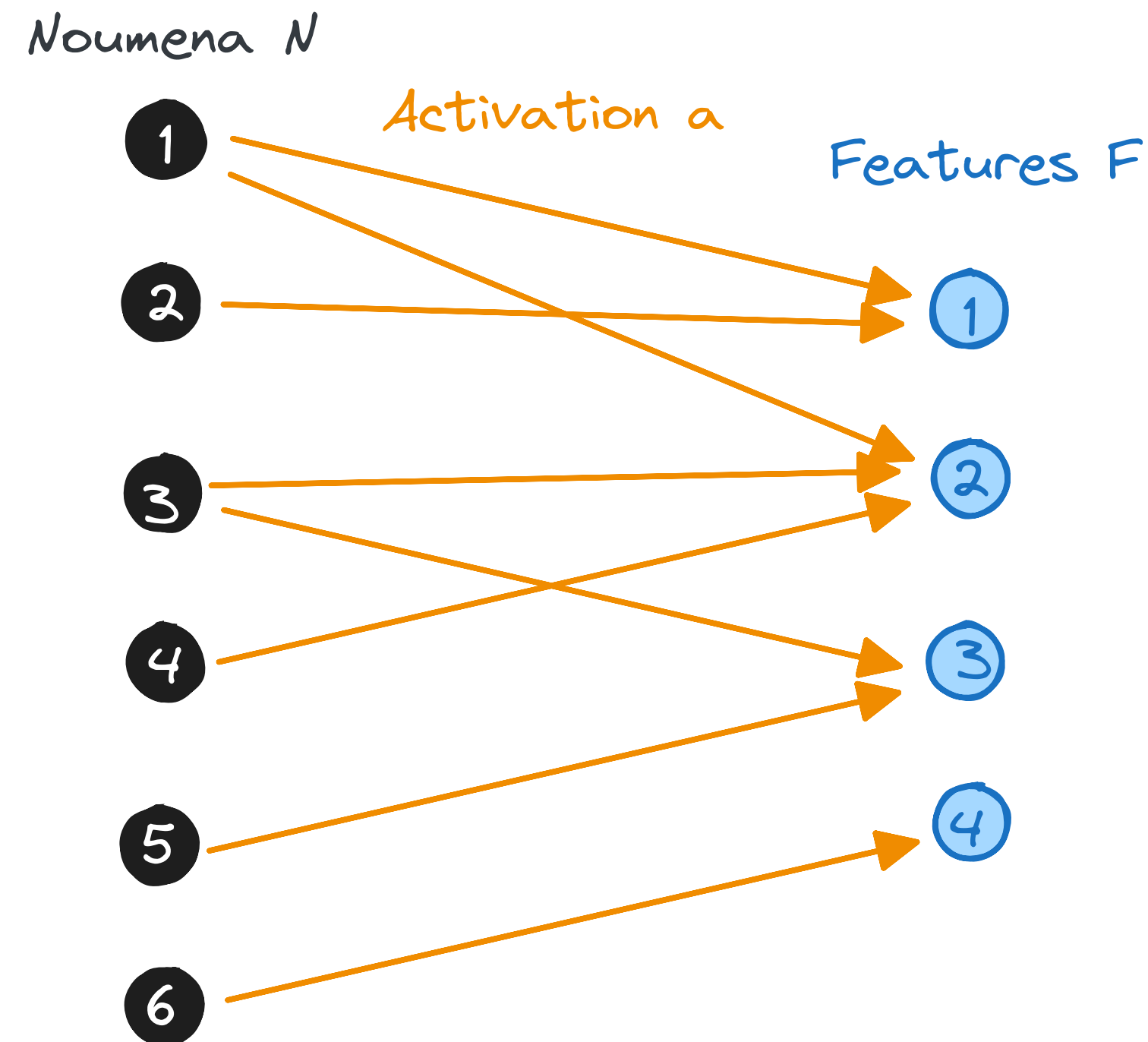
# Representational vs processing capacity

# Representational vs processing capacity

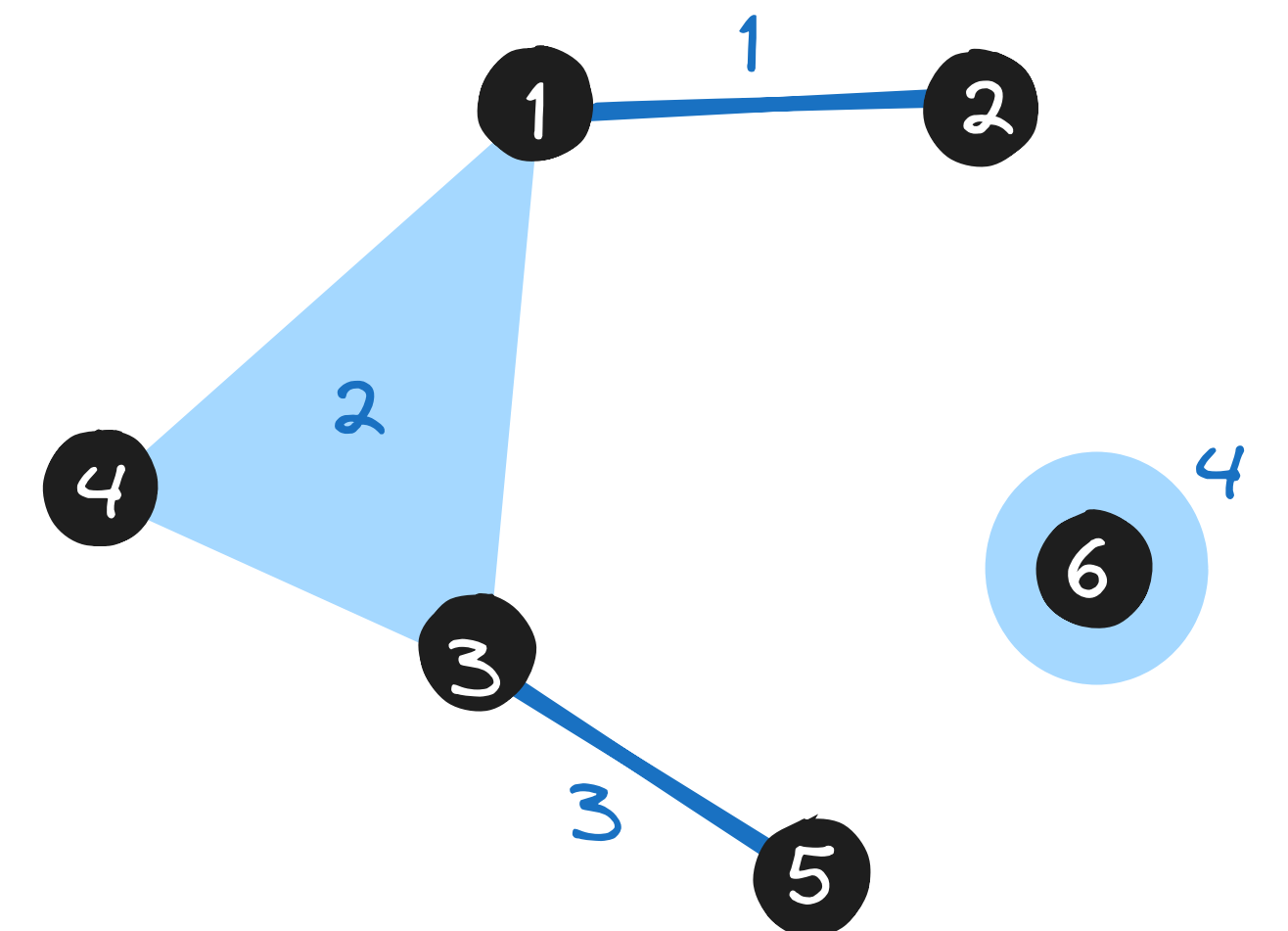
A code of  $N$  is a map

$$a : N \rightarrow 2^F$$

Bipartite graph representation



Hypergraph representation



# Representational vs processing capacity

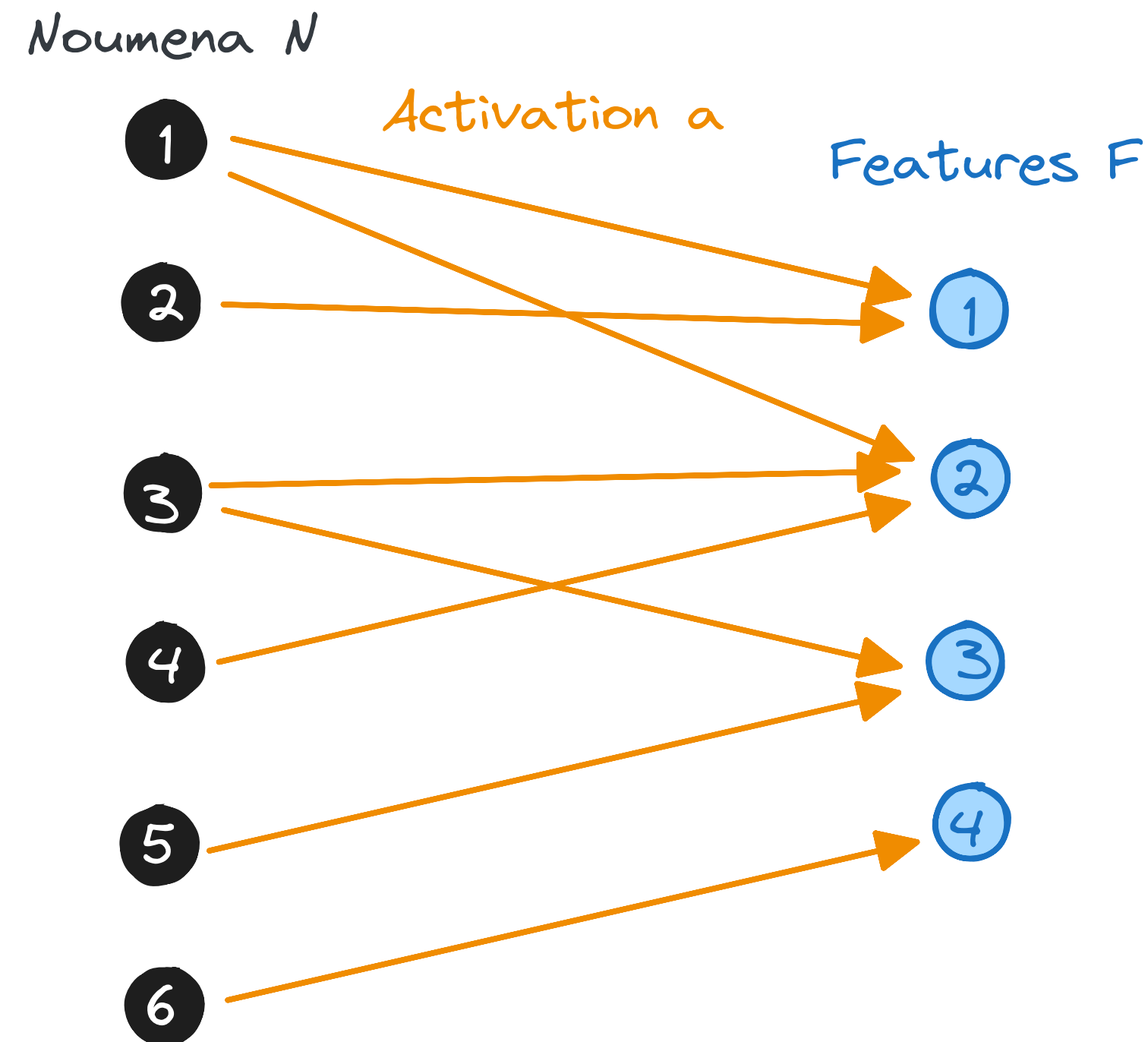
A code of  $N$  is a map

$$a : N \rightarrow 2^F$$

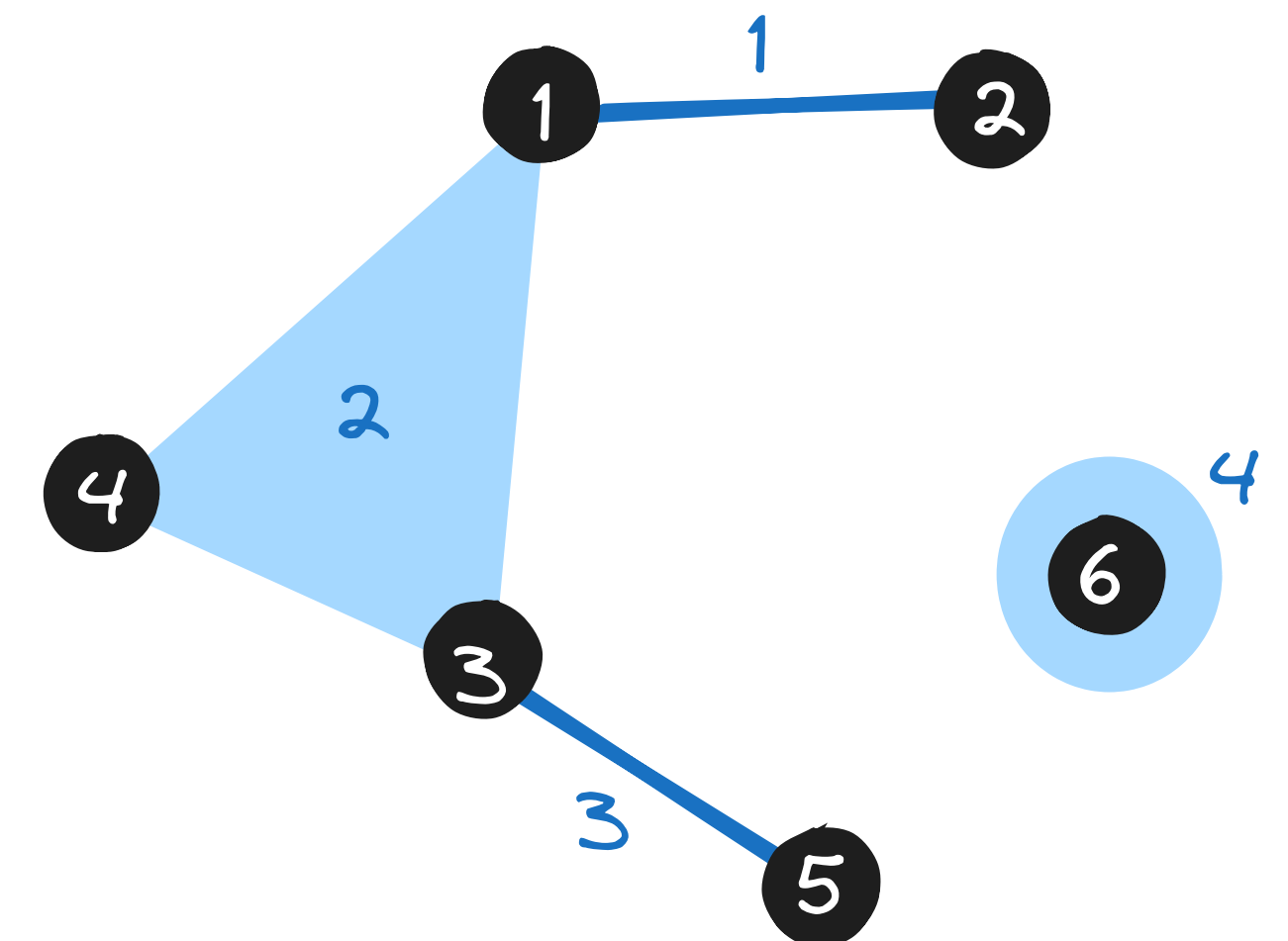
1-to-1 code (deterministically)

$$a(x) \neq a(y) \text{ if } x \neq y$$

Bipartite graph representation



Hypergraph representation



# Representational vs processing capacity

A code of  $N$  is a map

$$a : N \rightarrow 2^F$$

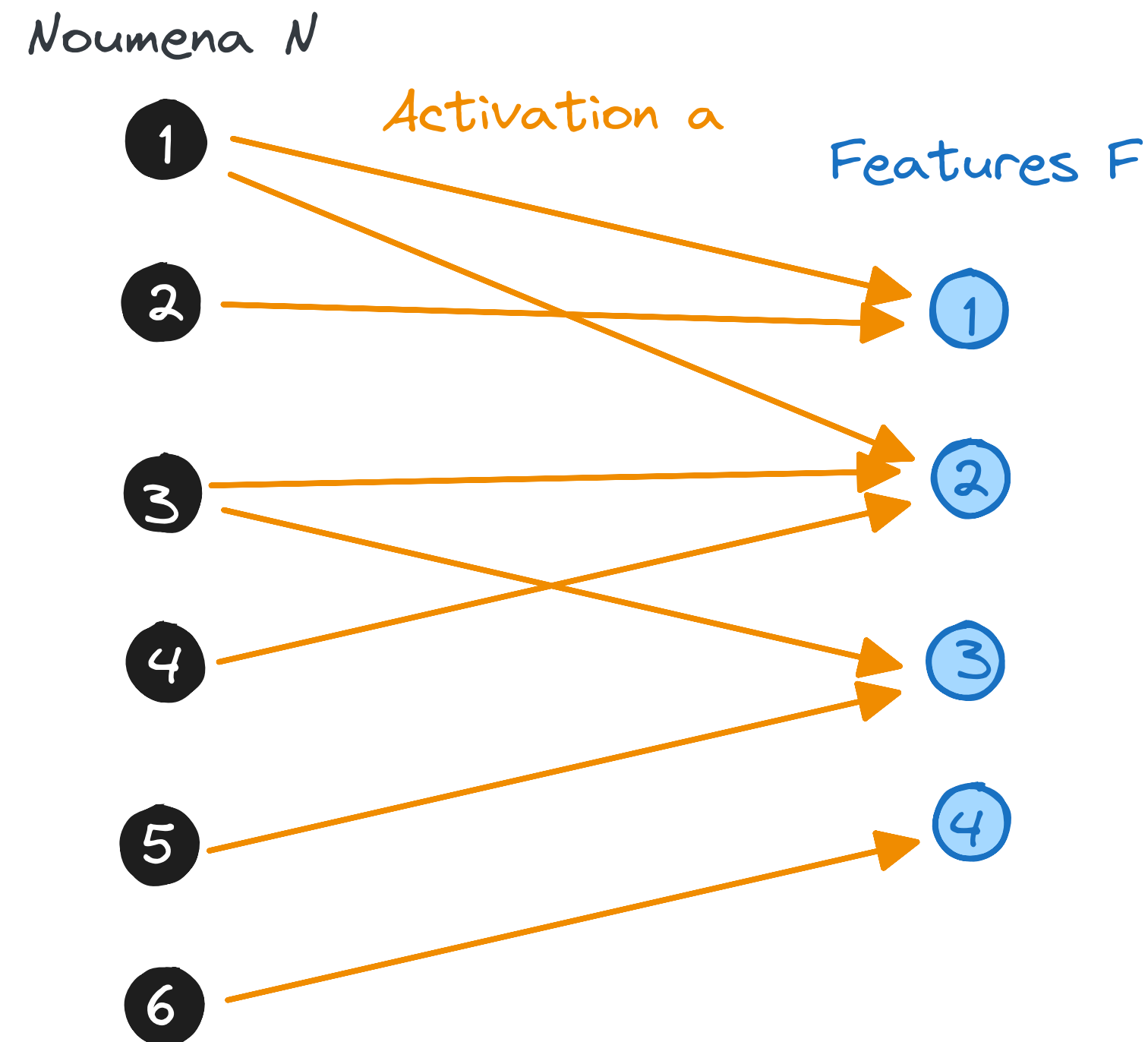
1-to-1 code (deterministically)

$$a(x) \neq a(y) \text{ if } x \neq y$$

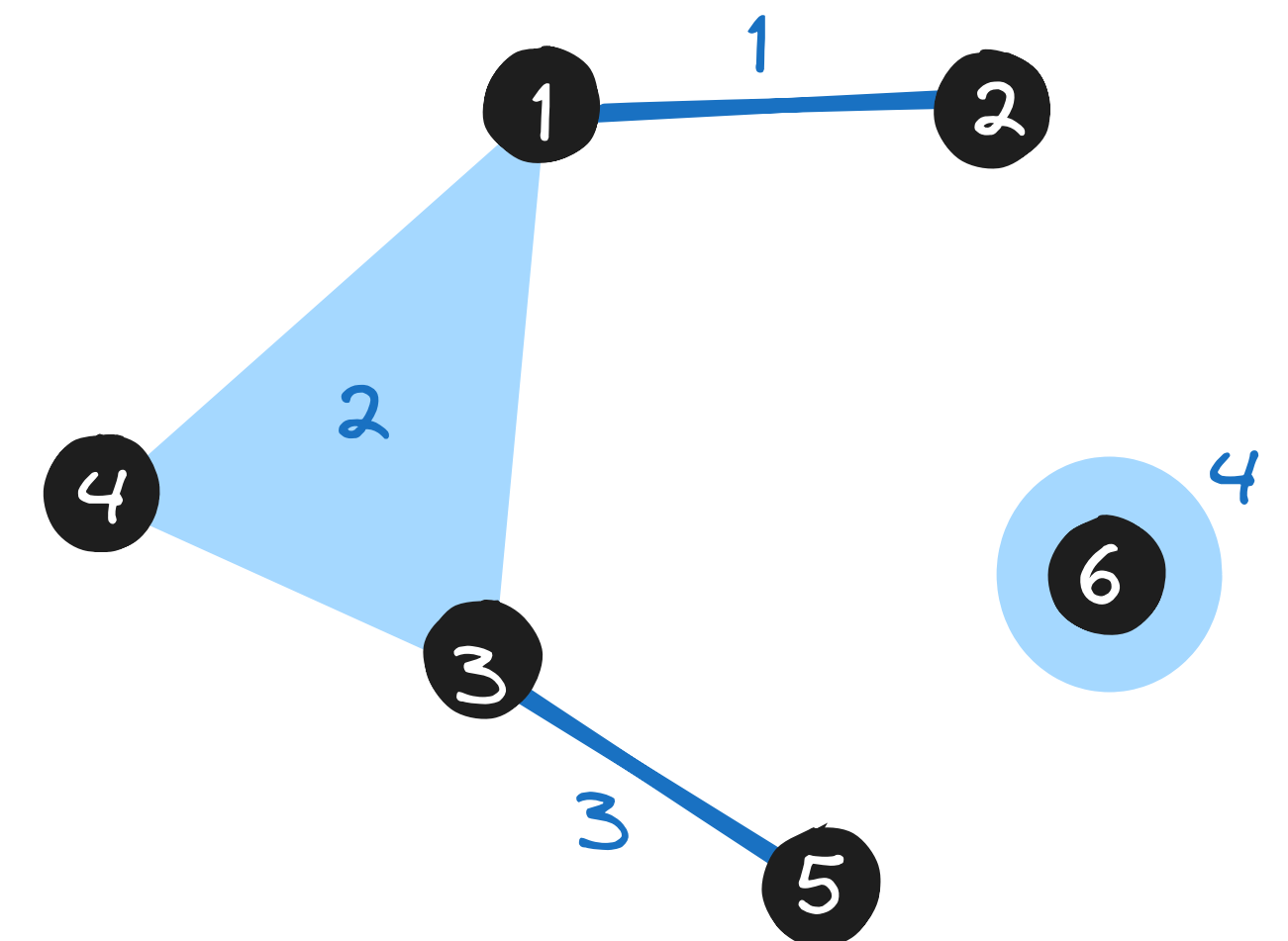
Dist  $p(x)$  on  $N$  induced dist  $\mathbb{P}(a(x))$

$$I(X; a(X)) = - \sum_{x \in N} p_x \log_2 \mathbb{P}(a(x)),$$

Bipartite graph representation



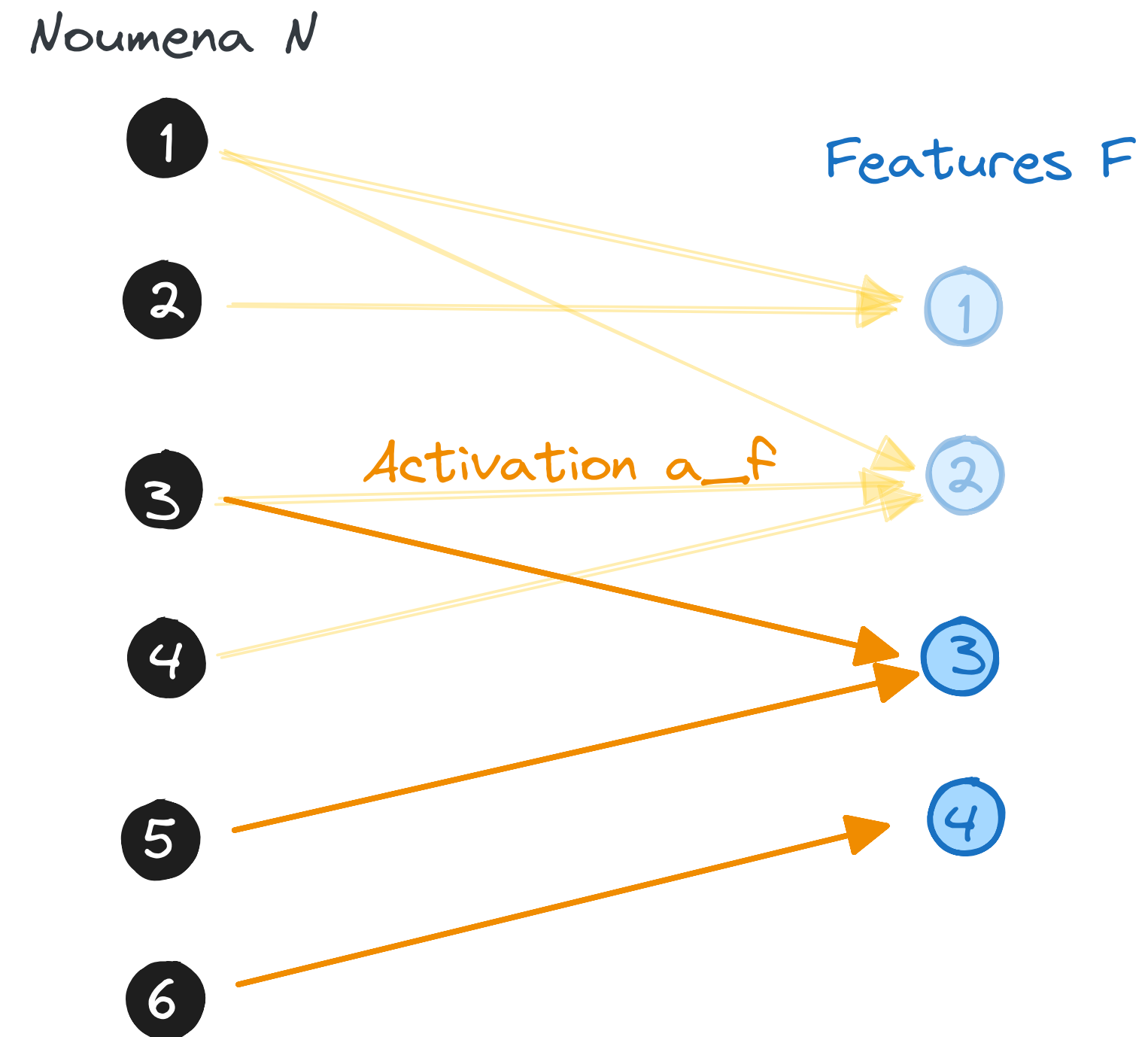
Hypergraph representation



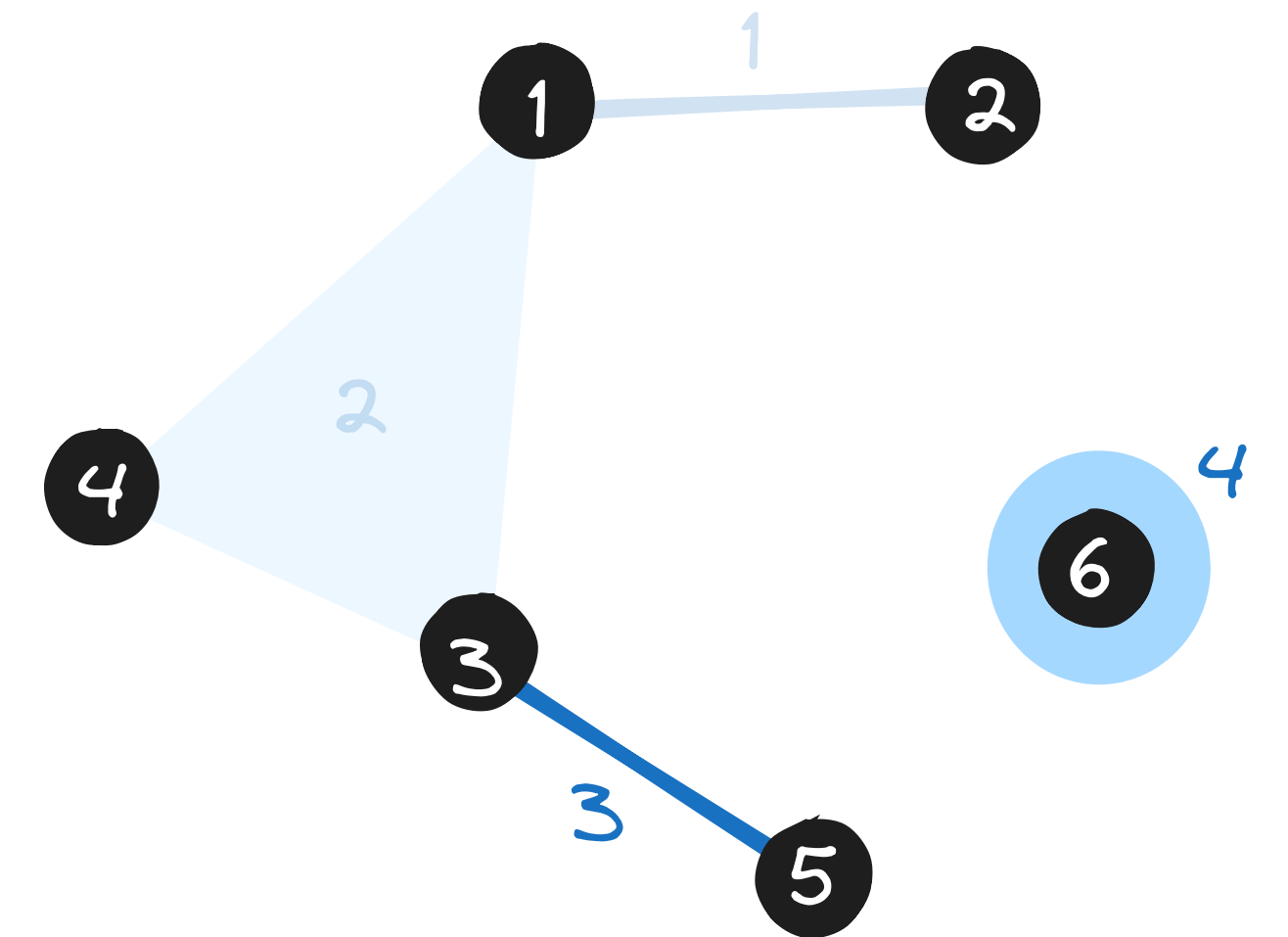
# Representational vs processing capacity

## Representation

Bipartite graph representation



Hypergraph representation



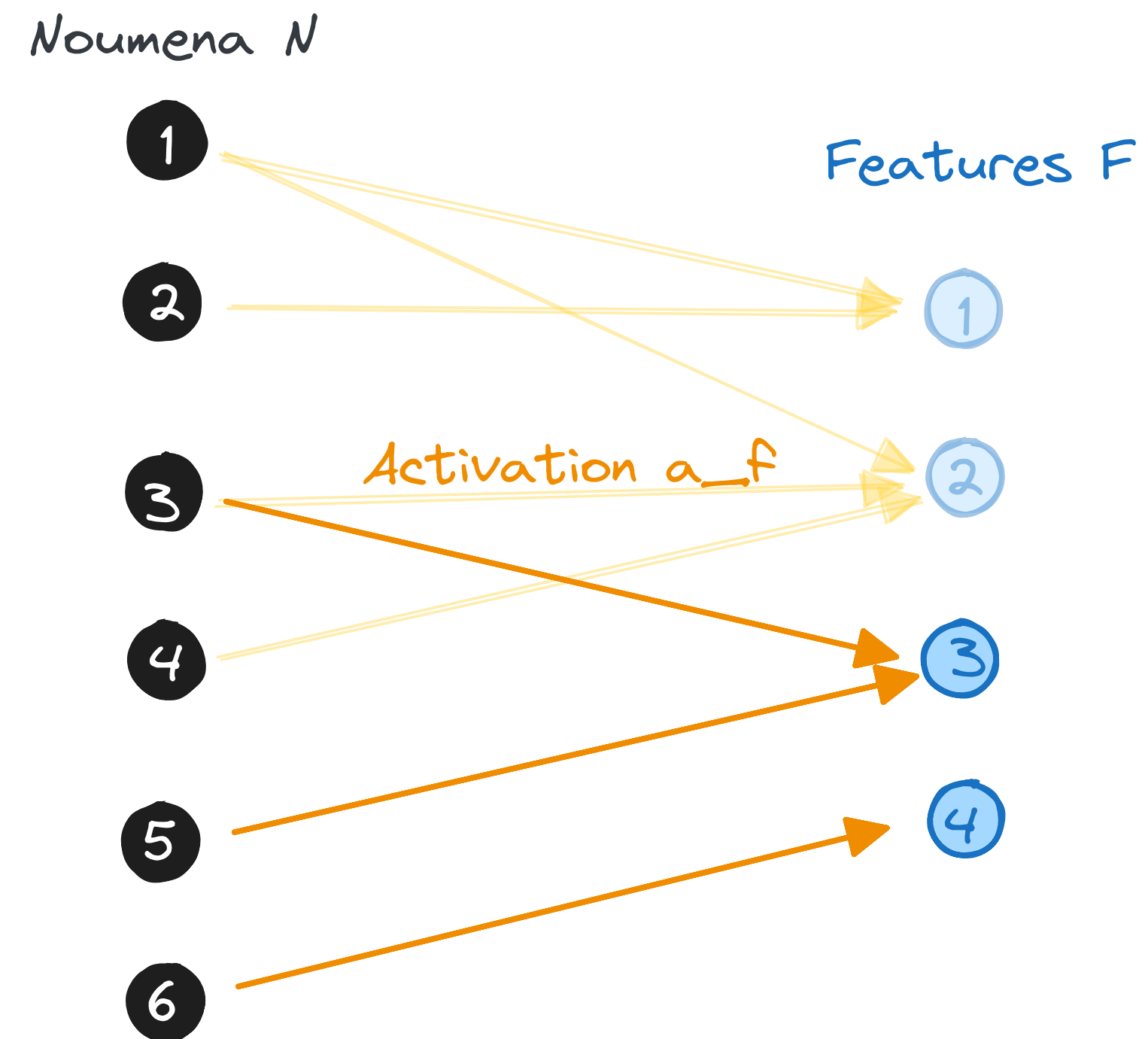
# Representational vs processing capacity

## Representation

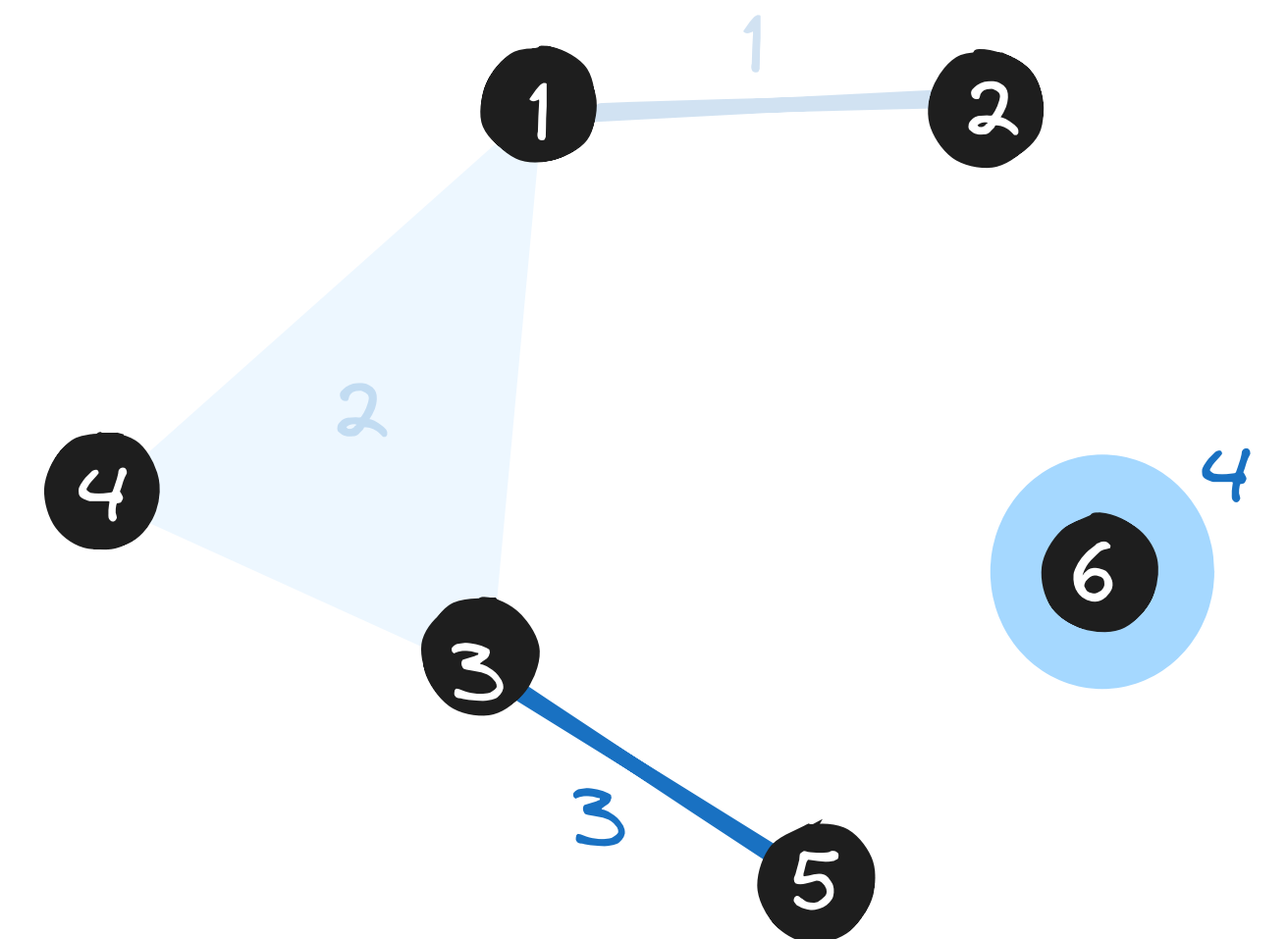
Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Bipartite graph representation



Hypergraph representation



# Representational vs processing capacity

## Representation

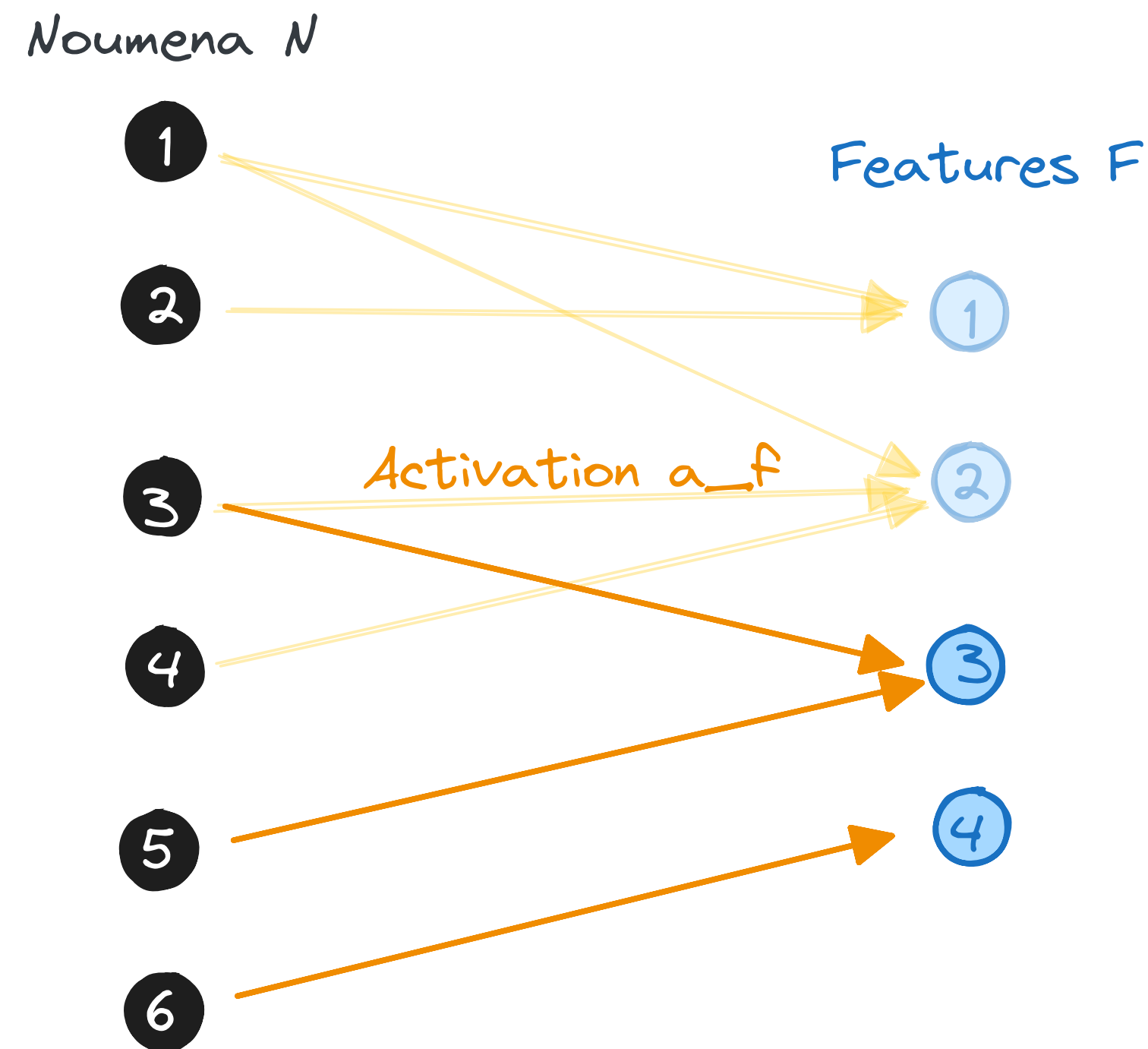
Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

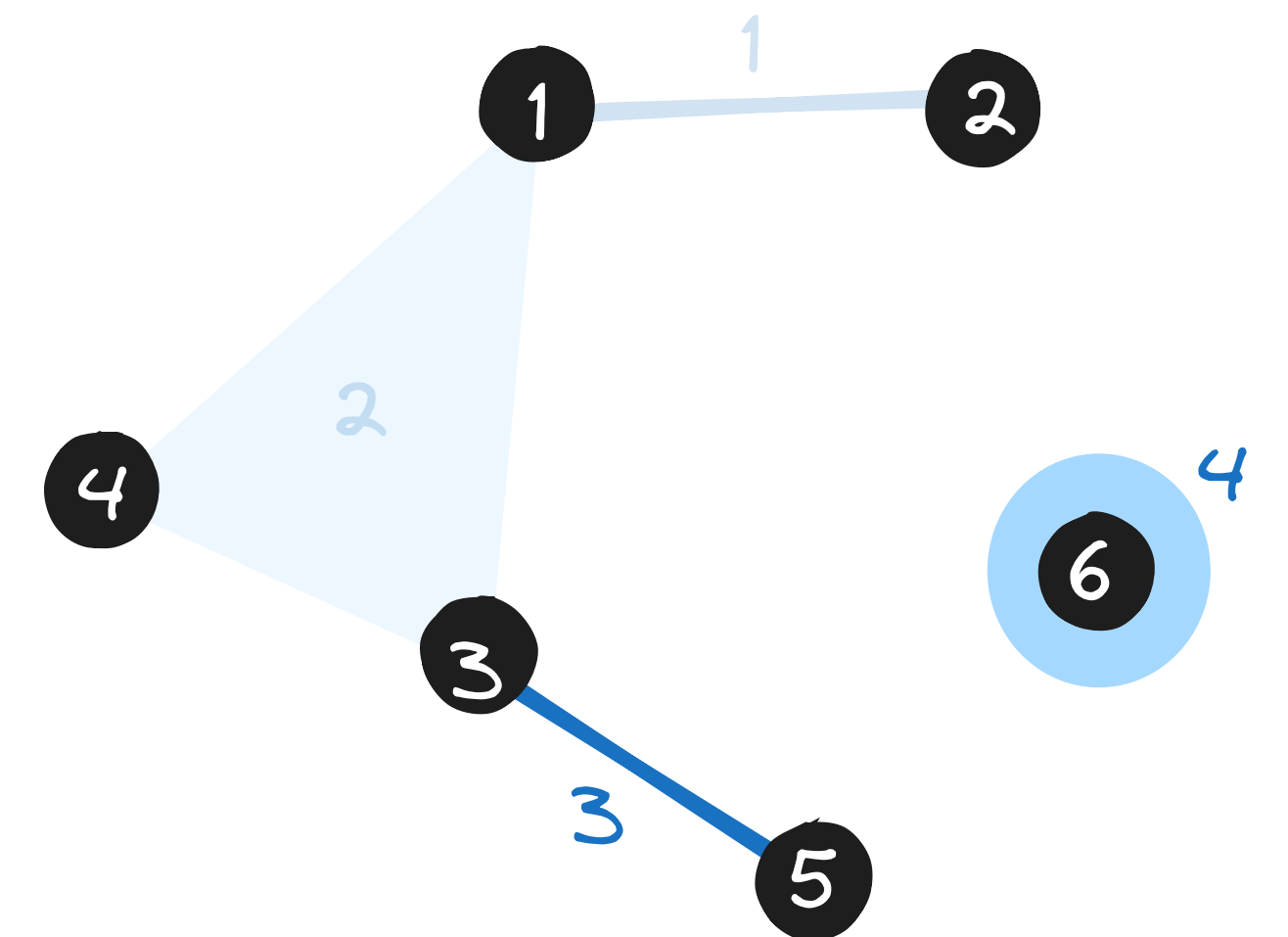
Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

Bipartite graph representation



Hypergraph representation



# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

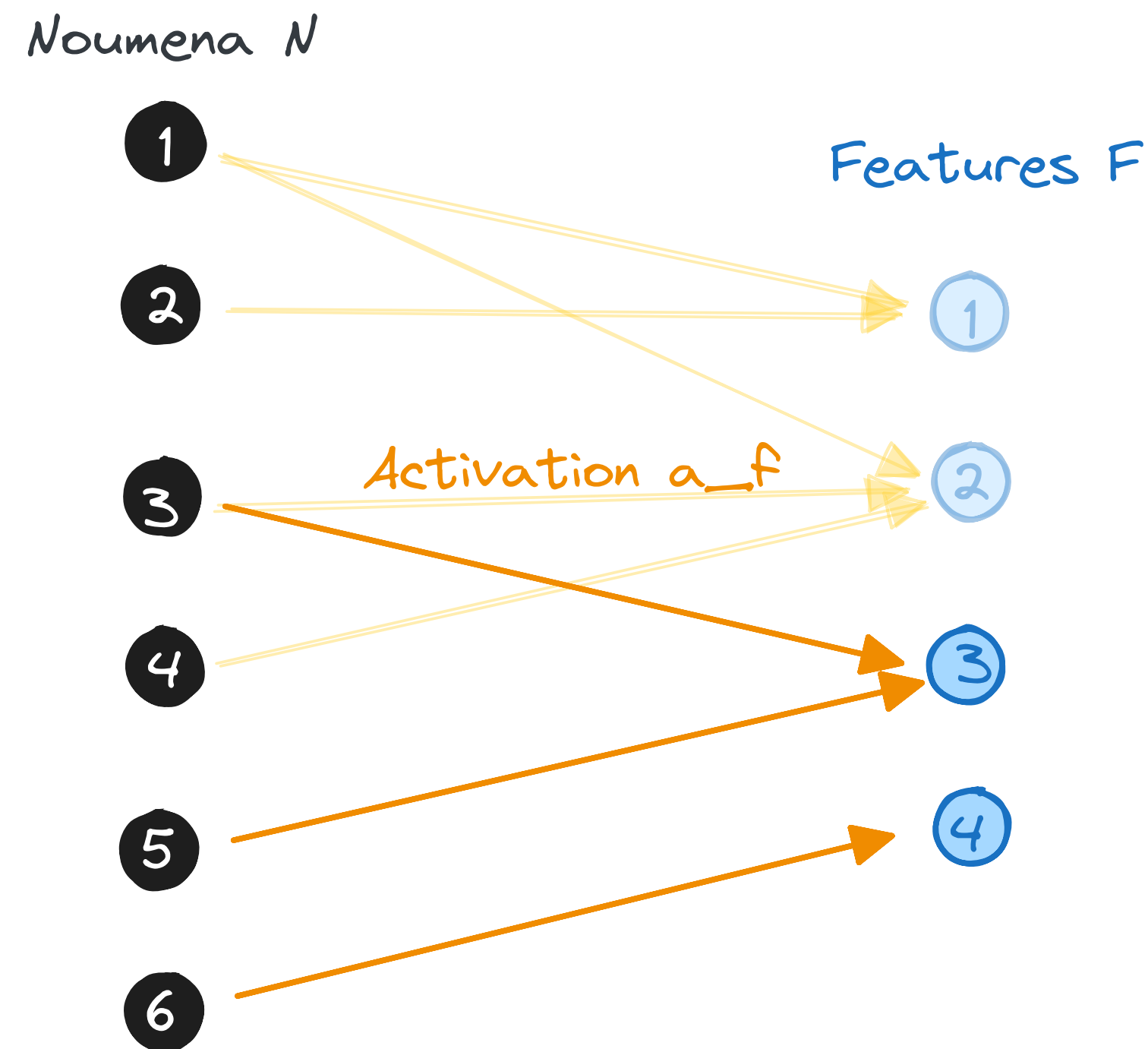
$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

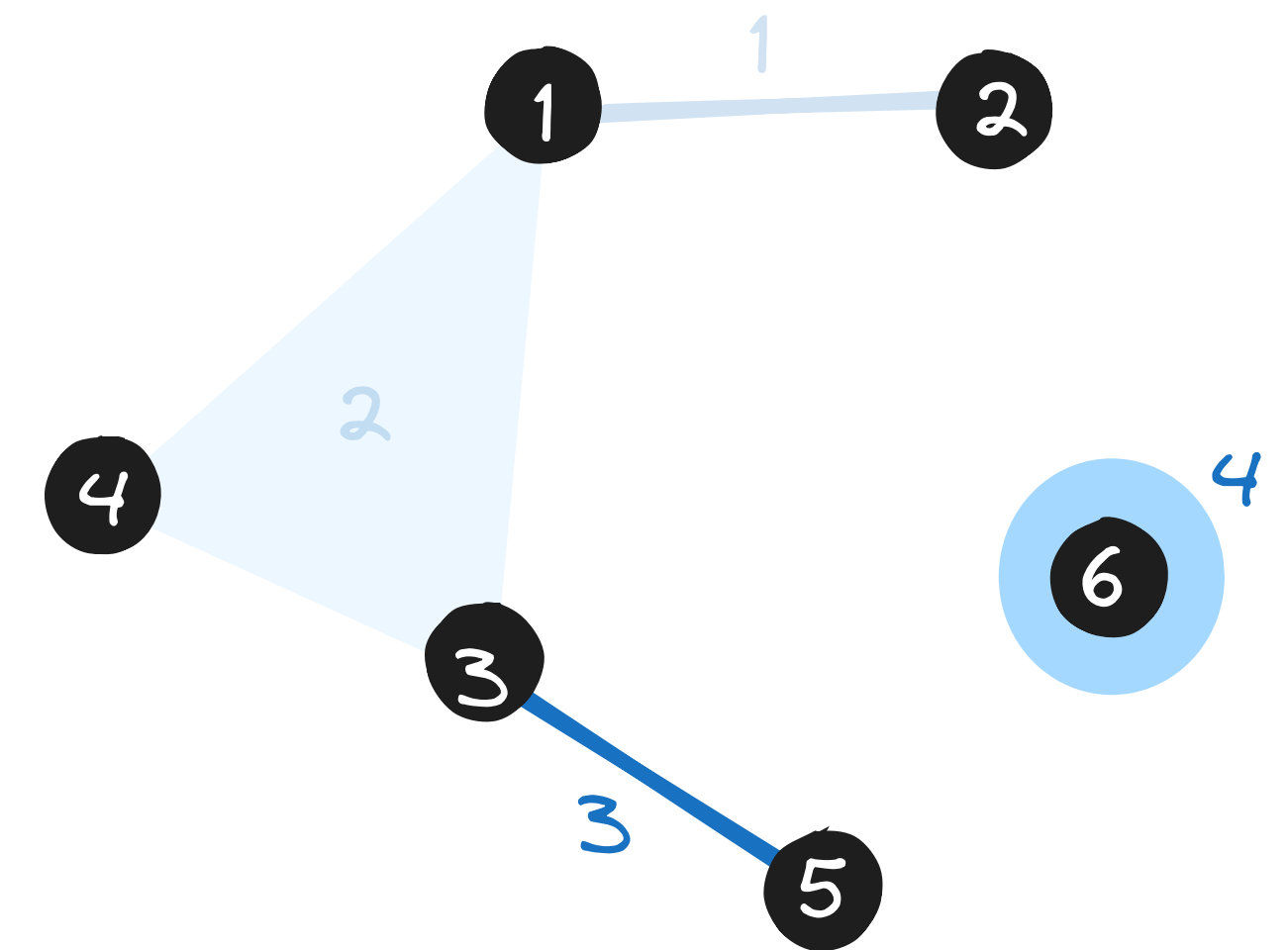
$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should increase  $I_r(X, a)$  rapidly with  $r$

Bipartite graph representation



Hypergraph representation



# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should  
increase  $I_r(X, a)$  rapidly with  $r$

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should  
increase  $I_r(X, a)$  rapidly with  $r$

## Processing

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should  
increase  $I_r(X, a)$  rapidly with  $r$

## Processing

Apple  $\rightarrow (1,0,0,0,1)$

Banana  $\rightarrow (1,1,0,0,1)$

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should  
increase  $I_r(X, a)$  rapidly with  $r$

## Processing

Apple  $\rightarrow (1,0,0,0,1)$

Banana  $\rightarrow (1,1,0,0,1)$

(Apple, banana)  $\rightarrow (1,1,0,0,1)$

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should  
increase  $I_r(X, a)$  rapidly with  $r$

## Processing

Apple  $\rightarrow (1,0,0,0,1)$

Banana  $\rightarrow (1,1,0,0,1)$       (Apple, banana)  $\rightarrow (1,1,0,0,1)$

$$a^k : \binom{N}{k} \rightarrow 2^F, \quad a^k((x_1, \dots, x_k)) = \bigcup_{i=1}^k a(x_i)$$

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^f, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should increase  $I_r(X, a)$  rapidly with  $r$

## Processing

Apple  $\rightarrow (1,0,0,0,1)$

Banana  $\rightarrow (1,1,0,0,1)$       (Apple, banana)  $\rightarrow (1,1,0,0,1)$

$$a^k : \binom{N}{k} \rightarrow 2^F, \quad a^k((x_1, \dots, x_k)) = \bigcup_{i=1}^k a(x_i)$$

$$I^k(X; a) = \sum_{S \in \binom{N}{k}} \mathbb{P}[S] \log_2 \mathbb{P}(a^k(S)).$$

# Representational vs processing capacity

## Representation

Feature restriction  $f$ :

$$a_f : N \rightarrow 2^F, \quad a_f(x) = a(x) \cap f,$$

Info coded by  $f$ :

$$\begin{aligned} I_r(X; a) &:= \mathbb{E}_{f \in \binom{N}{r}} [I(x, a_f(X))] \\ &= \sum_{f \in \binom{F}{r}} \mathbb{P}[f] I(X; a_f(X)), \end{aligned}$$

A generalizing code should increase  $I_r(X, a)$  rapidly with  $r$

## Processing

Apple  $\rightarrow (1,0,0,0,1)$

Banana  $\rightarrow (1,1,0,0,1)$       (Apple, banana)  $\rightarrow (1,1,0,0,1)$

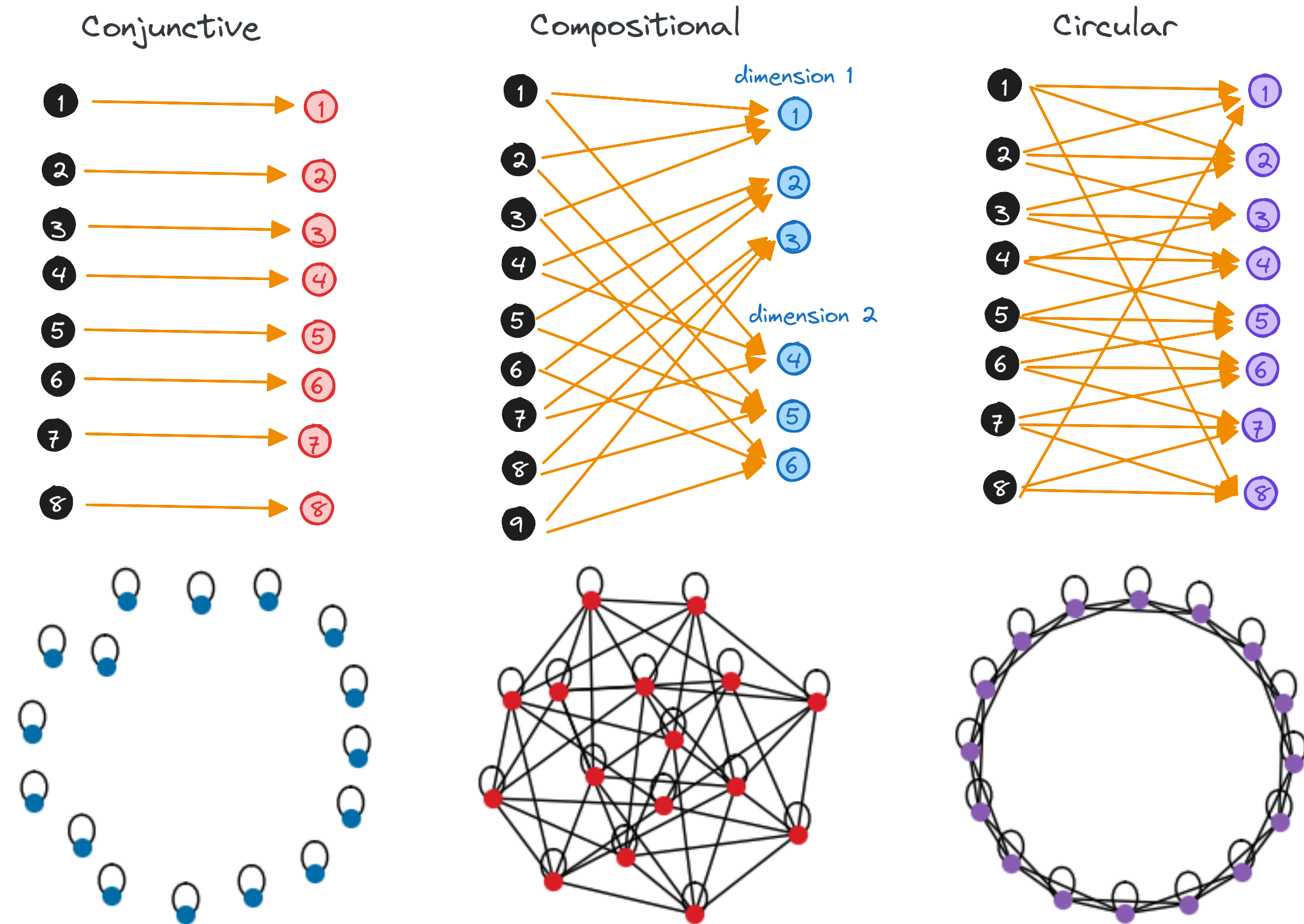
$$a^k : \binom{N}{k} \rightarrow 2^F, \quad a^k((x_1, \dots, x_k)) = \bigcup_{i=1}^k a(x_i)$$

$$I^k(X; a) = \sum_{S \in \binom{N}{k}} \mathbb{P}[S] \log_2 \mathbb{P}(a^k(S)).$$

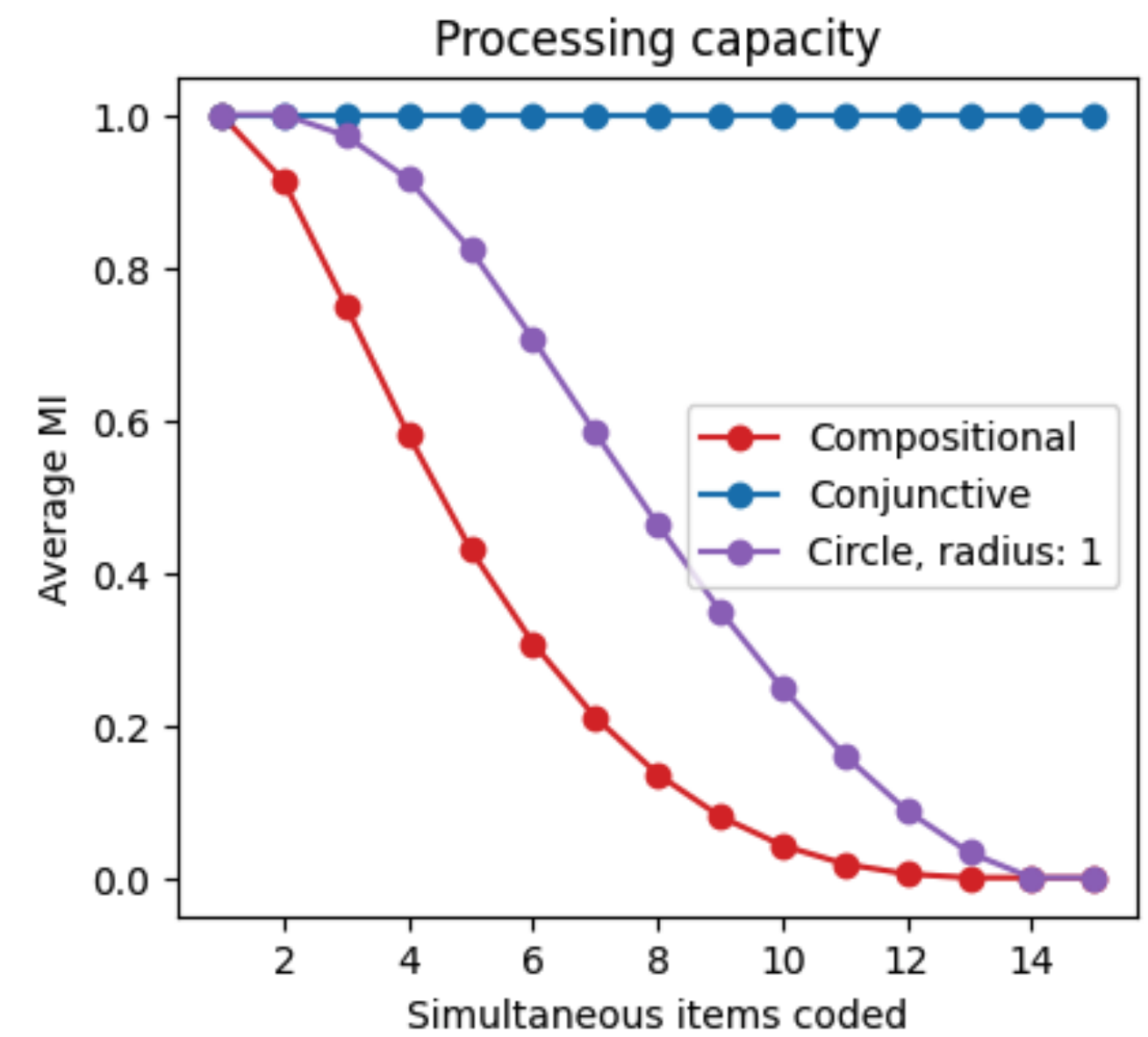
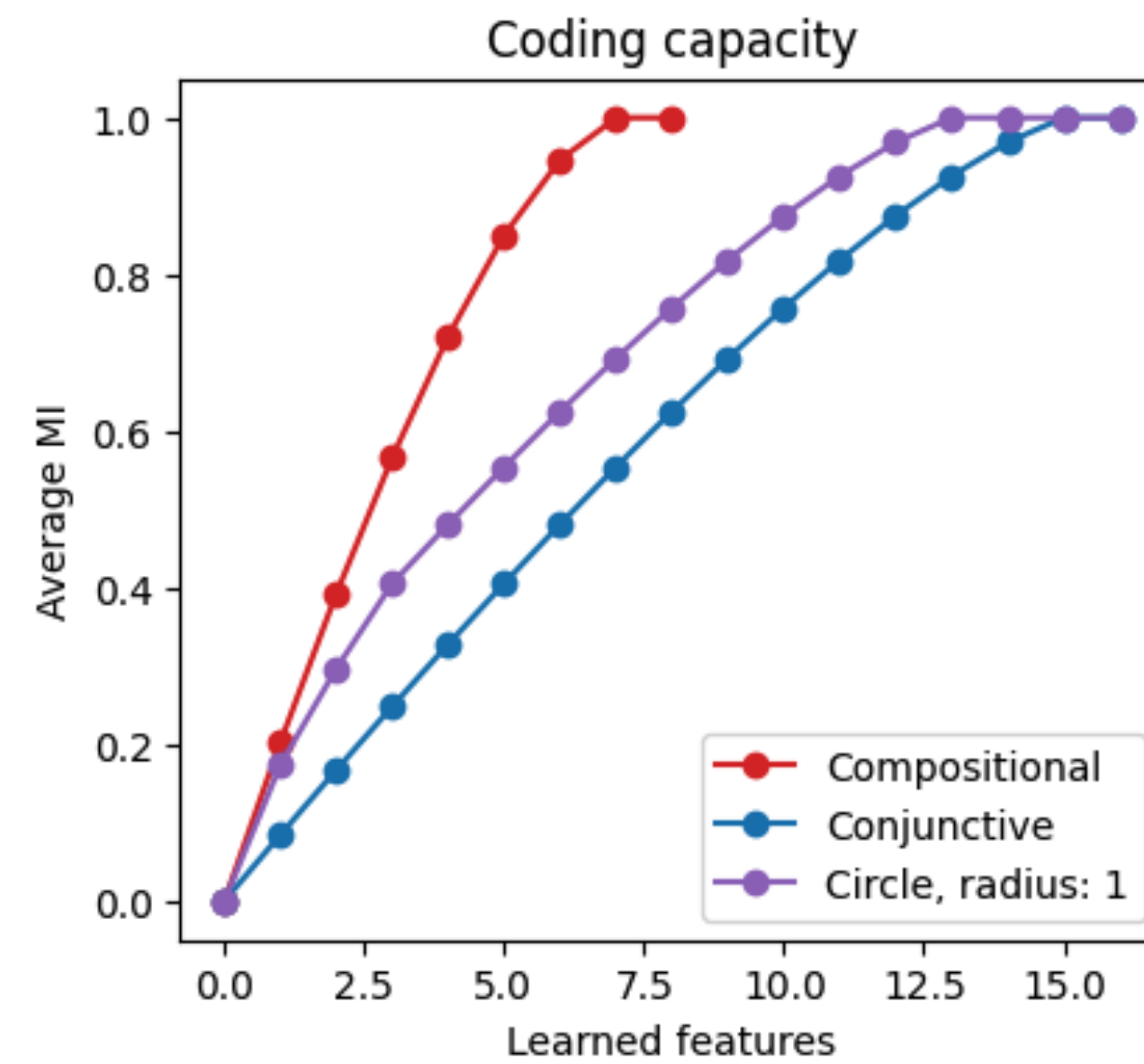
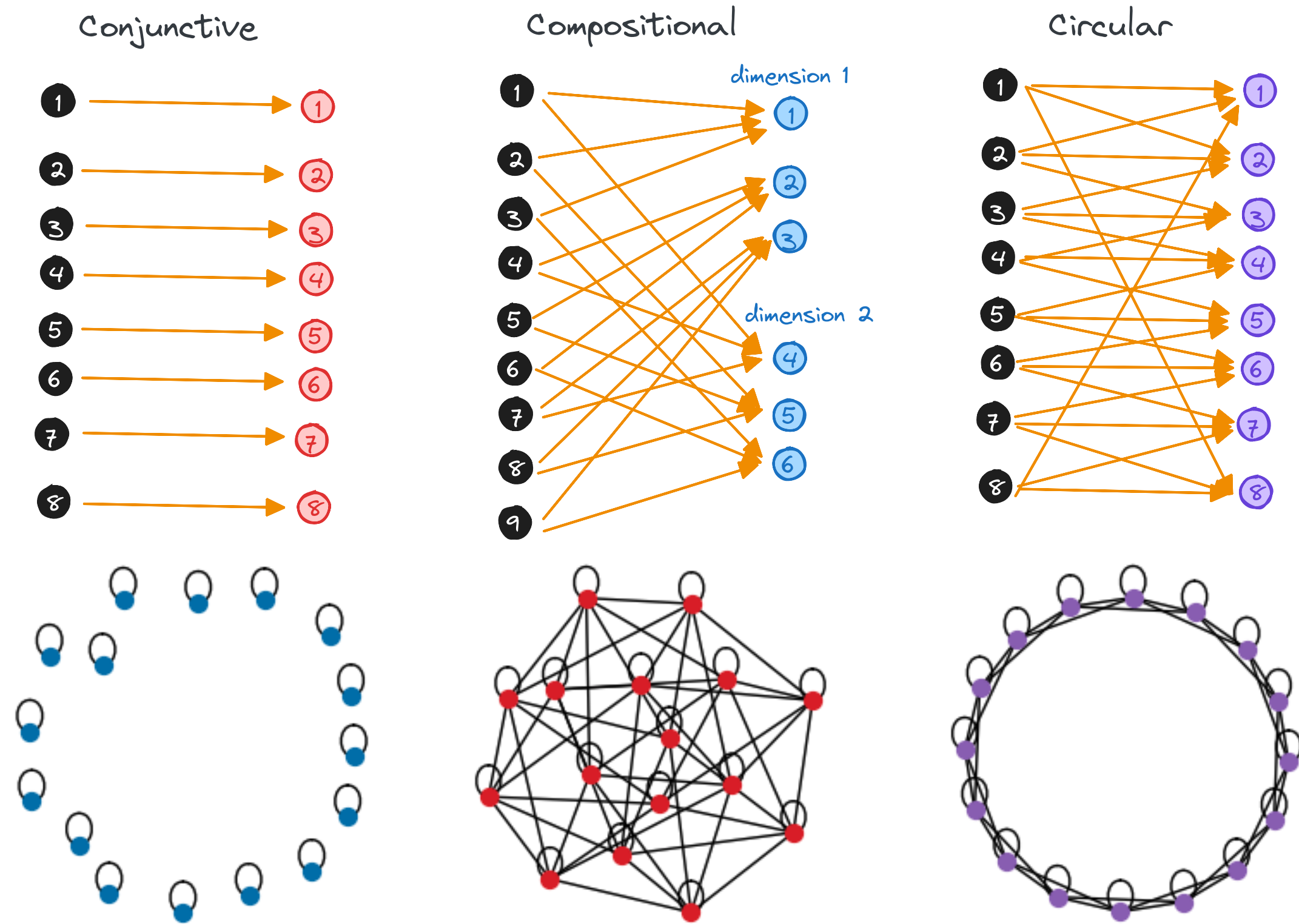
A good processing code should not decrease  $I^k(X, a)$  with  $k$

# Representational vs processing capacity

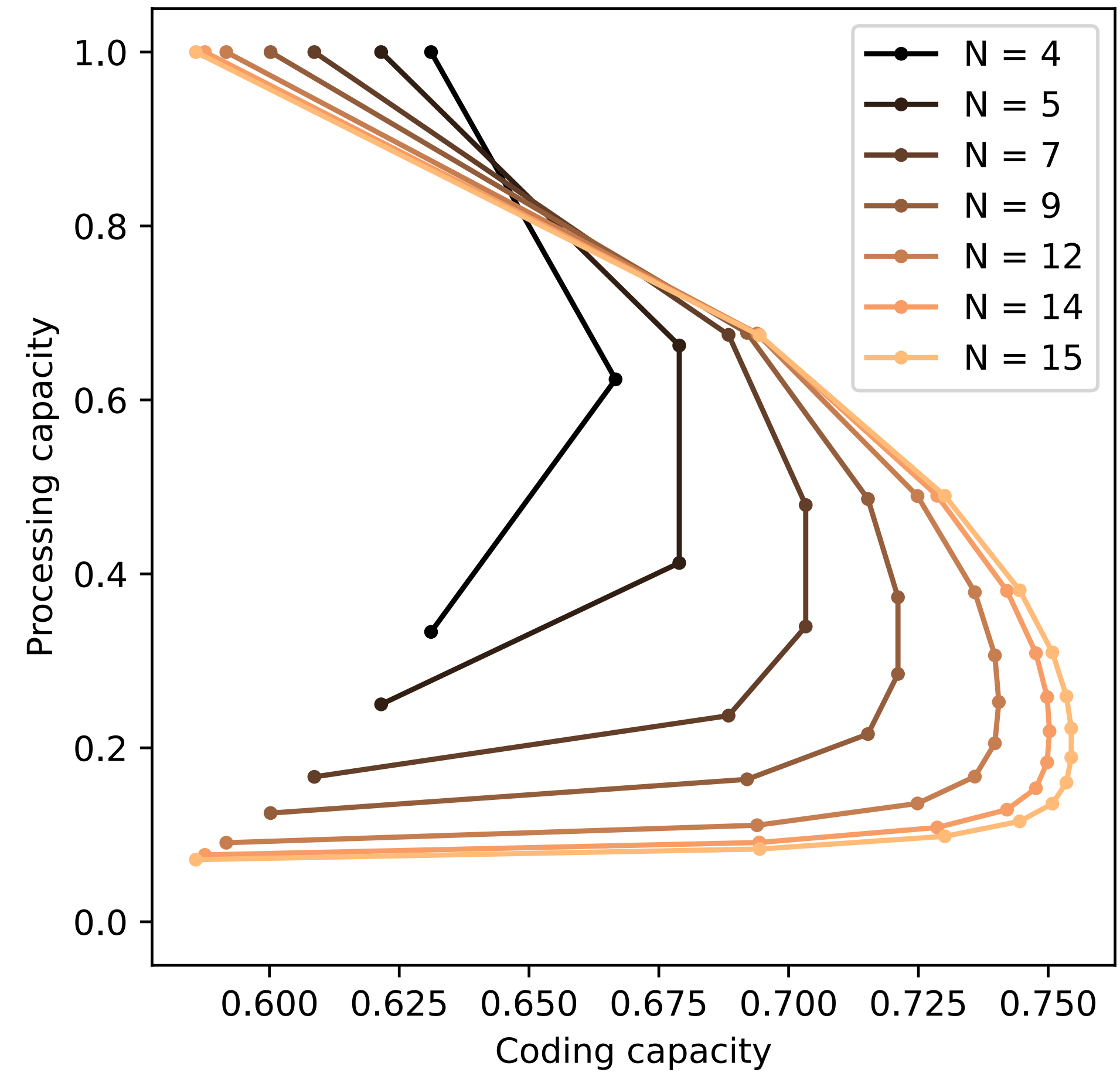
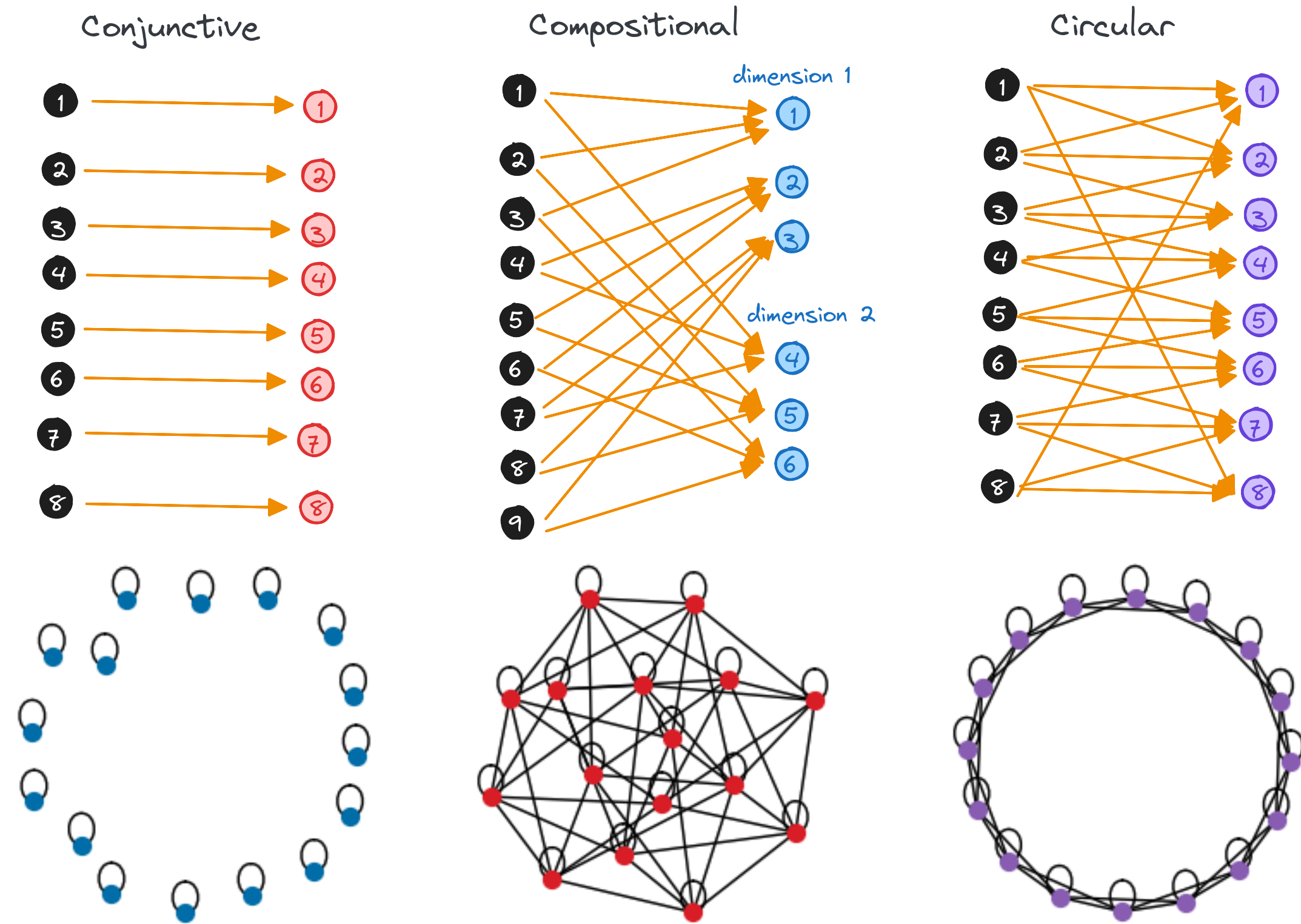
# Representational vs processing capacity



# Representational vs processing capacity



# Representational vs processing capacity



# Emergence of resolution in simple neural networks



# Implementation

Consider the *toy model of superposition* architecture  $f(x) = \sigma(W^\top Wx)$ .

The data is composed of  $n$  items (encoded as one-hots) + a distance matrix  $D$  (circle)

$i$ -th embedding  $z_i = We_i = W_i$

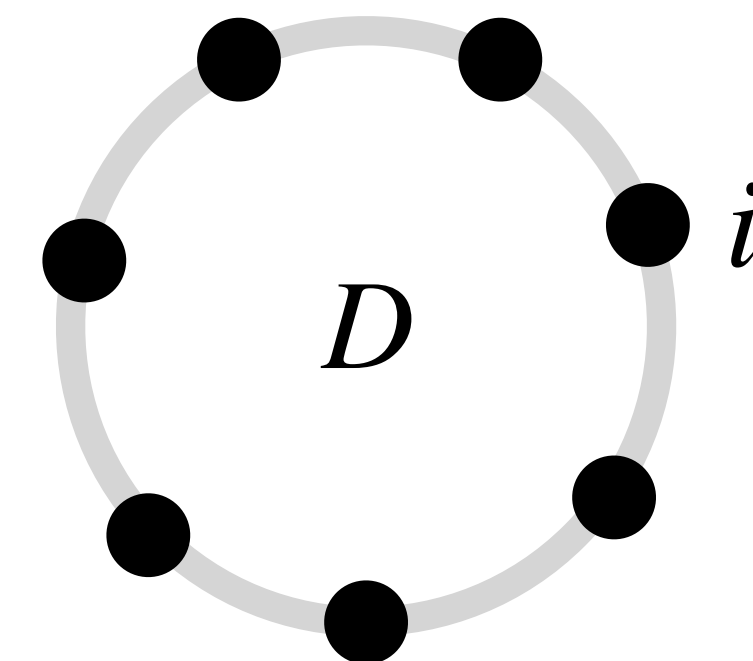
I train it using two losses:

1. **Reconstruction loss**  $L_{rec} = \frac{1}{n} \sum_{i=1}^n \|e_i - f(e_i)\|_2^2$

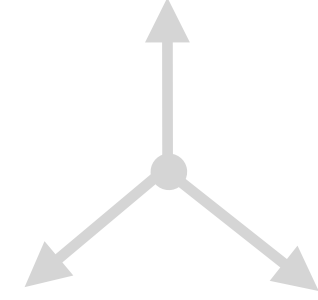
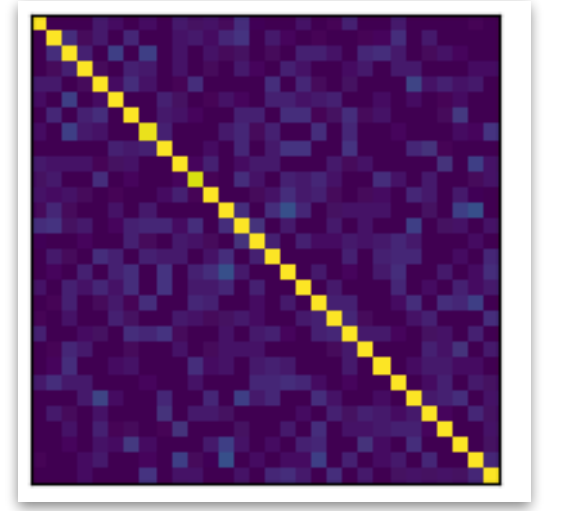
2. **Similarity test loss** = pick a triplet  $(i, j, k)$  and compute the probability

$$p_{ik} = \frac{\sigma(z_i^\top z_k)}{\sigma(z_i^\top z_k) + \sigma(z_j^\top z_k)}$$

and then compute Cross Entropy loss  $L_{sim}$  against the true index of the closest one w.r.t.  $D$

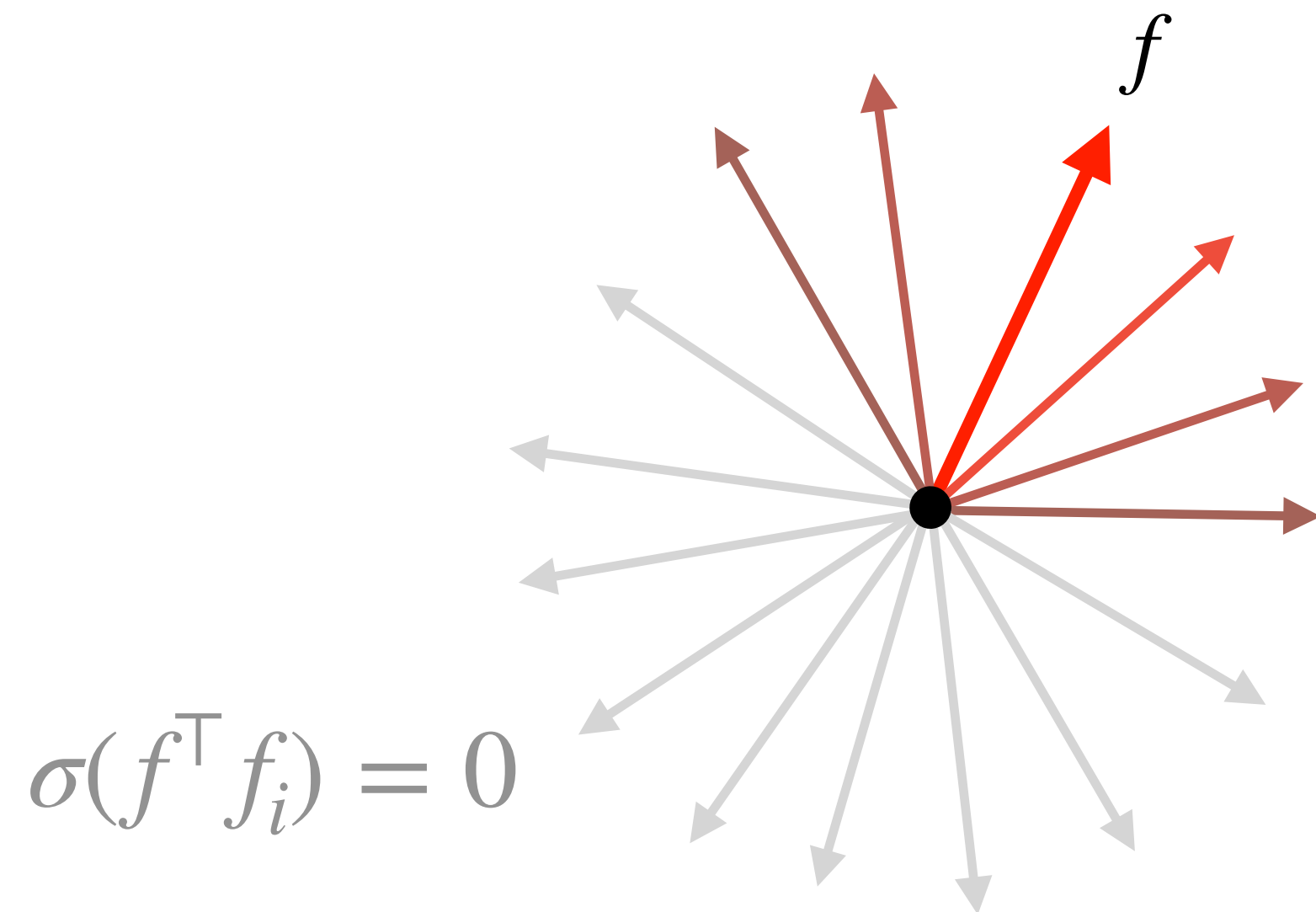


**Reconstruction**  $\cong$  Identification pushes for the embeddings to be “ReLU-orthogonal”

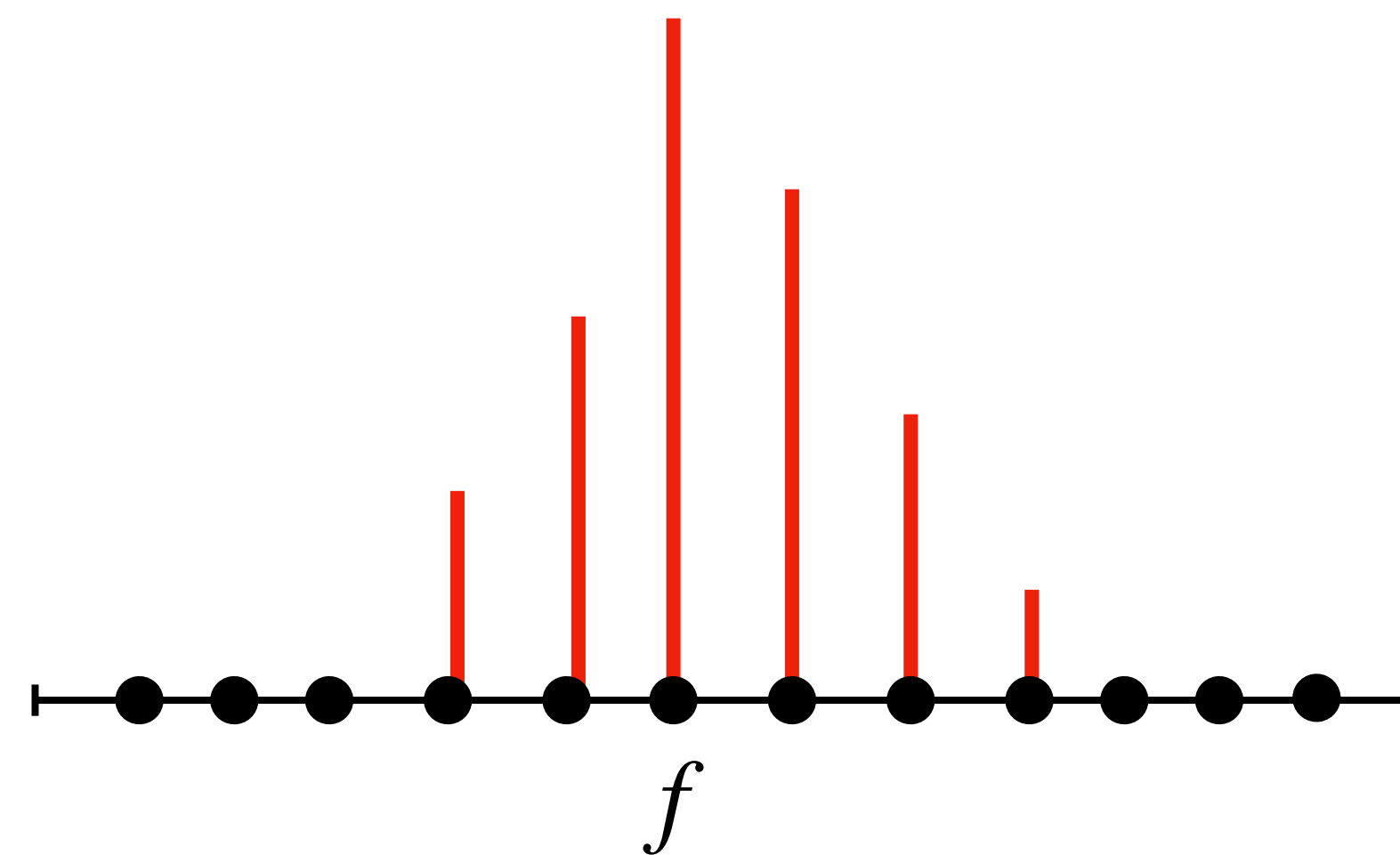


**Similarity** pushes for the embeddings to adapt to the metric space structure

Reconstruction drives the model to put many features in a quasi-orthogonal state  
 $\implies$  Similarity encourages a “metric” resolution  $\implies$  Miller’s law (?)

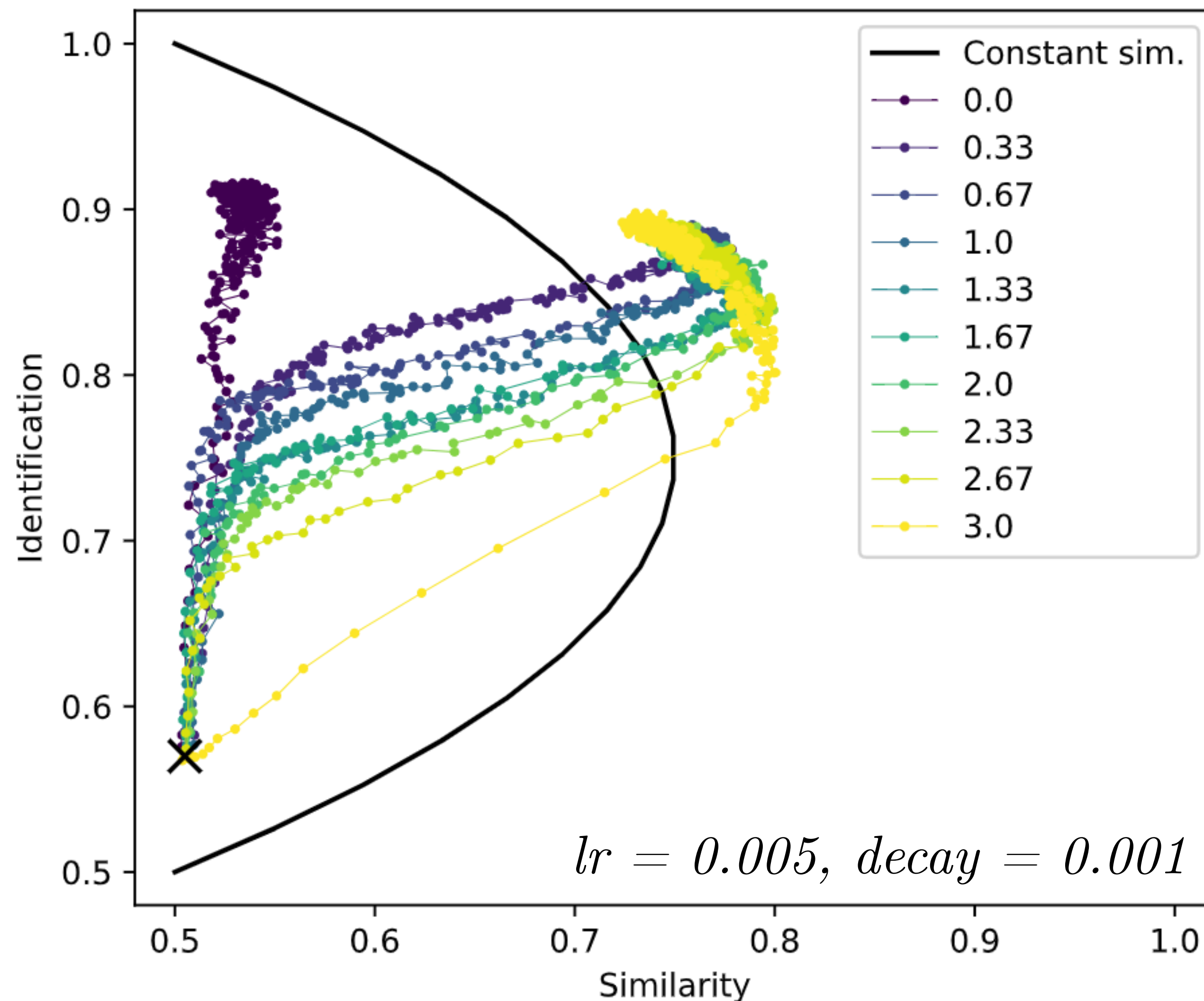


$\mathbb{R}$



$$\text{Total loss} = \lambda_{\text{rec}} L_{\text{rec}} + \lambda_{\text{sim}} L_{\text{sim}}$$

We can look at the training profiles when we change the values of  $\lambda_{\text{rec}}$ ,  $\lambda_{\text{sim}}$



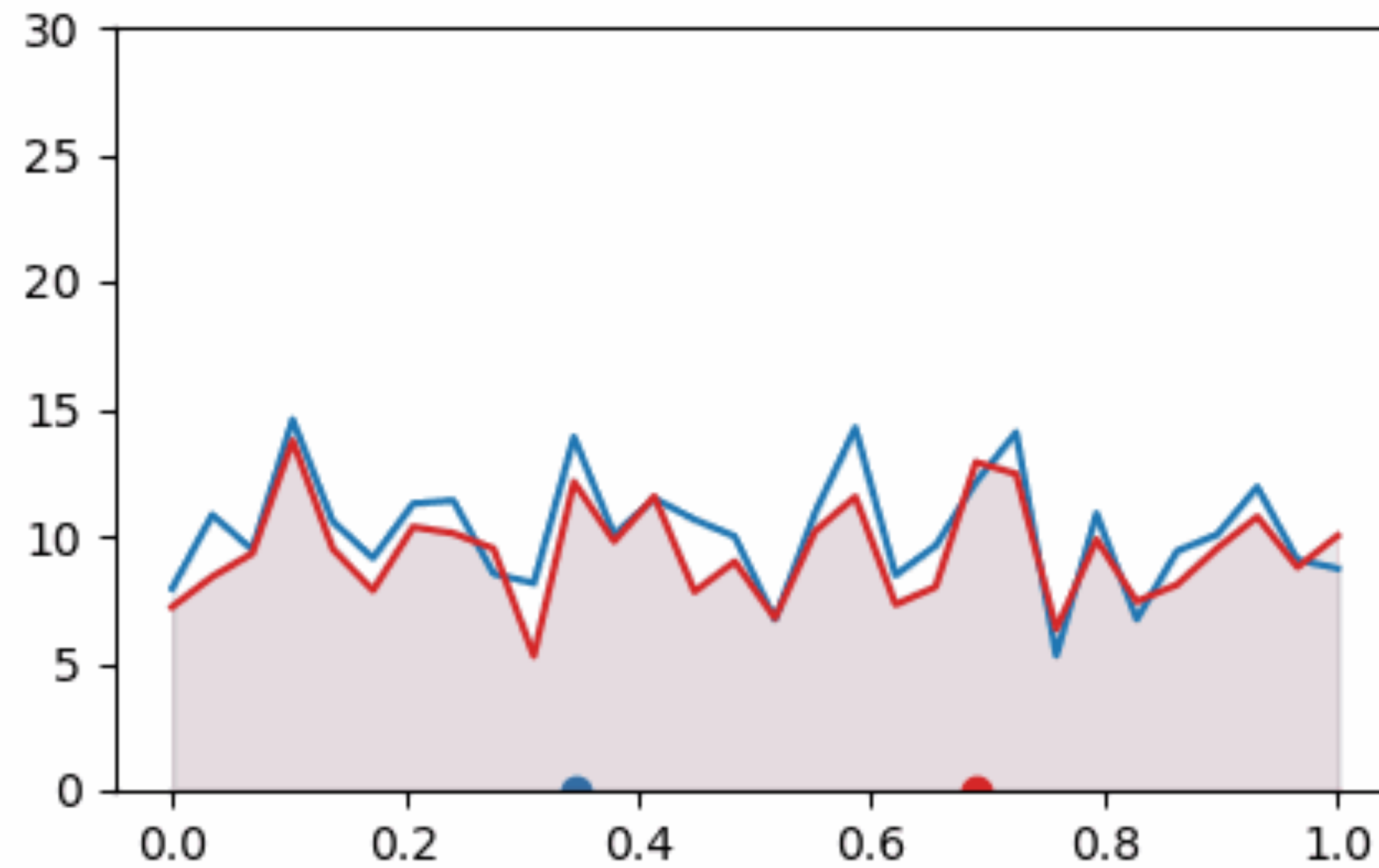
It seems that, even for similarity-dominated trainings, there is a force leading the model to better  $I$  at the expense of some  $G$  (*opt. issue?*).

This is an optimization issue but it reveals the law's curve

The max identification is due to the fact that we're squeezing 30 features in 10 dimensions.

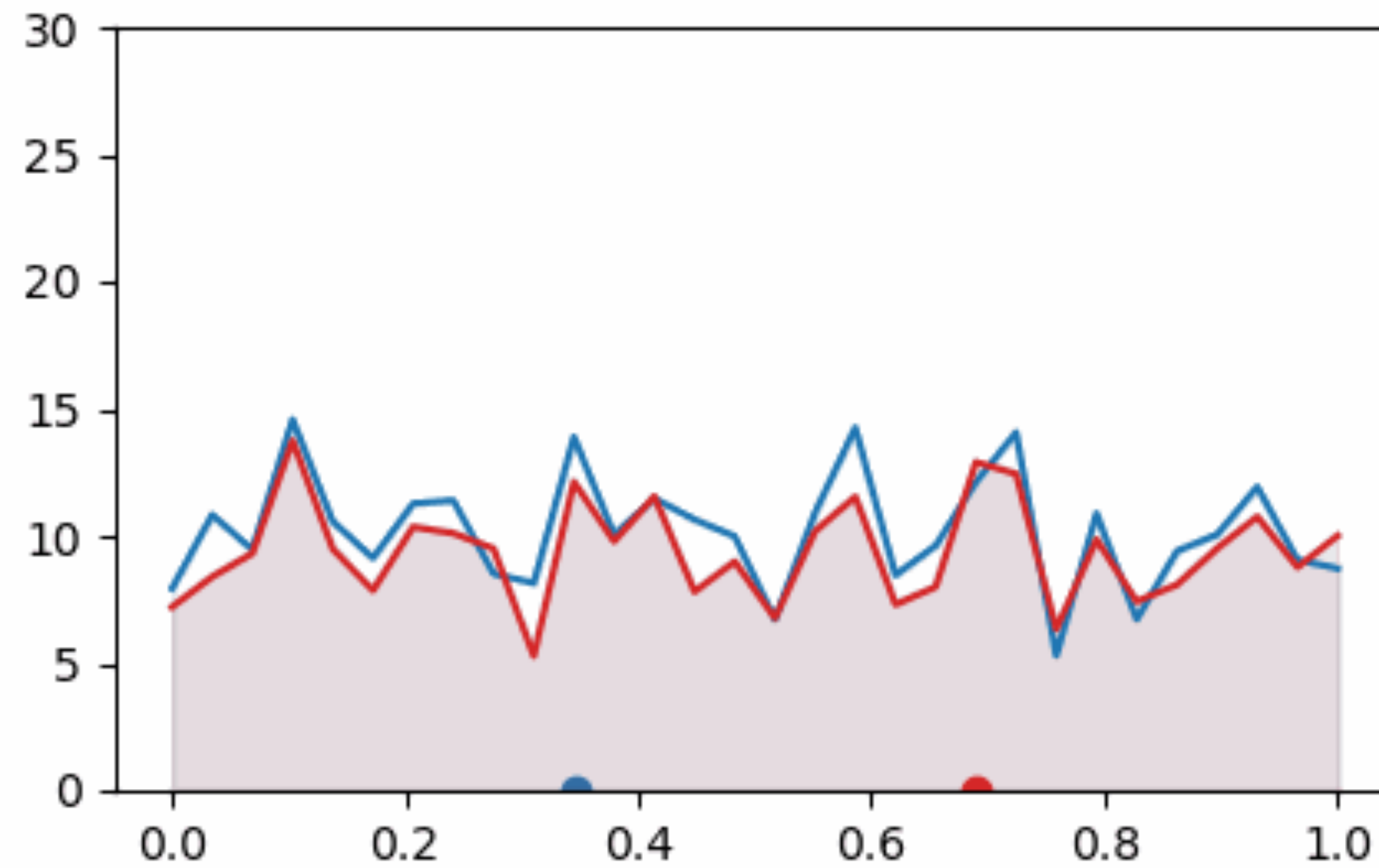
Train a single model only on similarity for *500 epochs* with *0.001 lr* and *0 decay*

## Similarity function over time



Train a single model only on similarity for *500 epochs* with *0.001 lr* and *0 decay*

## Similarity function over time



The shape of the similarity function suggests that the constant similarity of the theory may not be a good fit.

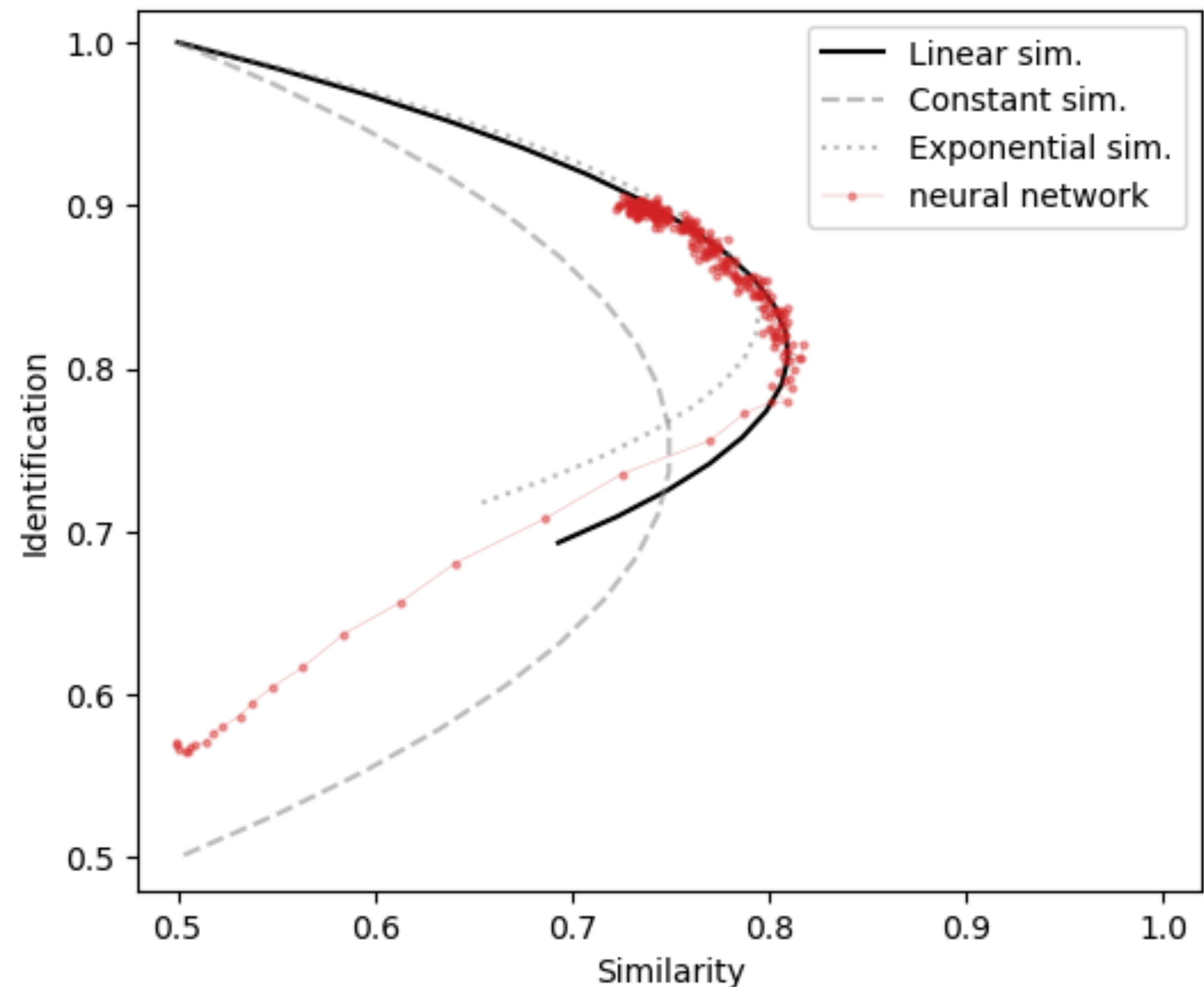
If we assume linear “hat” similarities  we can still find the formulae, at least for the circle

$$G(\varepsilon) = \frac{1}{2} + 2\varepsilon - 2(3 - 2 \log 2)\varepsilon^2$$

$$I(\varepsilon) = 1 - 2(1 - \log 2)\varepsilon$$

Linear decay gives us a great fit!

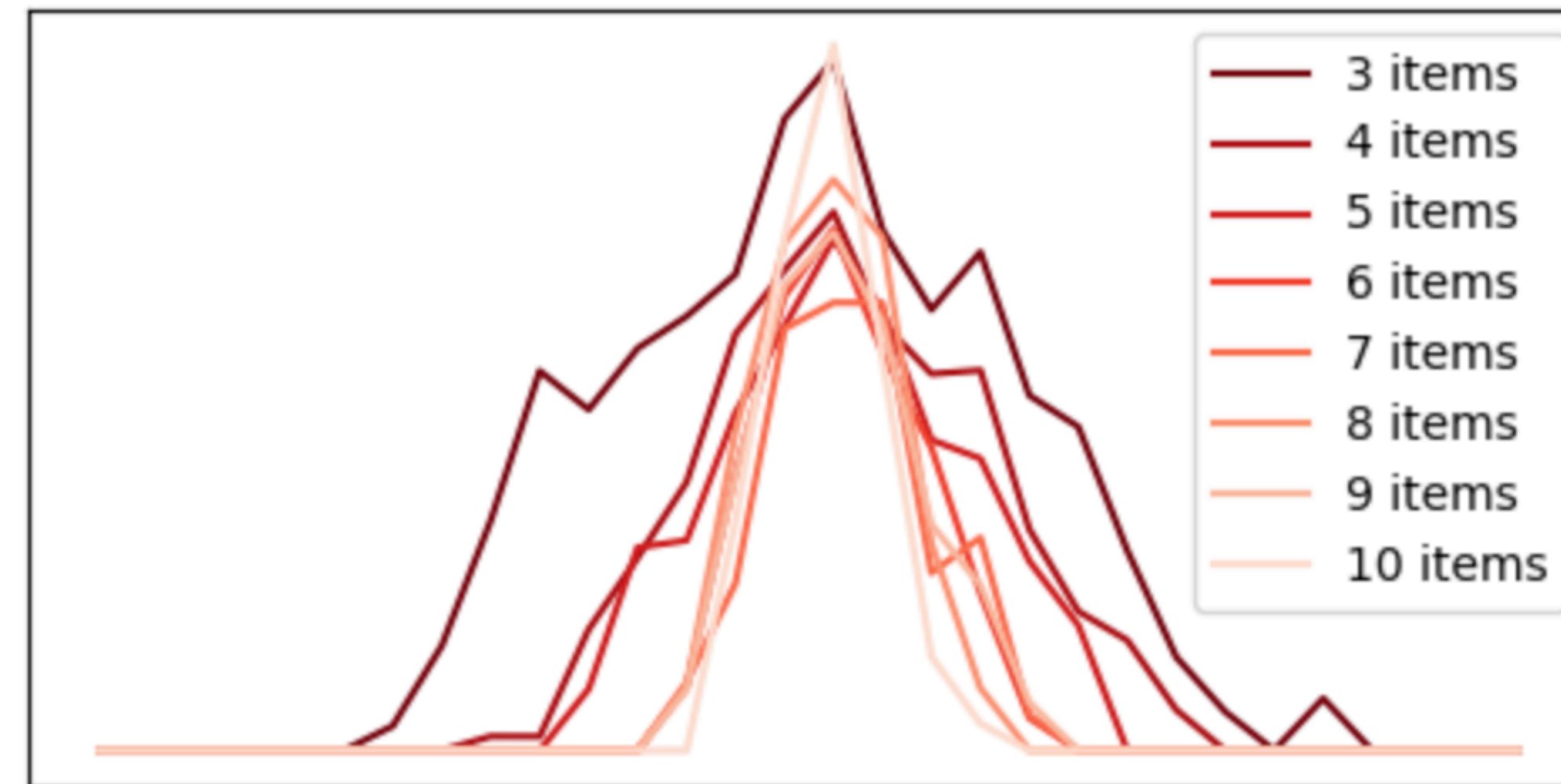
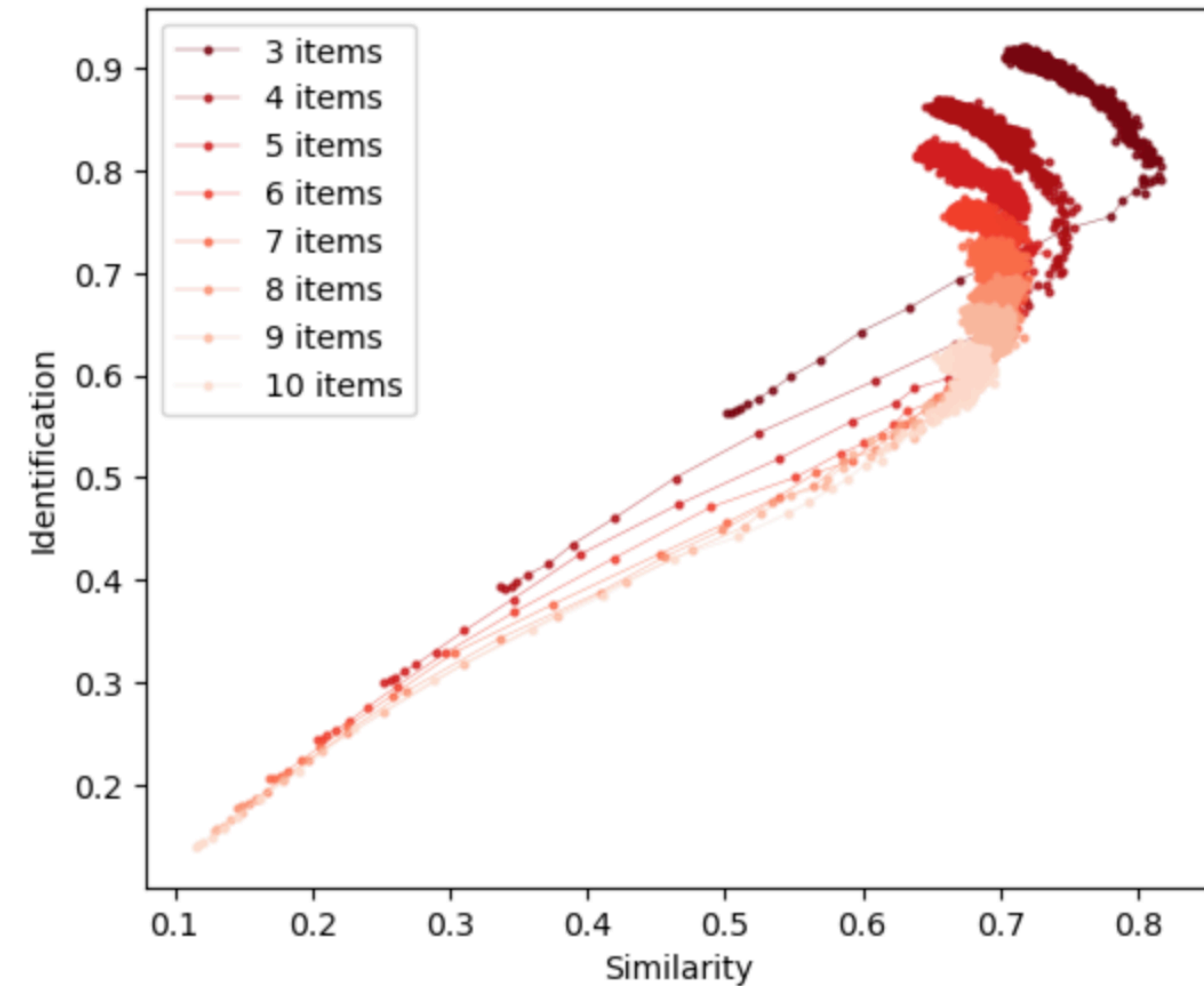
Notice that  $\varepsilon_* = \frac{1}{2(3 - 2 \log 2)} \approx 0.31$  is the radius maximizing  $G$



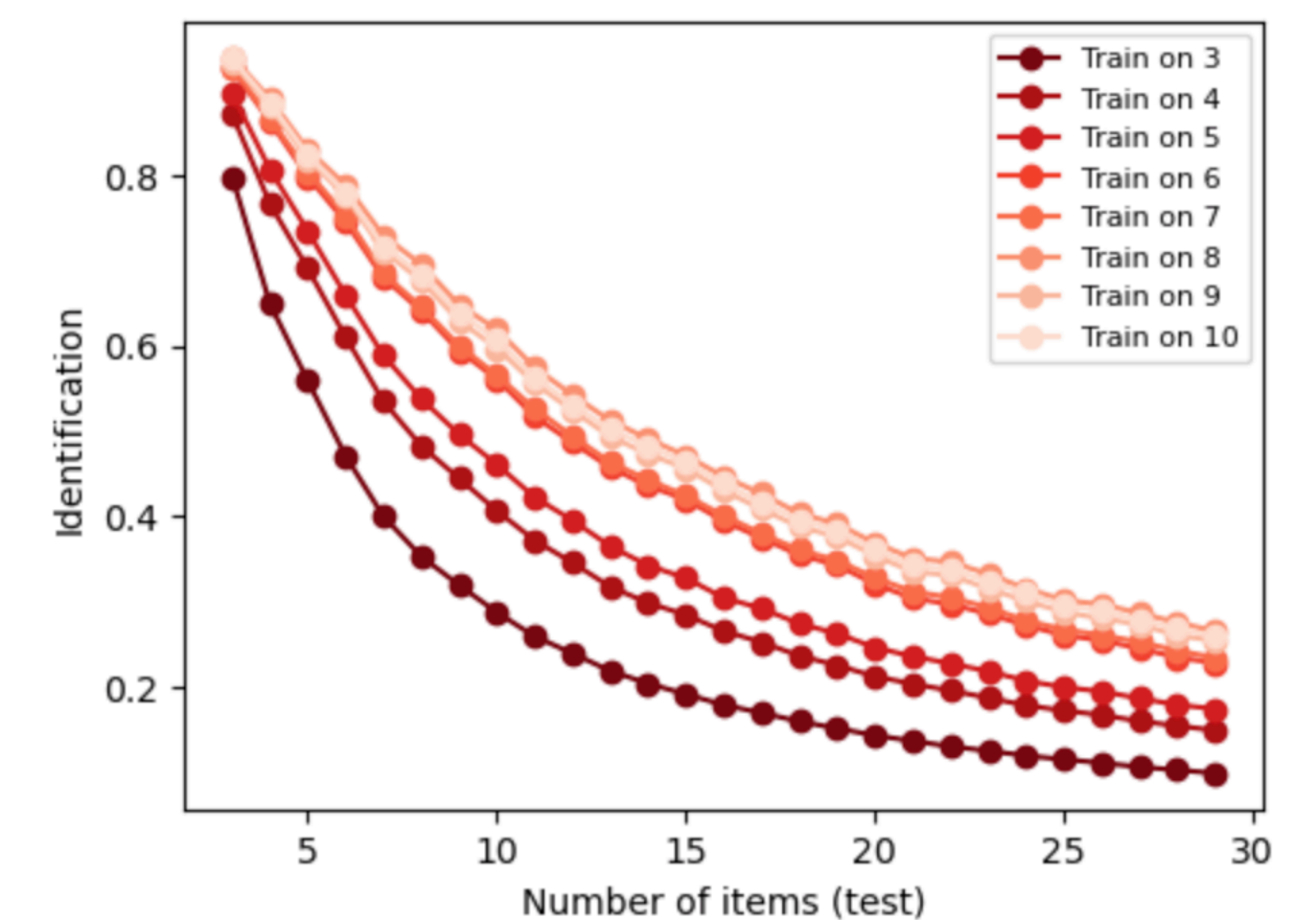
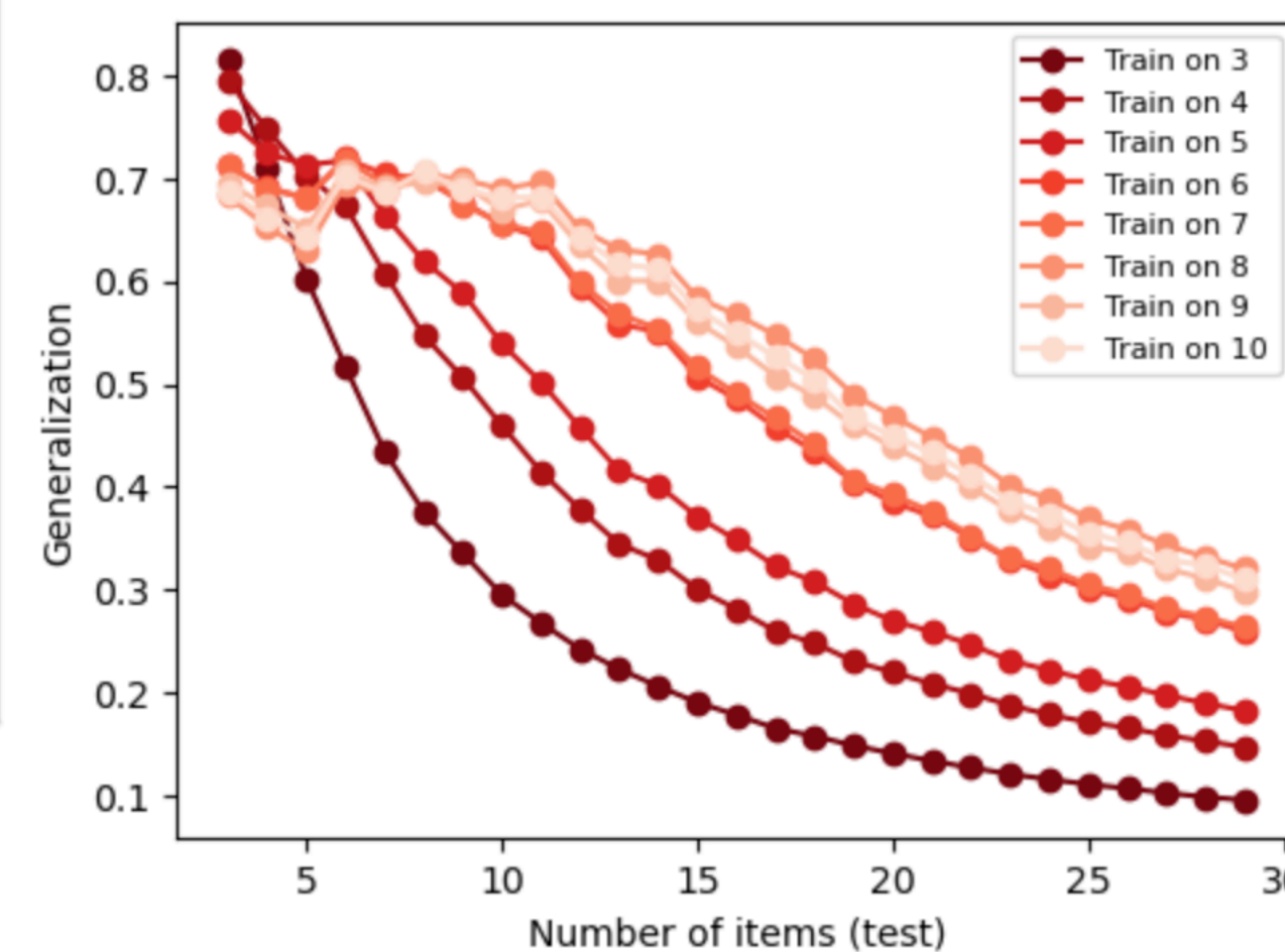
# Multi-item analysis

30 inputs, circle structure, 10 latents, 2000 samples, 1000 epochs, 128 batch size, init scale [0,2] uniform, 0.007 lr, 0 wd,  $\lambda_{sim} = 0.1$ ,  $\lambda_{rec} = 0$

**Note:** higher number of items means that the model sees proportionally more data in the training process. *Probably not the best idea*



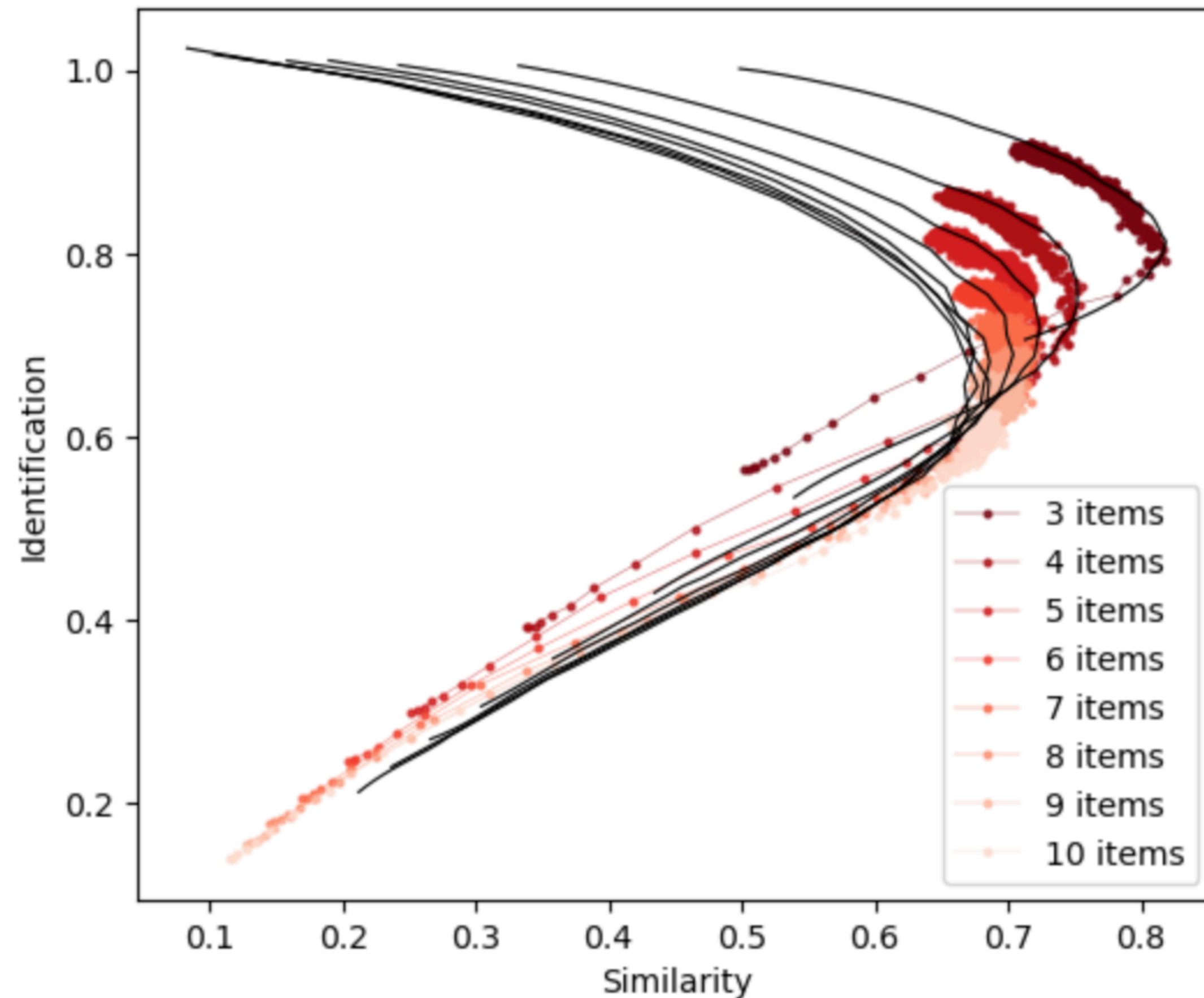
Smaller sim. functions as the number of items gets larger



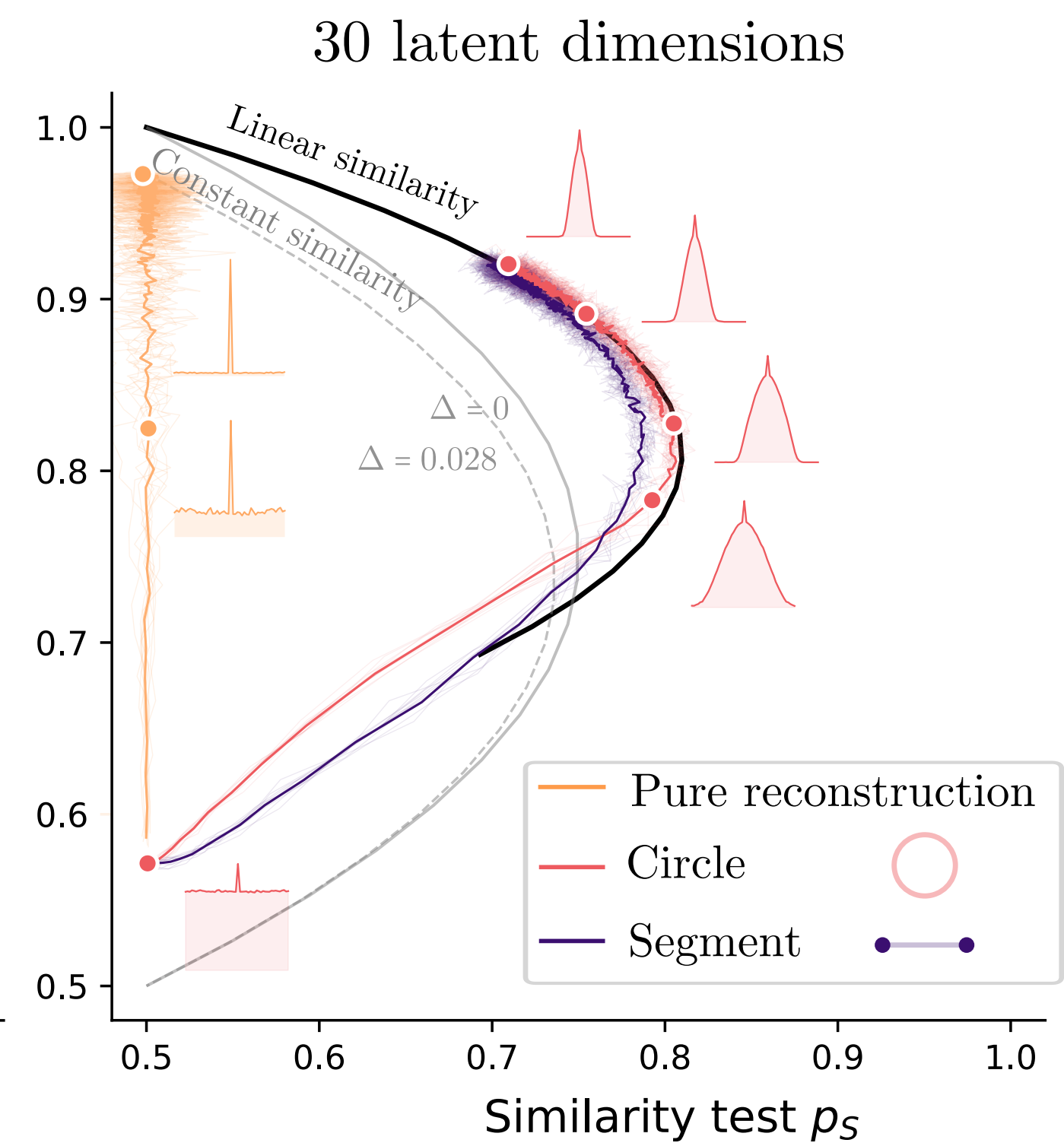
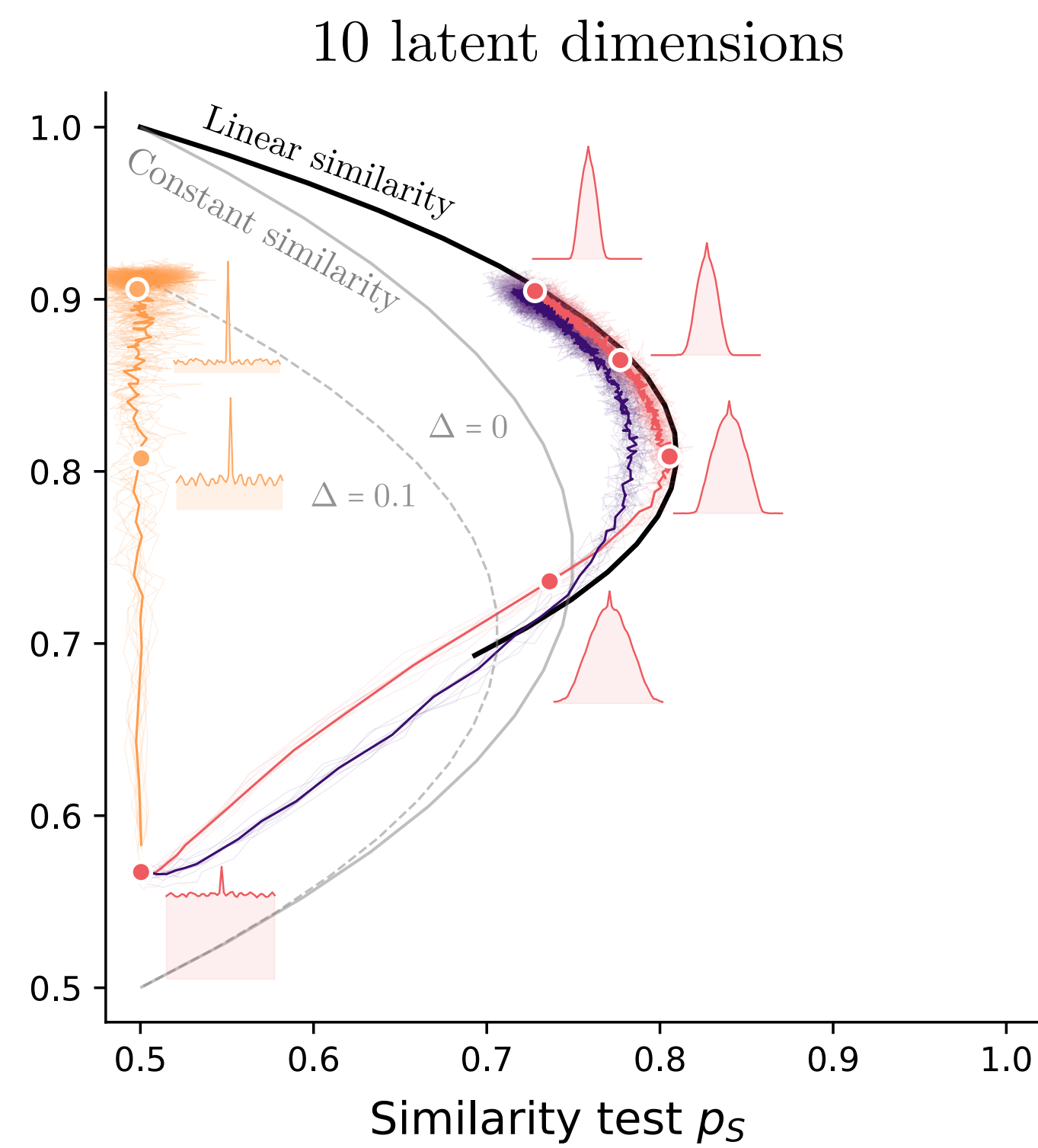
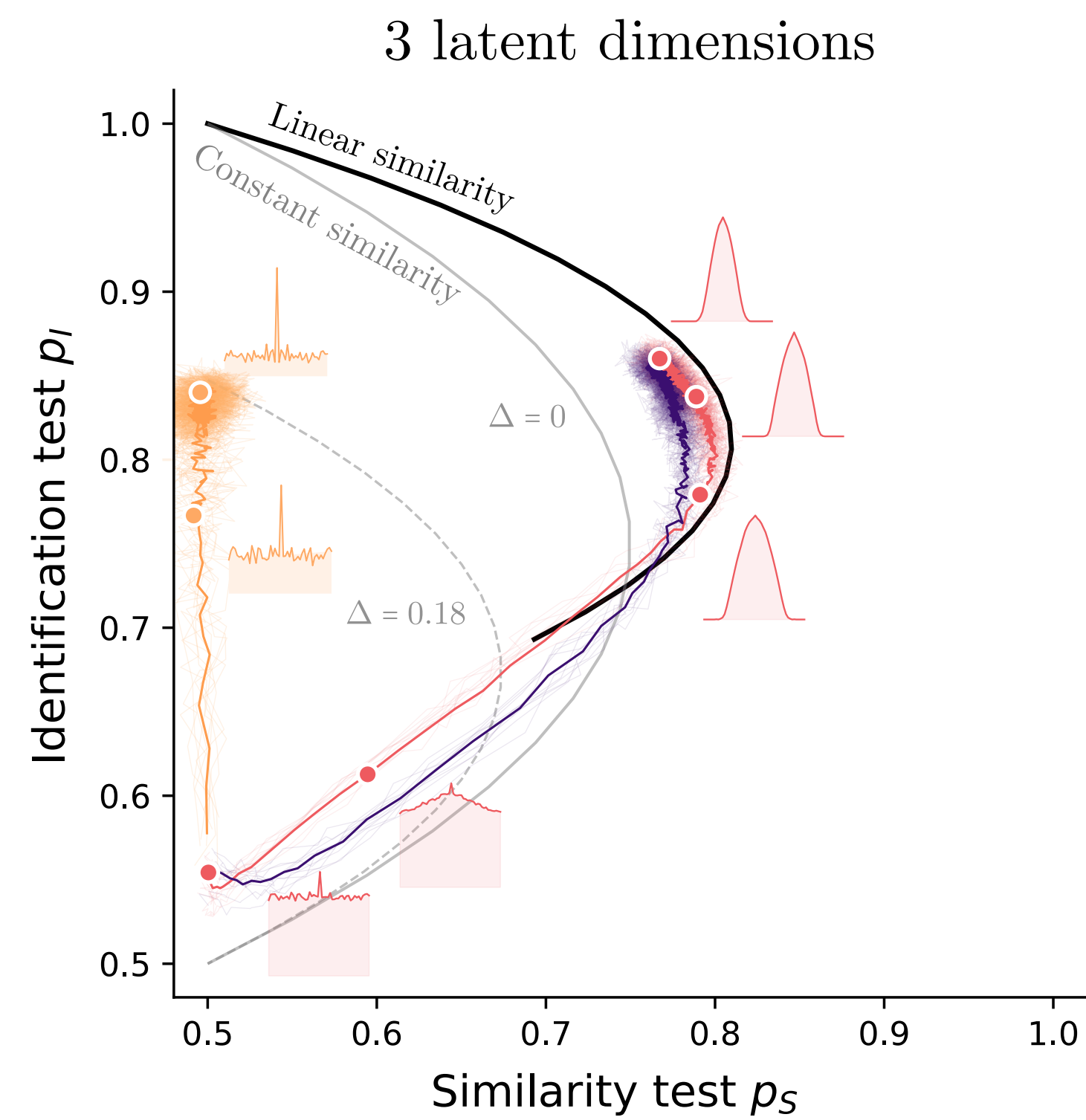
# Multi-item analysis

If I plot the multi-item curves for the linear decaying similarity function (by simulation + smoothing), I get the following

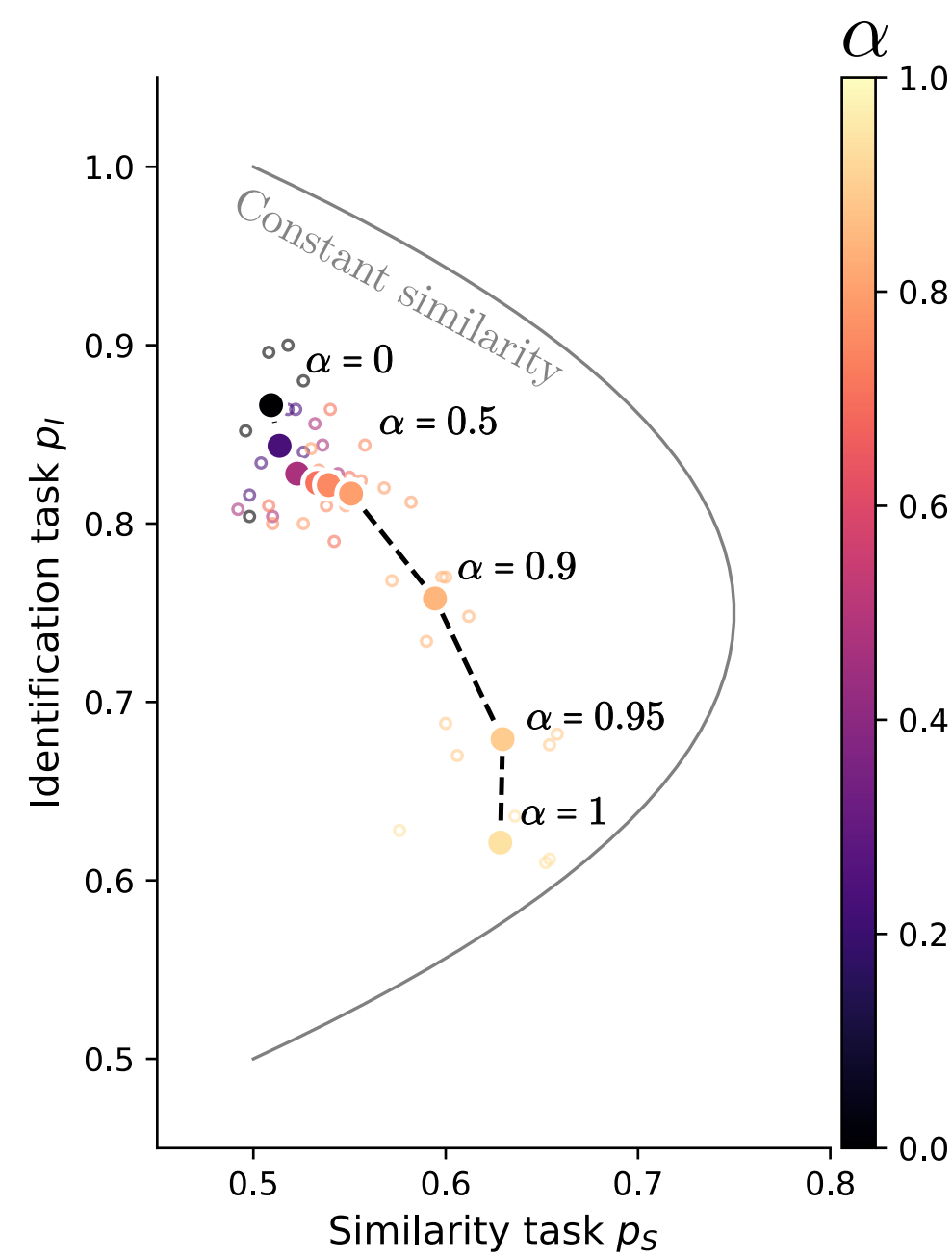
Not perfect but still quite good



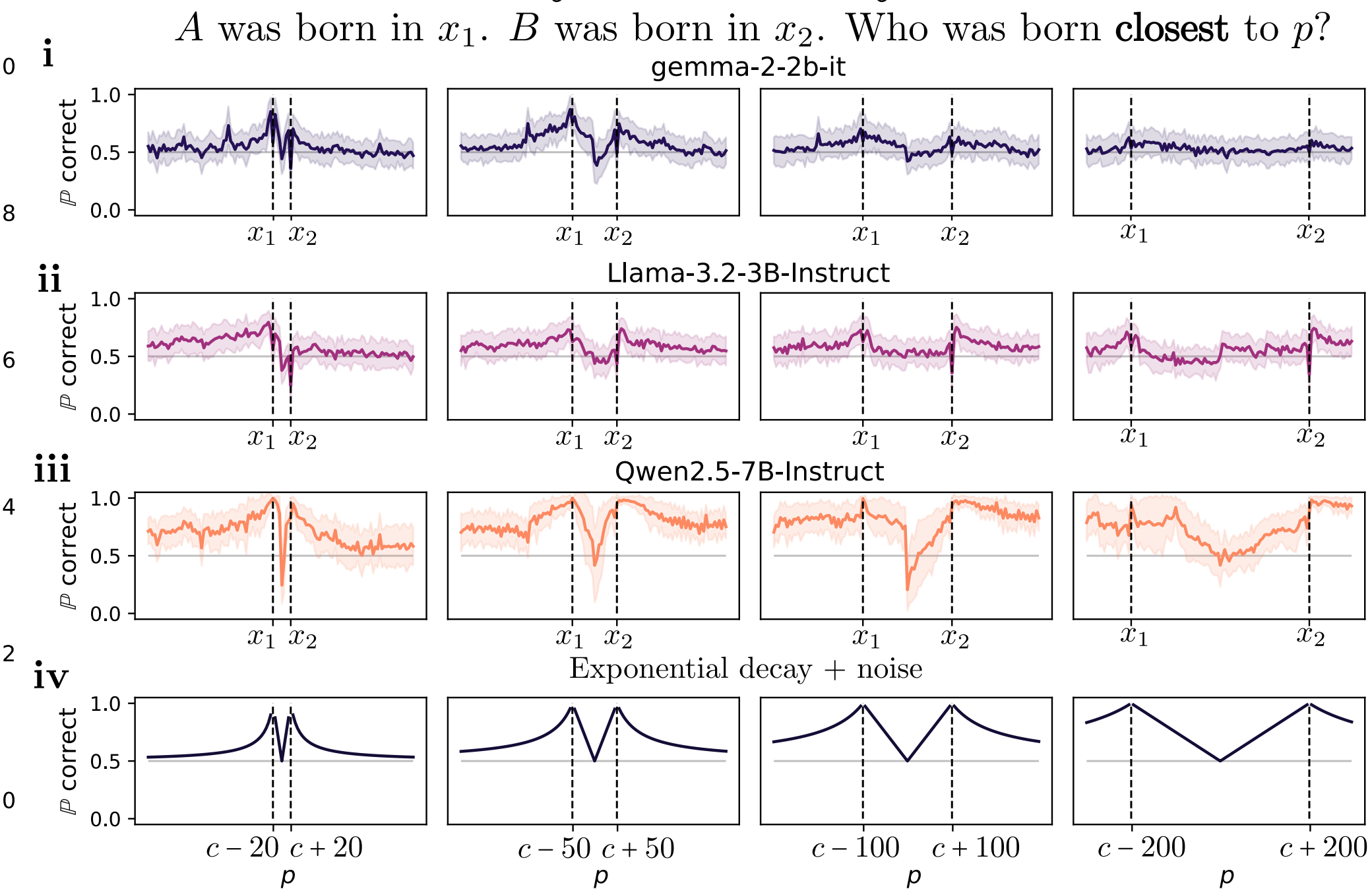
This means that the similarity function learned can be approximated by linear for small  $n$  but changes when we increase.



### a CNN finetuning



### b LLM year similarity task



### c VLM spatial similarity task

