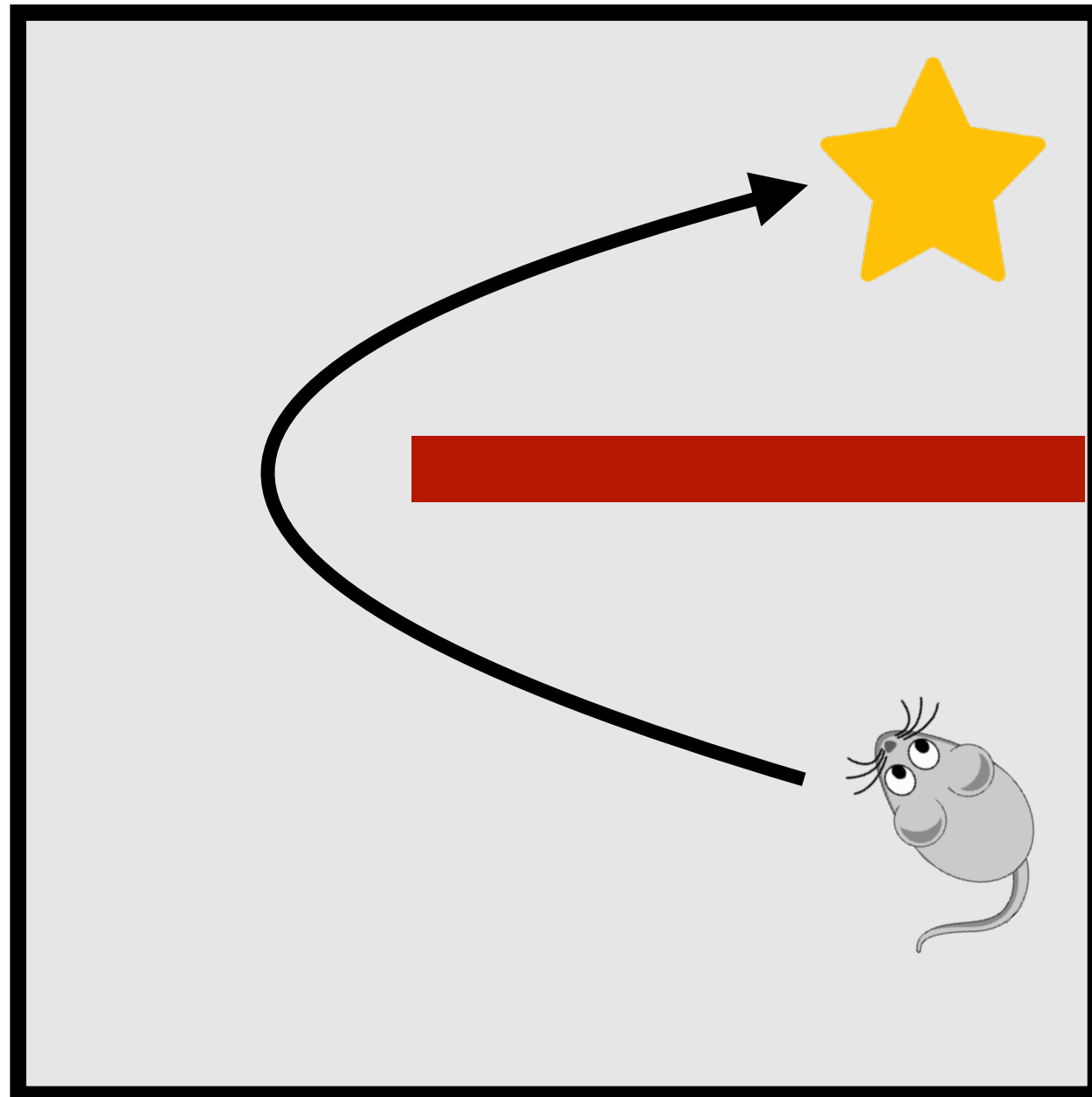


# How to build cognitive maps

James Whittington

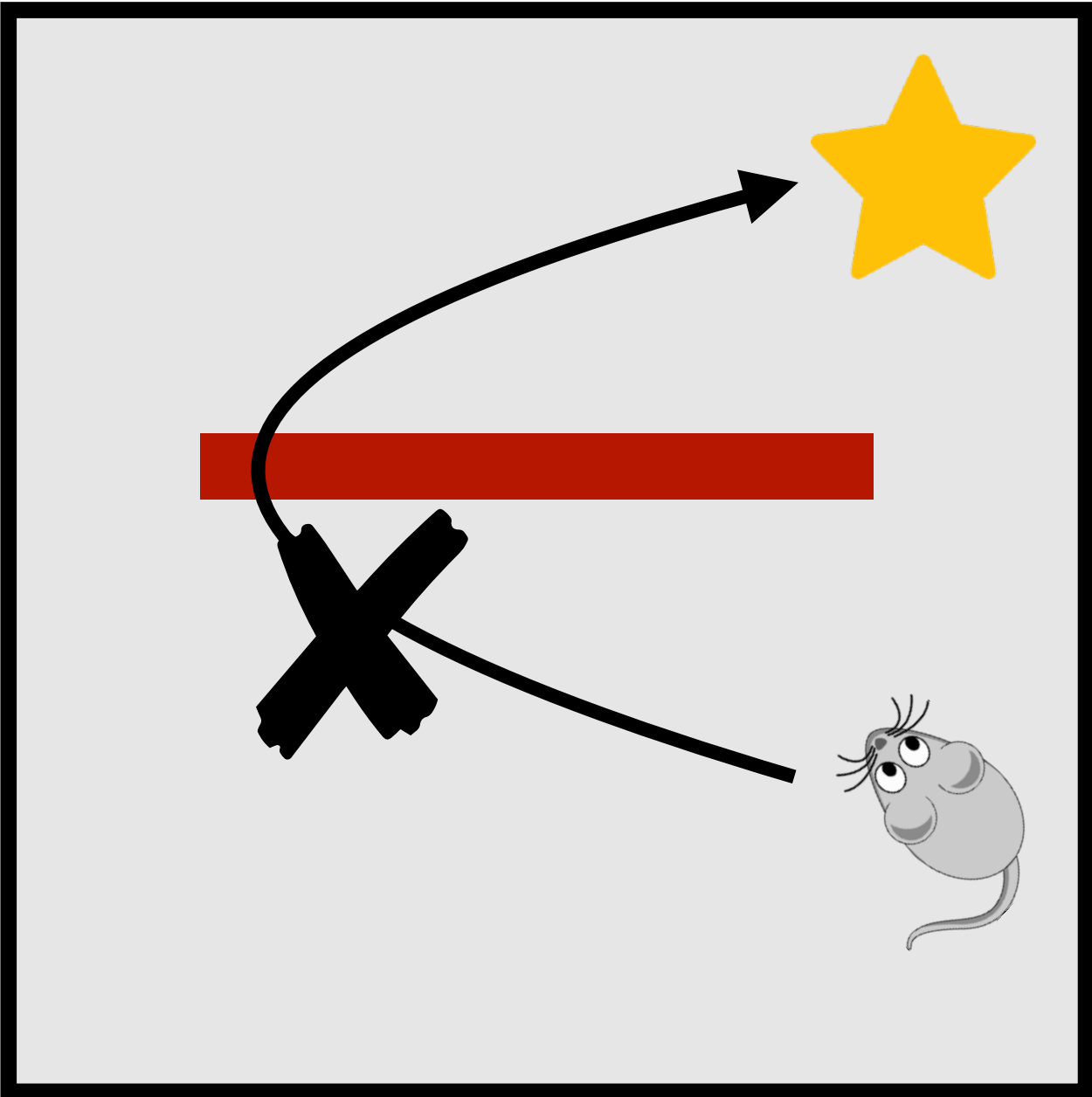
**What is a cognitive map? Representating knowledge for flexible behaviour**

# What is a cognitive map? Representating knowledge for flexible behaviour



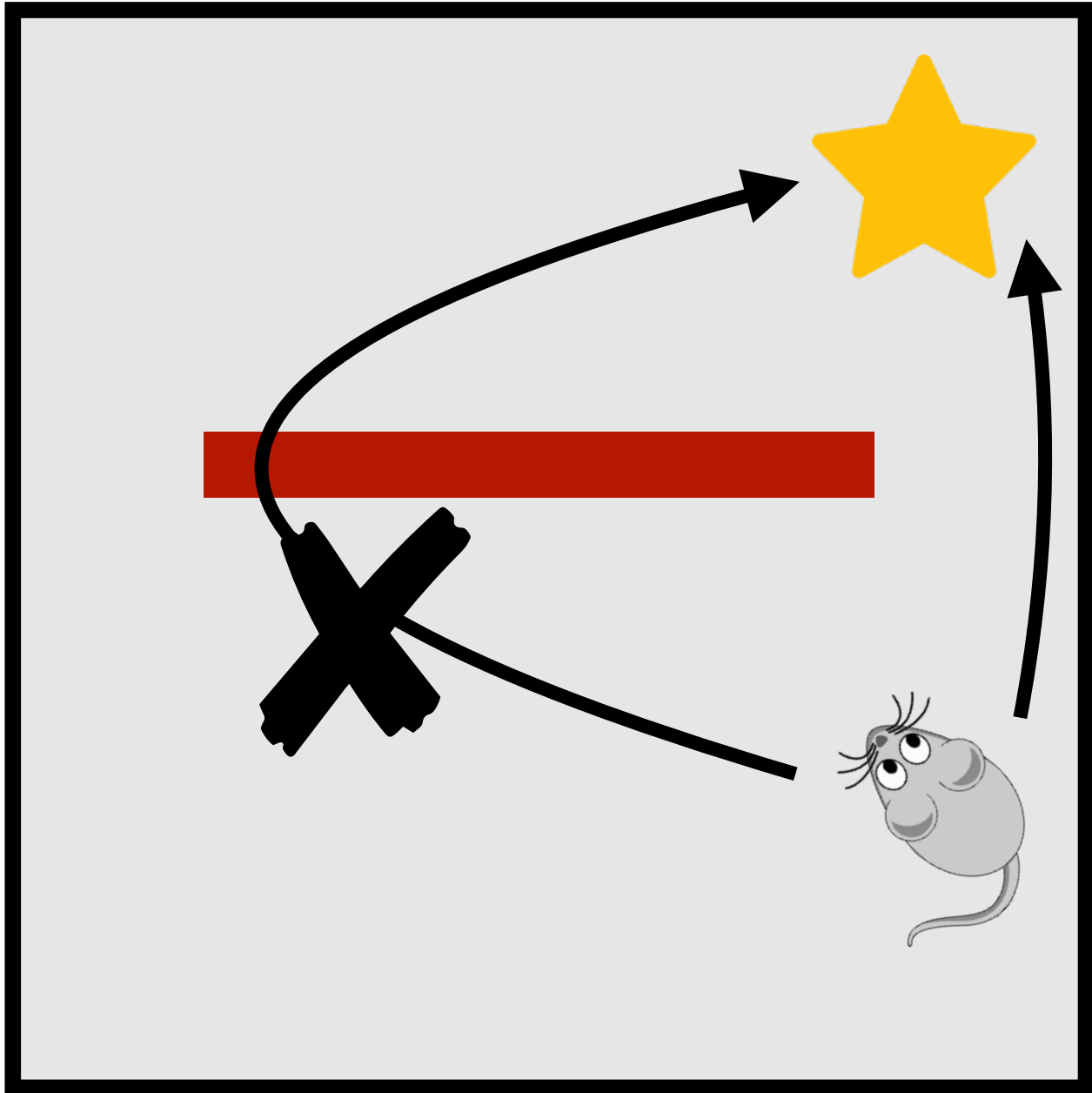
# What is a cognitive map? Representating knowledge for flexible behaviour

Repeating  
old actions  
isn't optimal



# What is a cognitive map? Representating knowledge for flexible behaviour

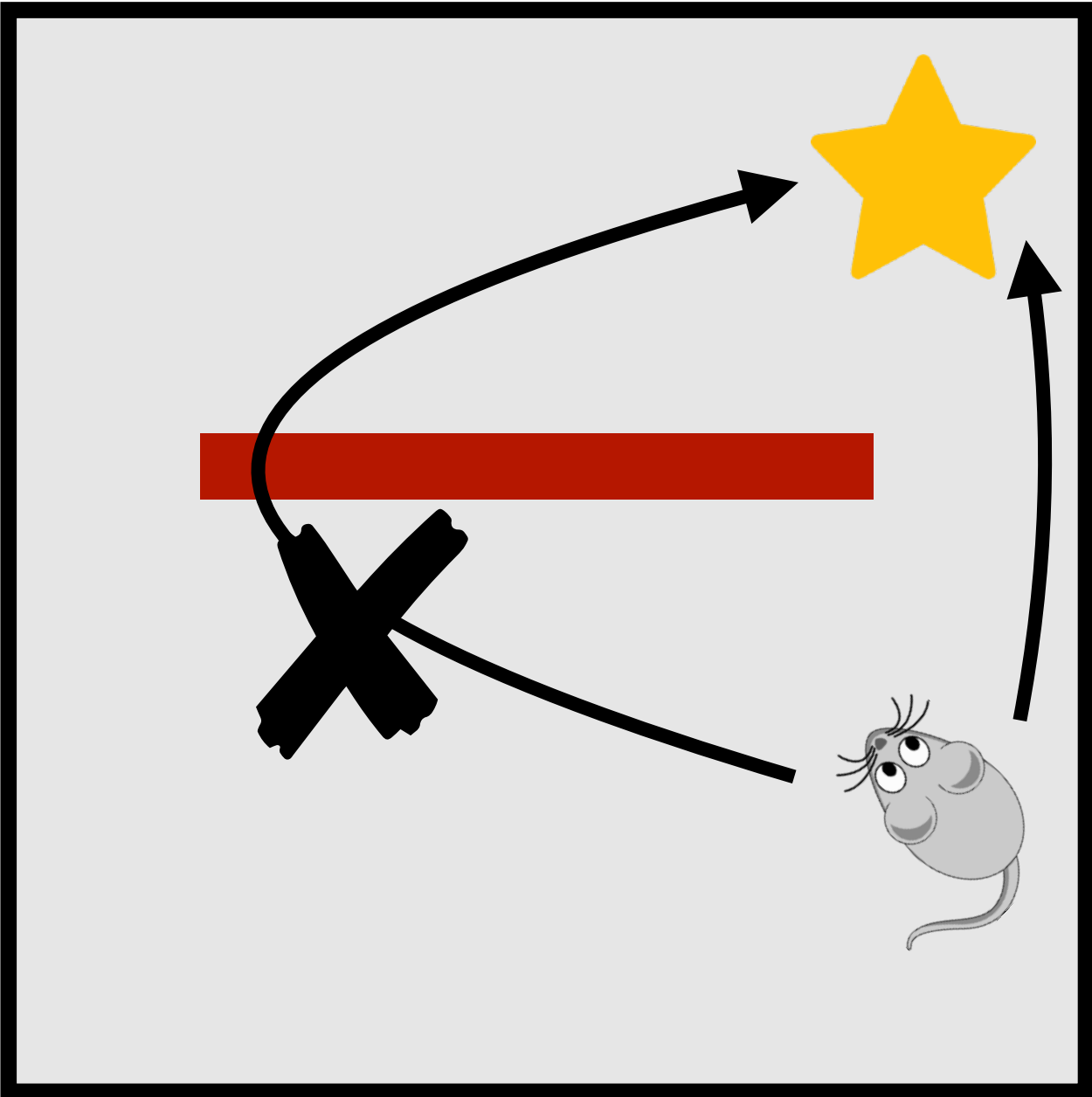
Repeating old actions isn't optimal



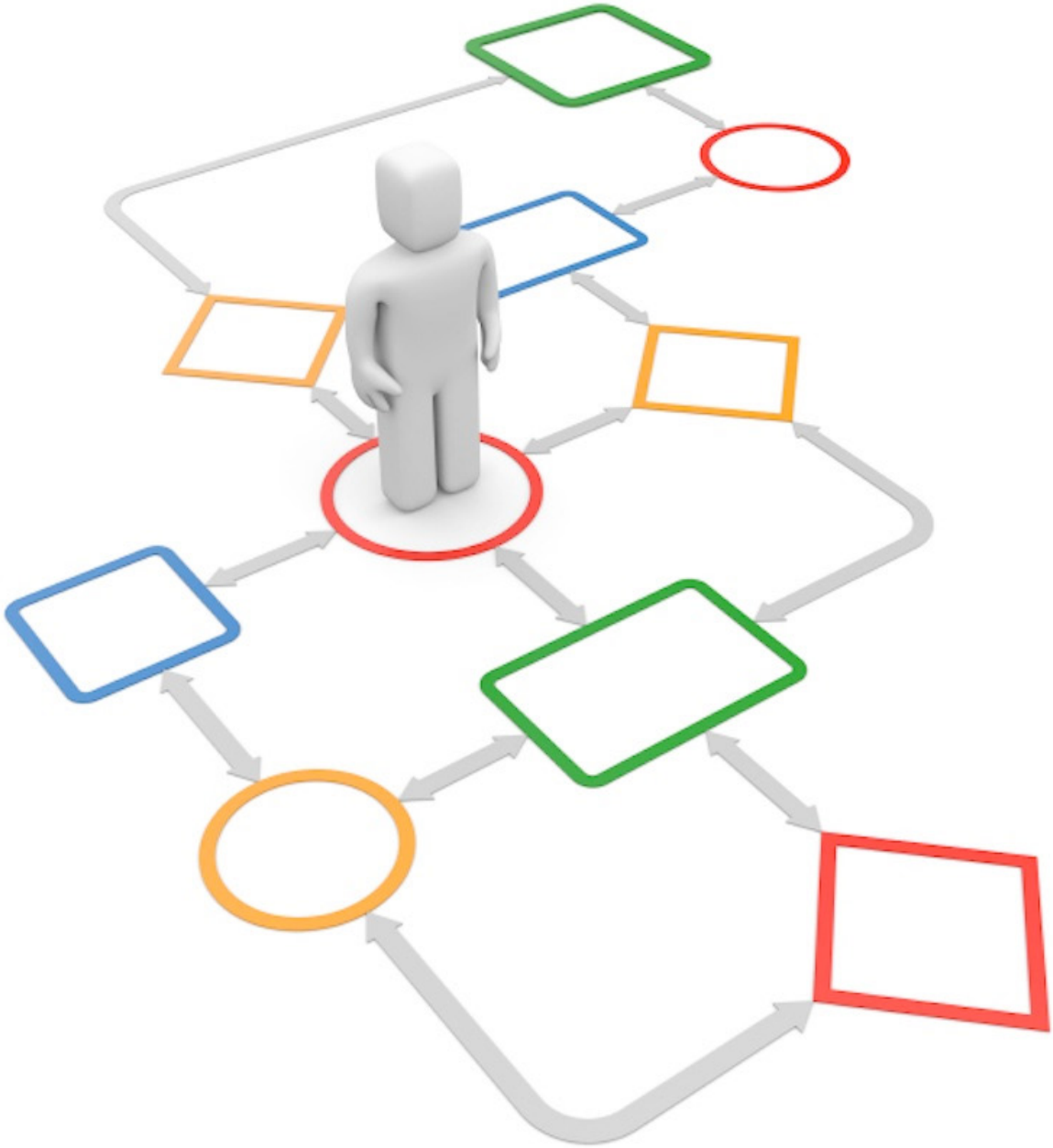
Can take shortcuts with a world model

# What is a cognitive map? Representating knowledge for flexible behaviour

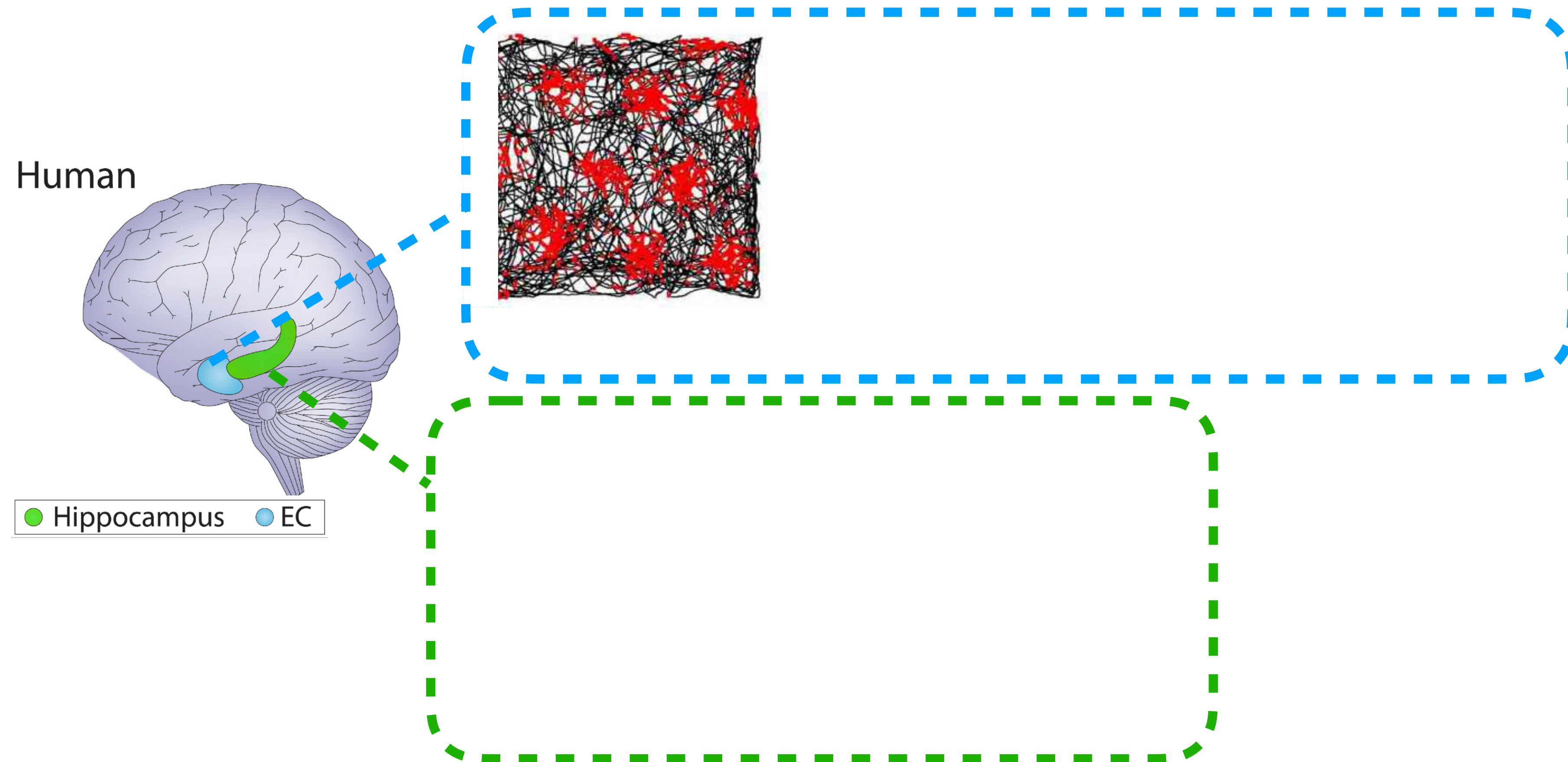
Repeating old actions isn't optimal



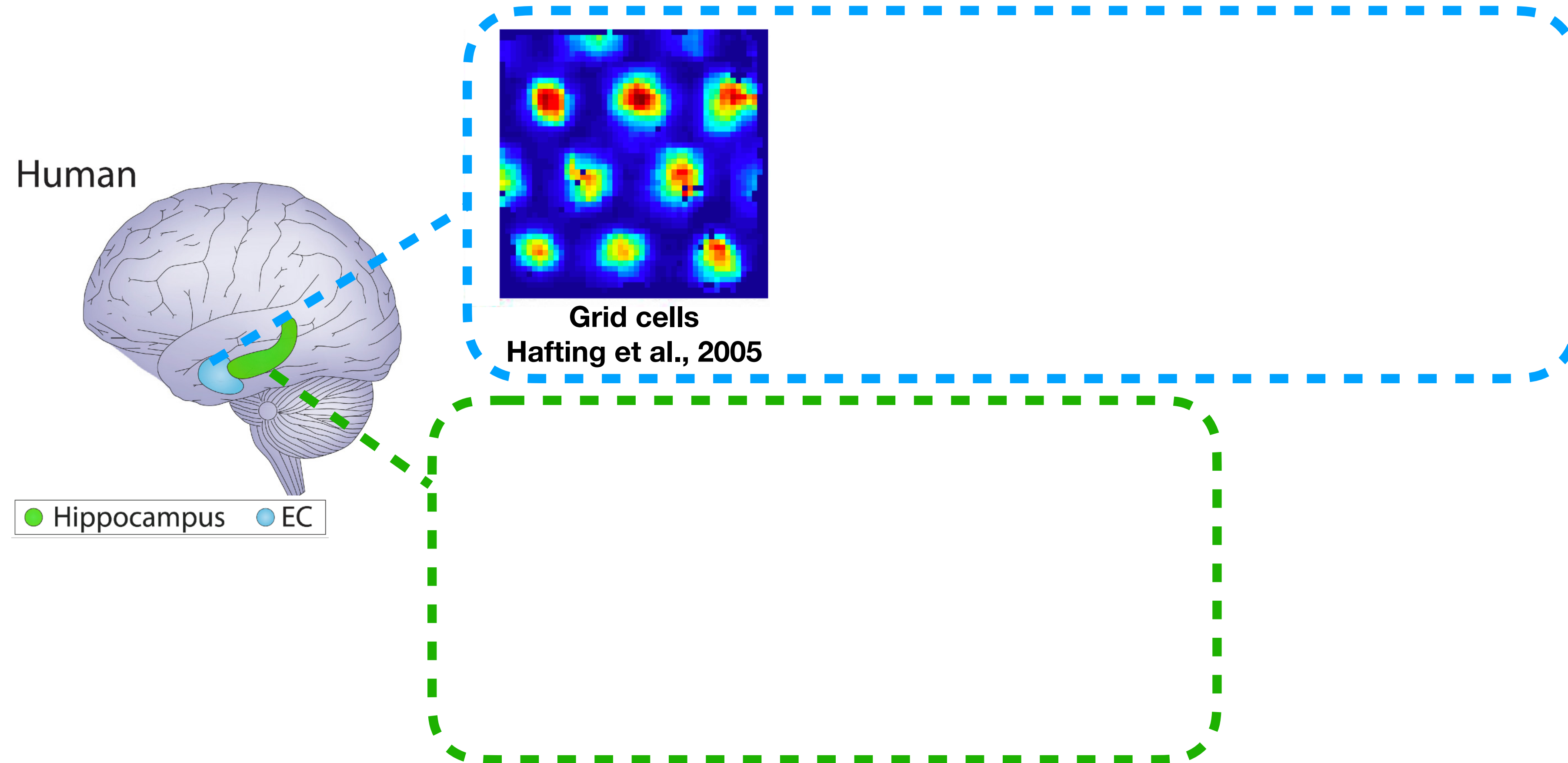
Can take shortcuts with a world model



# The hippocampal and entorhinal cognitive map

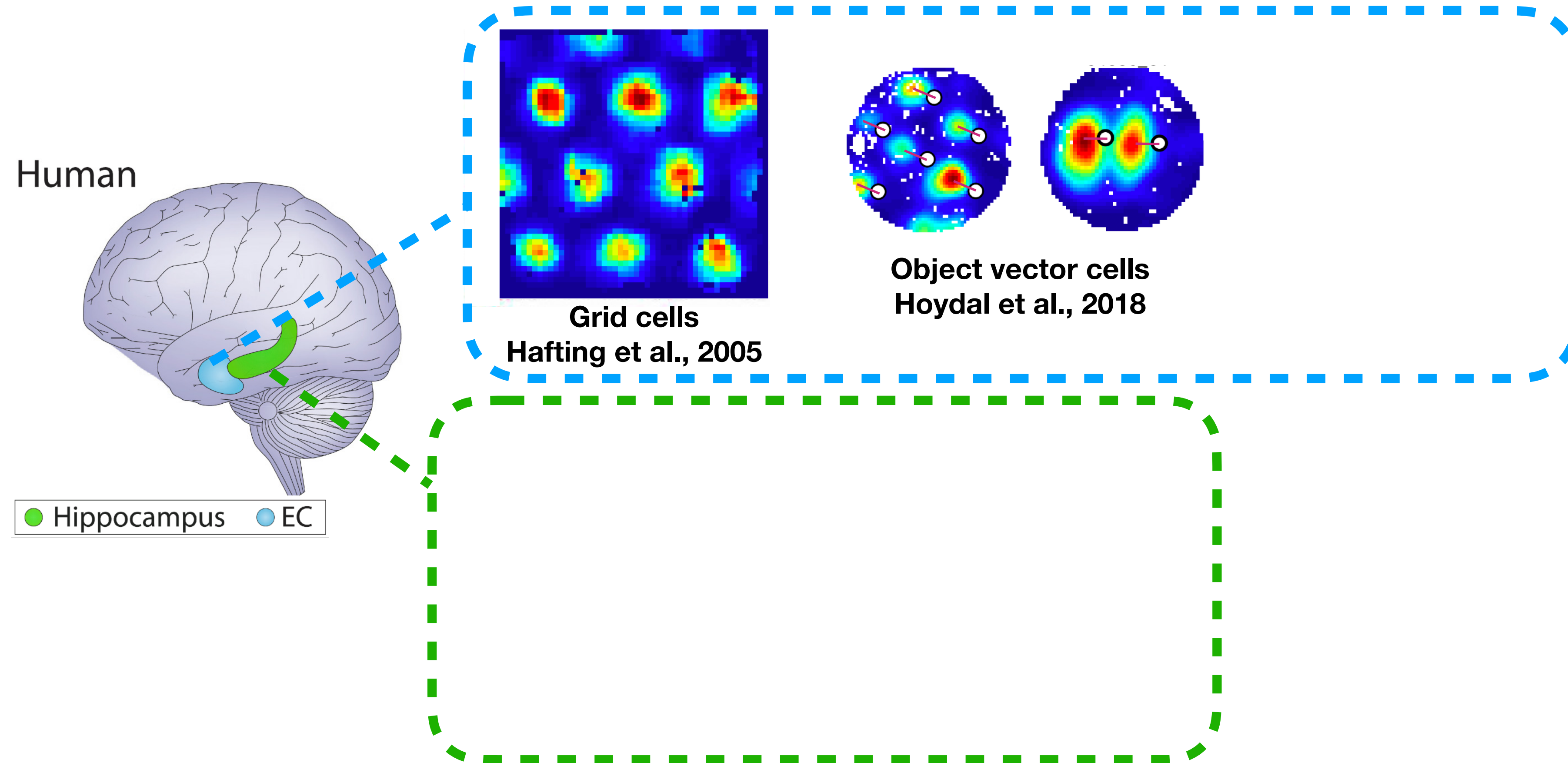


# The hippocampal and entorhinal cognitive map

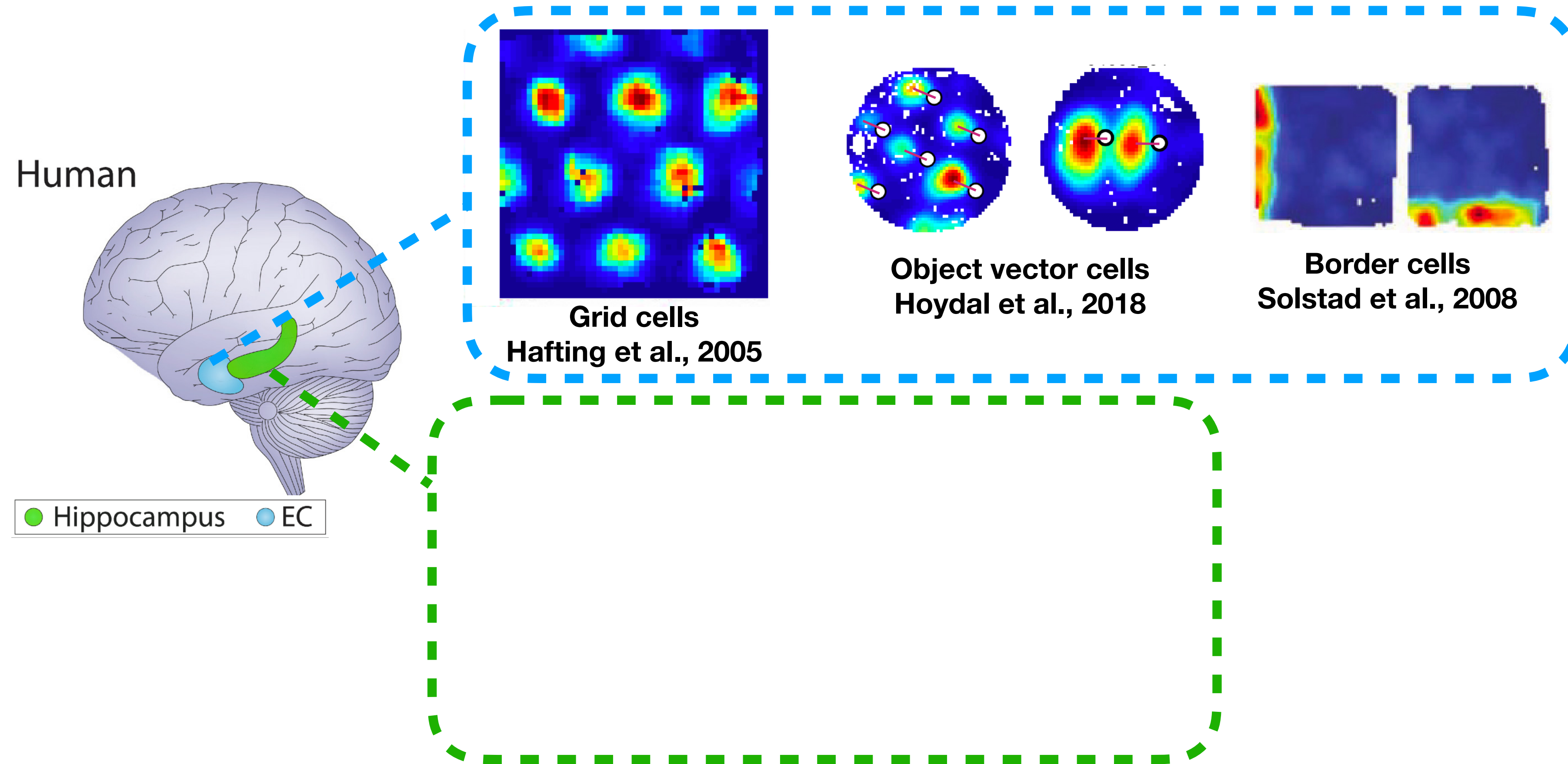




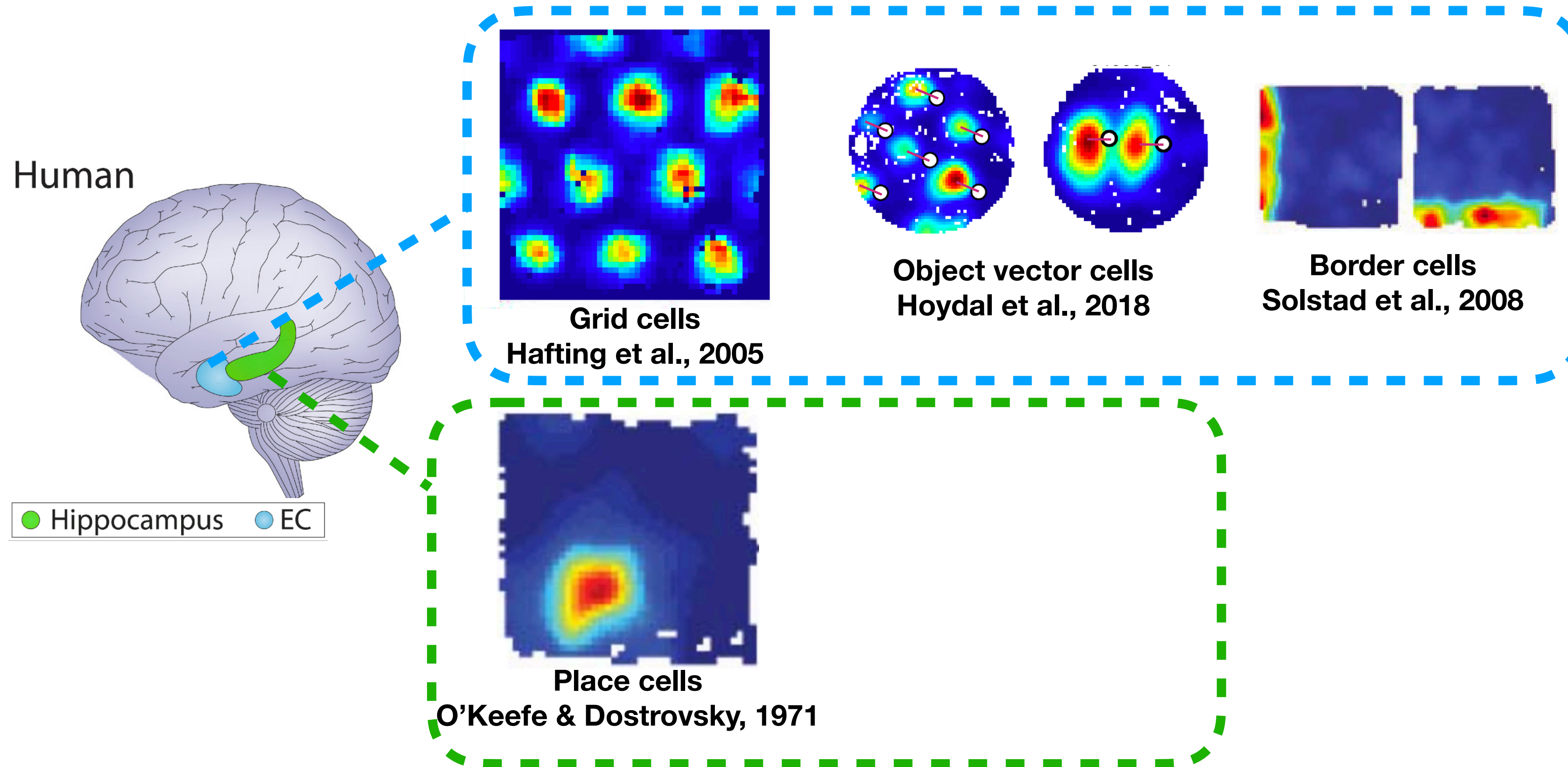
# The hippocampal and entorhinal cognitive map



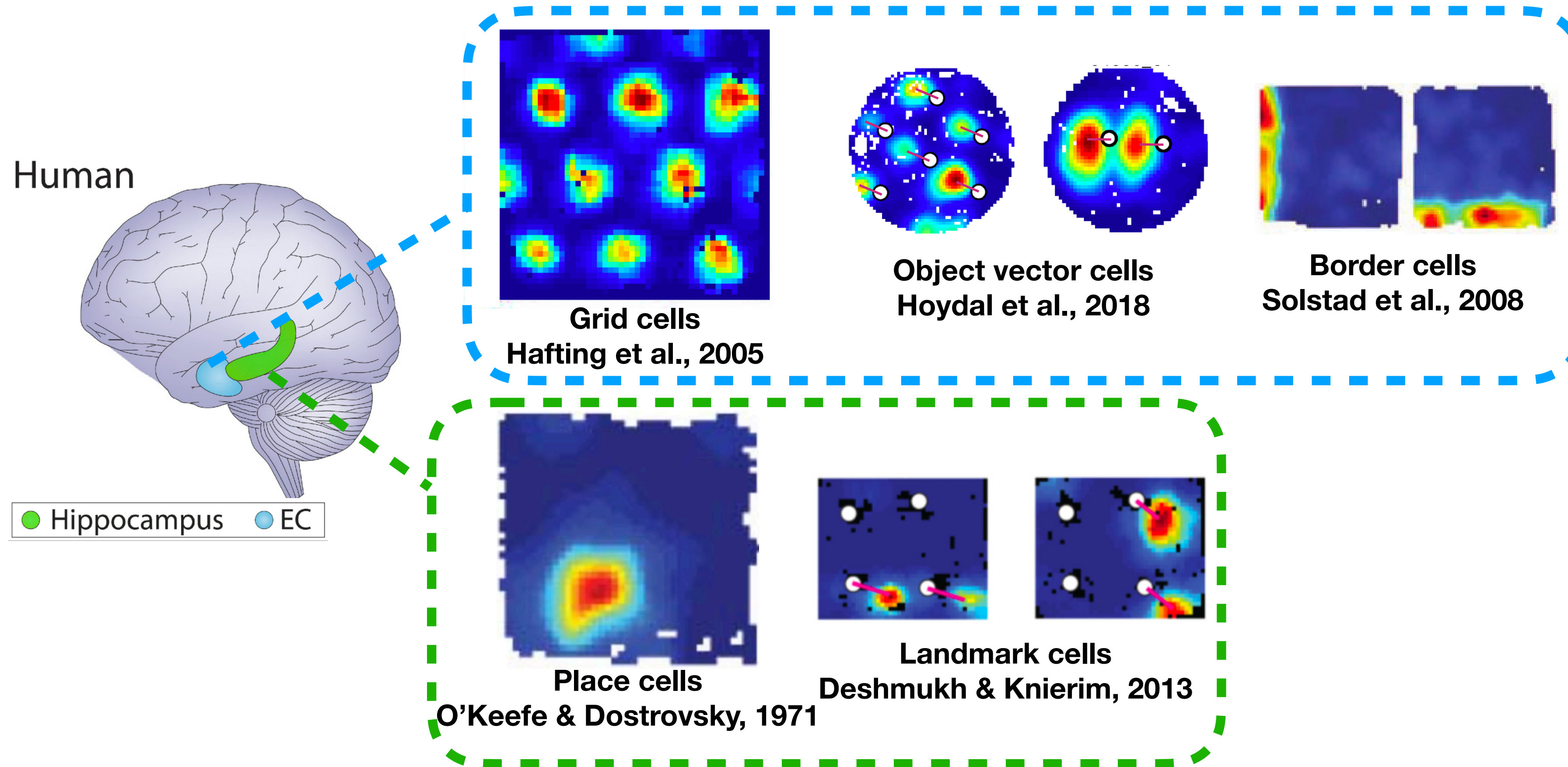
# The hippocampal and entorhinal cognitive map



# The hippocampal and entorhinal cognitive map

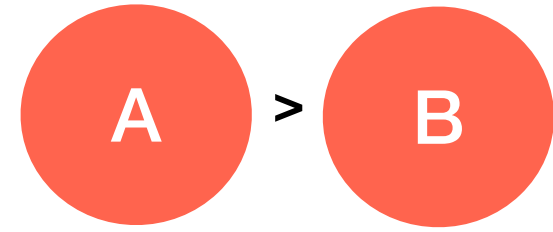


# The hippocampal and entorhinal cognitive map

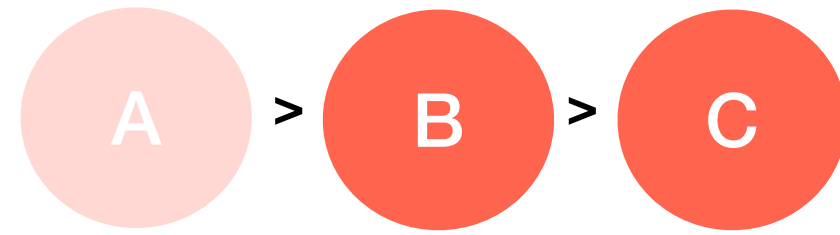


**These brain areas are critical for non-spatial 'shortcuts'**

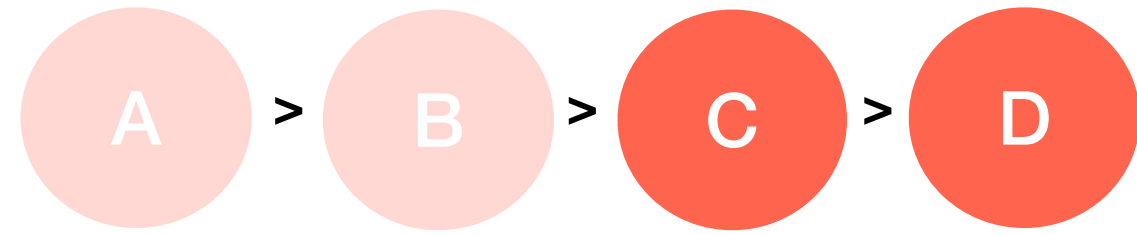
# These brain areas are critical for non-spatial 'shortcuts'



# These brain areas are critical for non-spatial 'shortcuts'

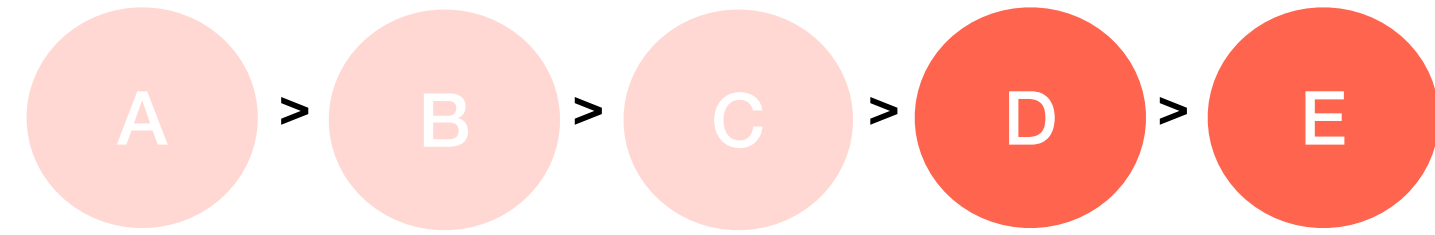


# These brain areas are critical for non-spatial 'shortcuts'

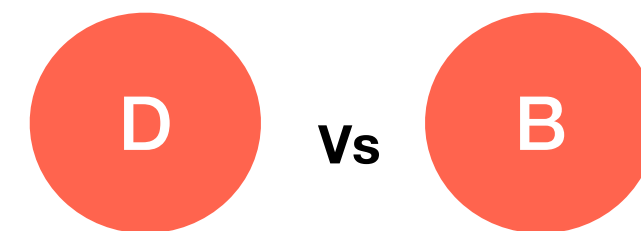
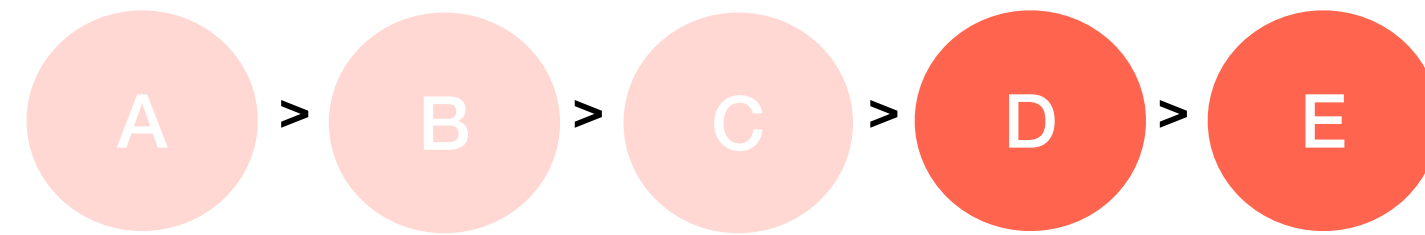




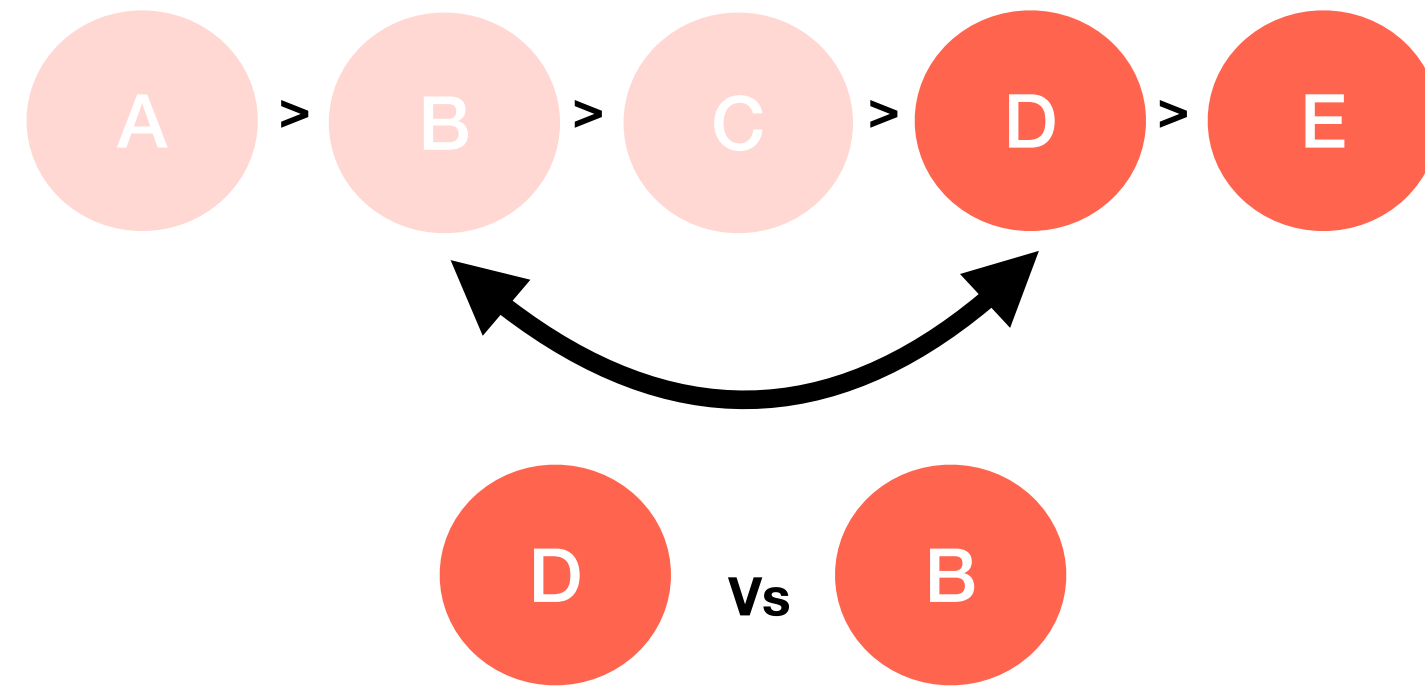
# These brain areas are critical for non-spatial 'shortcuts'



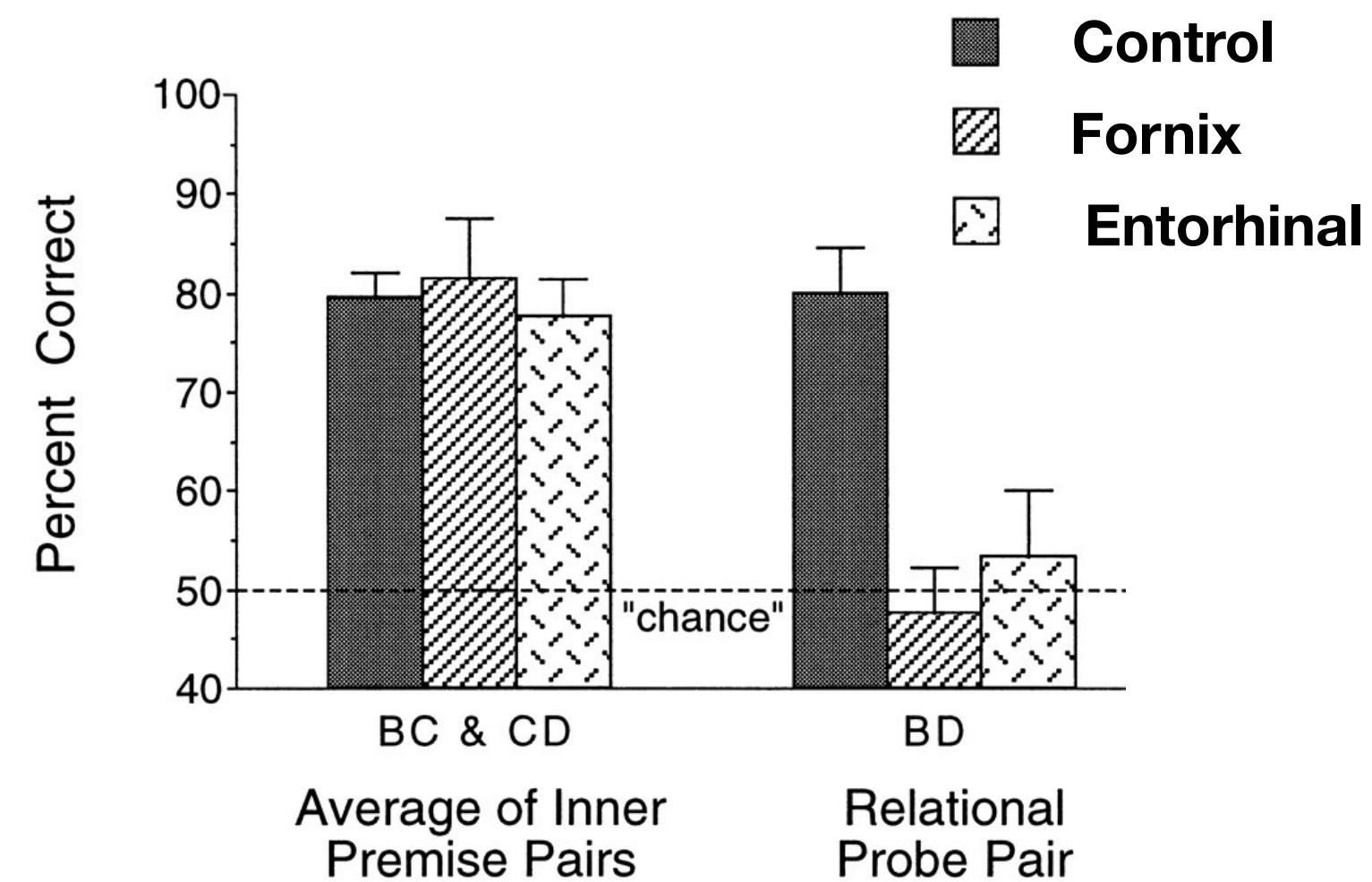
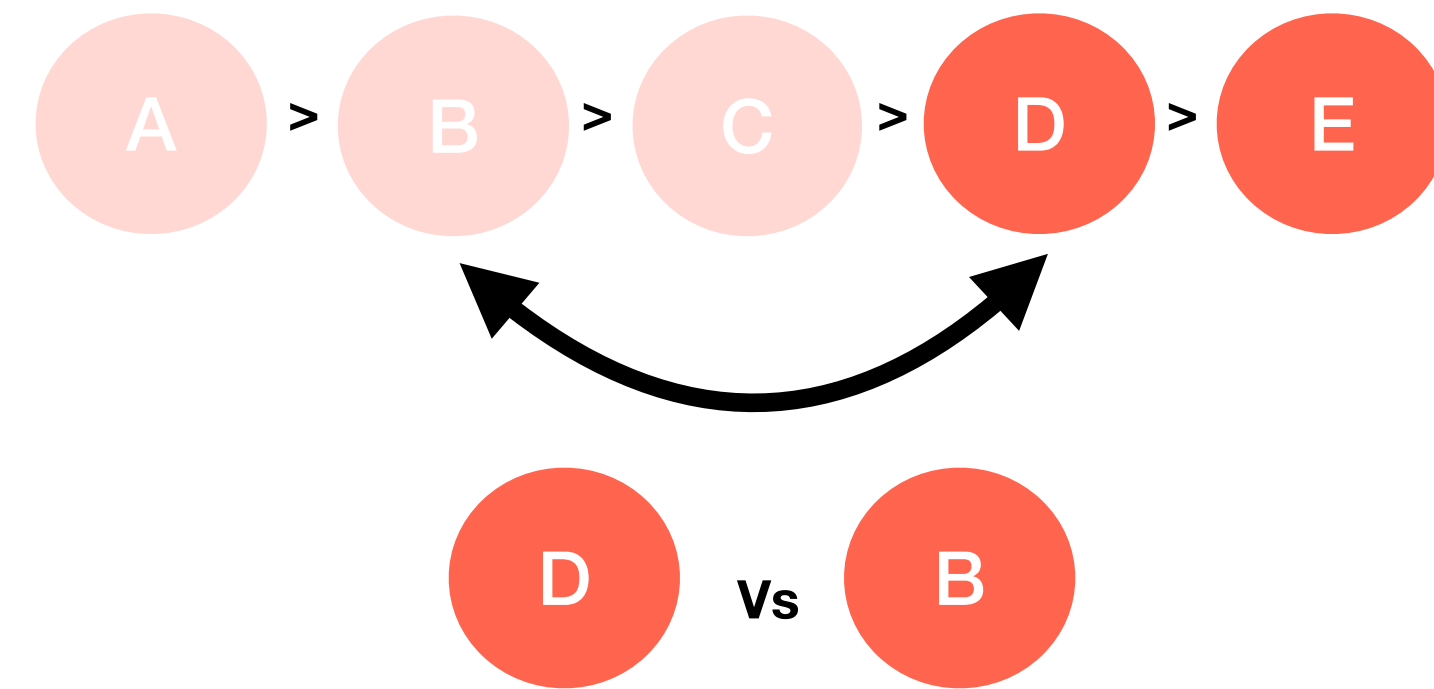
# These brain areas are critical for non-spatial 'shortcuts'



# These brain areas are critical for non-spatial 'shortcuts'



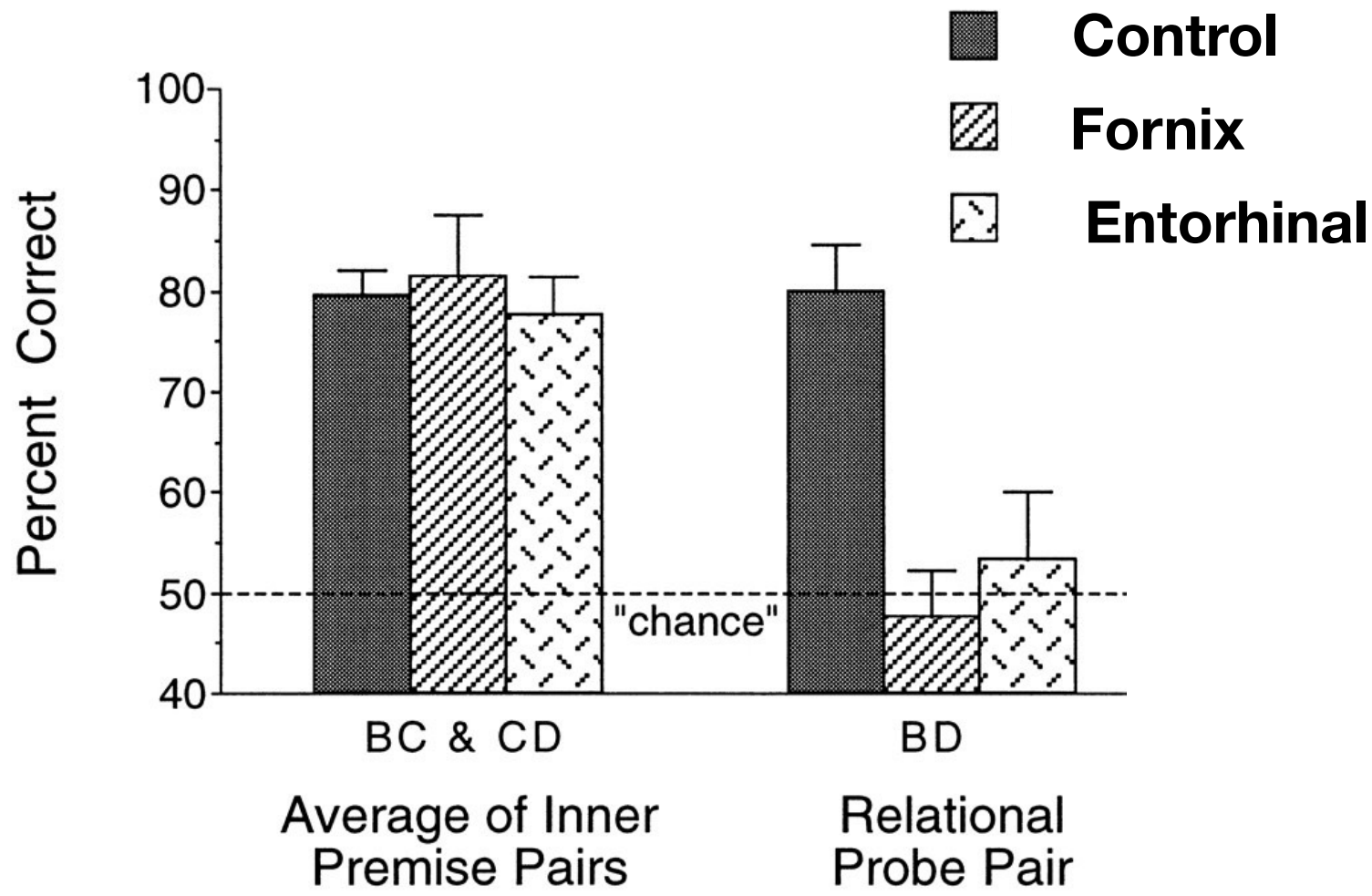
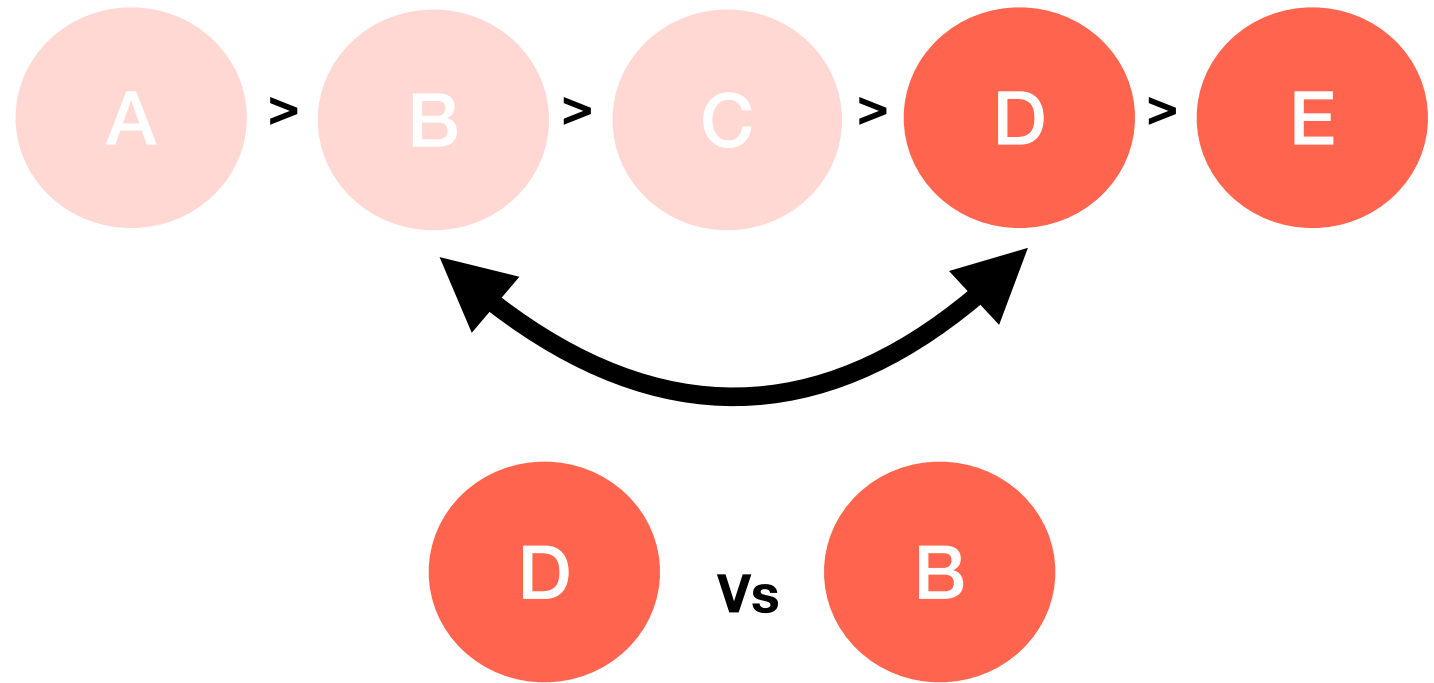
# These brain areas are critical for non-spatial 'shortcuts'



Dusek et al. PNAS 1997

# These brain areas are critical for non-spatial 'shortcuts'

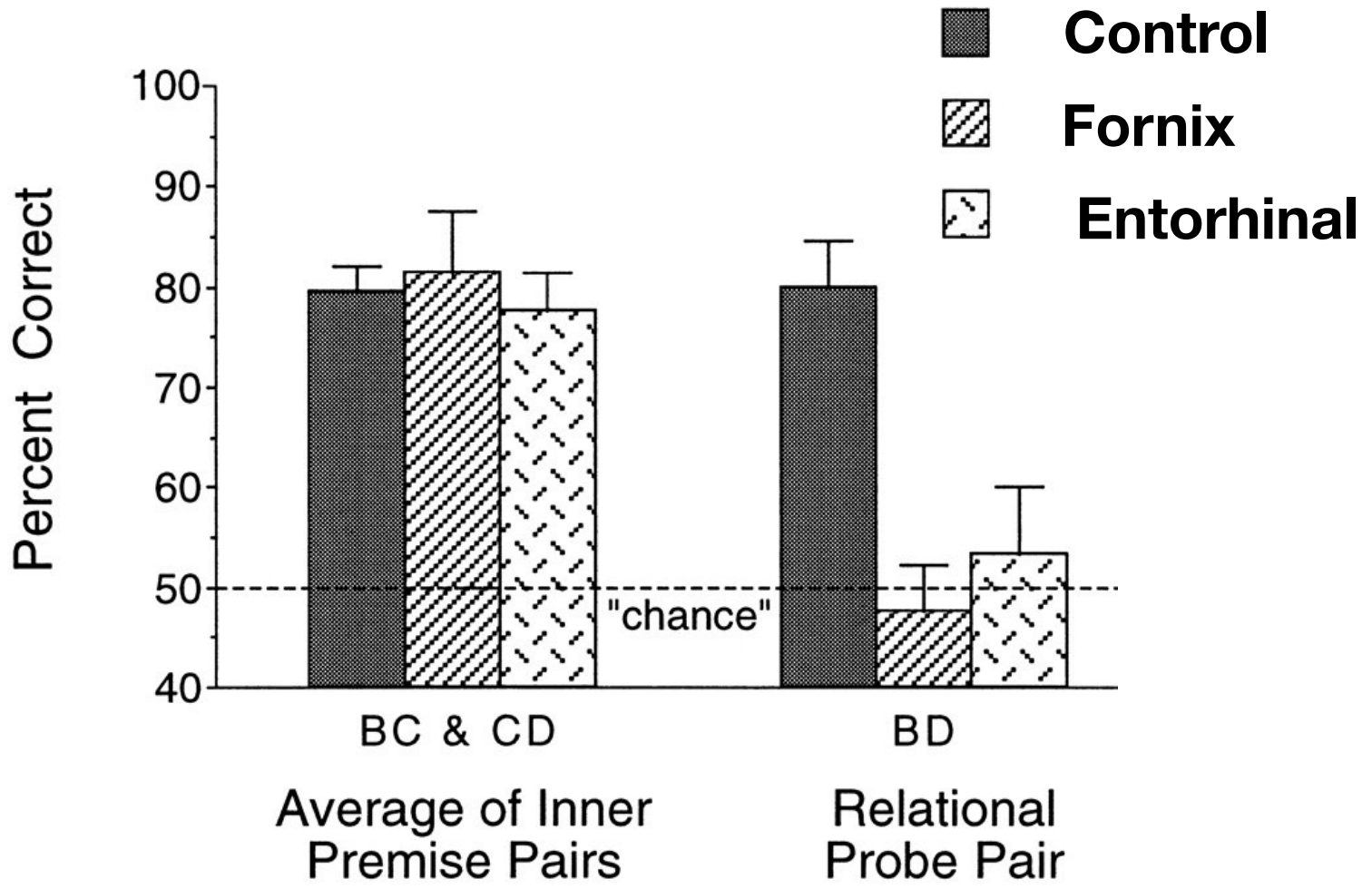
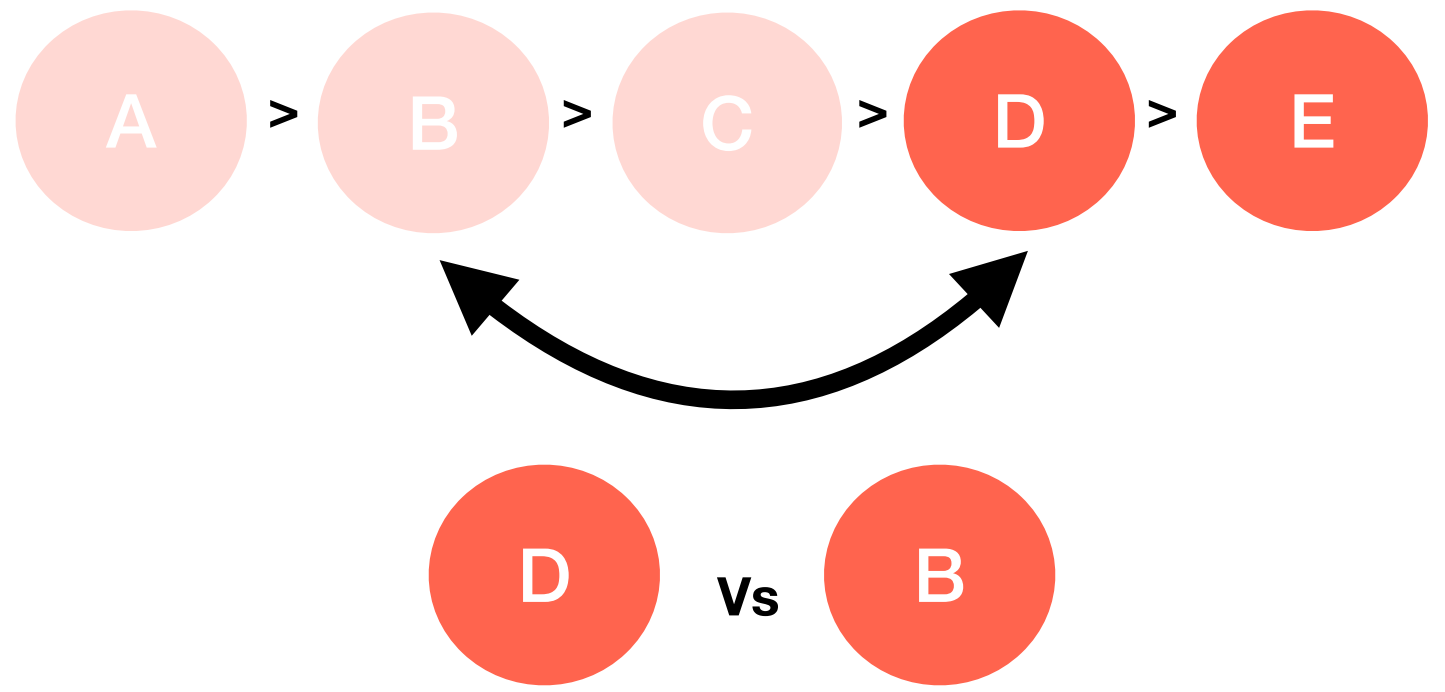
**“Spatial” cells seem to do similar things in non-spatial problems.**



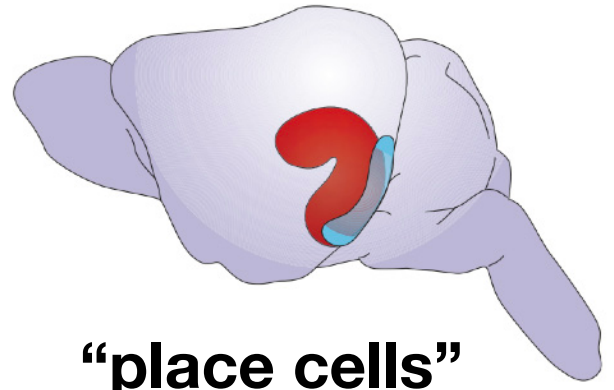
Dusek et al. PNAS 1997

# These brain areas are critical for non-spatial 'shortcuts'

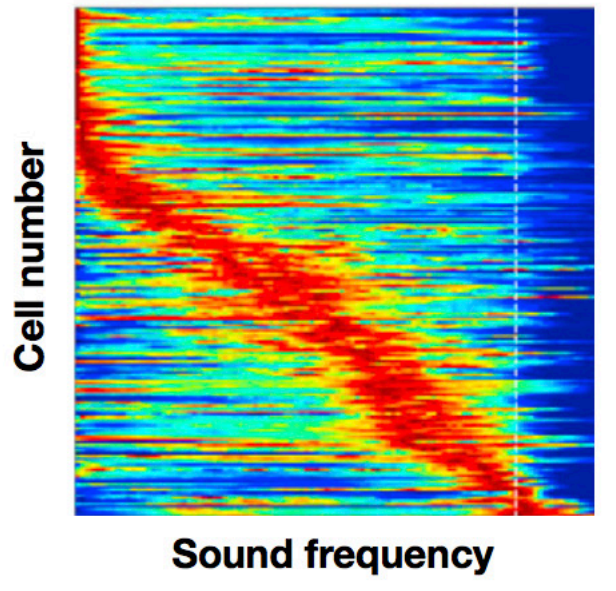
**“Spatial” cells seem to do similar things in non-spatial problems.**



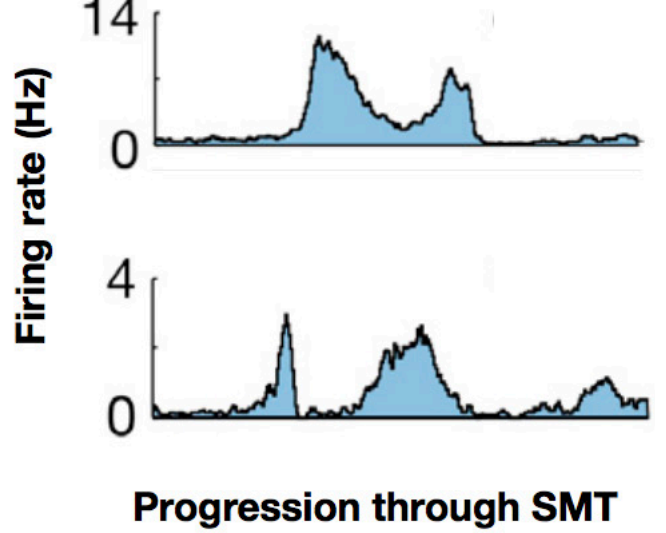
Dusek et al. PNAS 1997



“place cells”

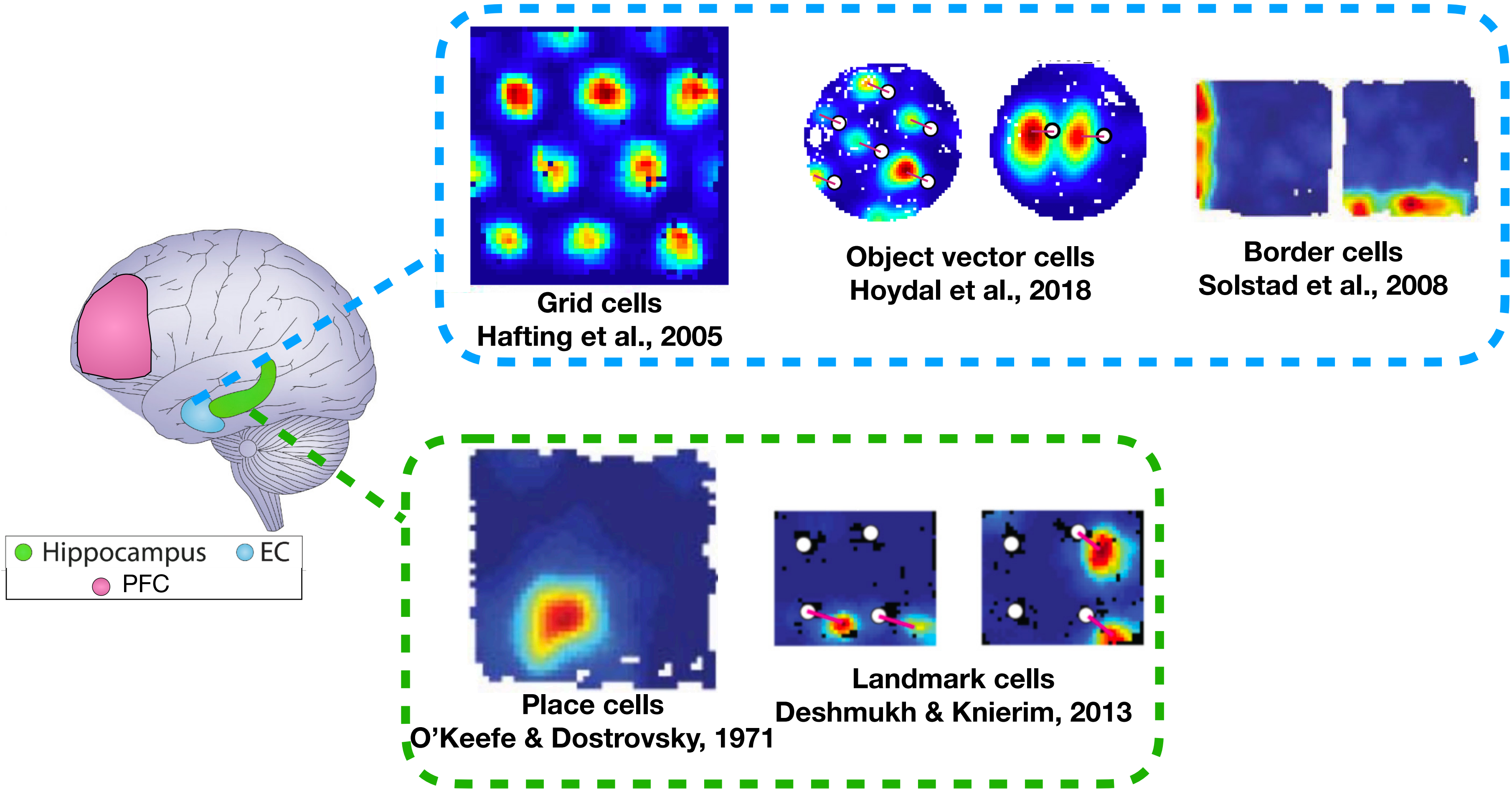


“linear grid cells”

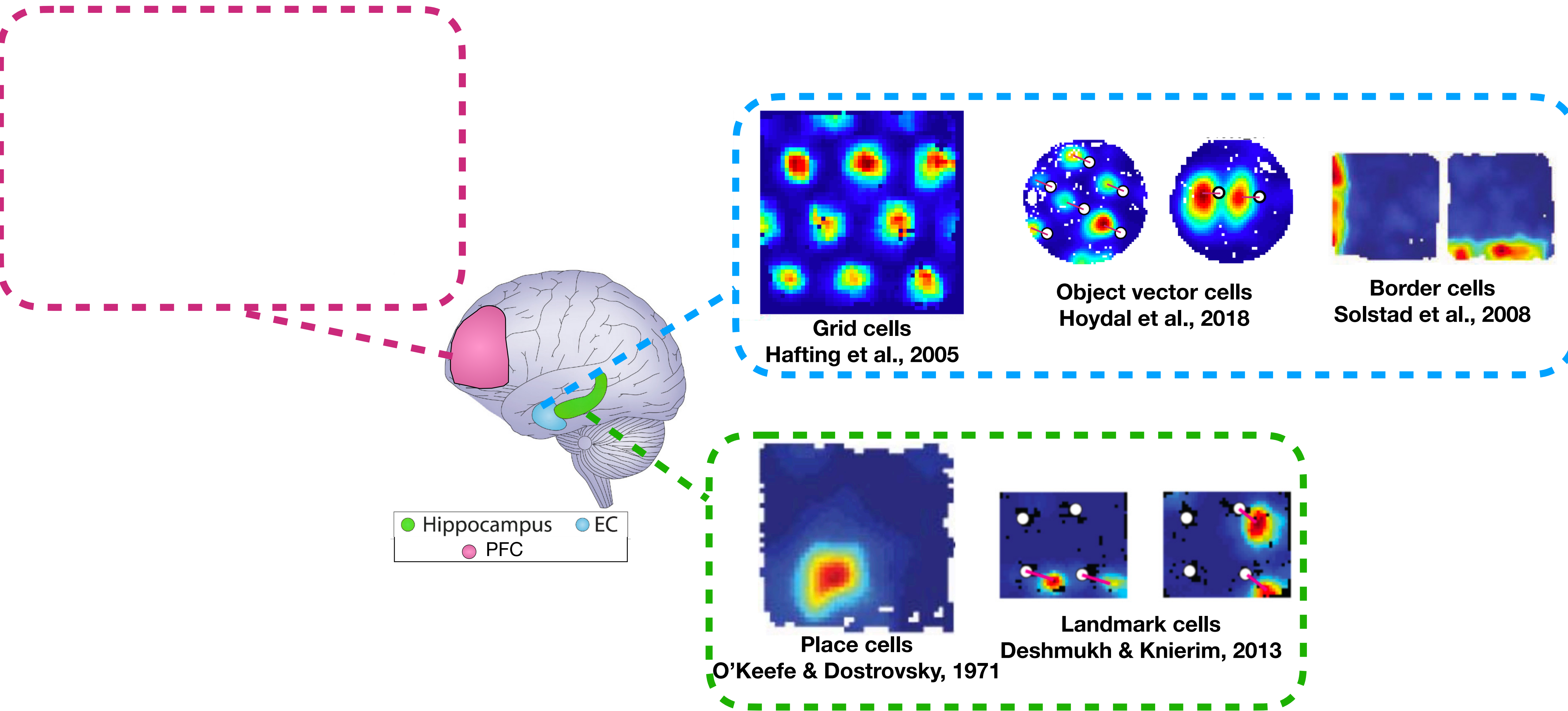


Aronov et al 2017

# But prefrontal cortex solves this problem in a different way

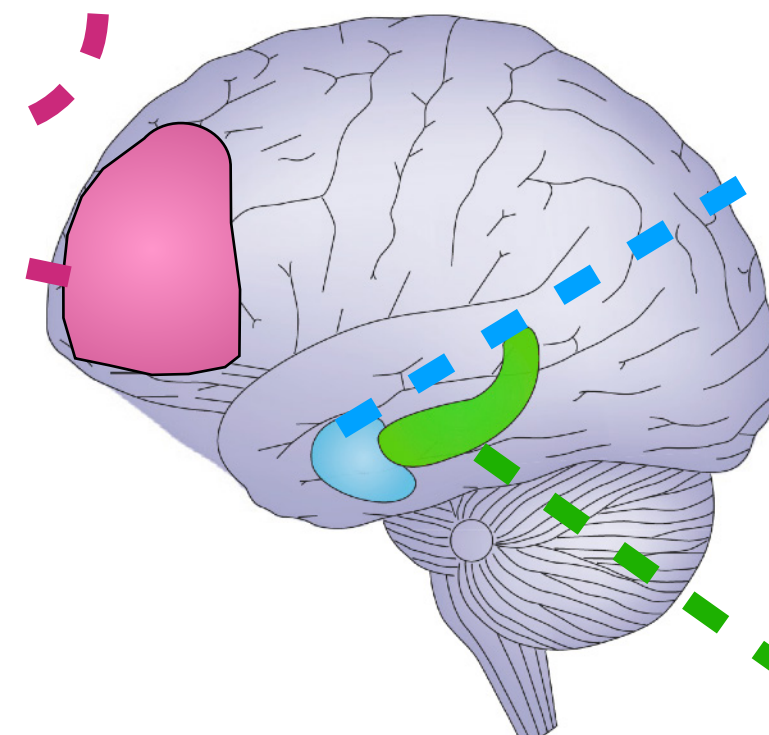
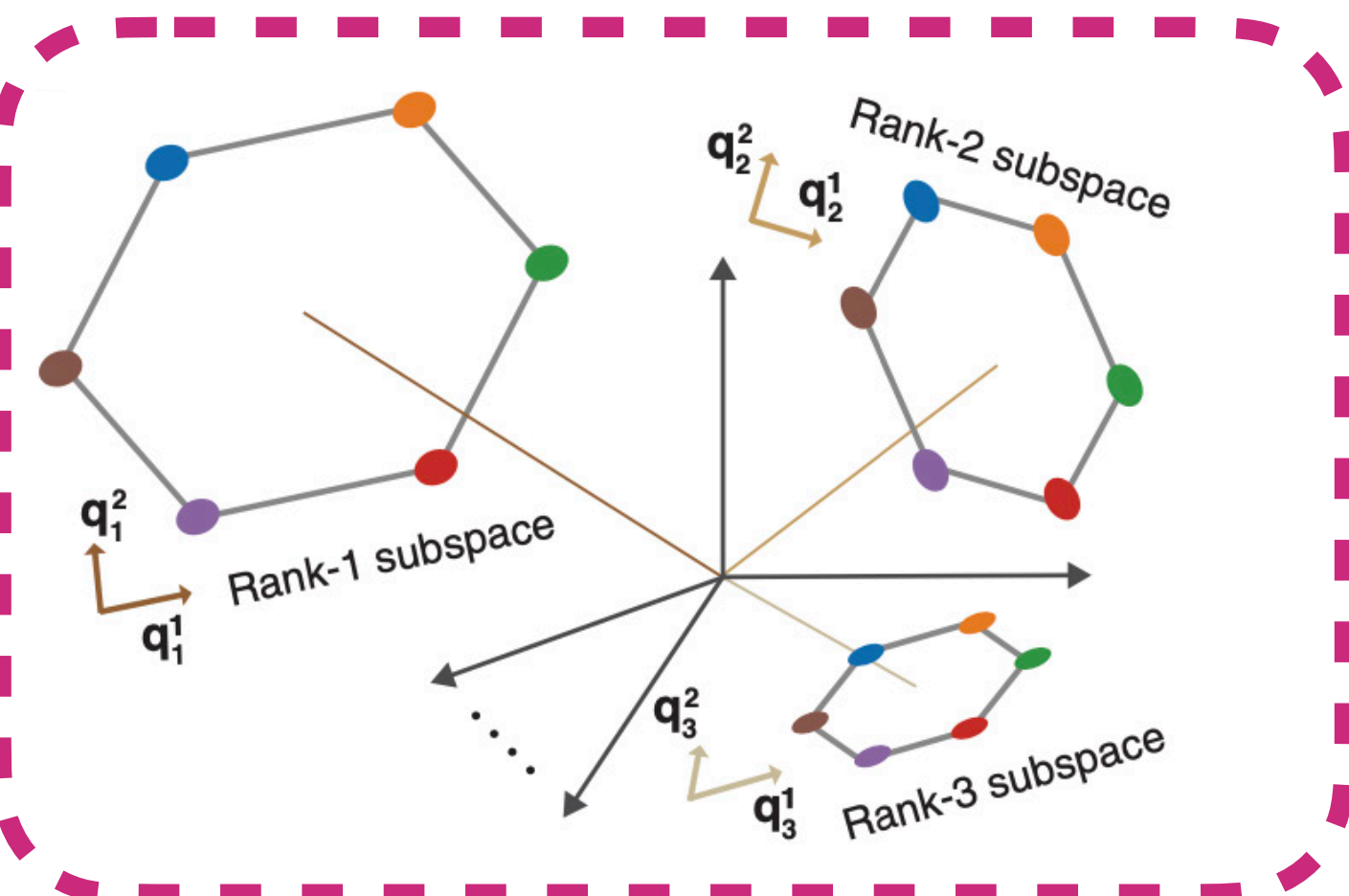


# But prefrontal cortex solves this problem in a different way

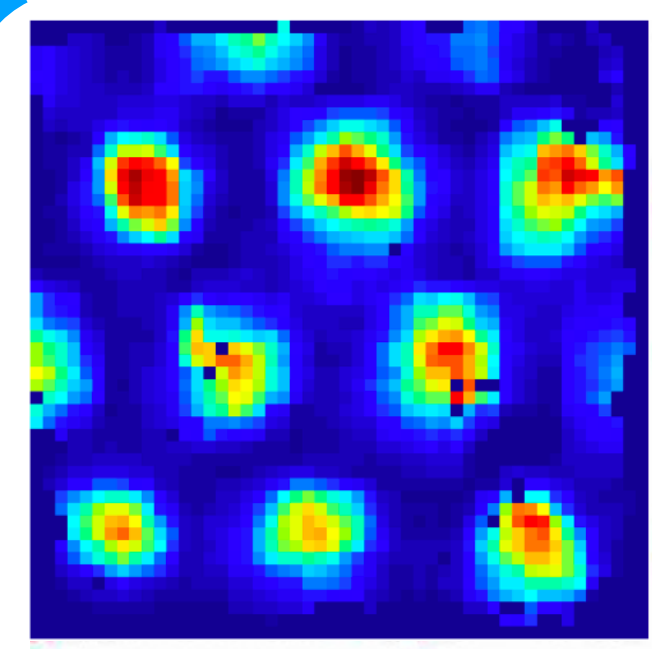




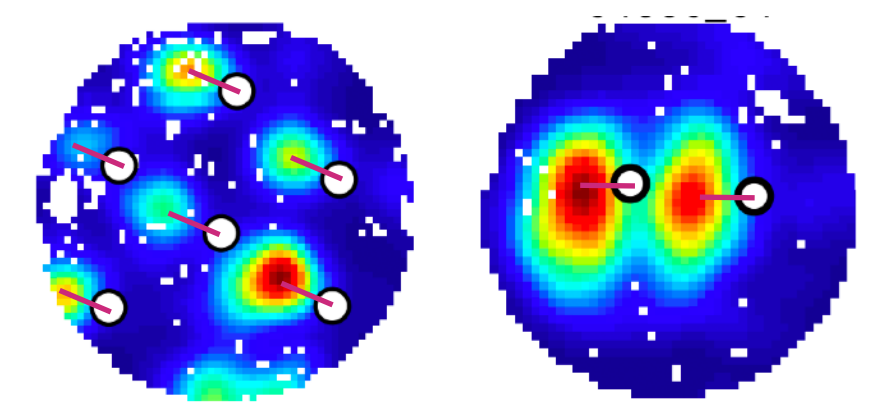
# But prefrontal cortex solves this problem in a different way



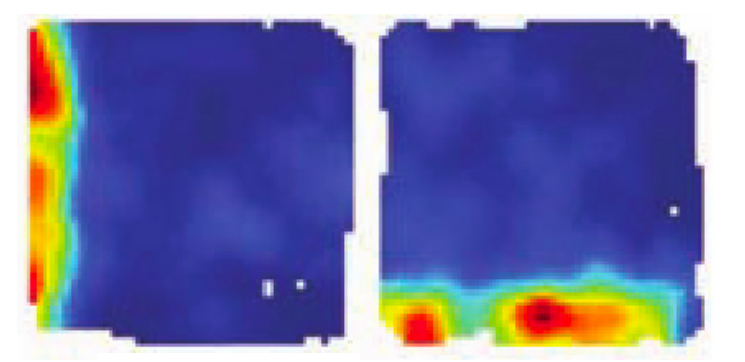
● Hippocampus ● EC  
● PFC



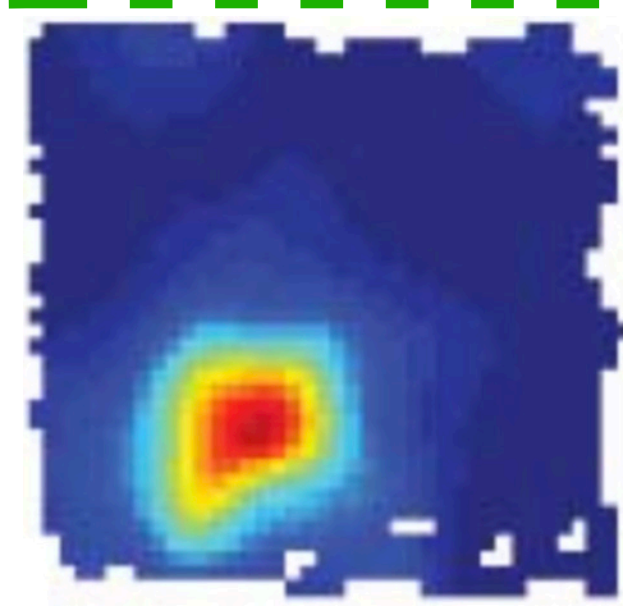
Grid cells  
Hafting et al., 2005



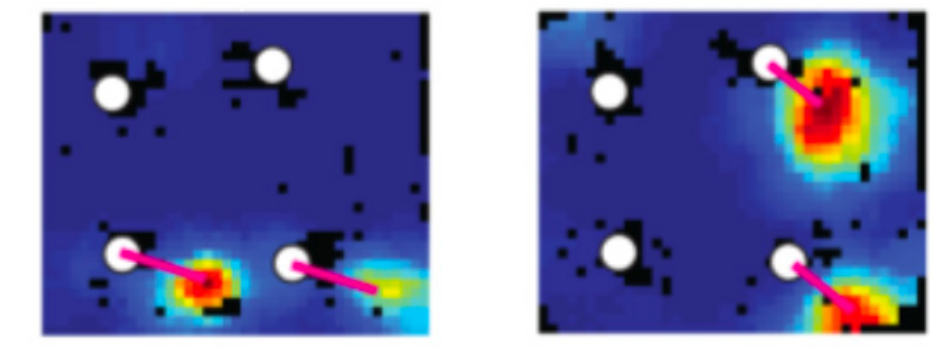
Object vector cells  
Hoydal et al., 2018



Border cells  
Solstad et al., 2008



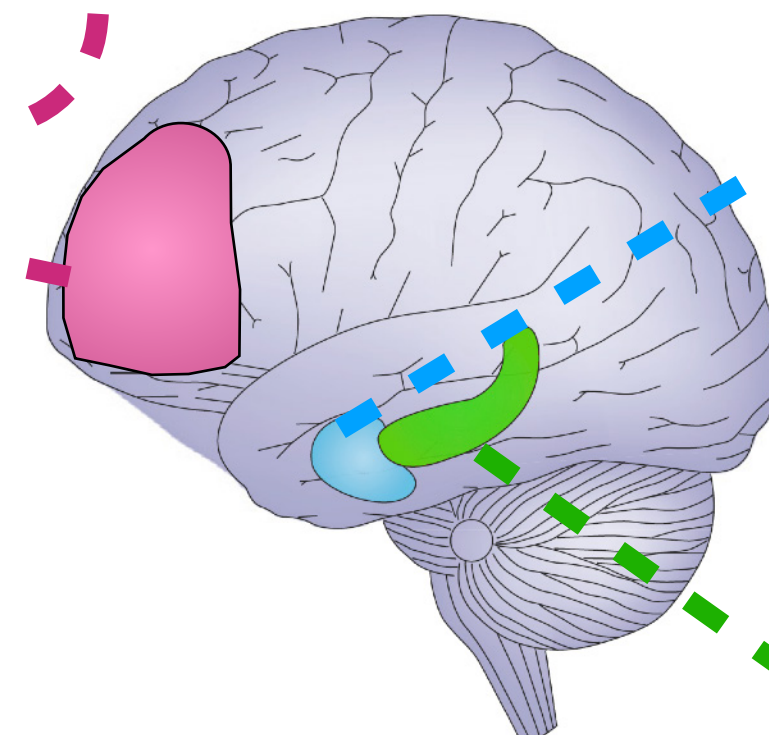
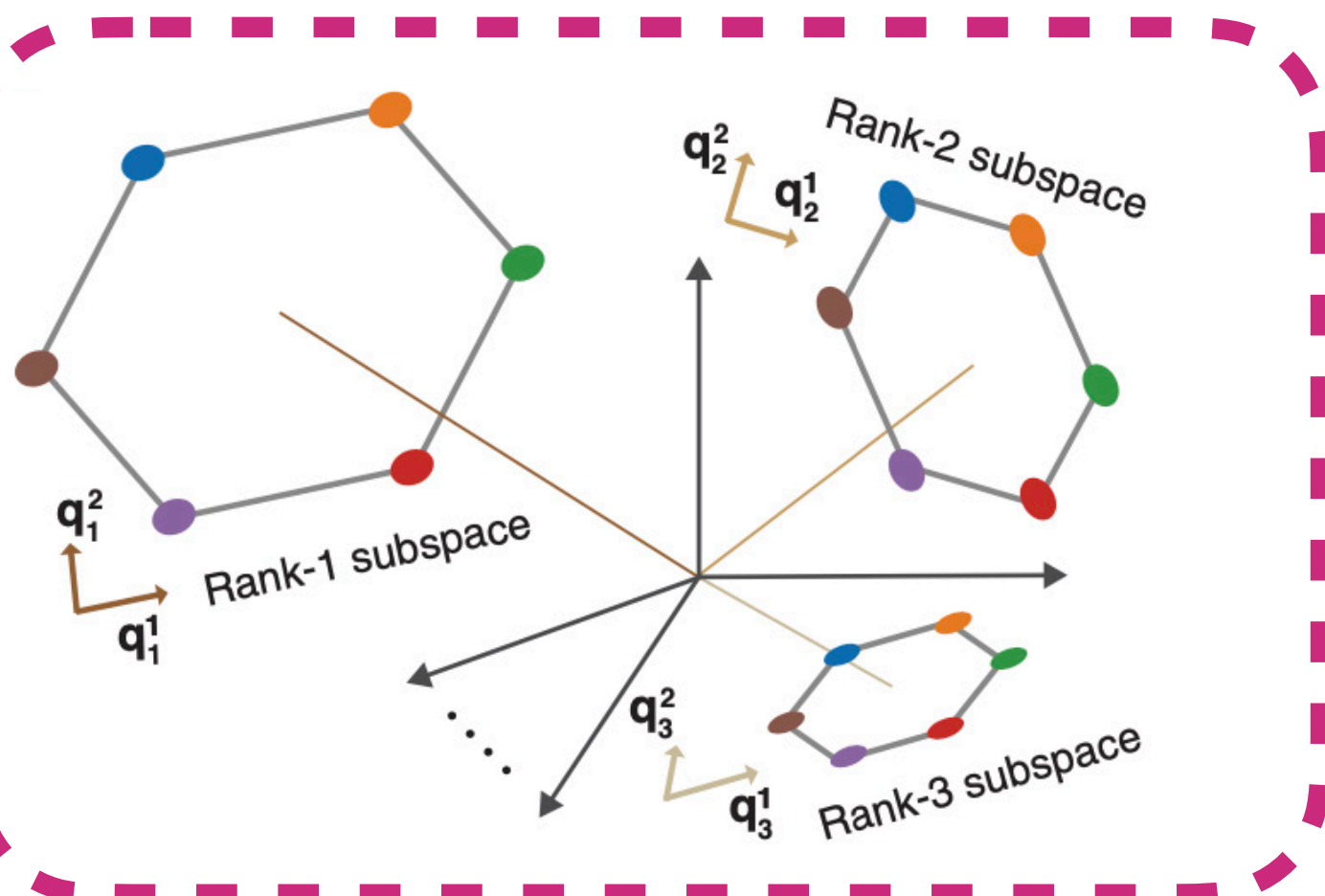
Place cells  
O'Keefe & Dostrovsky, 1971



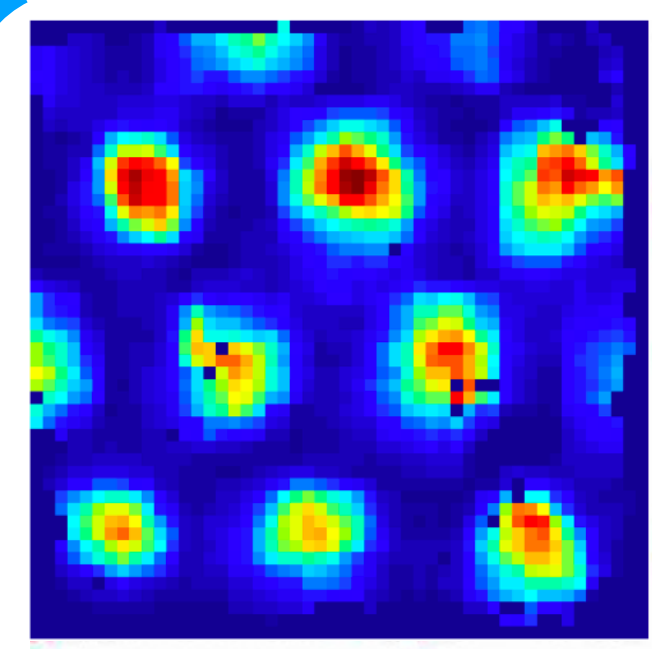
Landmark cells  
Deshmukh & Knierim, 2013

# But prefrontal cortex solves this problem in a different way

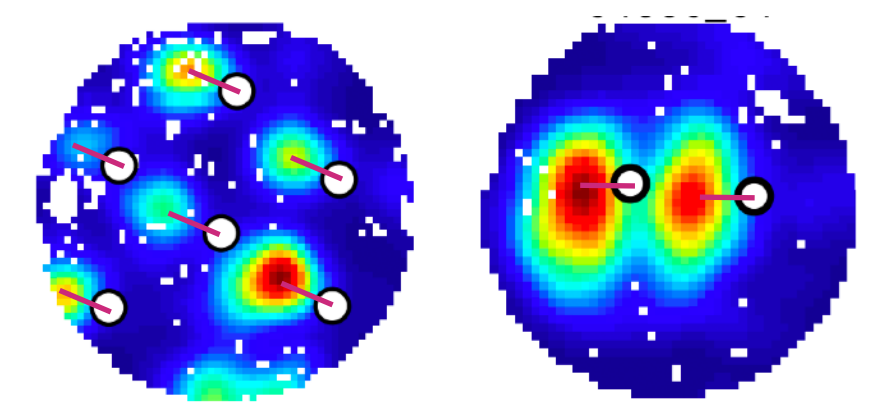
These cells tell you where you are or what you're seeing right now



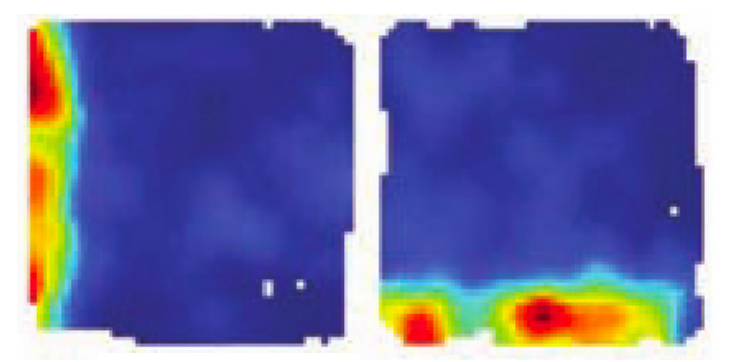
- Hippocampus
- EC
- PFC



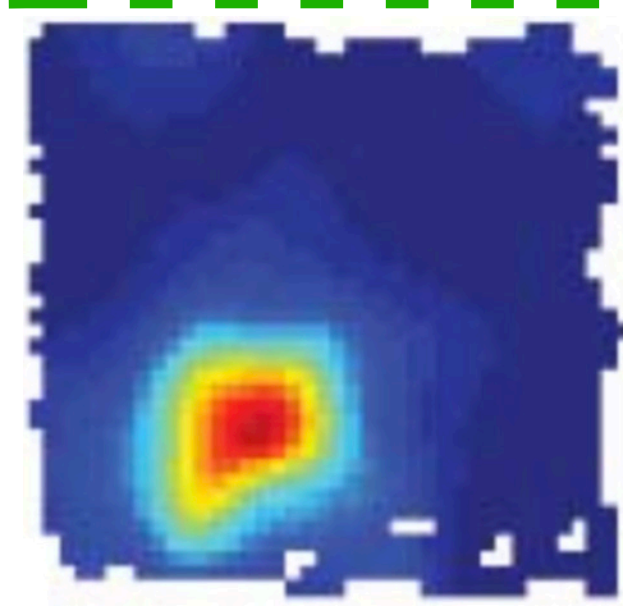
Grid cells  
Hafting et al., 2005



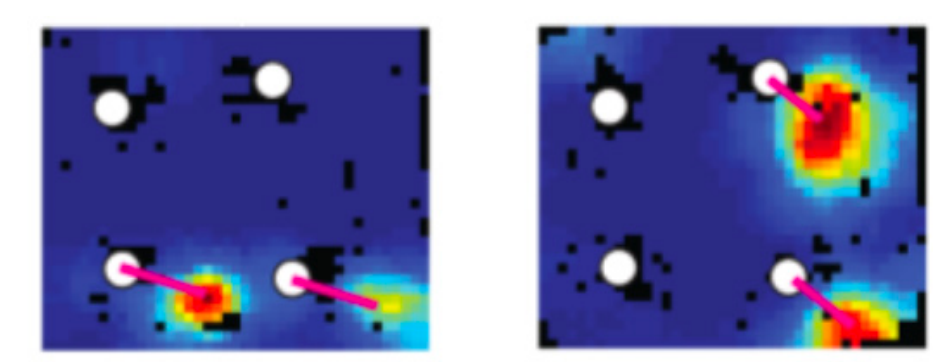
Object vector cells  
Hoydal et al., 2018



Border cells  
Solstad et al., 2008



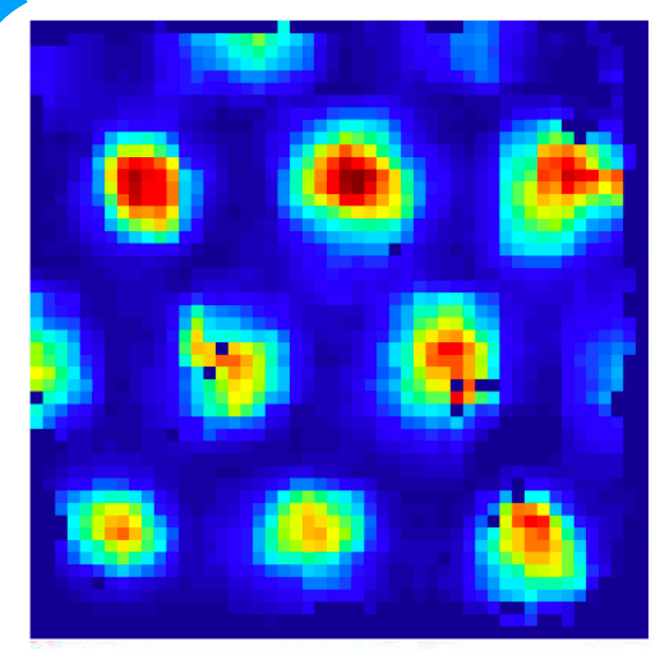
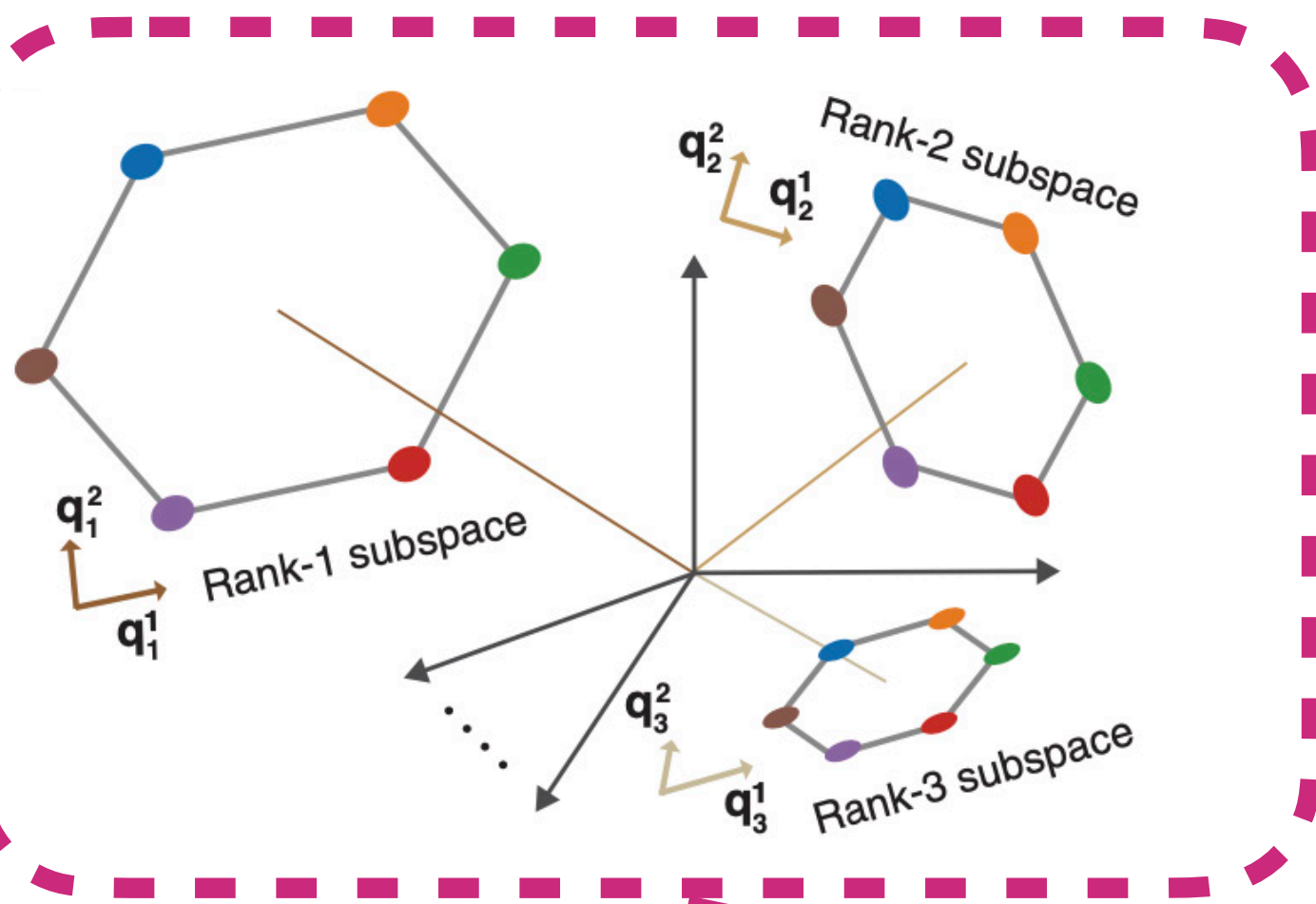
Place cells  
O'Keefe & Dostrovsky, 1971



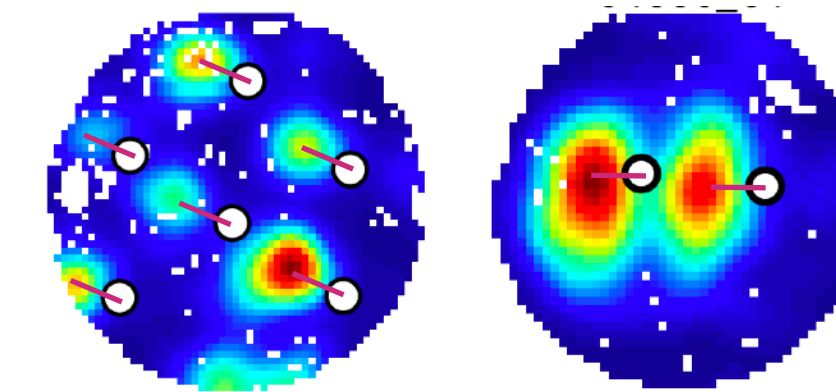
Landmark cells  
Deshmukh & Knierim, 2013

# But prefrontal cortex solves this problem in a different way

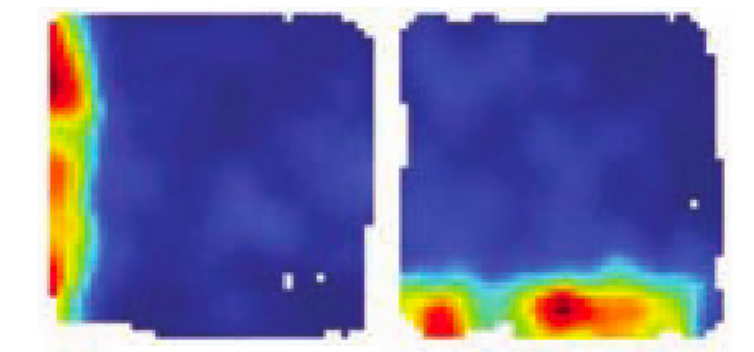
These cells tell you where you are or what you're seeing right now



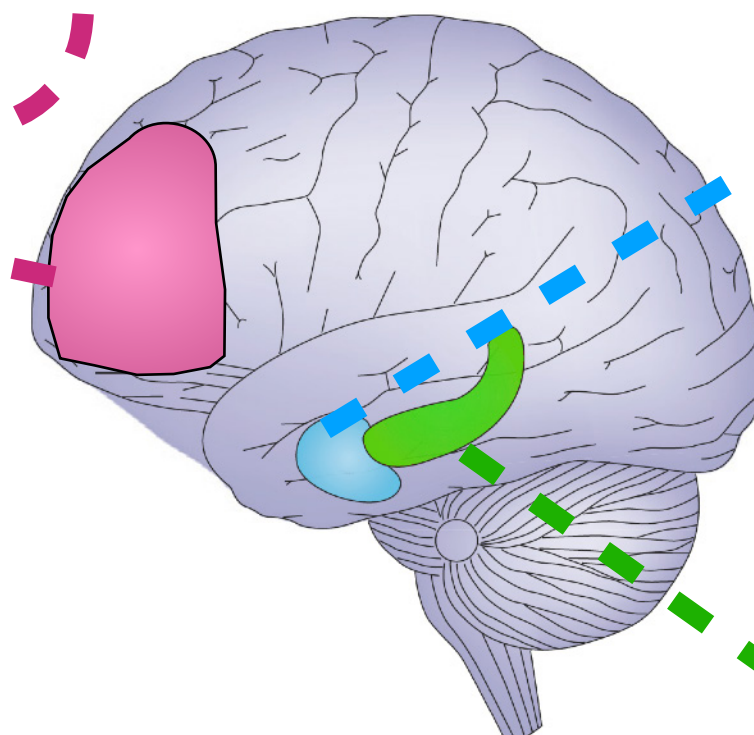
Grid cells  
Hafting et al., 2005



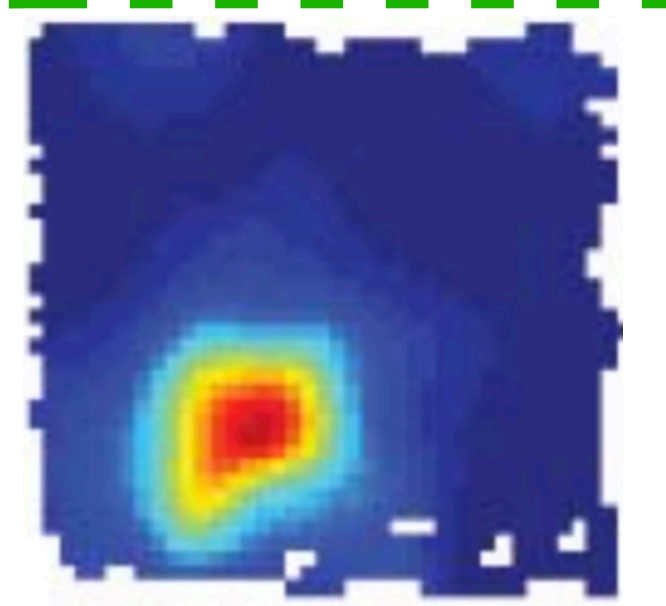
Object vector cells  
Hoydal et al., 2018



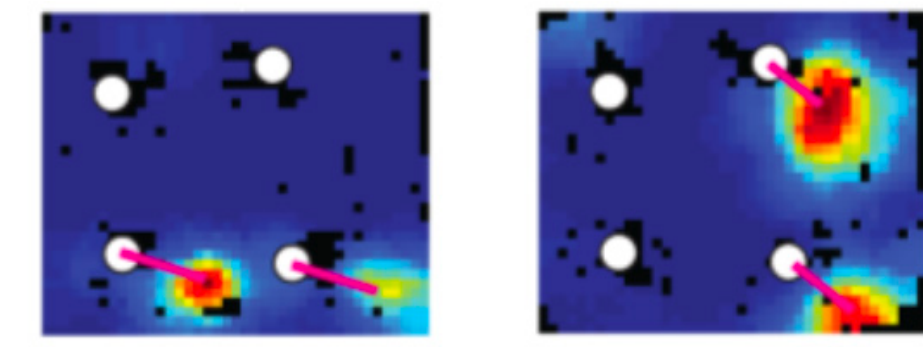
Border cells  
Solstad et al., 2008



- Hippocampus
- EC
- PFC



Place cells  
O'Keefe & Dostrovsky, 1971



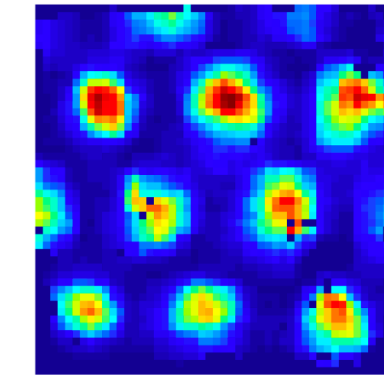
Landmark cells  
Deshmukh & Knierim, 2013

These cells include the past and future

# **Puzzles of cognitive maps in the brain**

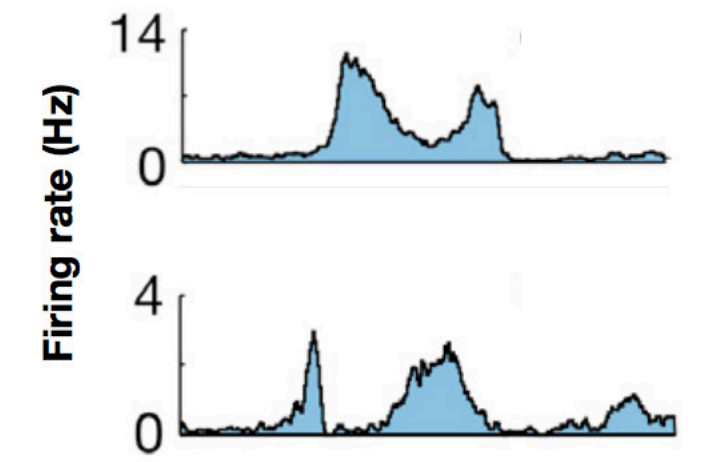
# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?



?

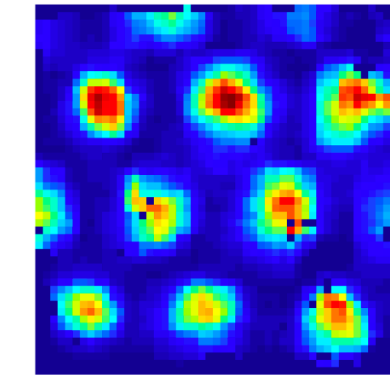
=



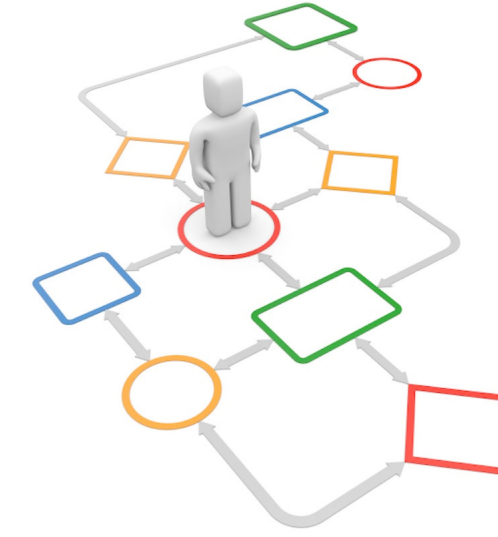
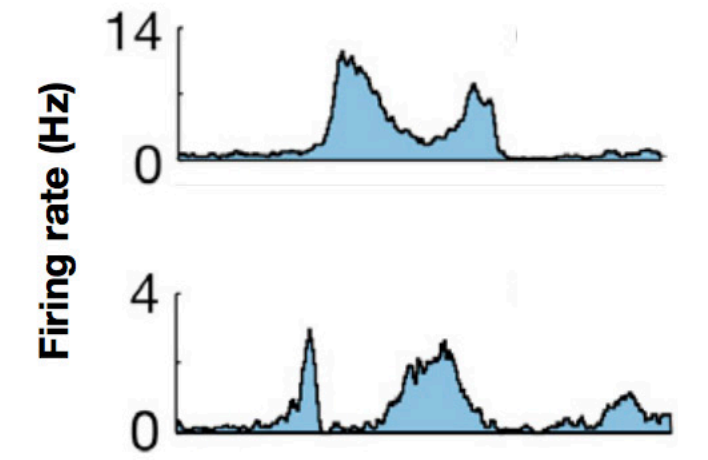
# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?

How can brains learn these maps?

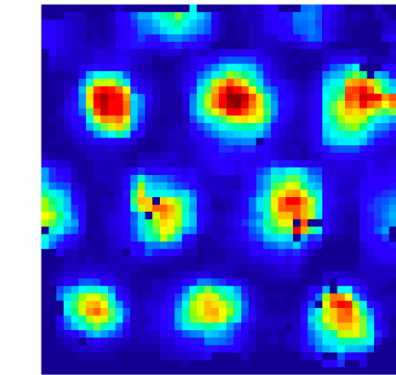


?  
=

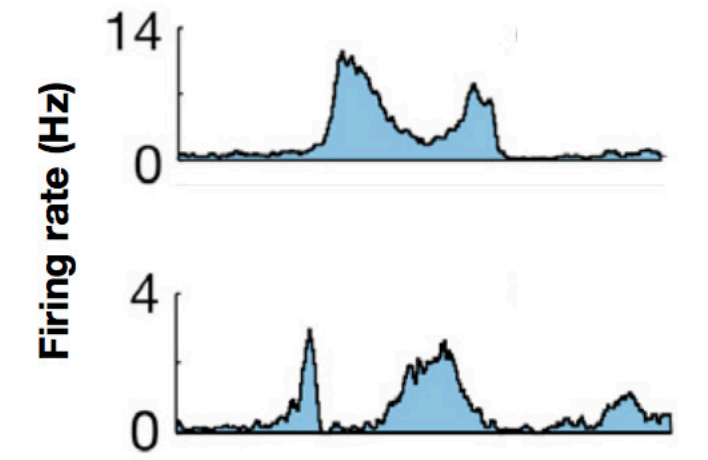


# Puzzles of cognitive maps in the brain

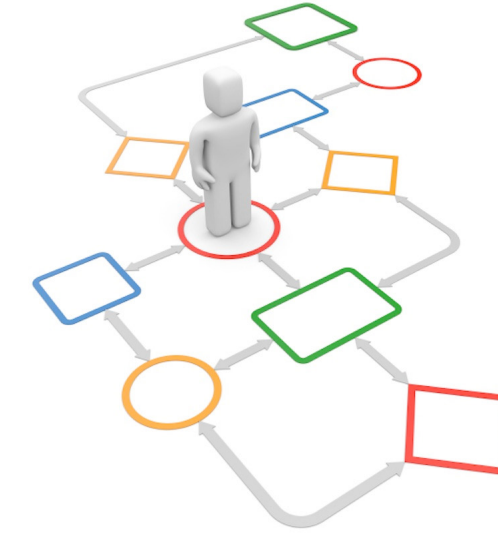
How does the same system do space and non-space?



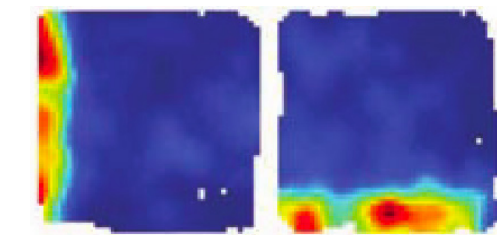
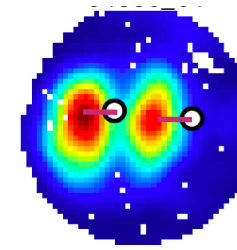
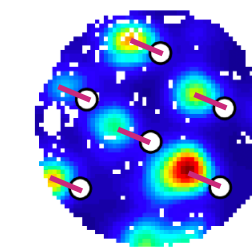
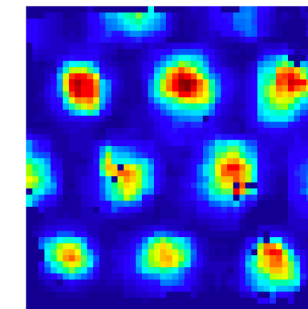
?  
=



How can brains learn these maps?

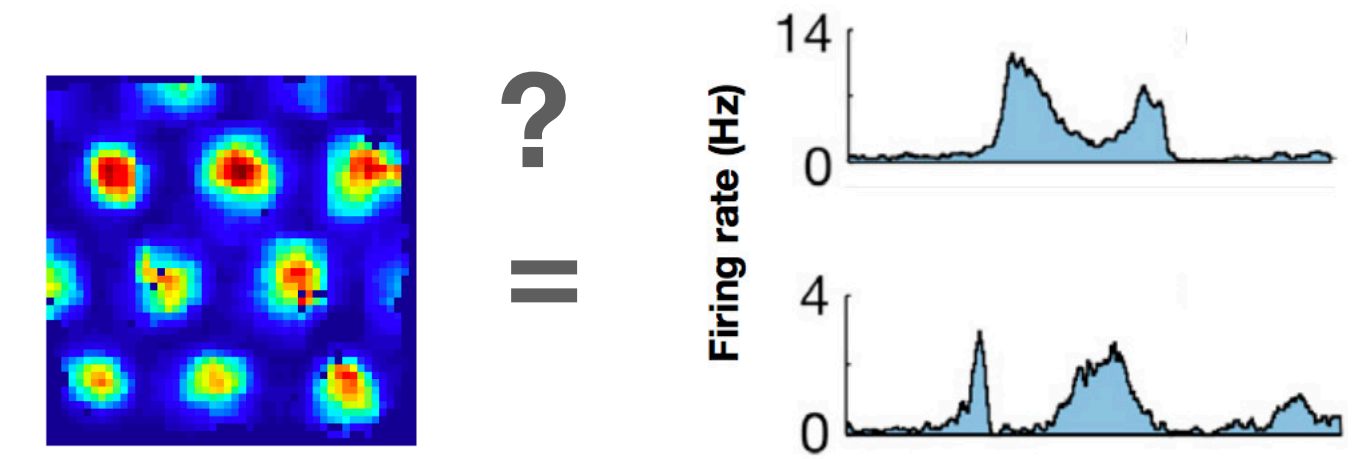


Why do the neurons look the ways they do?

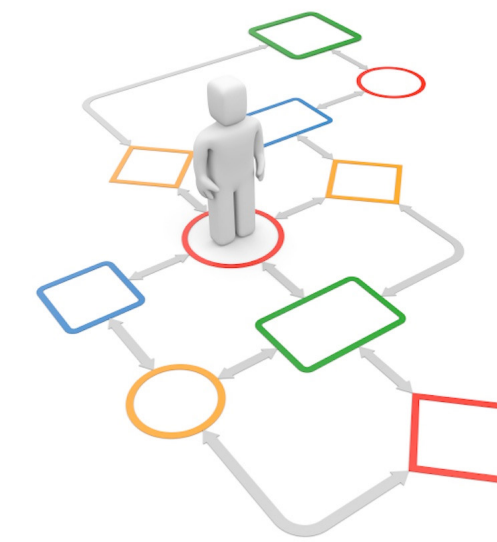


# Puzzles of cognitive maps in the brain

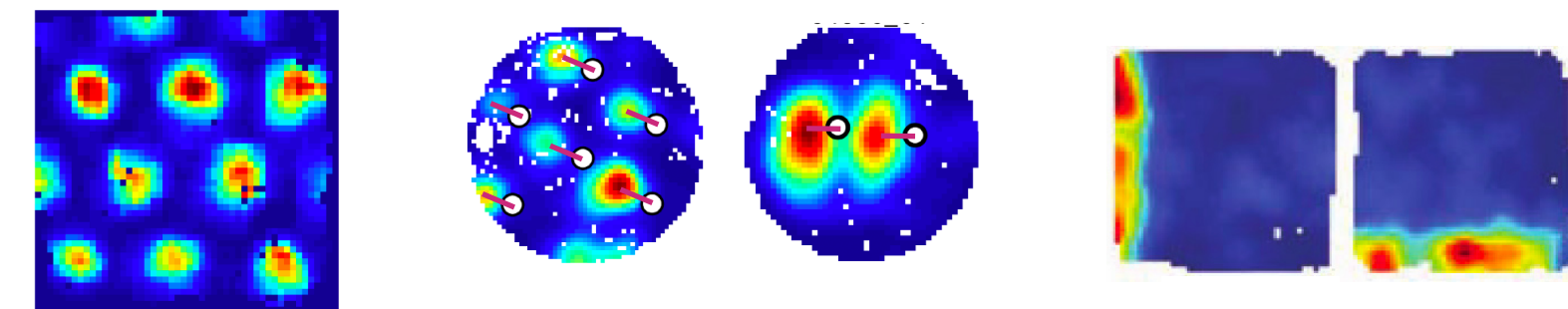
How does the same system do space and non-space?



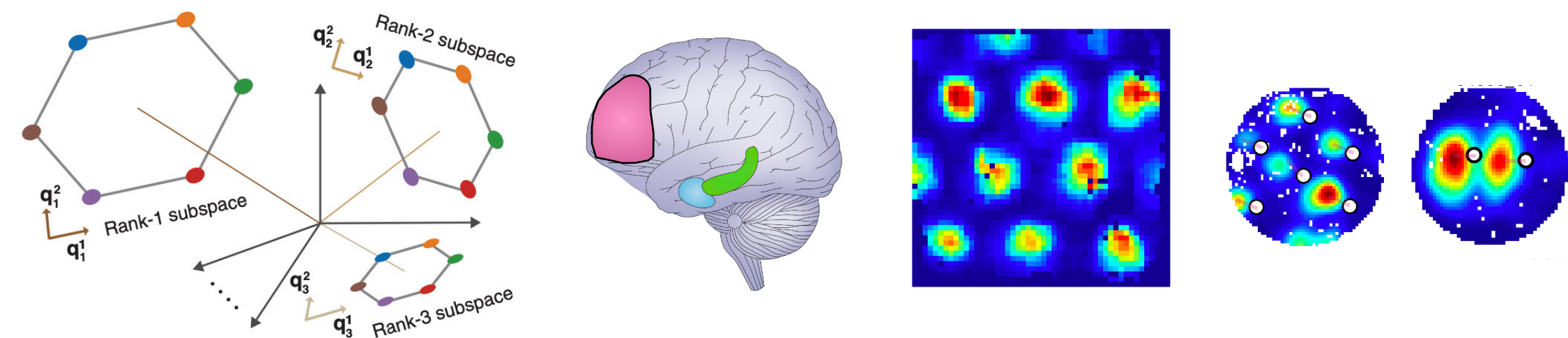
How can brains learn these maps?



Why do the neurons look the ways they do?



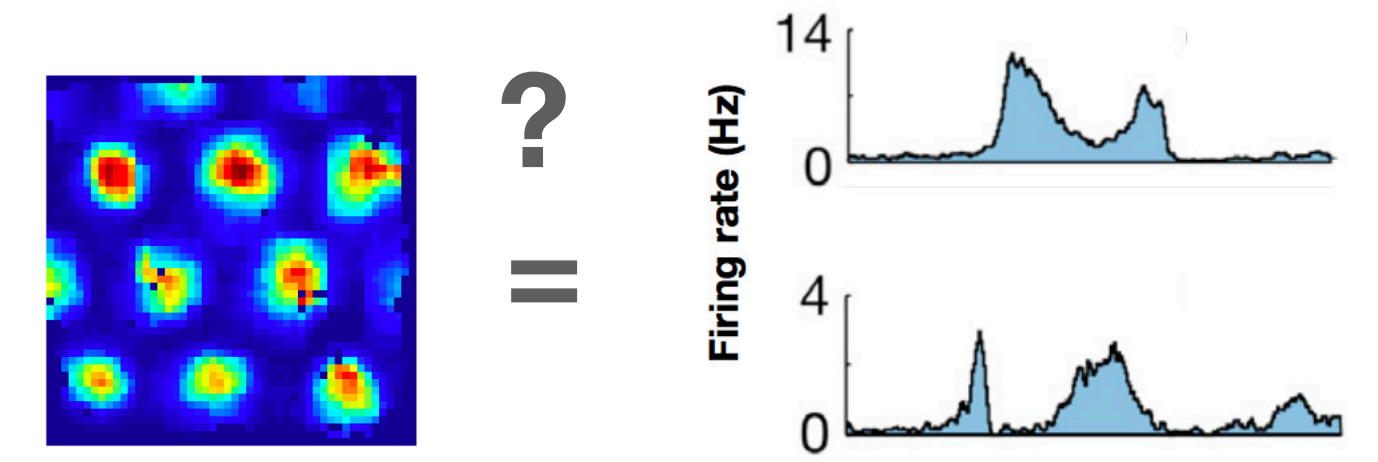
How do different brain regions solve the same problem in different ways?



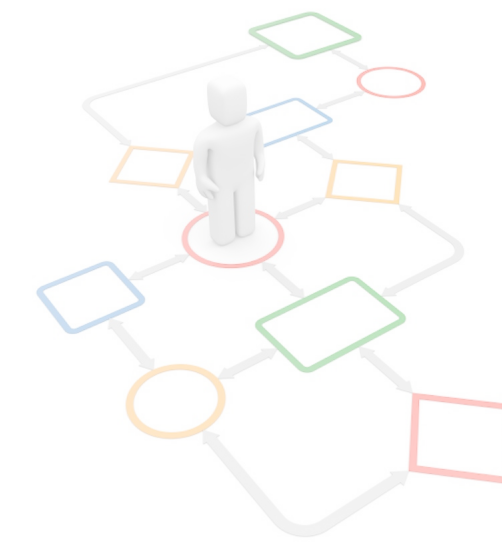


# Puzzles of cognitive maps in the brain

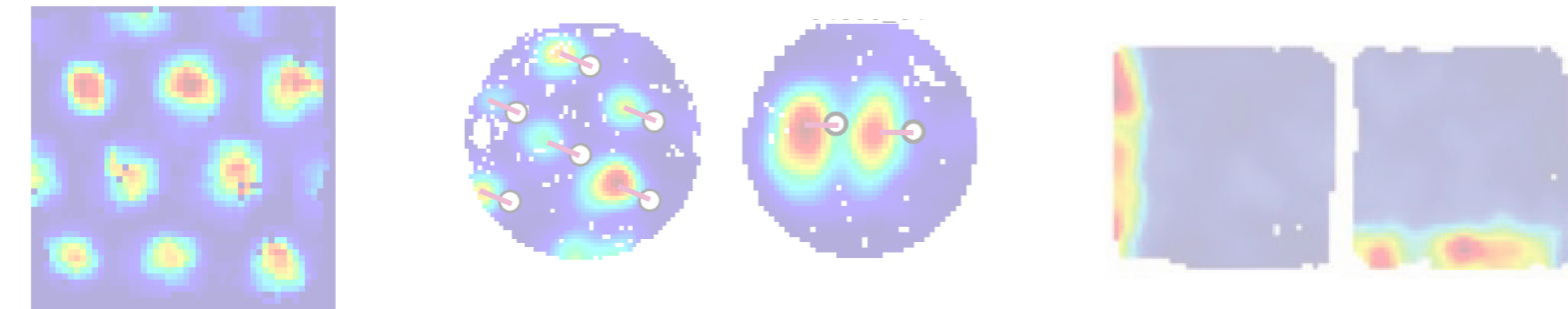
How does the same system do space and non-space?



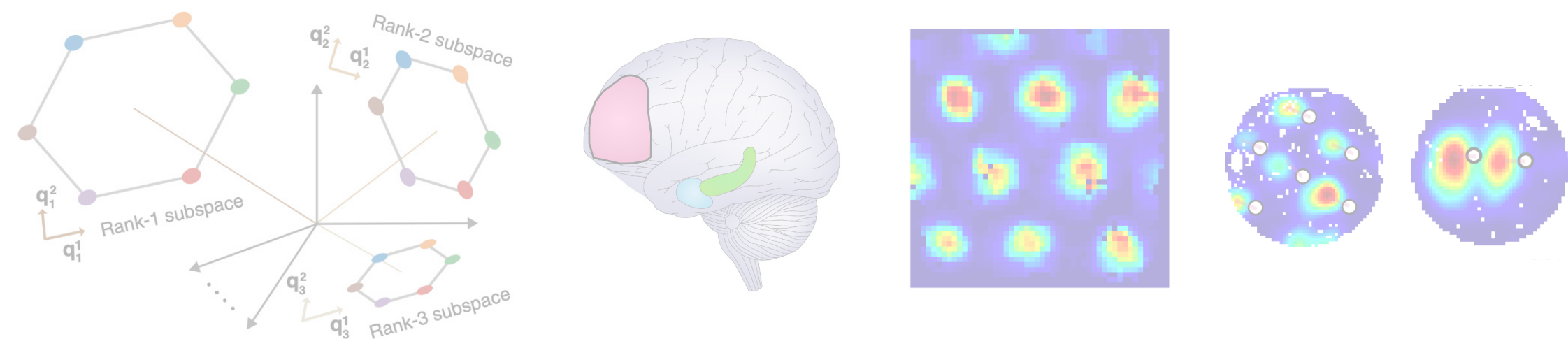
How can brains learn these maps?



Why do the neurons look the ways they do?

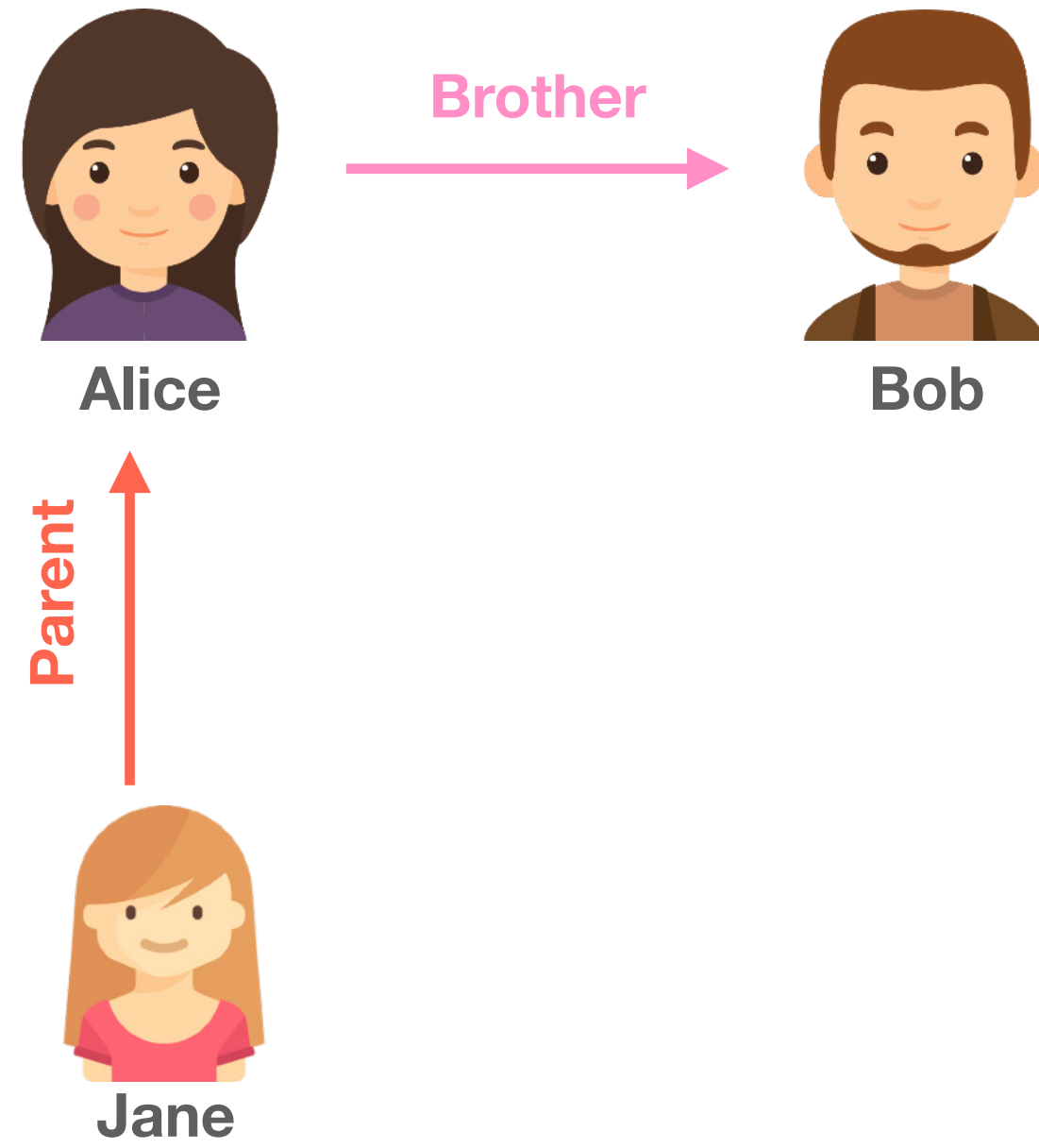


How do different brain regions solve the same problem in different ways?

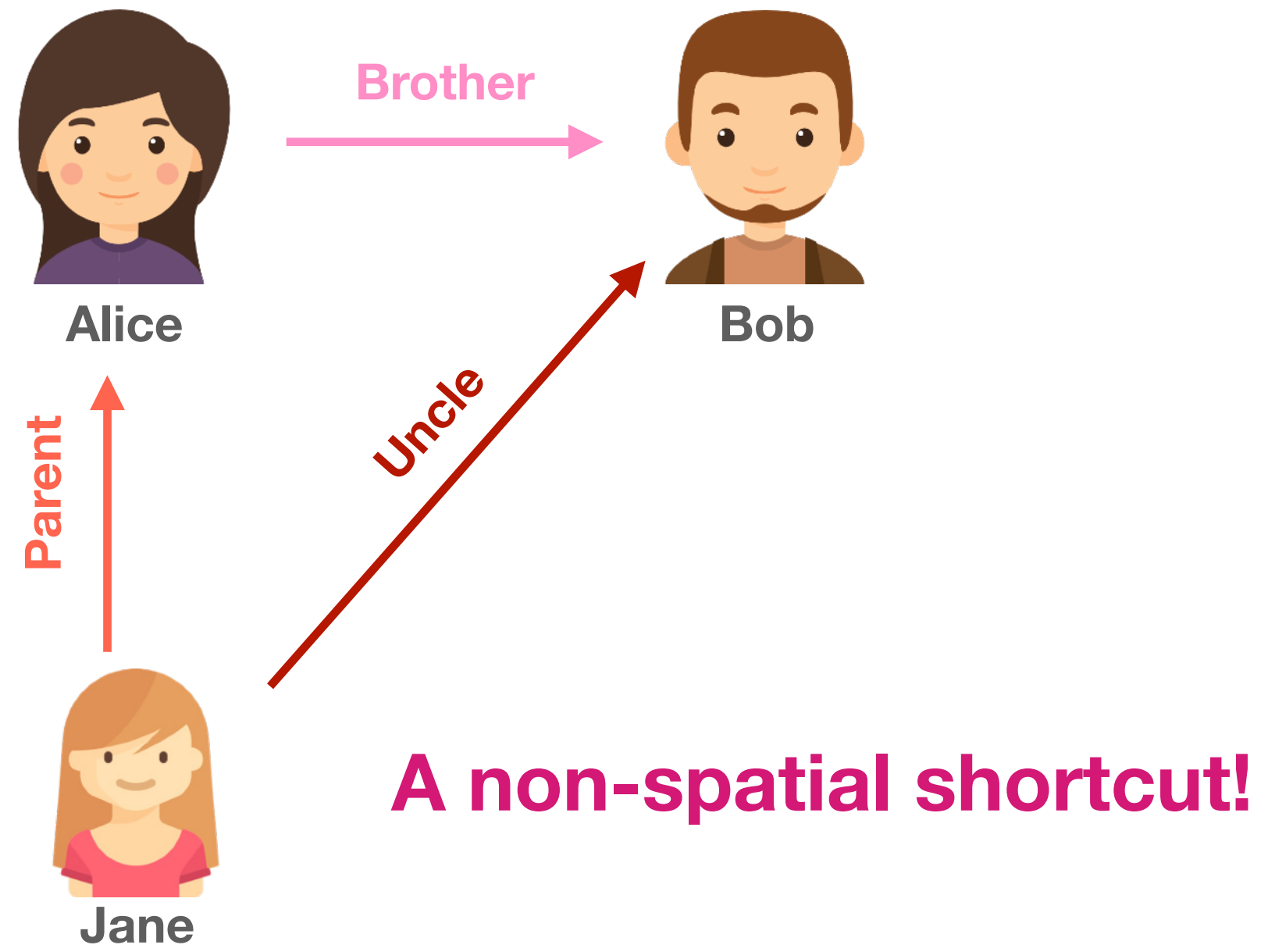


**Space and non-space unified by abstracting relationships**

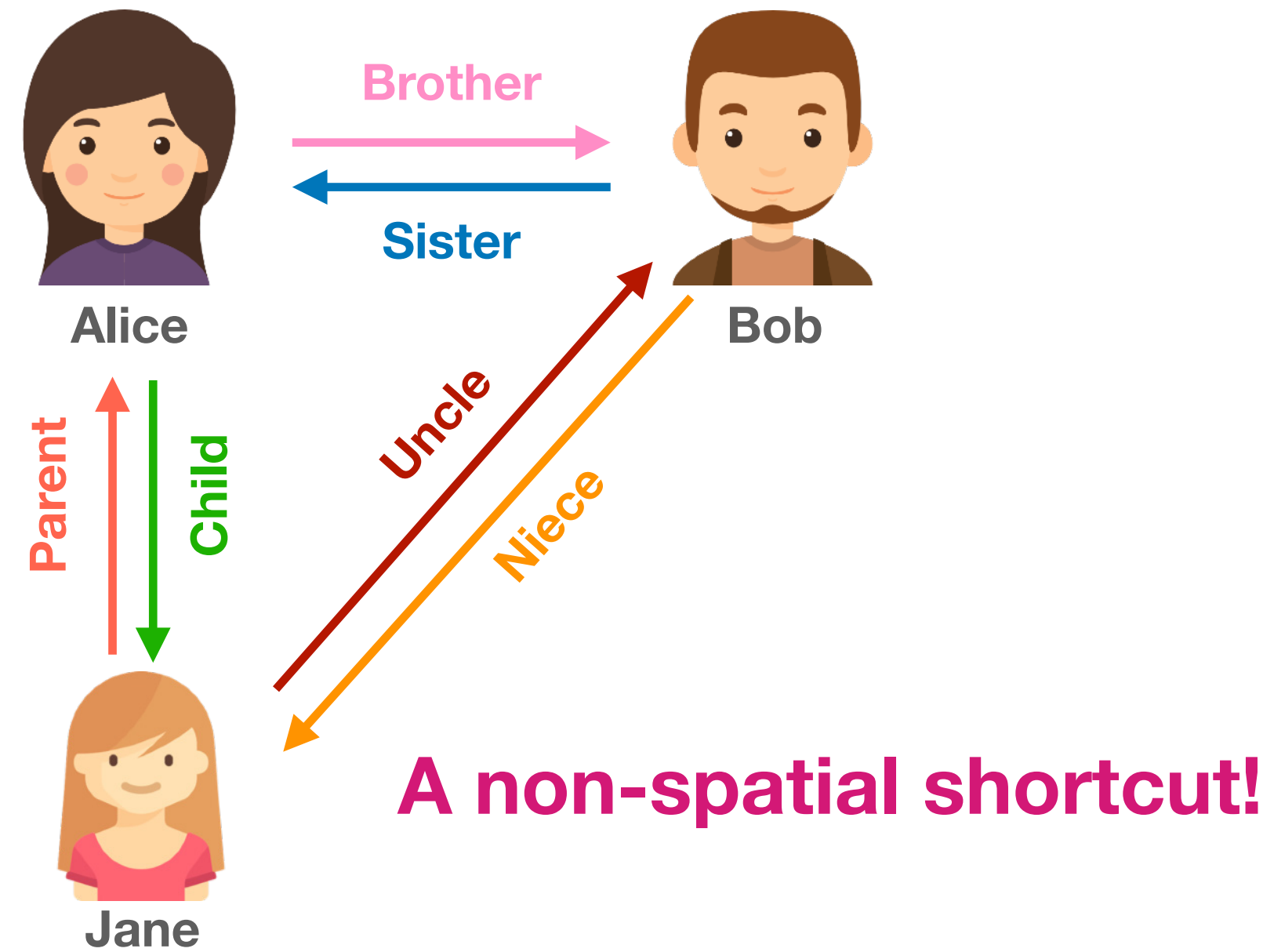
# Space and non-space unified by abstracting relationships



# Space and non-space unified by abstracting relationships



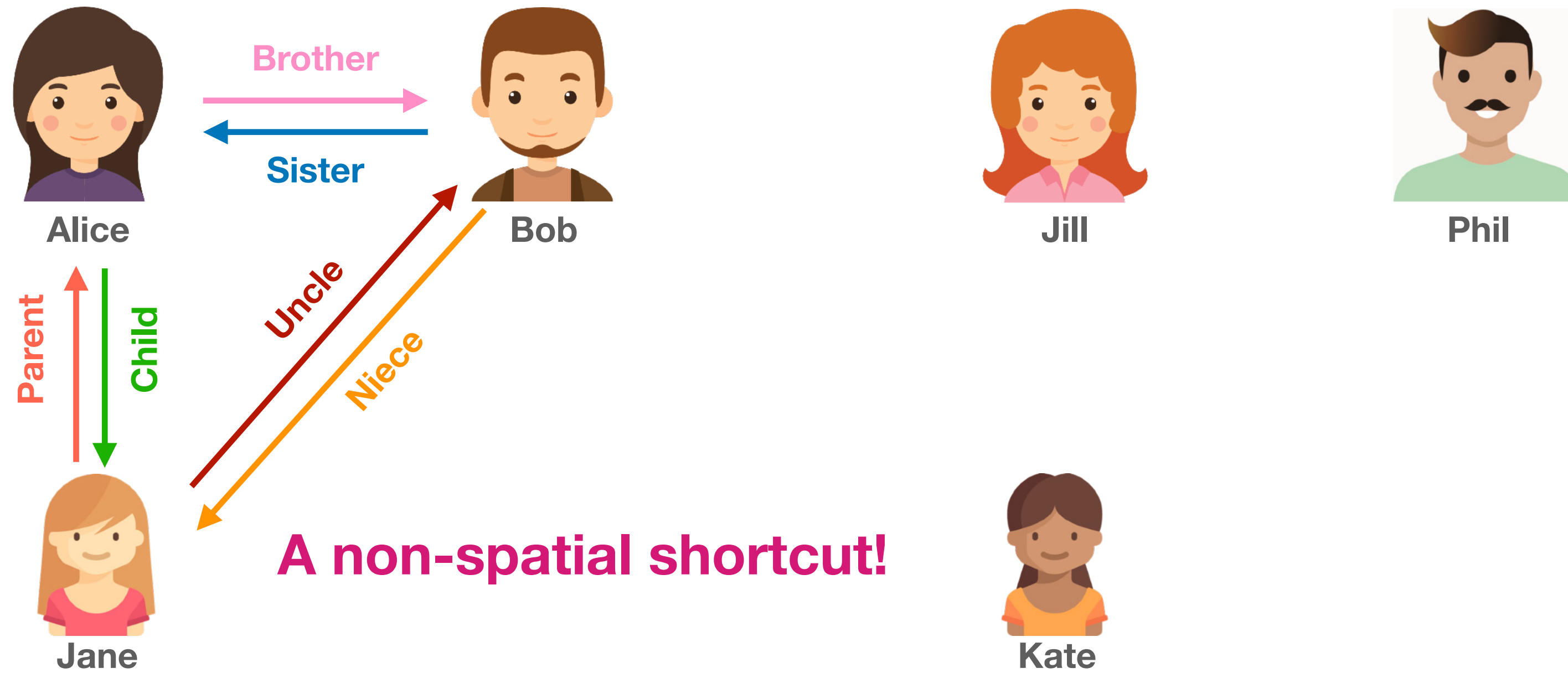
# Space and non-space unified by abstracting relationships



**Parent + Brother = Uncle**

Know which relationships add up to  
the same thing = path integration

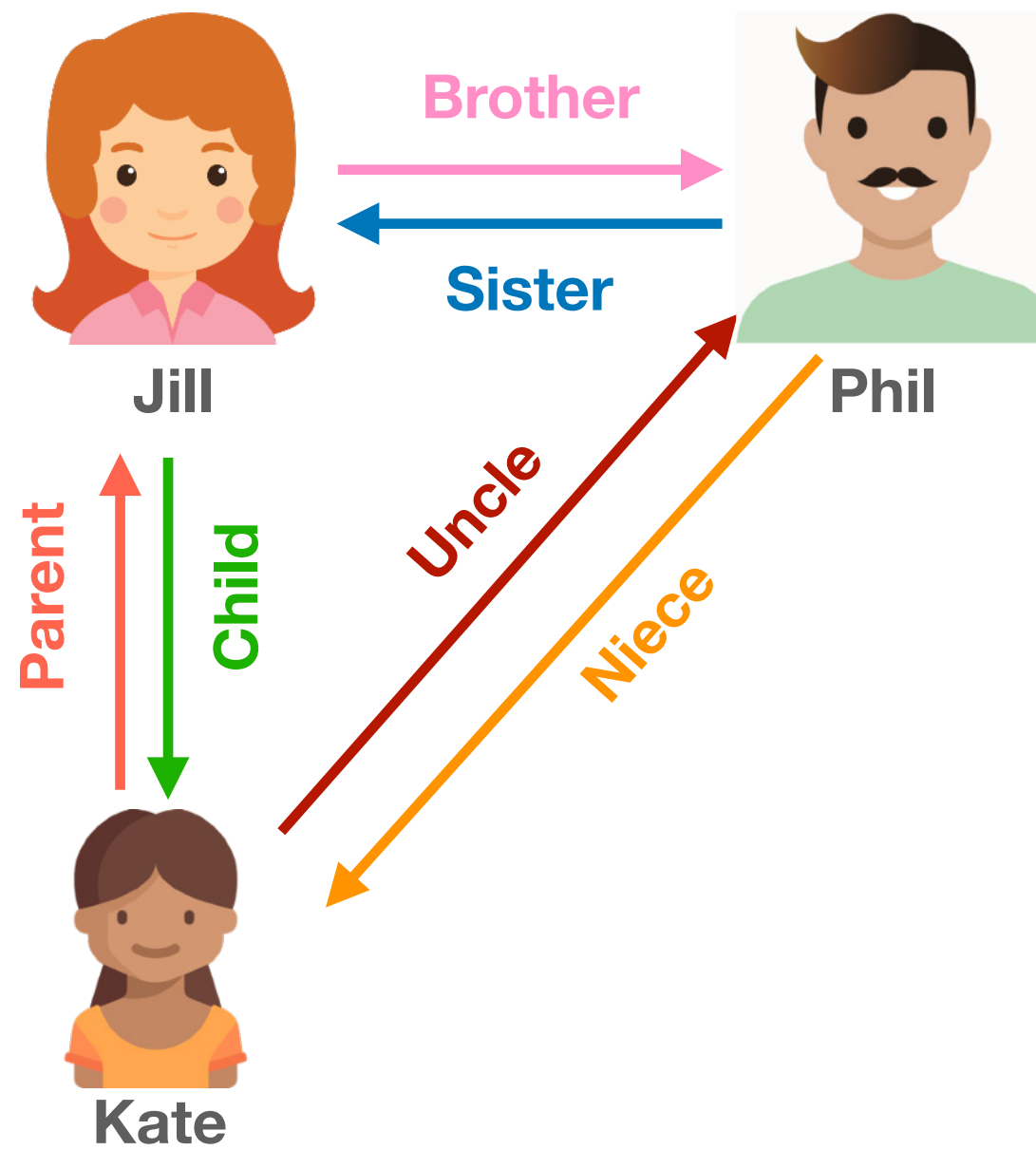
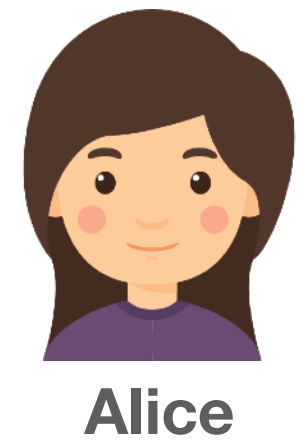
# Space and non-space unified by abstracting relationships



**Parent + Brother = Uncle**

Know which relationships add up to  
the same thing = path integration

# Space and non-space unified by abstracting relationships



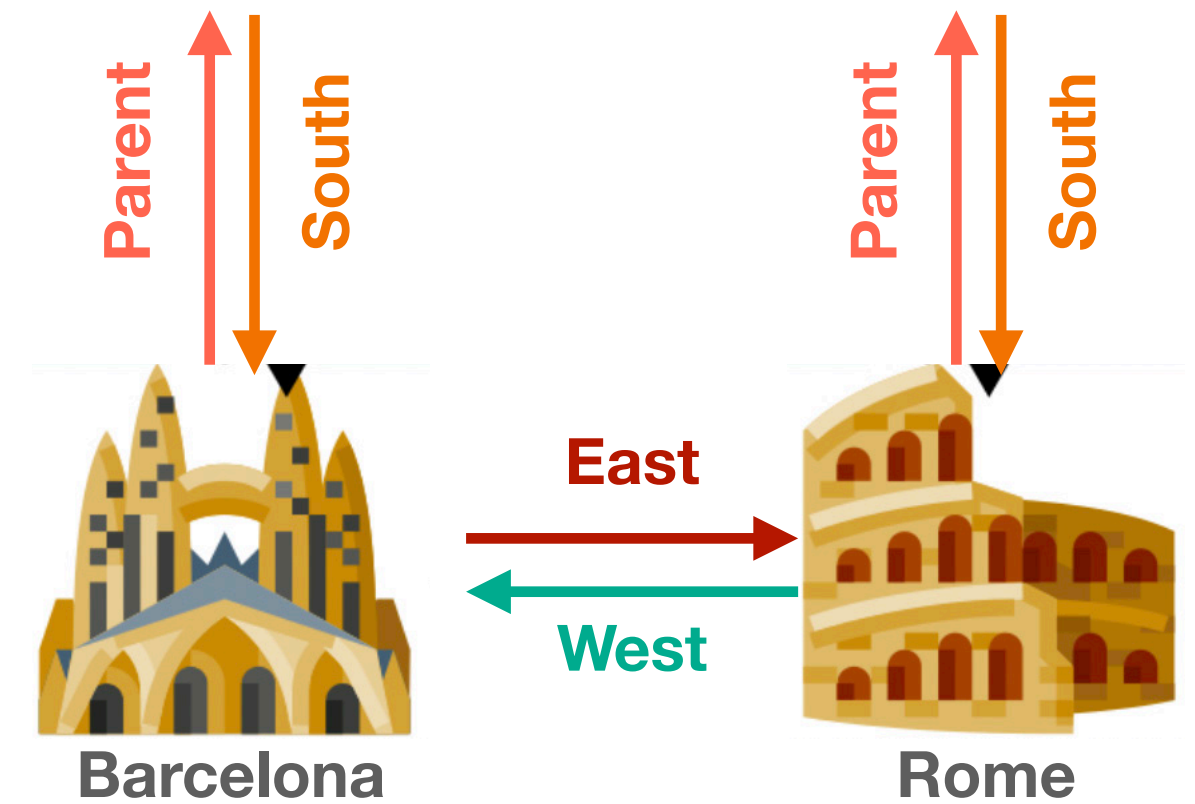
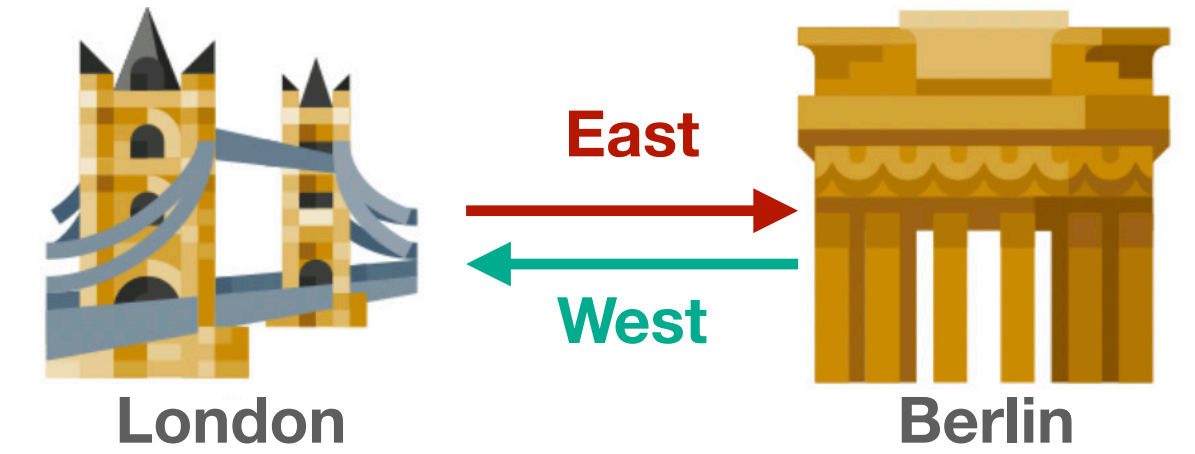
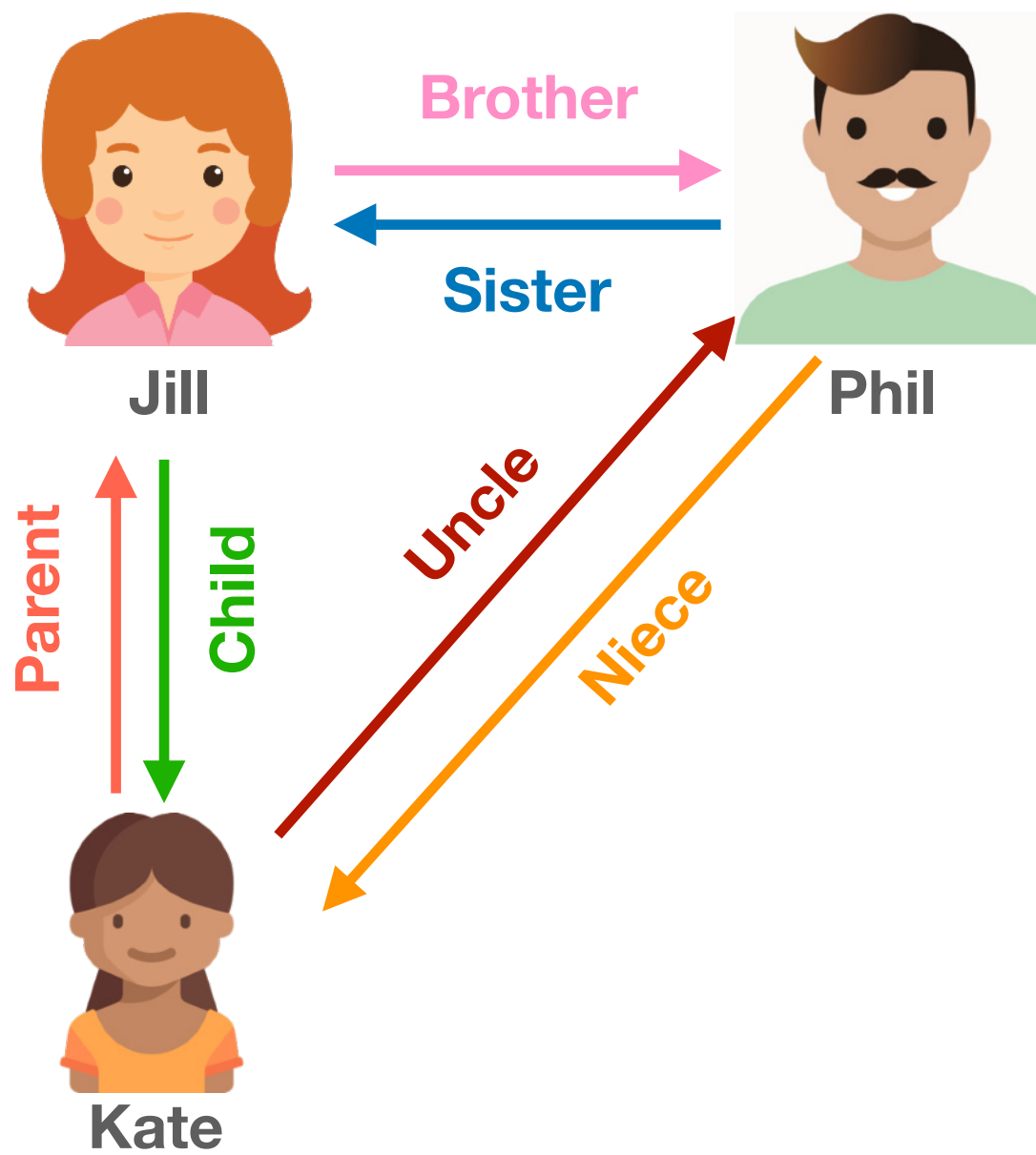
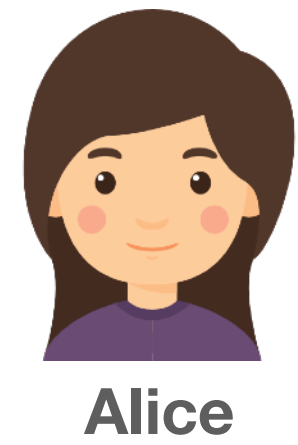
A non-spatial shortcut!

Parent + Brother = Uncle

Know which relationships add up to the same thing = path integration

Generalising abstract relationships to different scenarios

# Space and non-space unified by abstracting relationships



**A non-spatial shortcut!**

**Parent + Brother = Uncle**

**North + East + South + West = 0**

Know which relationships add up to the same thing = path integration

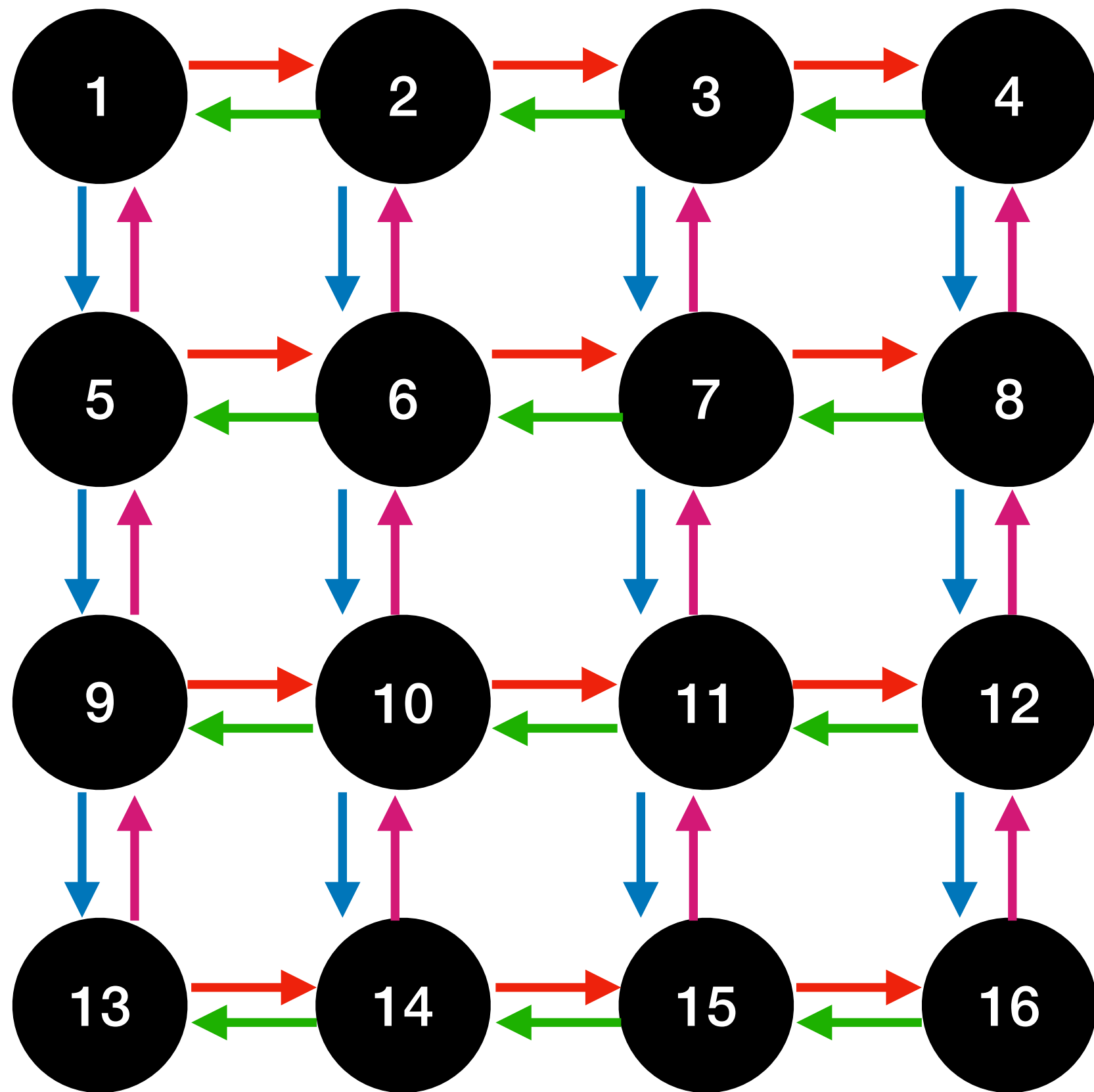
Generalising abstract relationships to different scenarios

This is just like traditional path integration in space



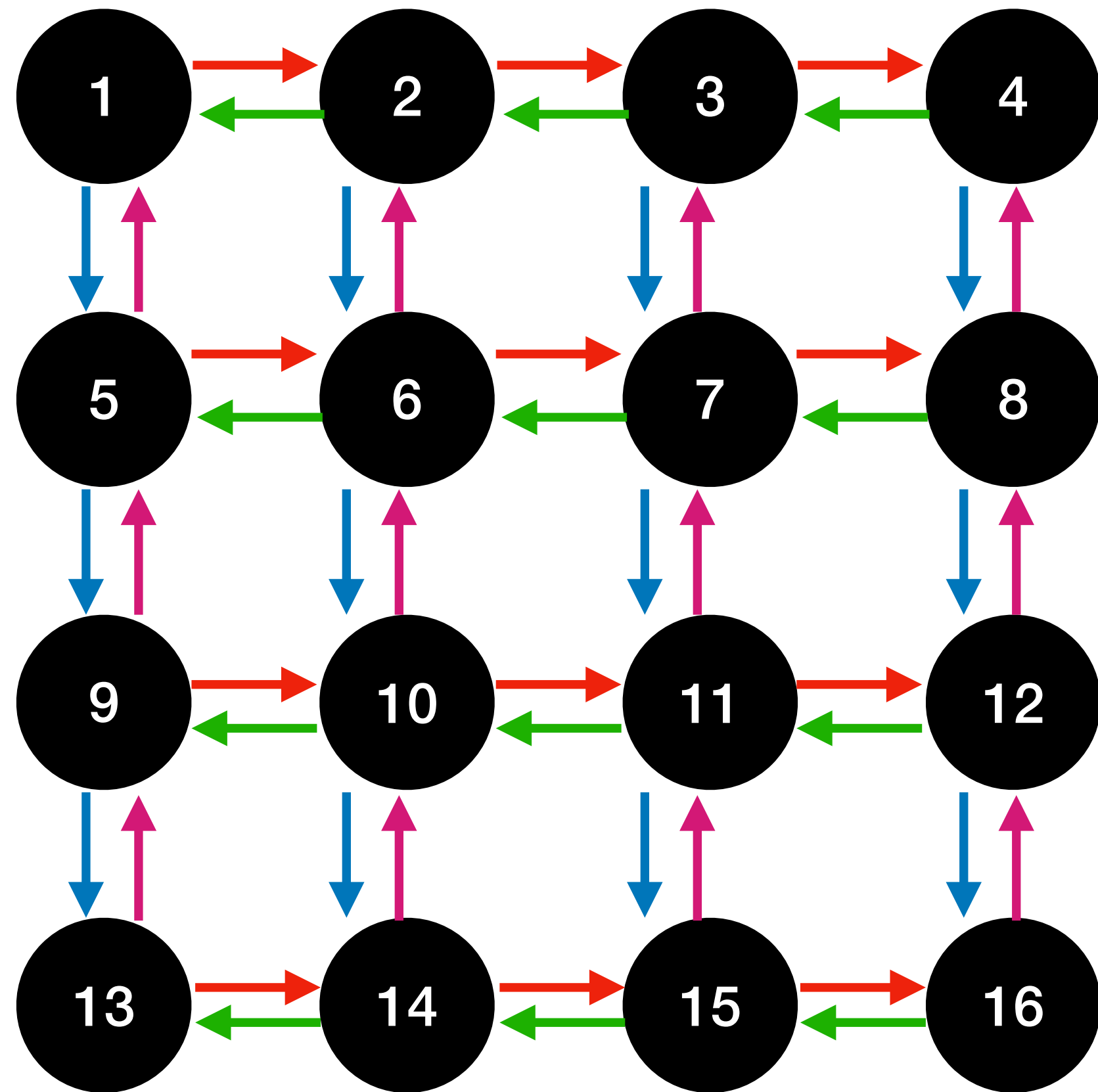
# Formalising relationships using graphs

# Formalising relationships using graphs

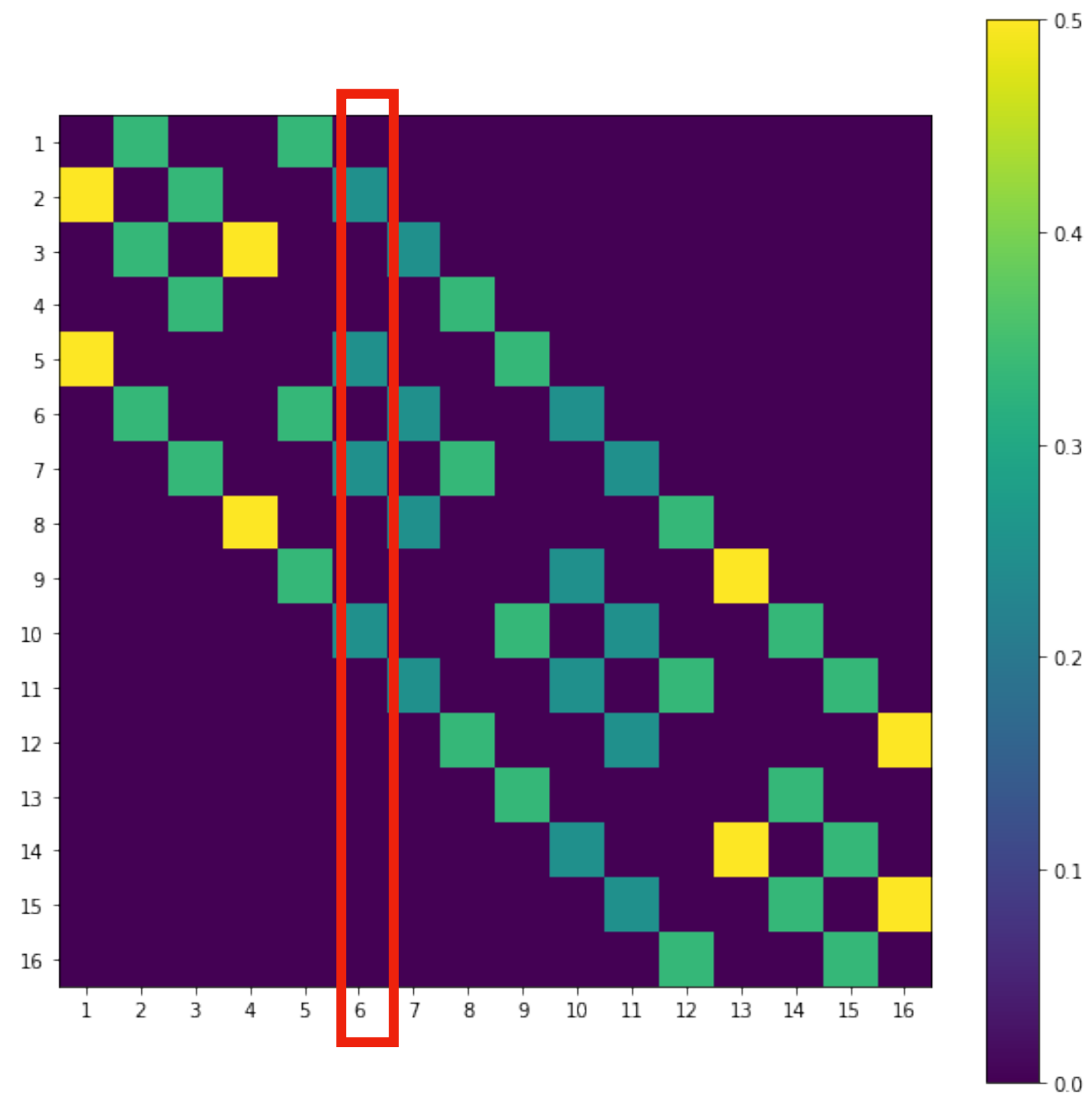




# Formalising relationships using graphs

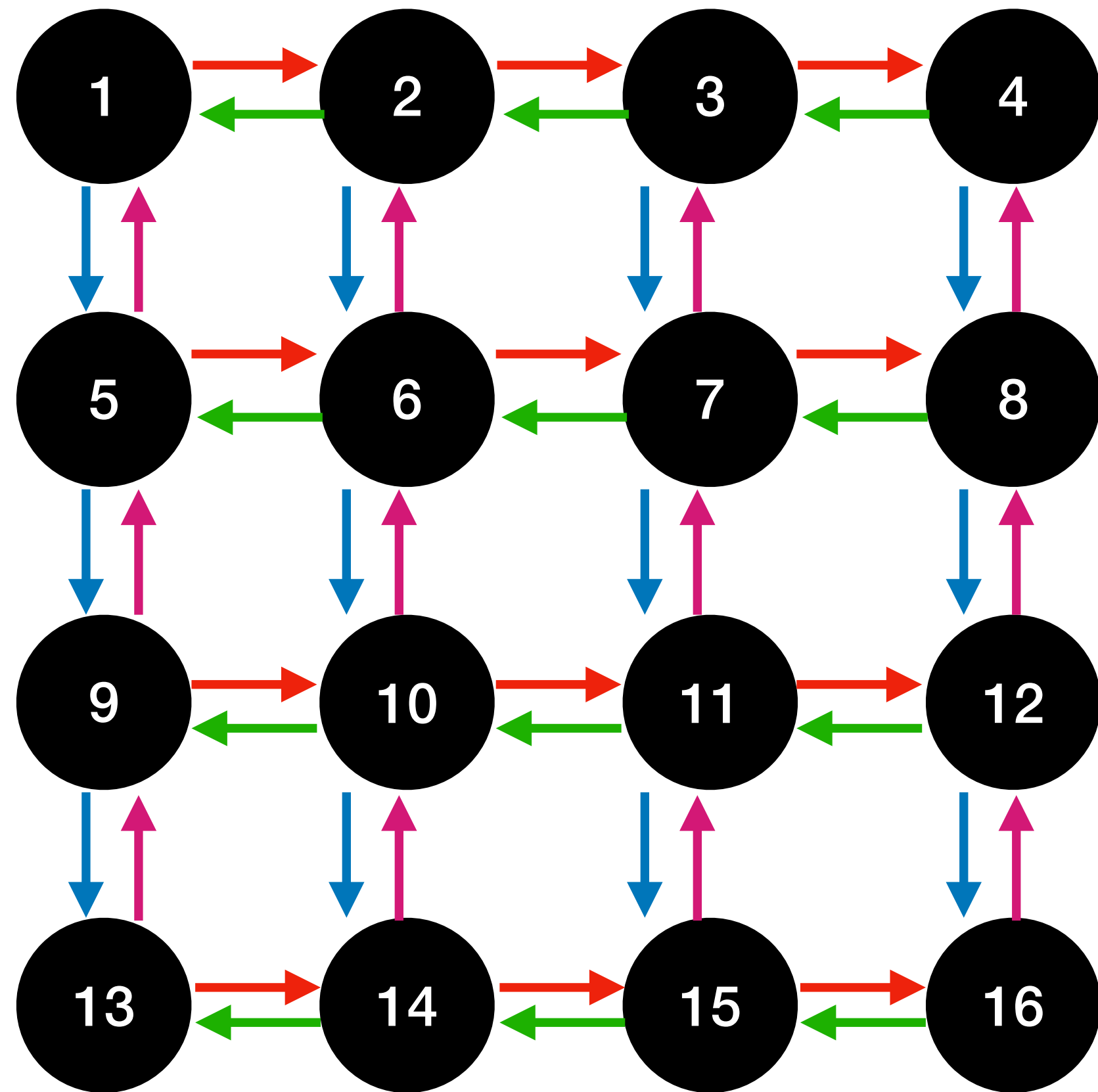


Transition matrix  $T_{ij}$   
= probability of going from state  $j$  to state  $i$

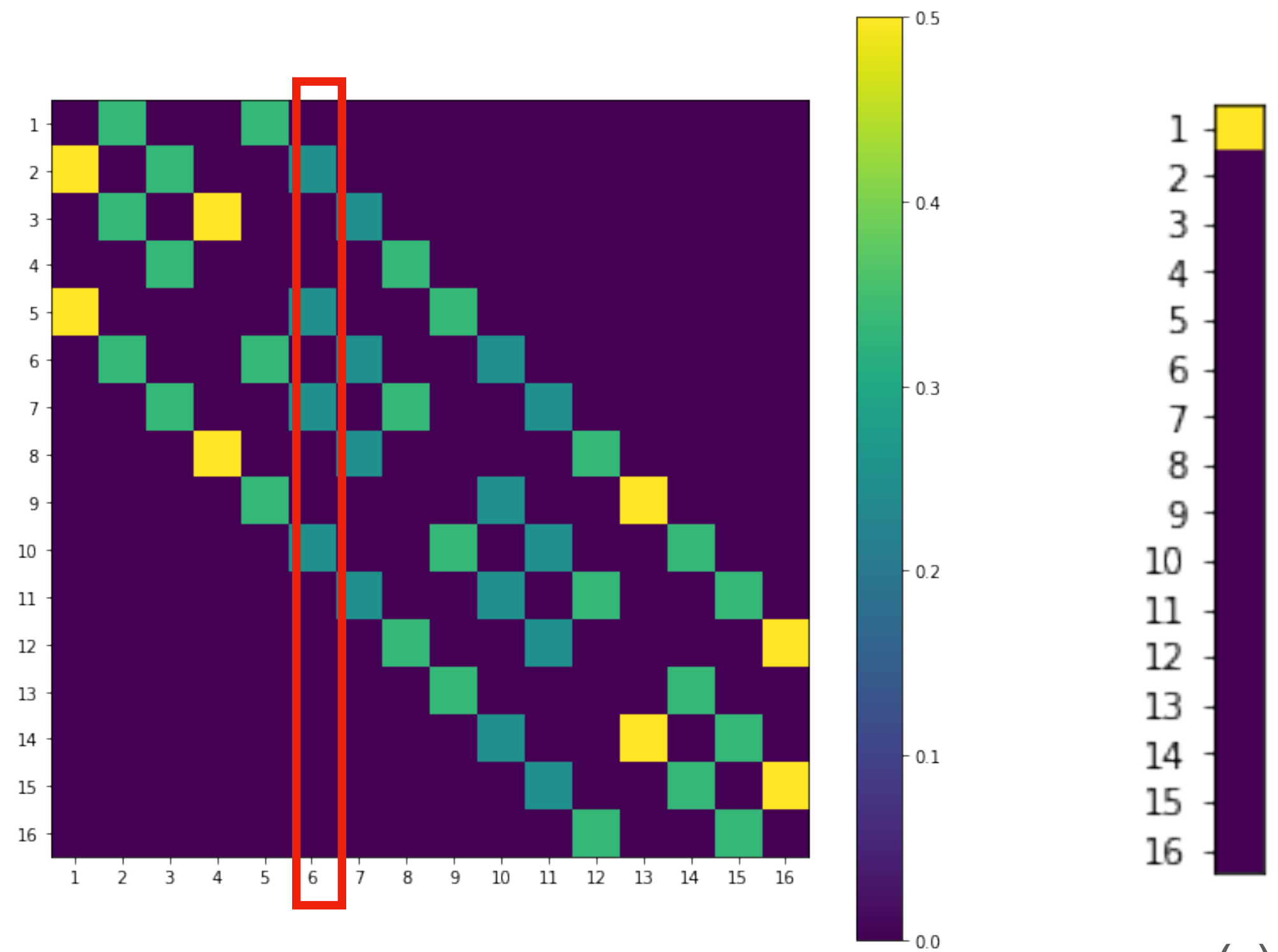


If I am in state 6, where will I be next?

# Formalising relationships using graphs



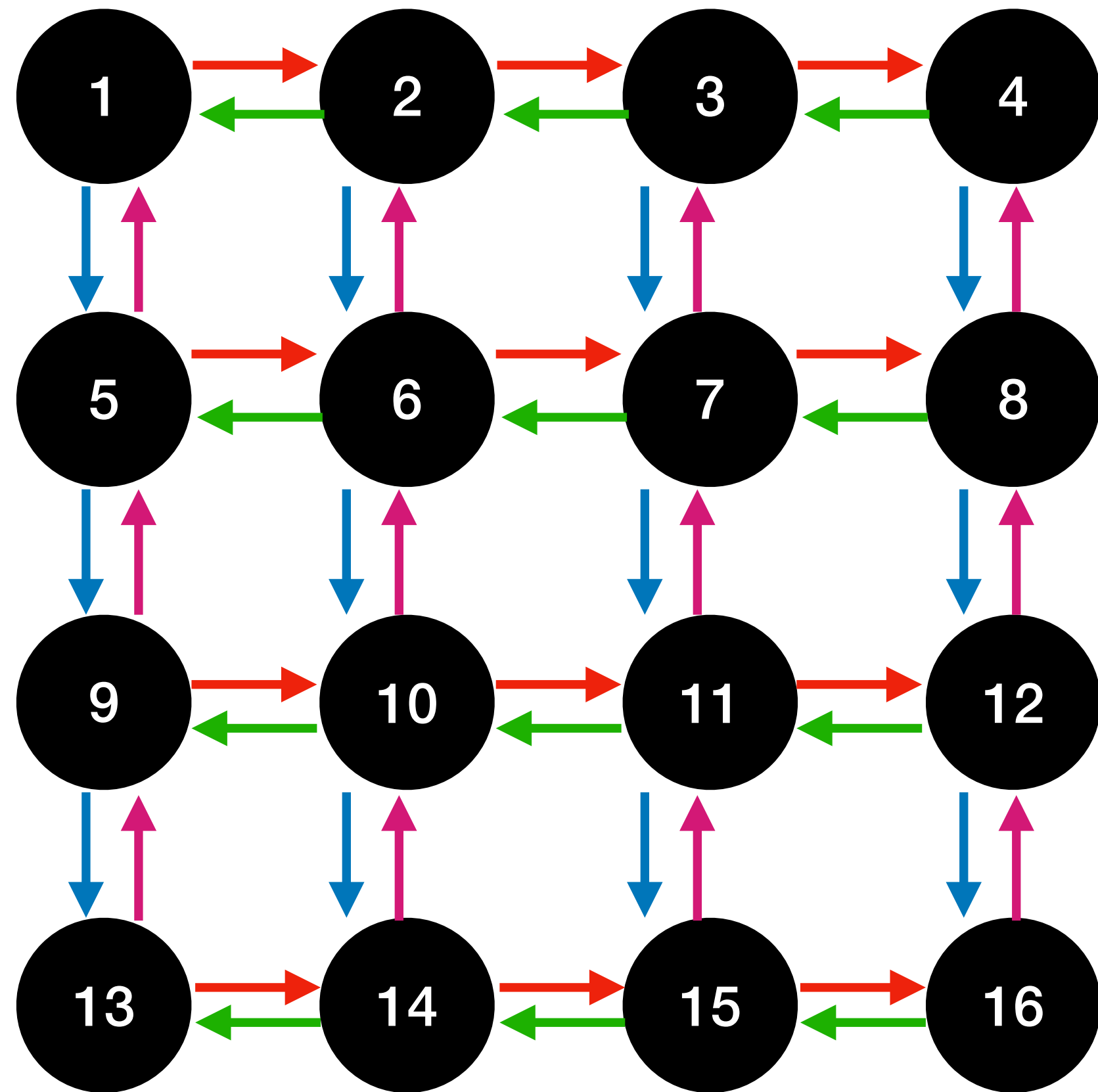
Transition matrix  $T_{ij}$   
 = probability of going from state  $j$  to state  $i$



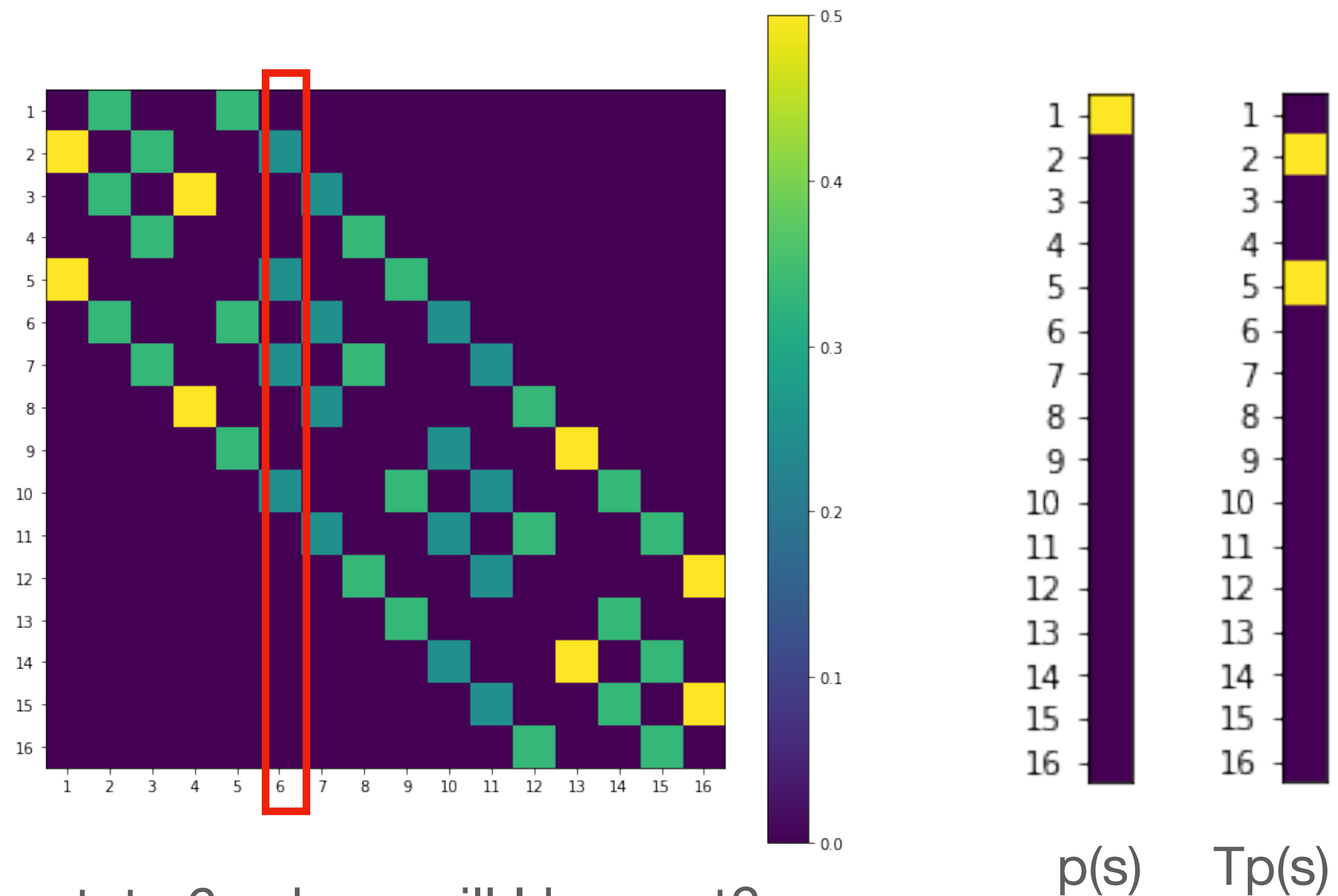
If I am in state 6, where will I be next?

$p(s)$

# Formalising relationships using graphs

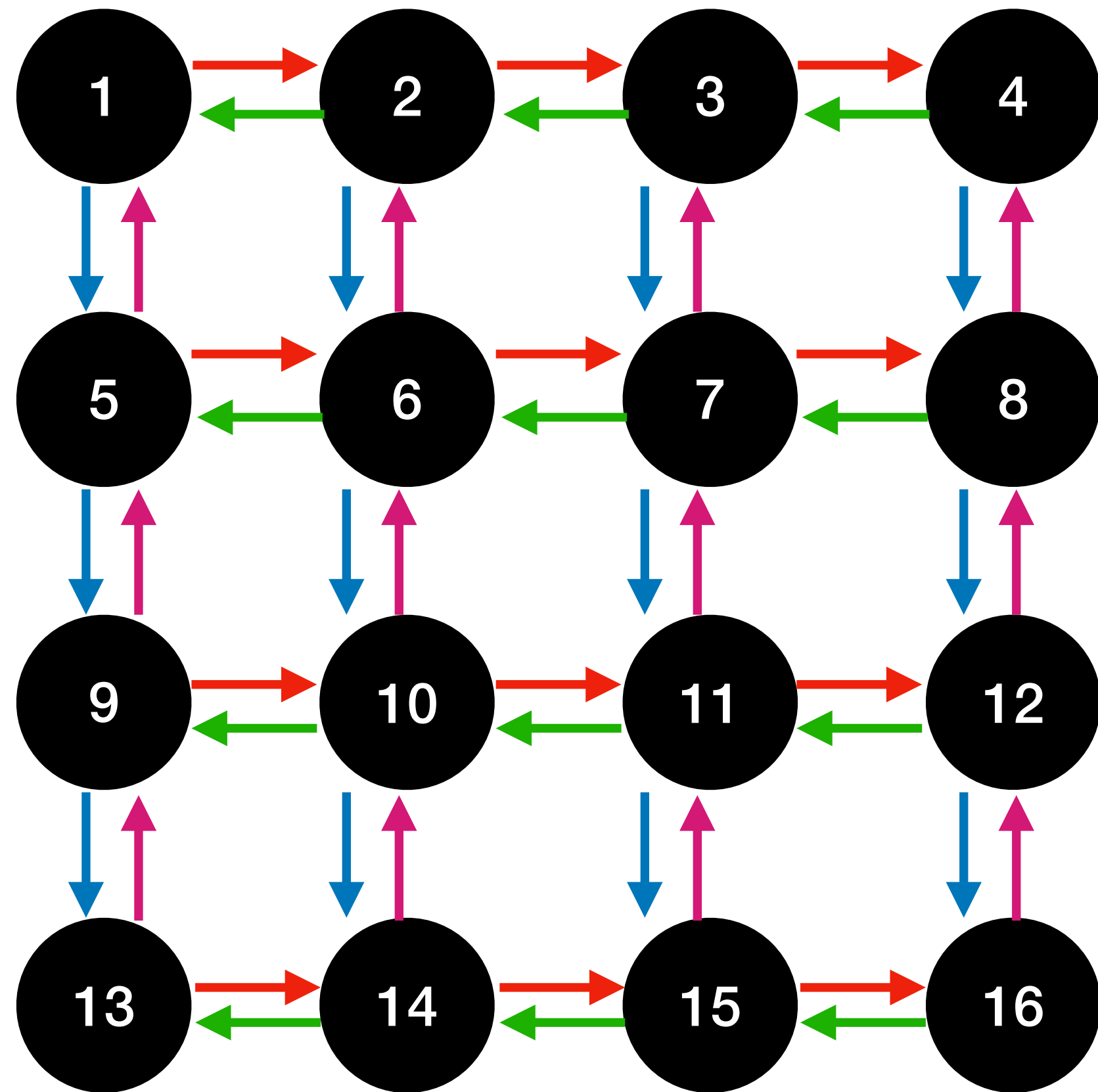


Transition matrix  $T_{ij}$   
 = probability of going from state  $j$  to state  $i$

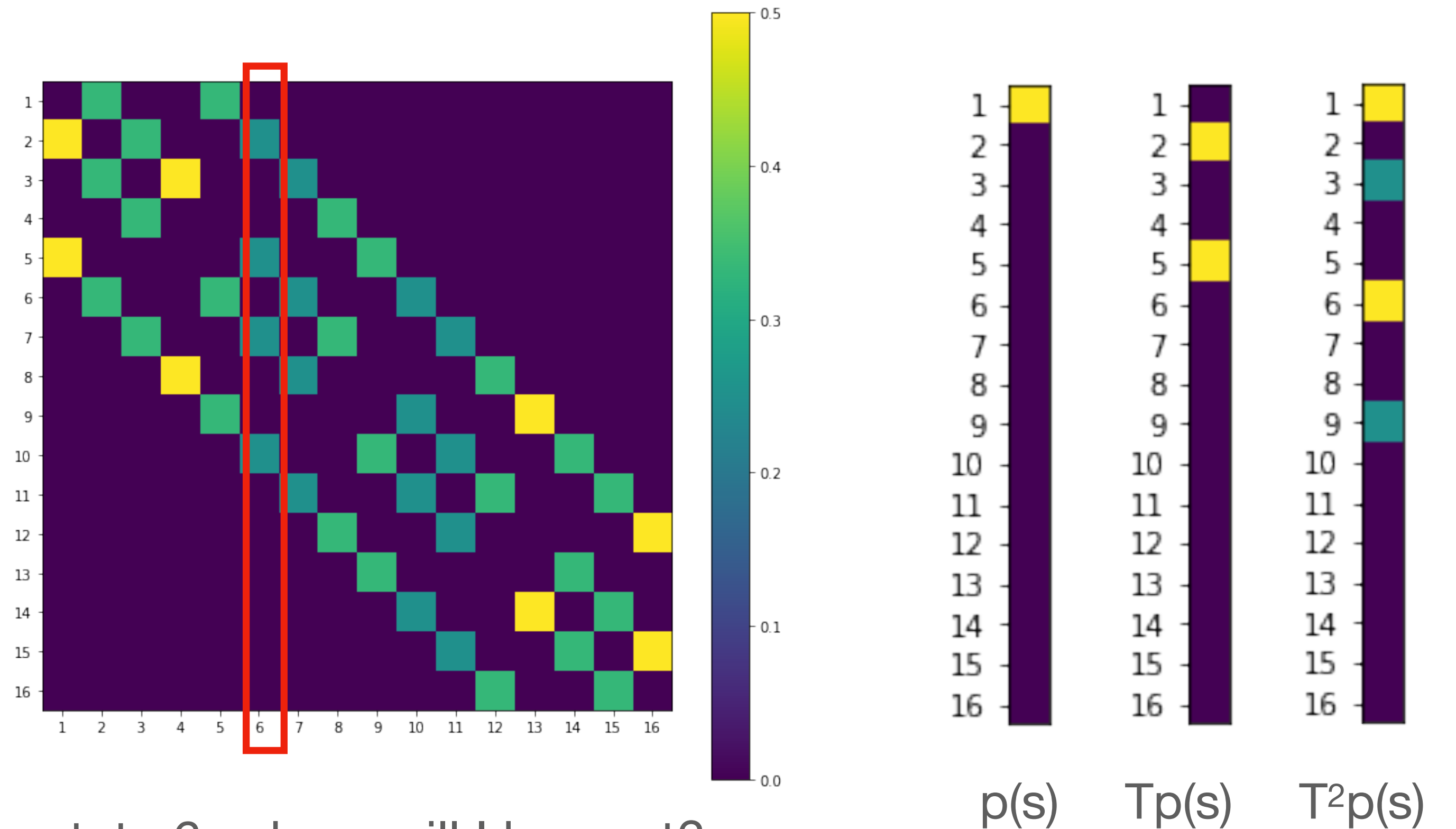


If I am in state 6, where will I be next?

# Formalising relationships using graphs



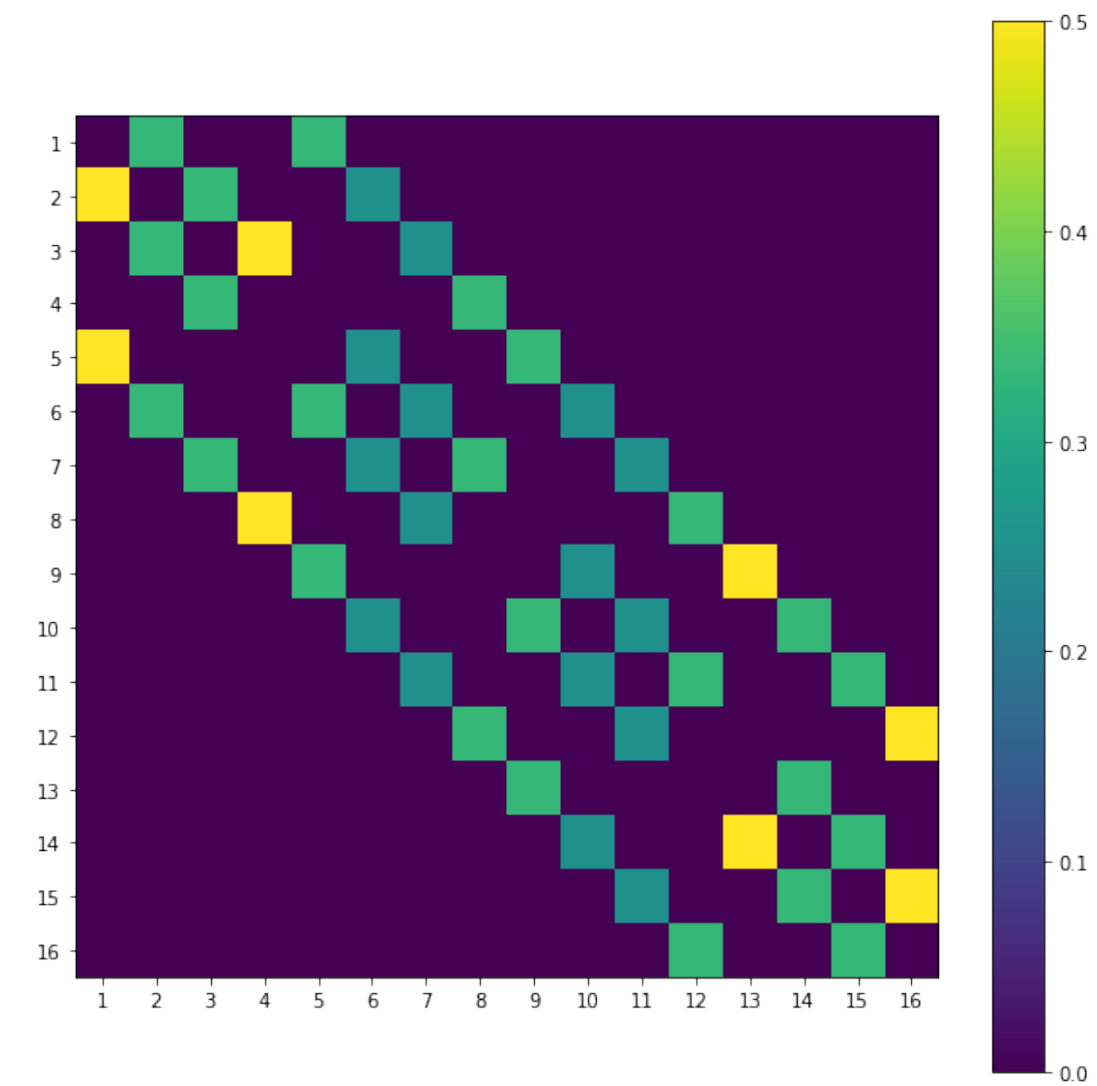
Transition matrix  $T_{ij}$   
 = probability of going from state  $j$  to state  $i$



If I am in state 6, where will I be next?

# The successor representation

The transition matrix says where you'll likely be in 1 step's time



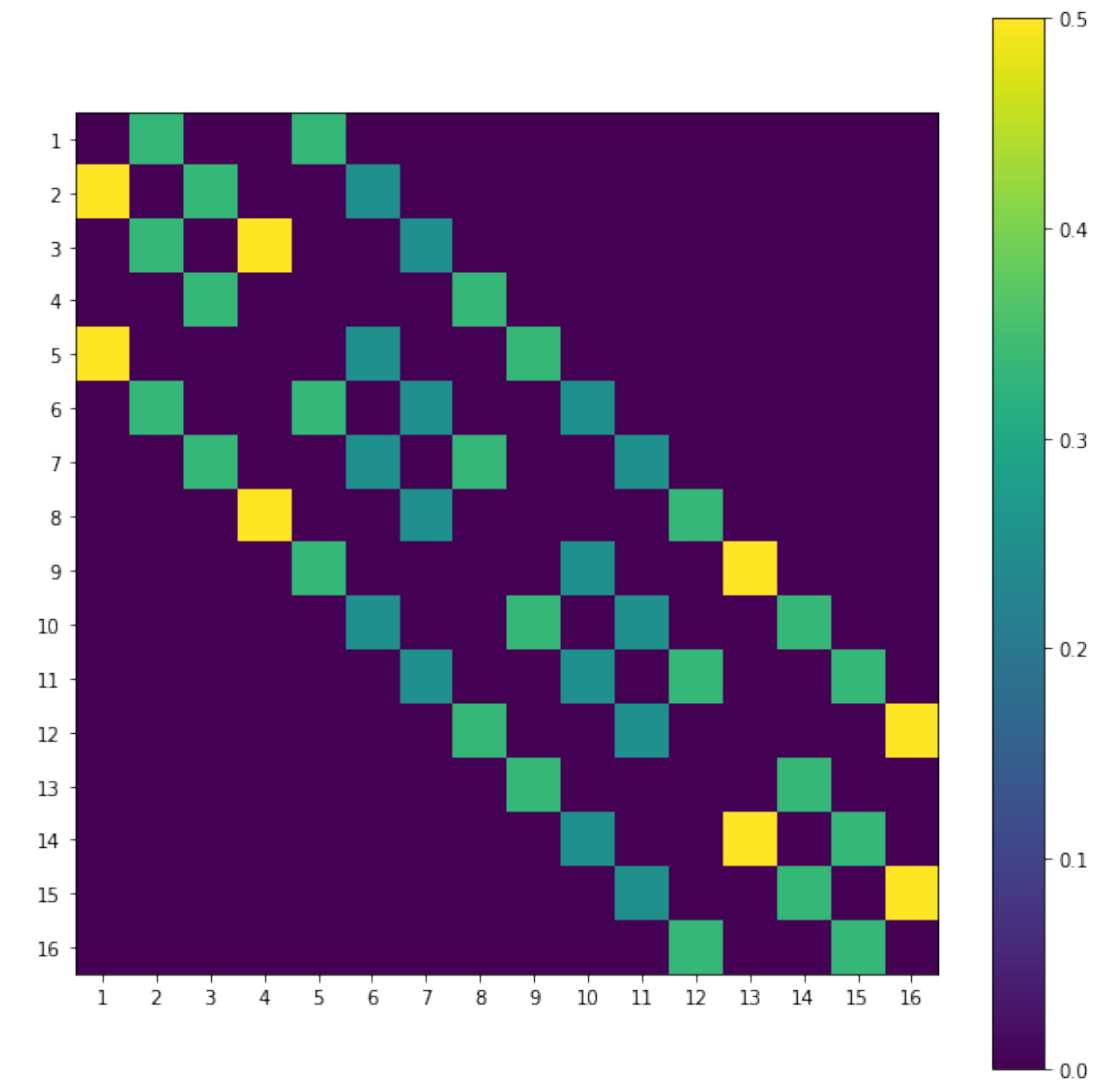


# The successor representation

The transition matrix says where you'll likely be in 1 step's time

The successor representation says where you'll likely be over all future time steps

$$SR = I + T + T^2 + T^3 + \dots$$

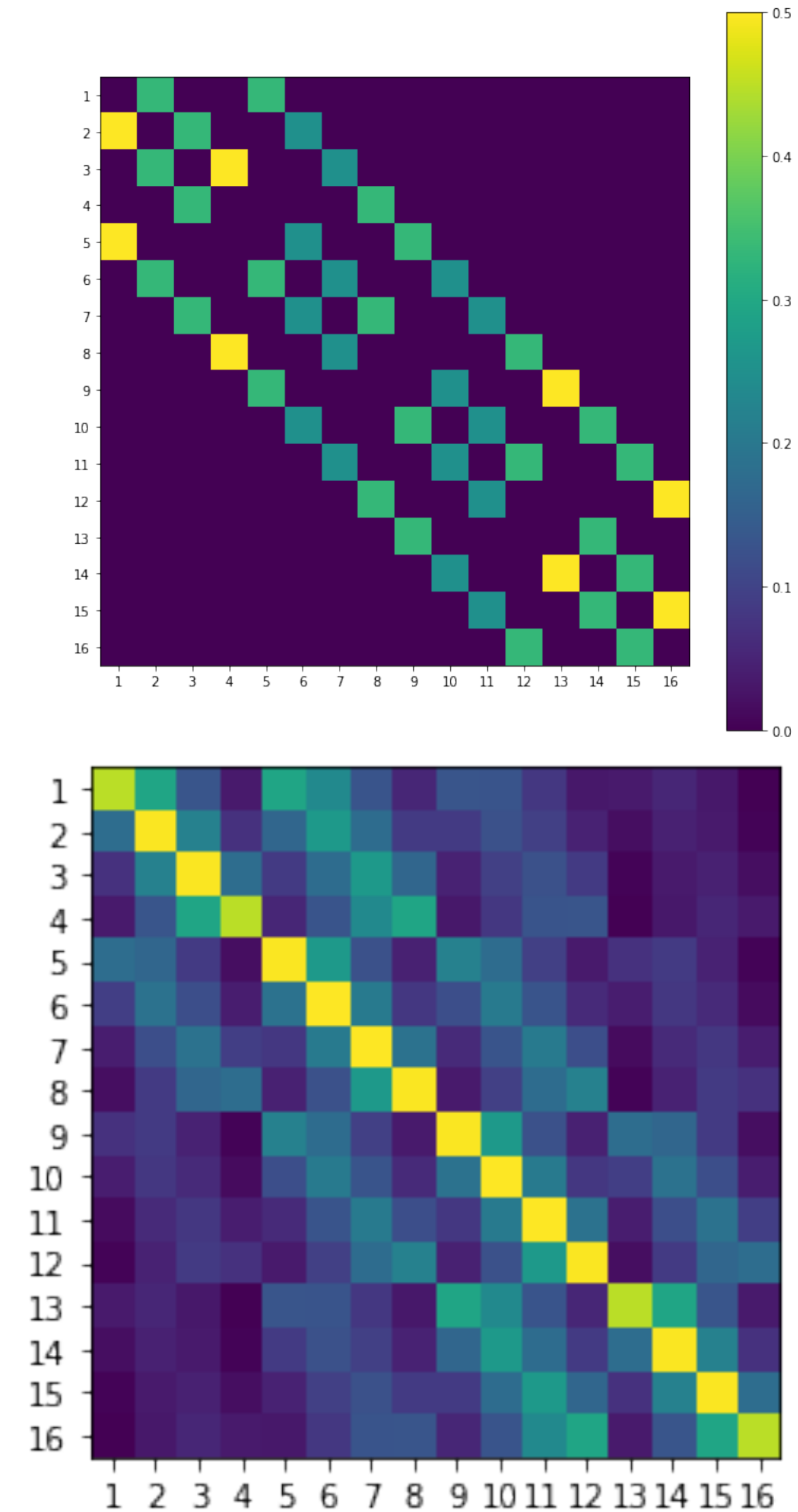


# The successor representation

The transition matrix says where you'll likely be in 1 step's time

The successor representation says where you'll likely be over all future time steps

$$SR = I + T + T^2 + T^3 + \dots$$



# The successor representation

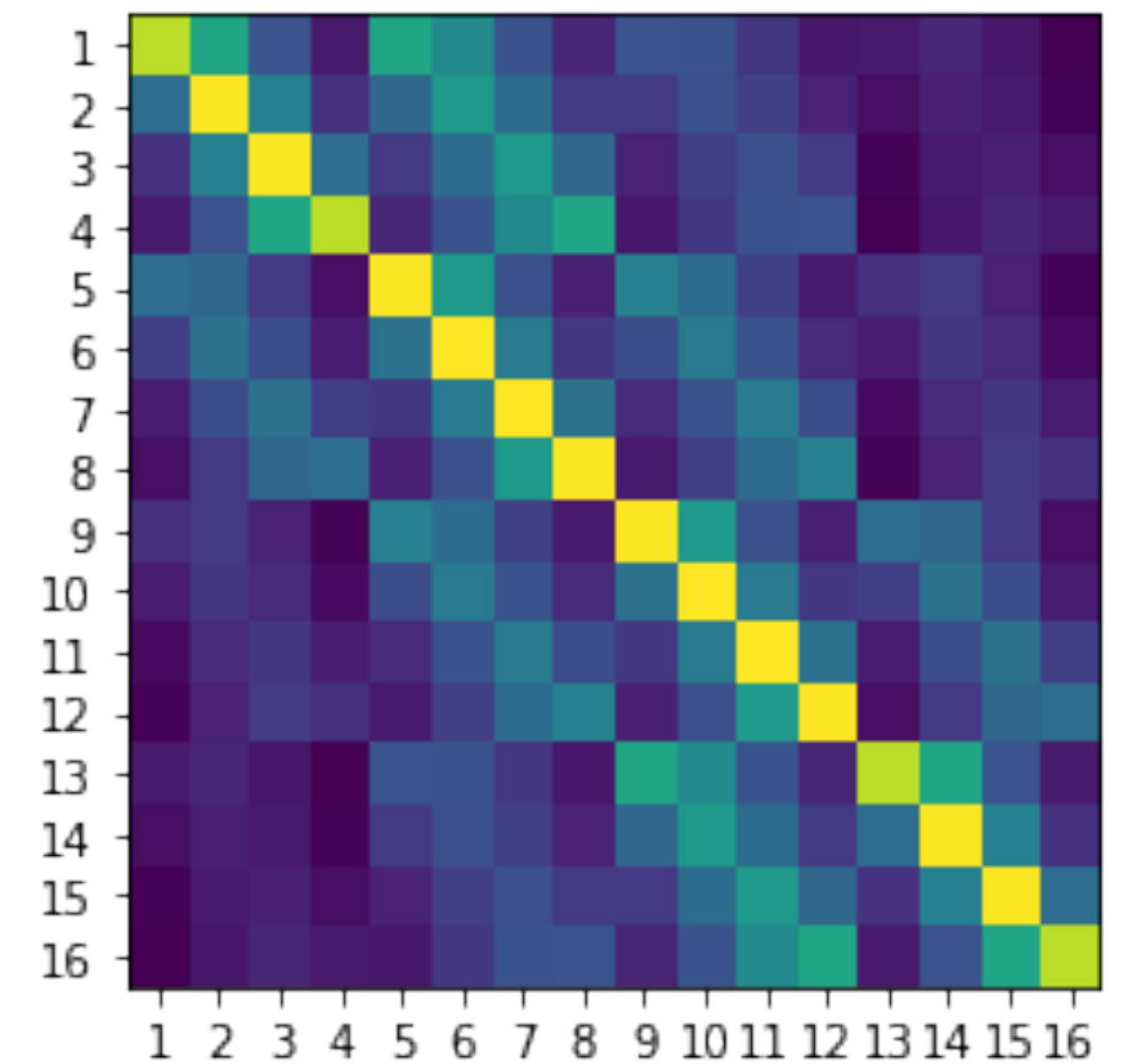
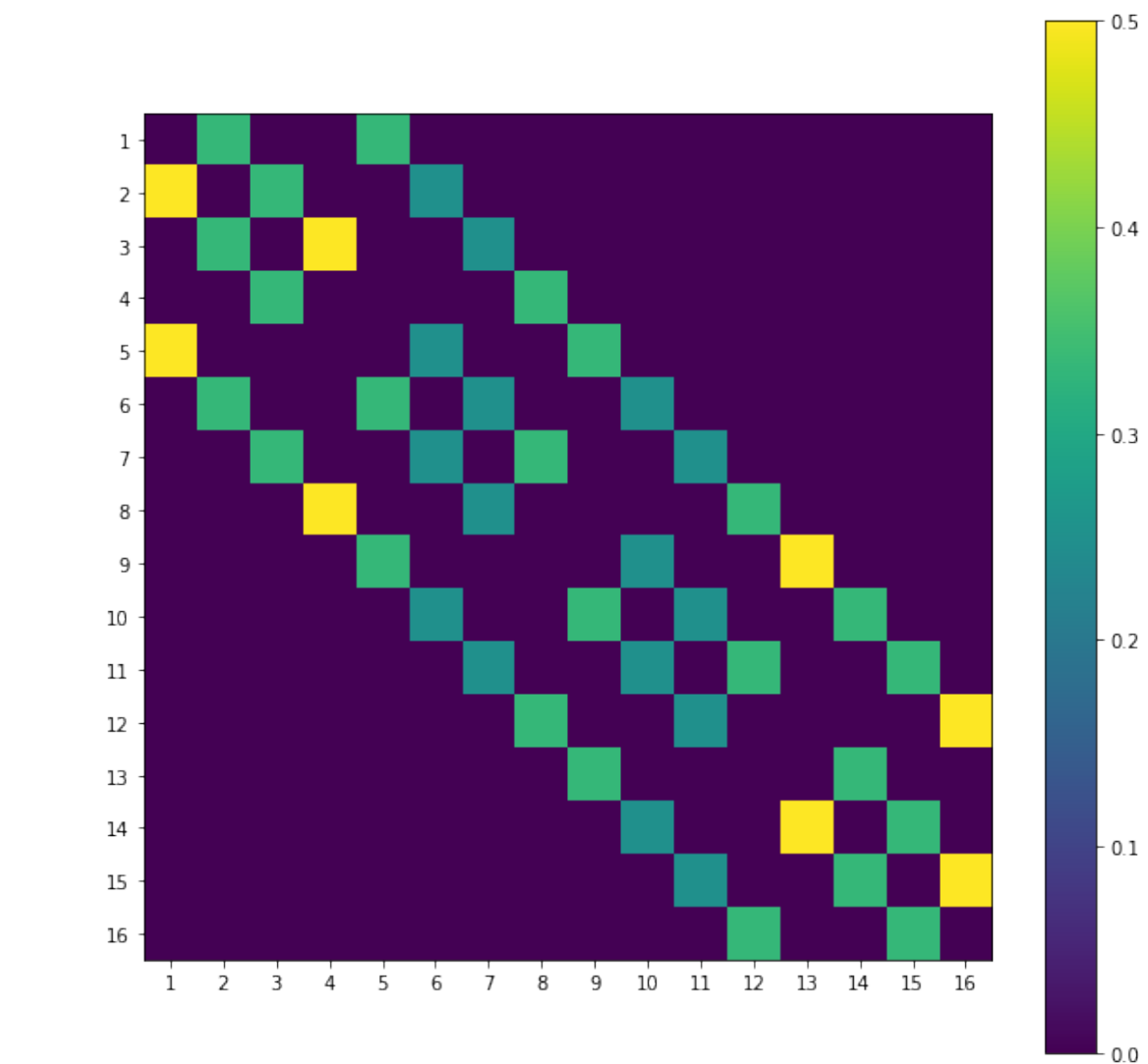
The transition matrix says where you'll likely be in 1 step's time

The successor representation says where you'll likely be over all future time steps

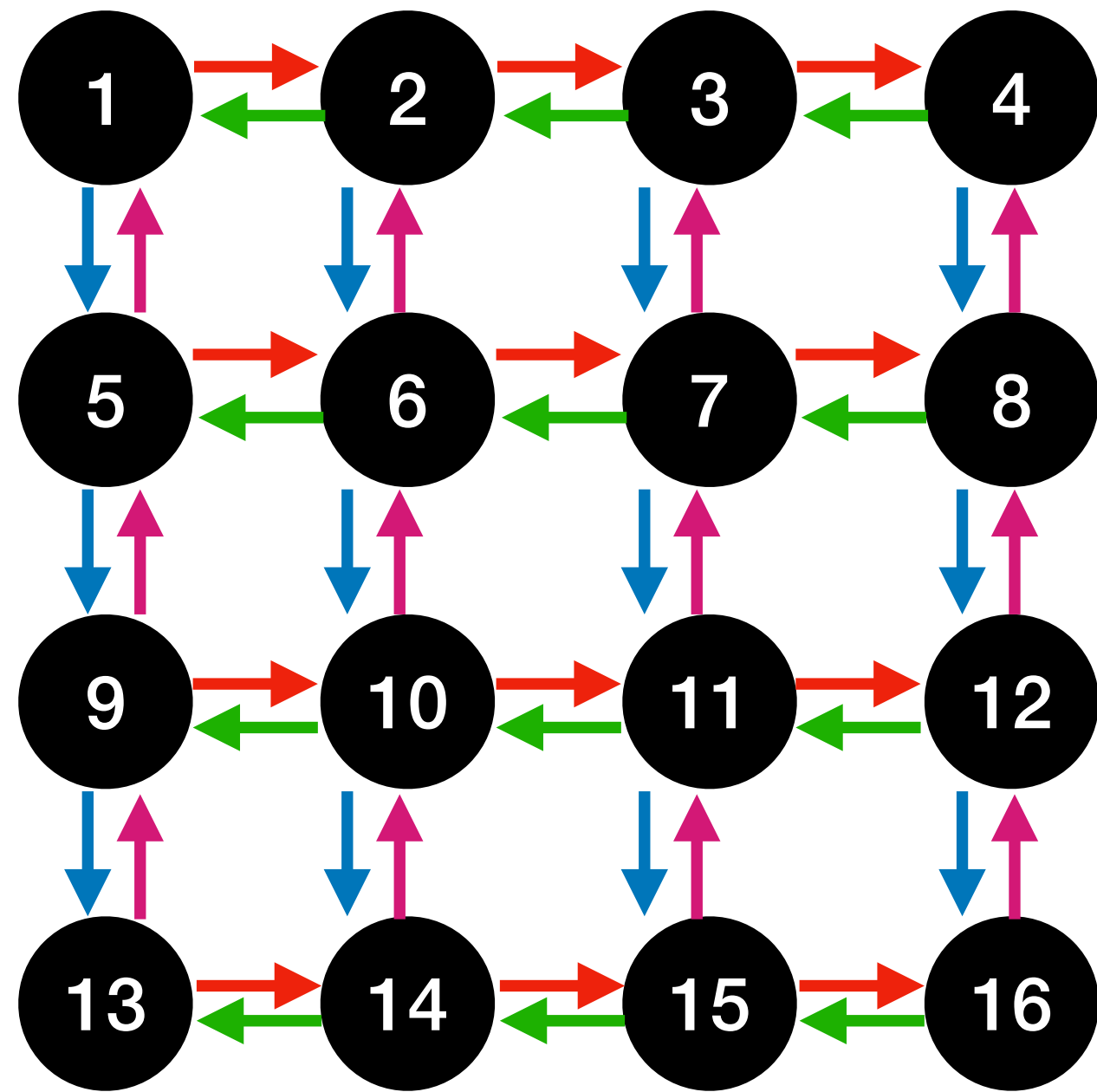
$$SR = I + T + T^2 + T^3 + \dots$$

Often discounted so the sum converges

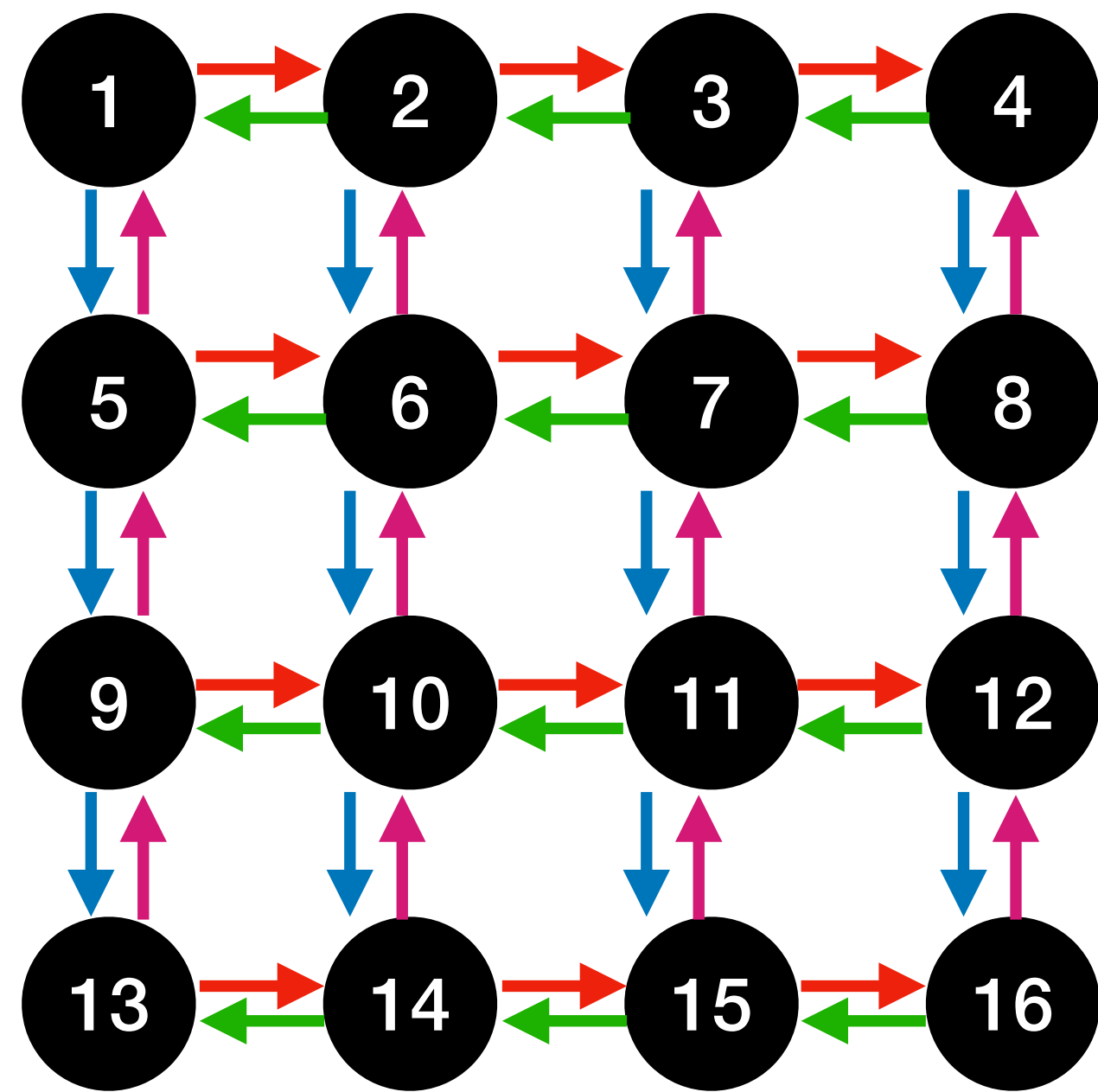
$$SR = I + \gamma T + \gamma^2 T^2 + \gamma^3 T^3 + \dots$$



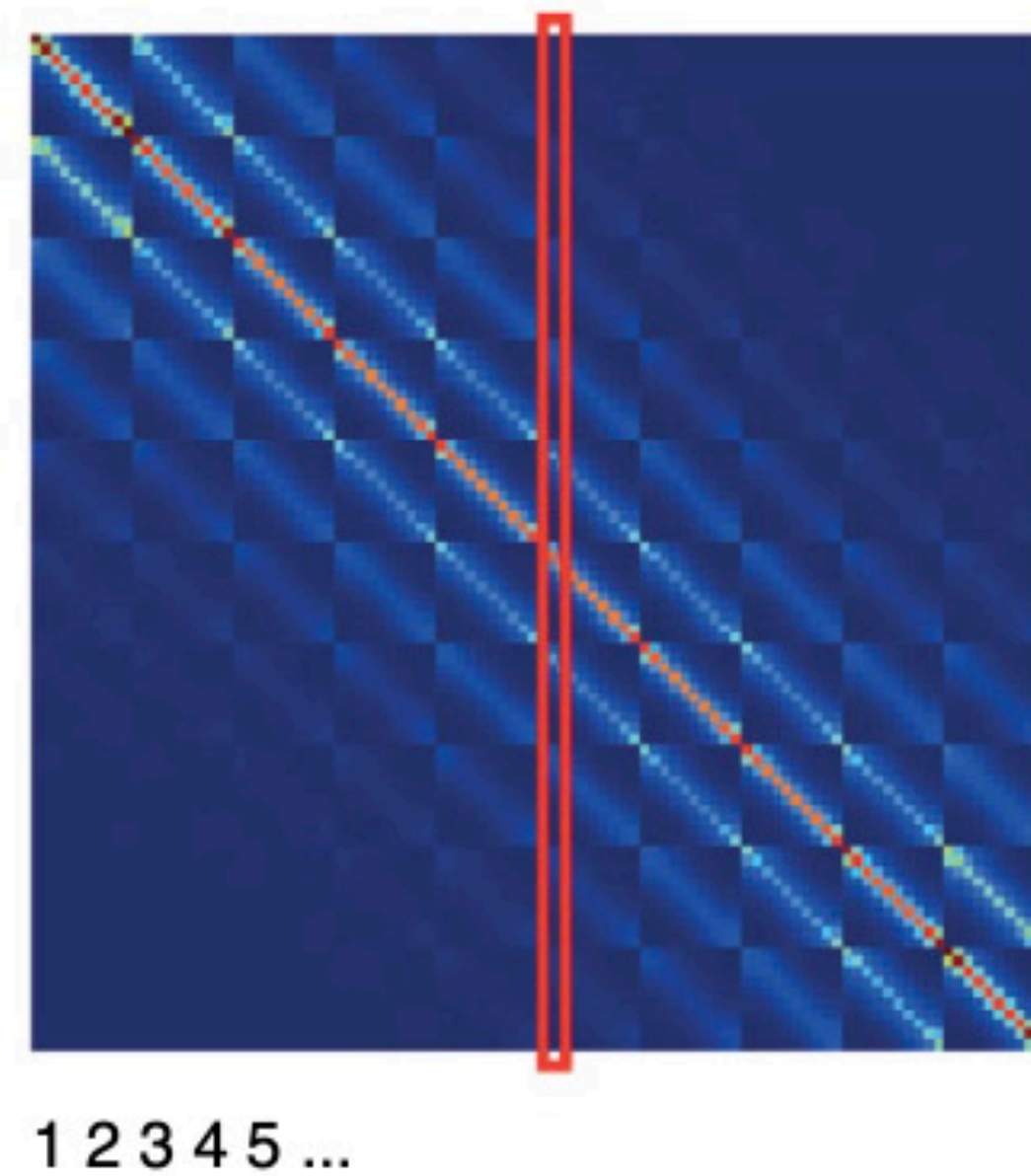
**It looks a bit like place cells...**



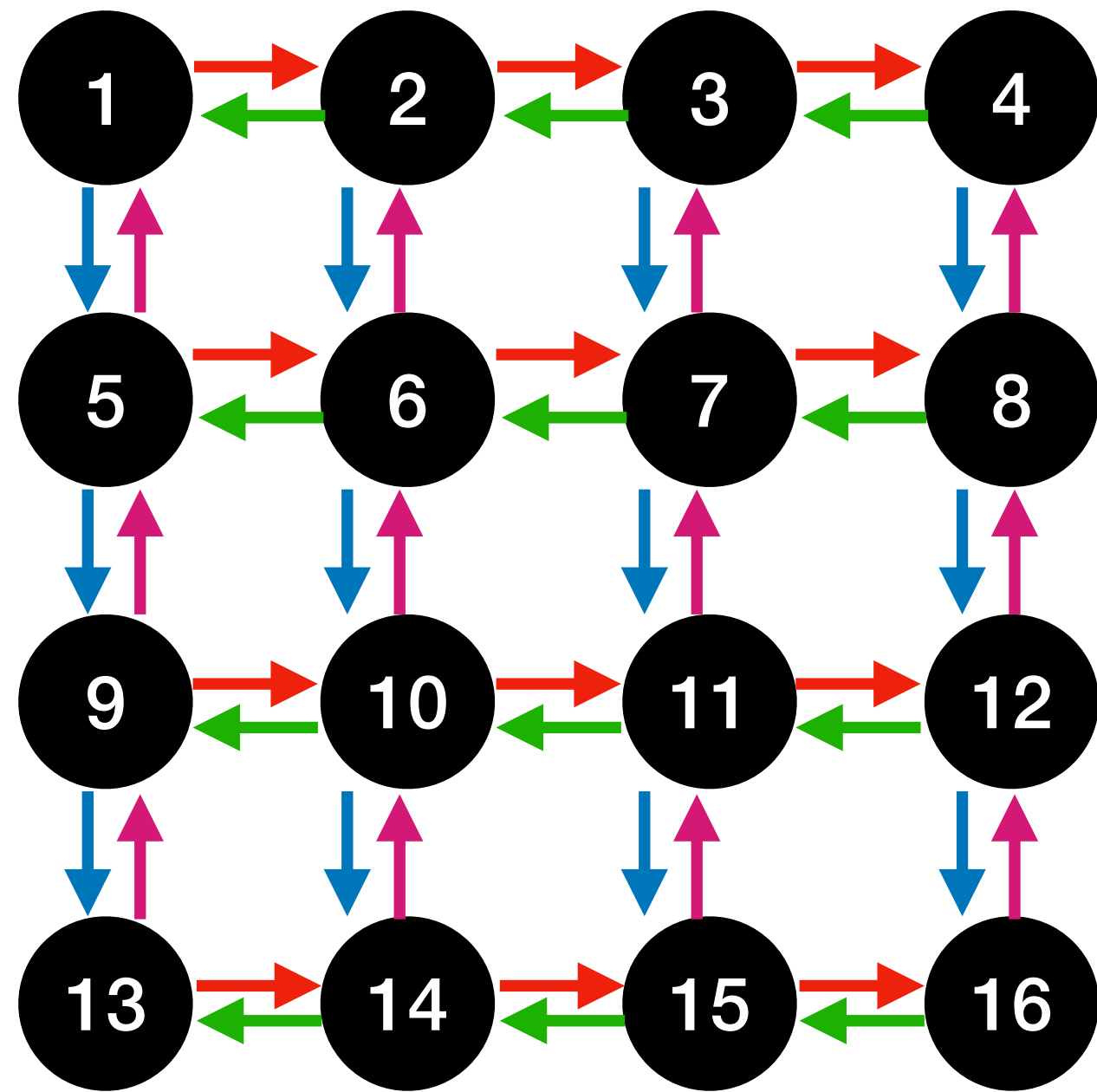
# It looks a bit like place cells...



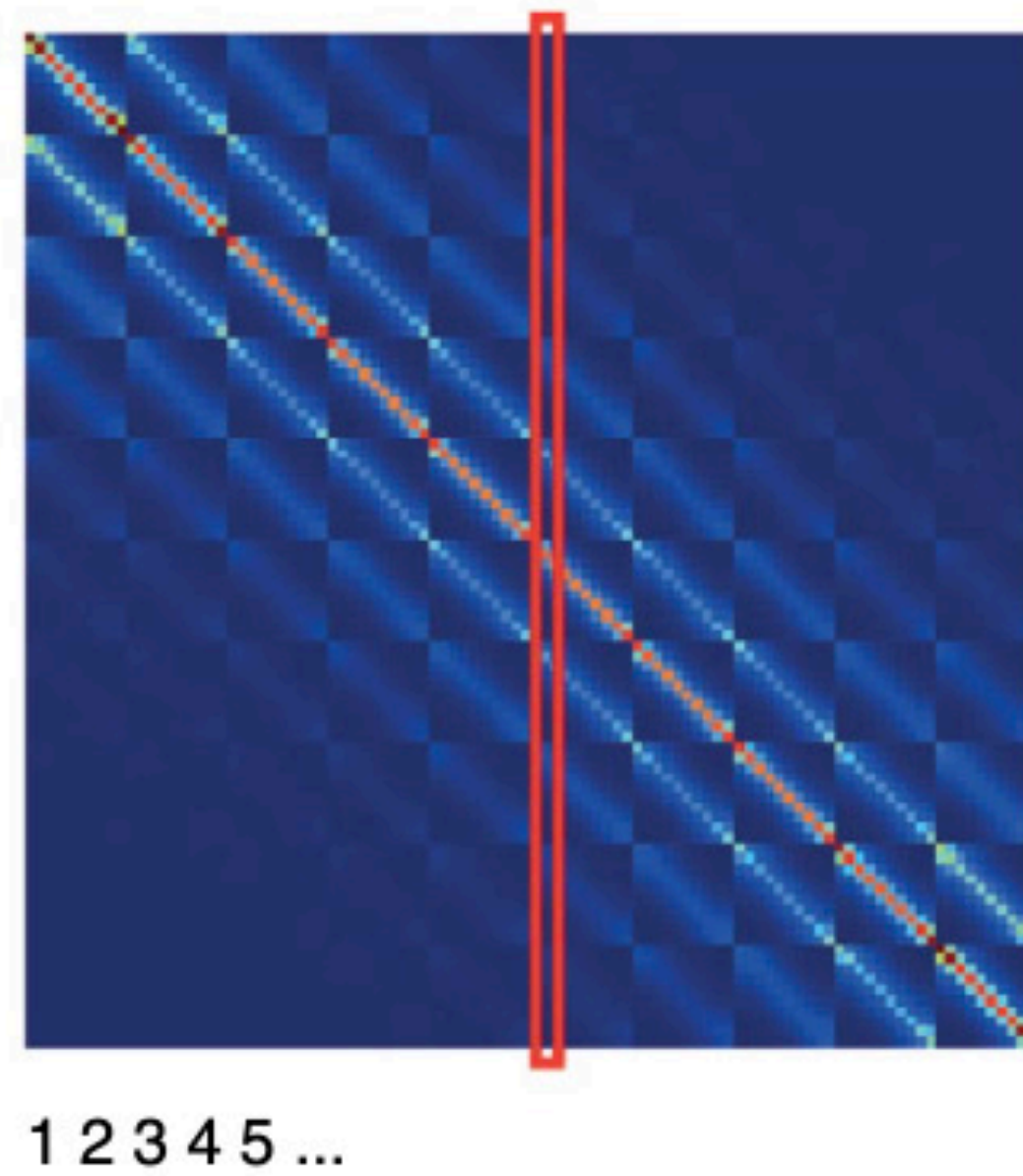
The successor representation



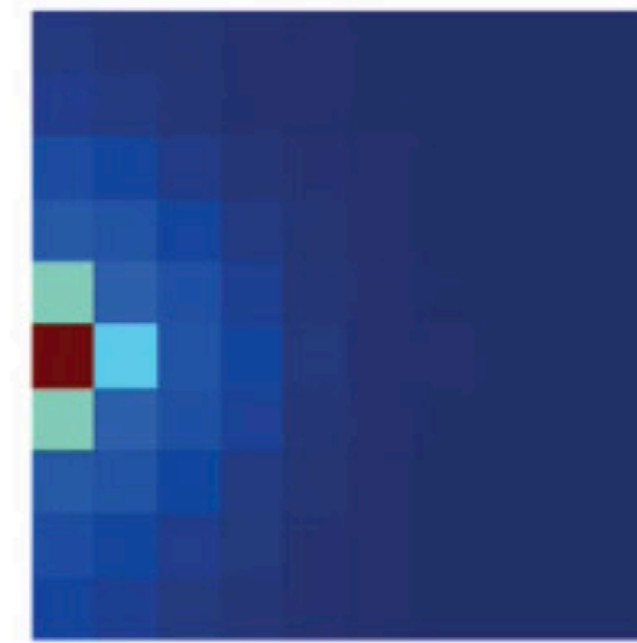
# It looks a bit like place cells...



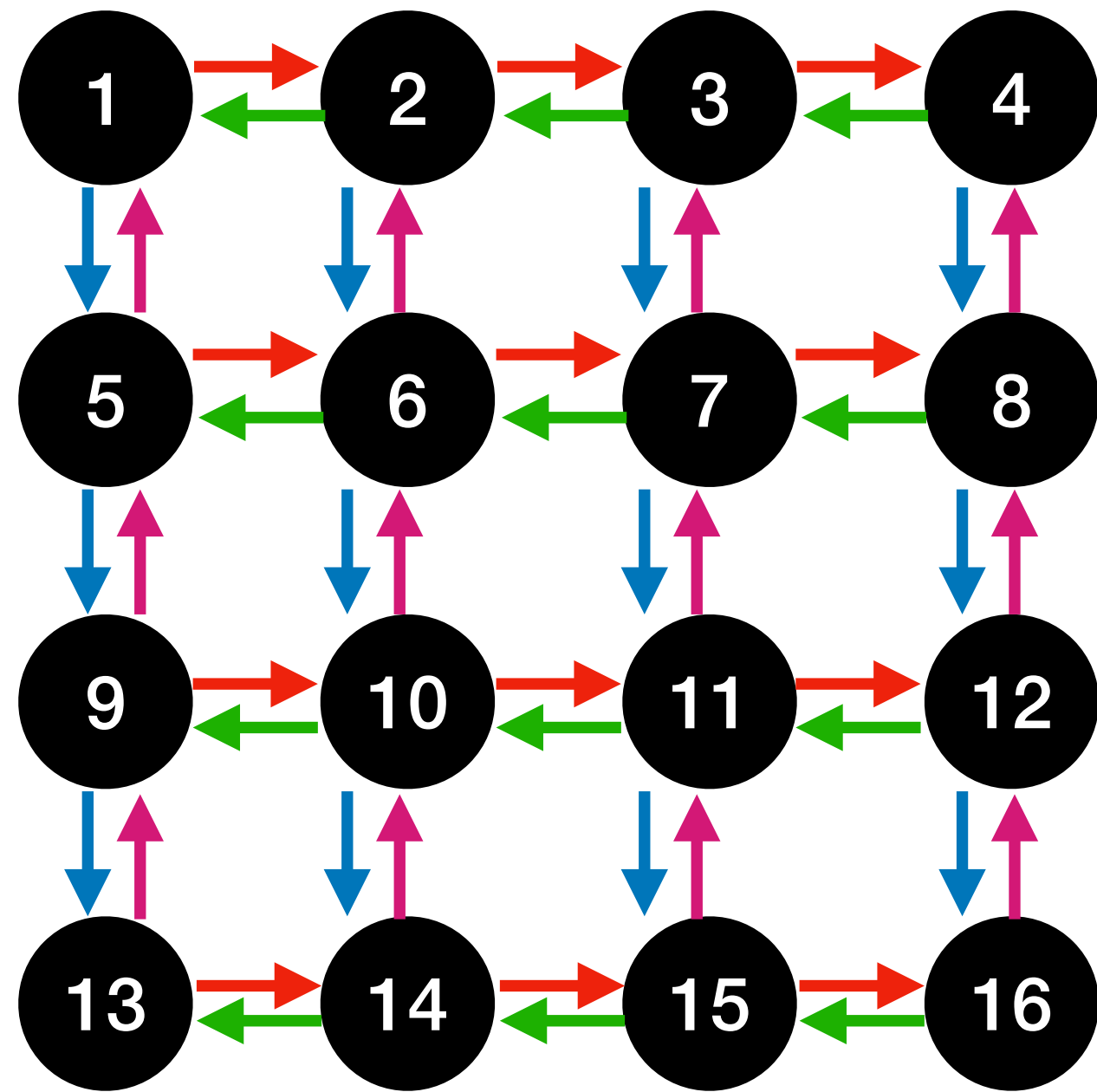
The successor representation



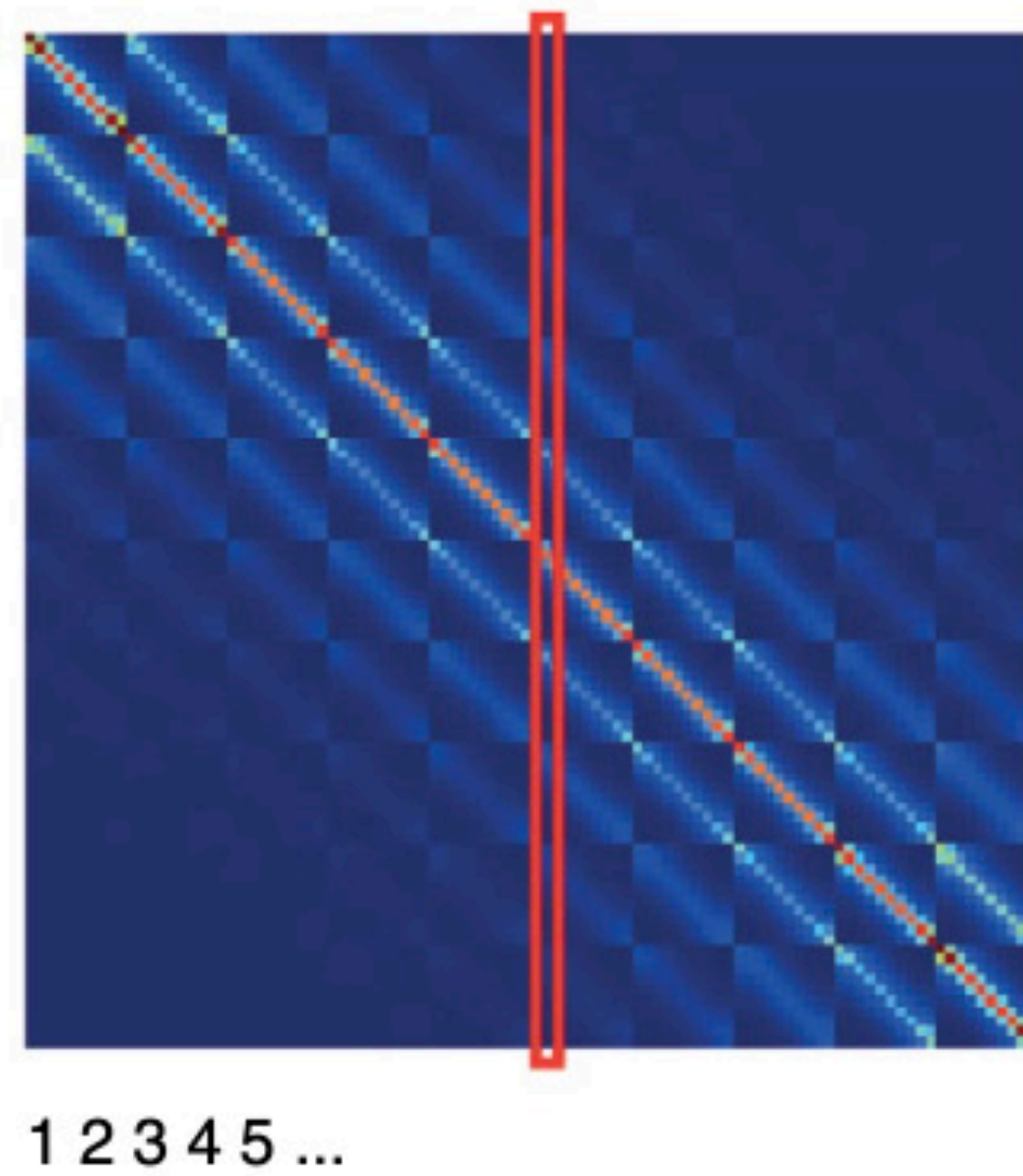
The column looks like a place cell



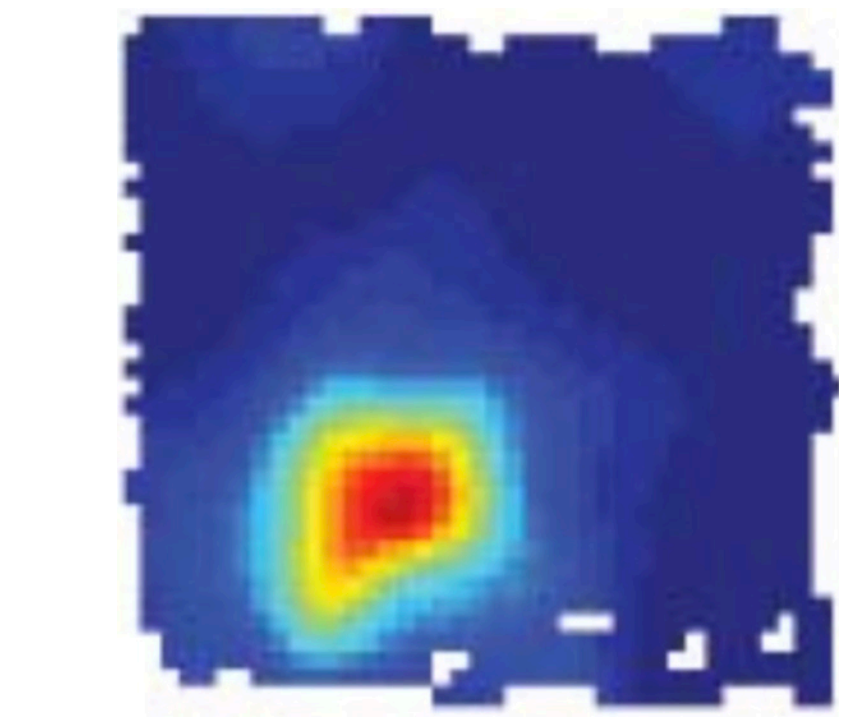
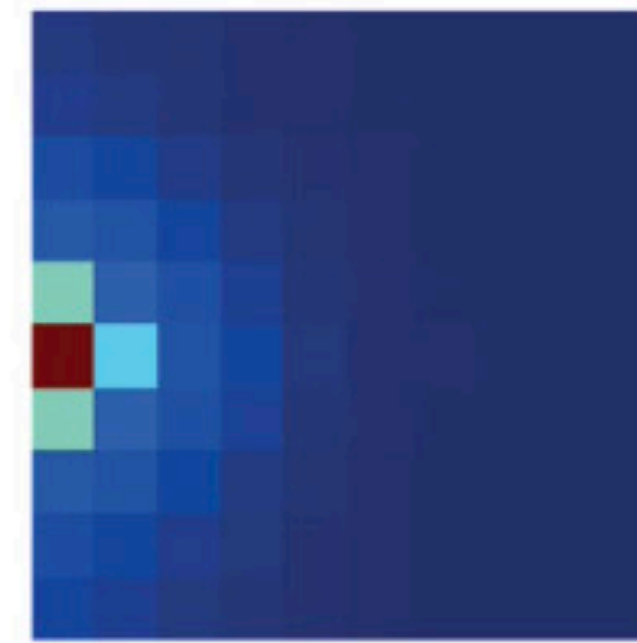
# It looks a bit like place cells...



The successor representation



The column looks like a place cell

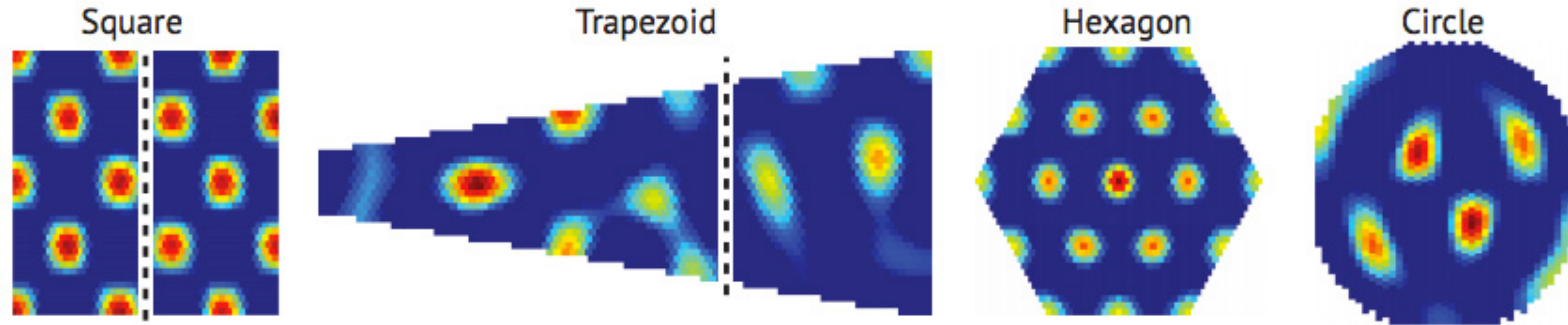


Place cells  
O'Keefe & Dostrovsky, 1971

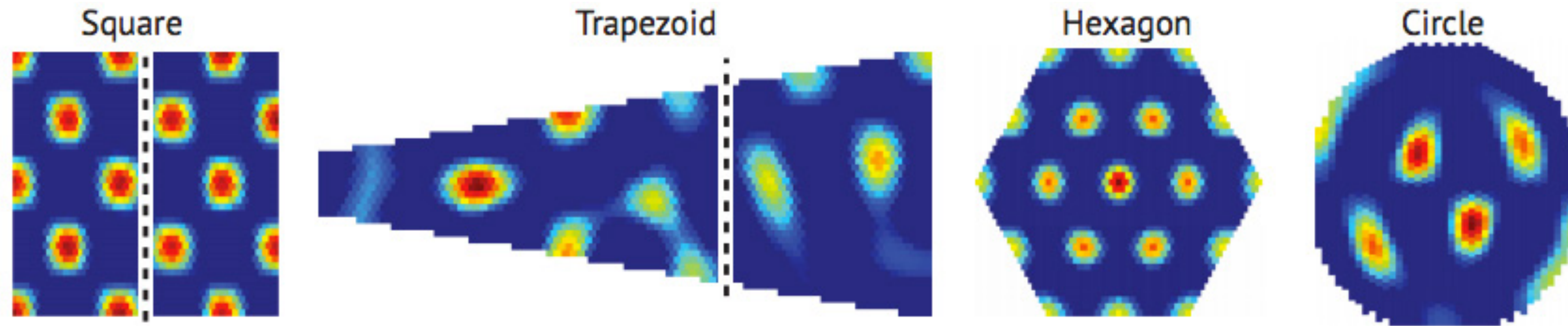
# Eigenvectors of the successor representation look grid cells



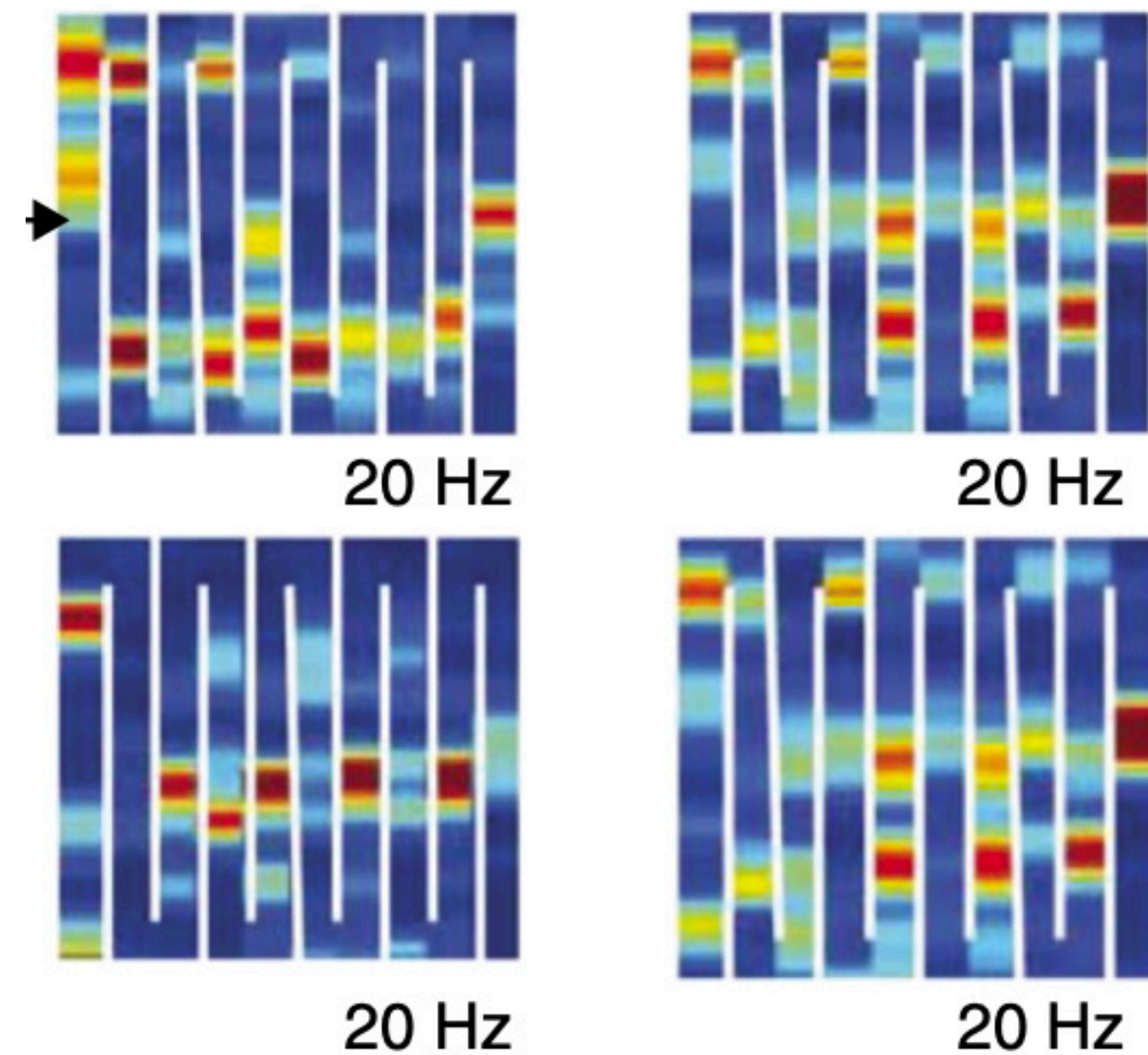
# Eigenvectors of the successor representation look grid cells



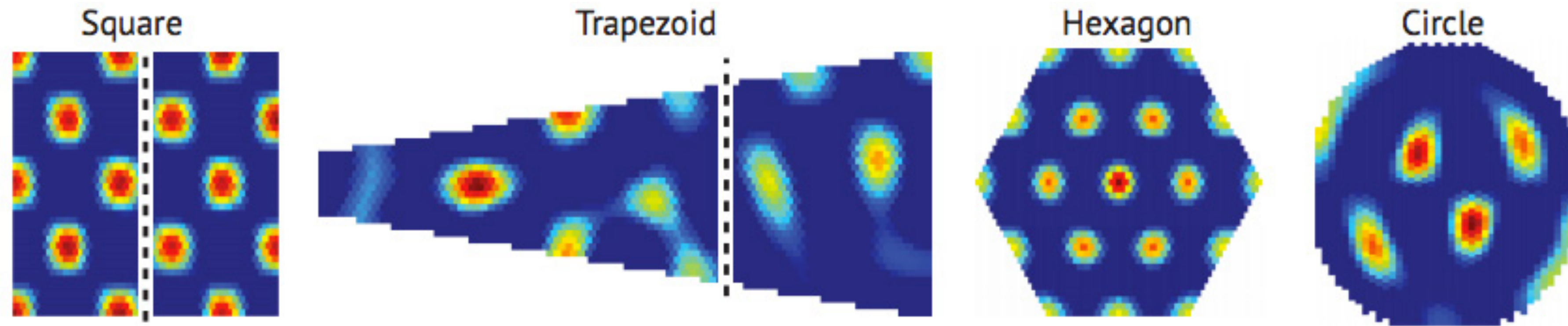
# Eigenvectors of the successor representation look grid cells



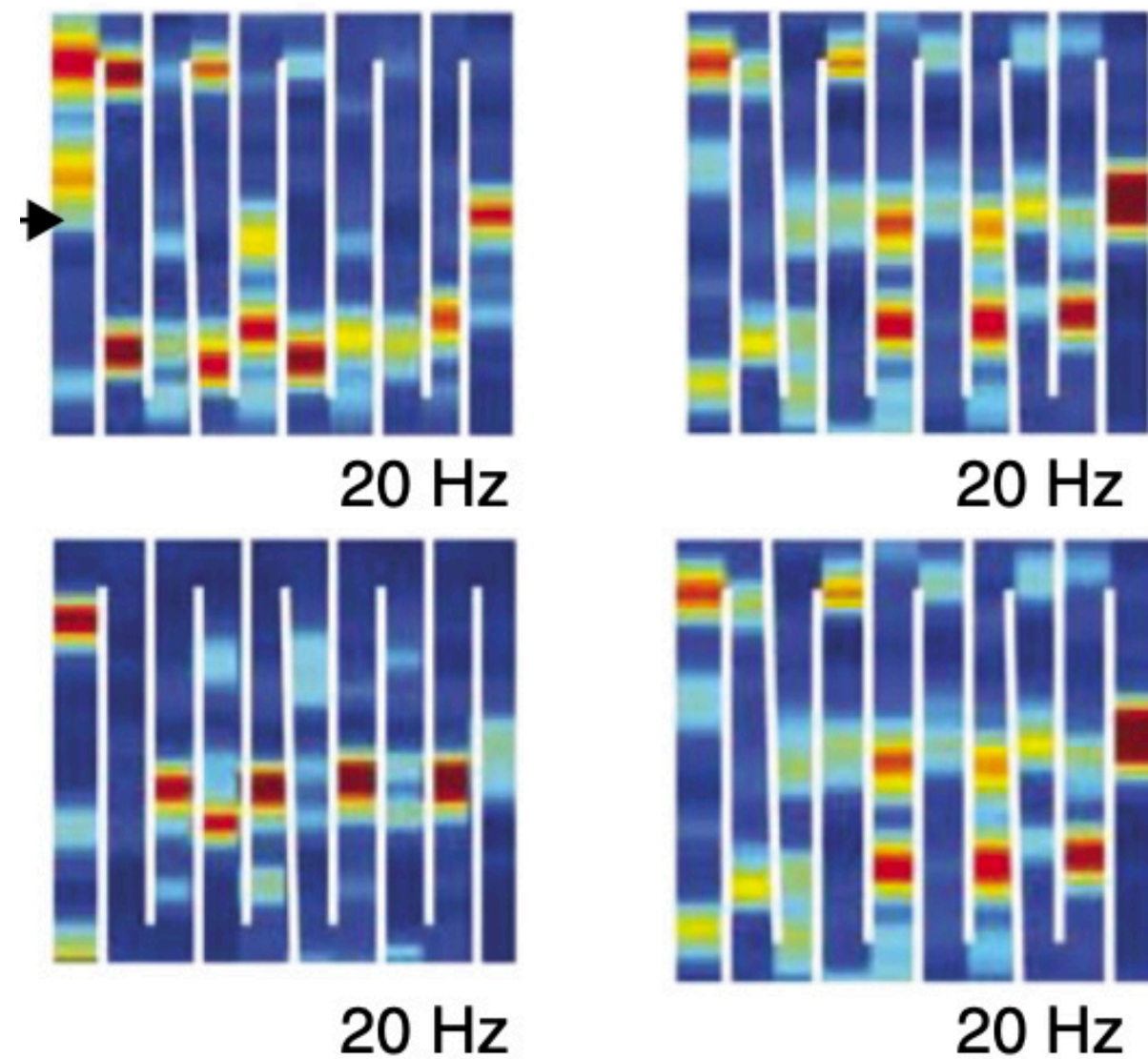
Grid cells in hairpin mazes 'fragment'



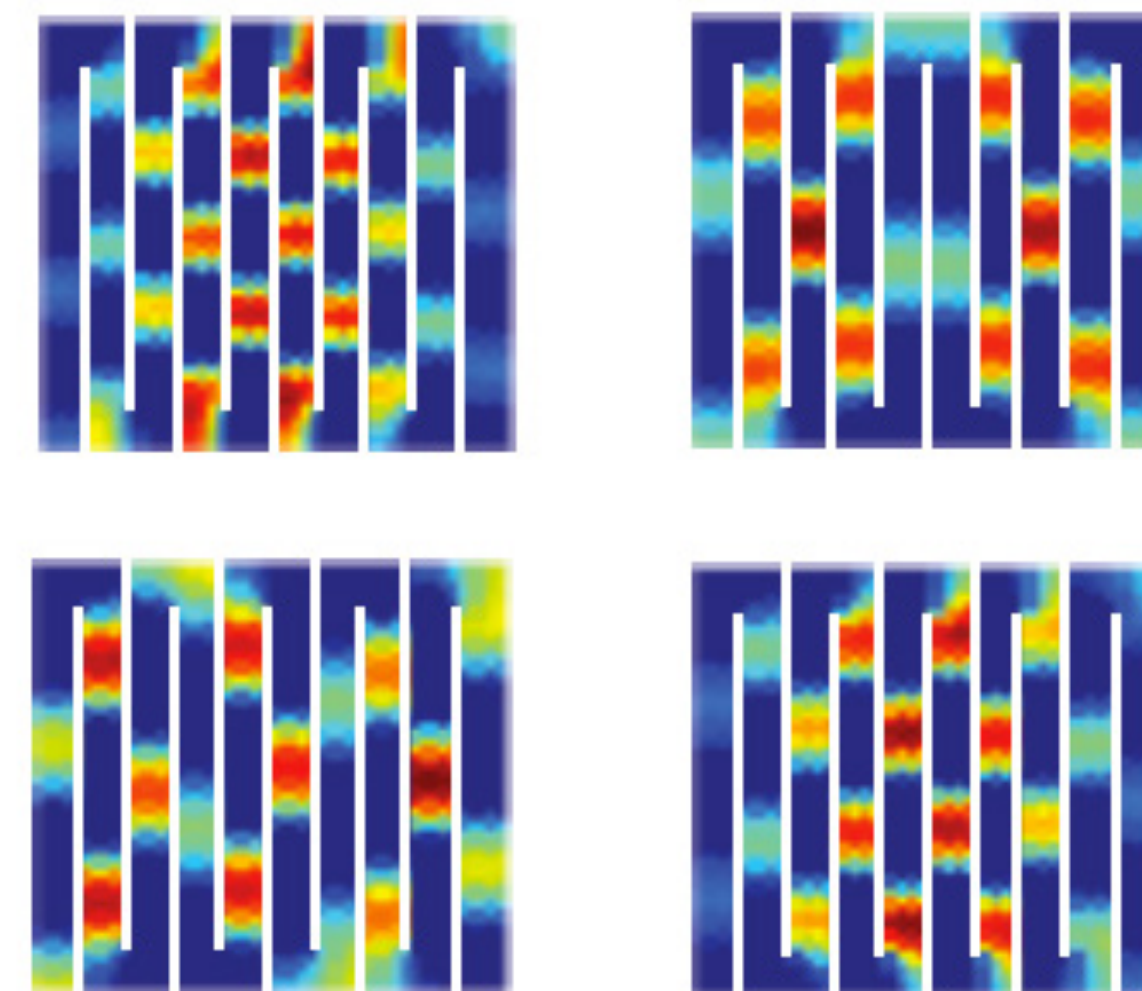
# Eigenvectors of the successor representation look grid cells



Grid cells in hairpin mazes 'fragment'



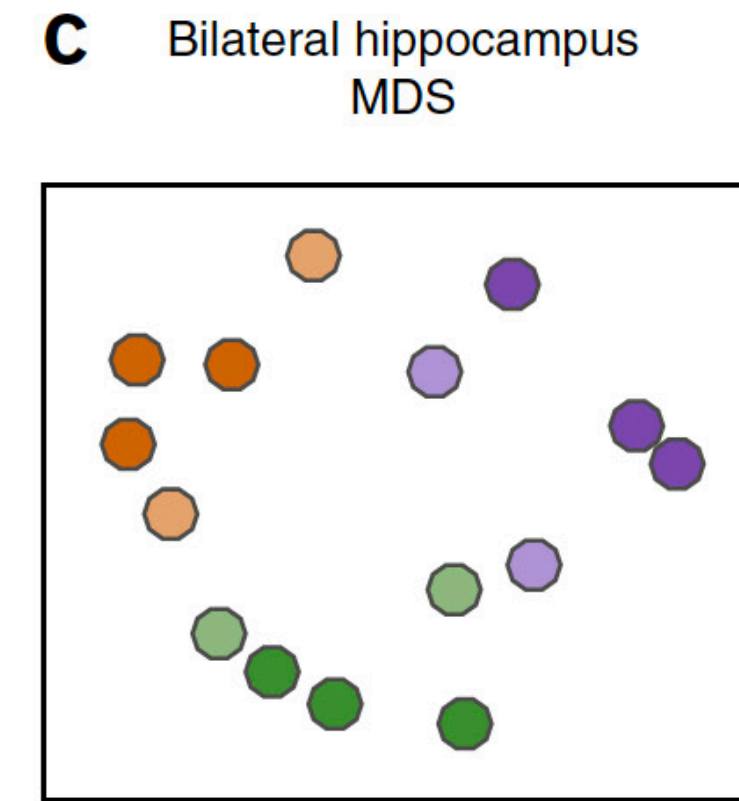
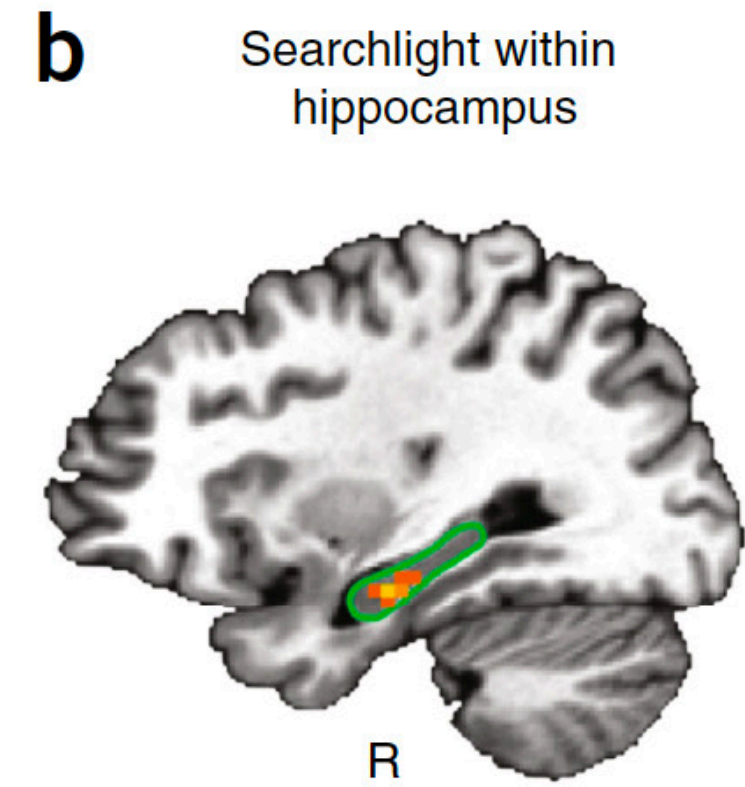
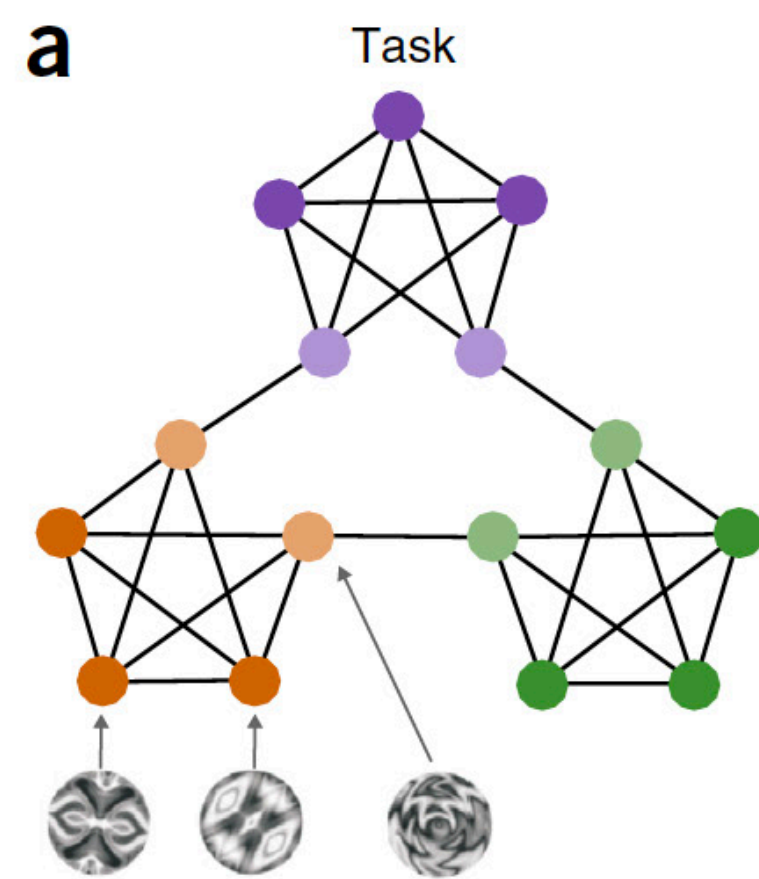
So too do the SR grid cells



# Accounts for non-spatial brain representations too

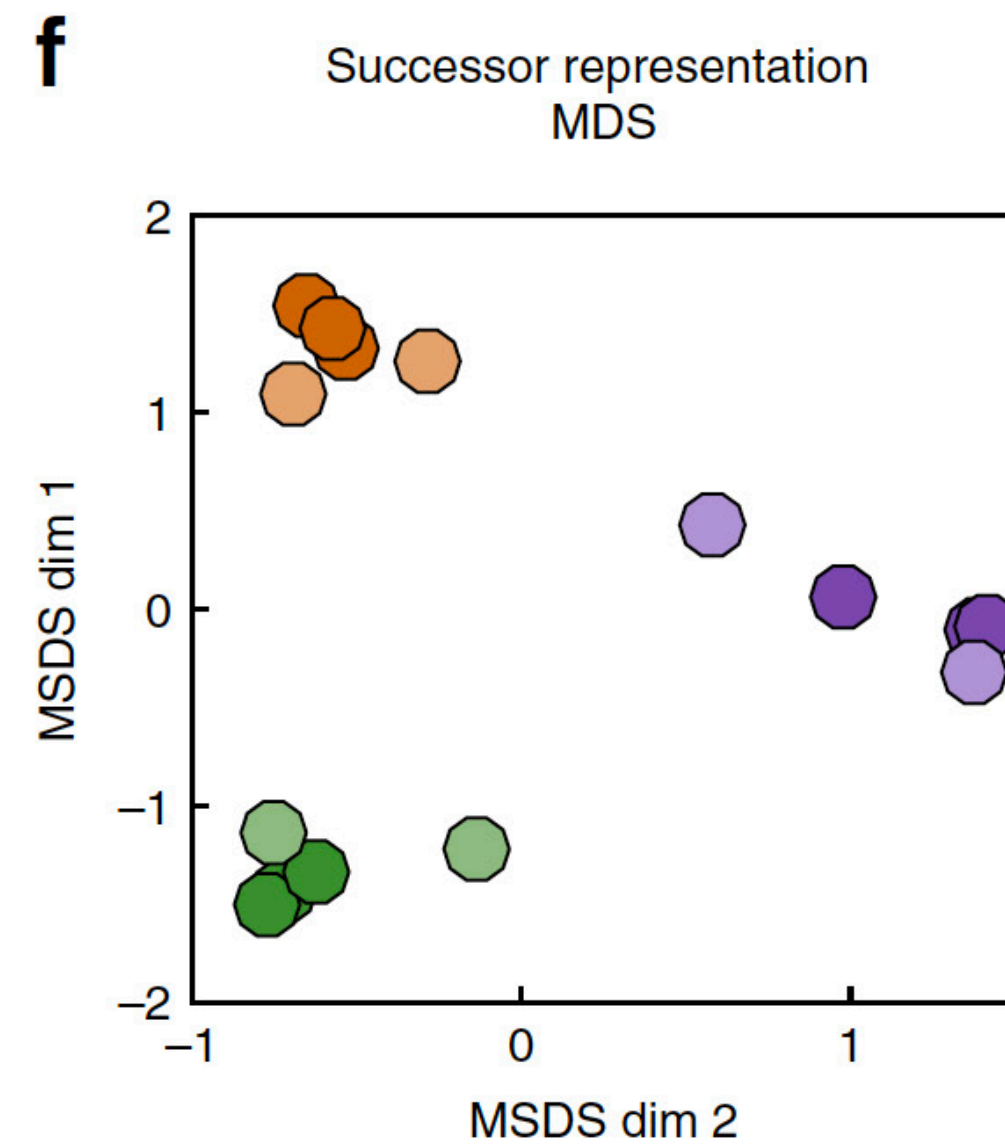
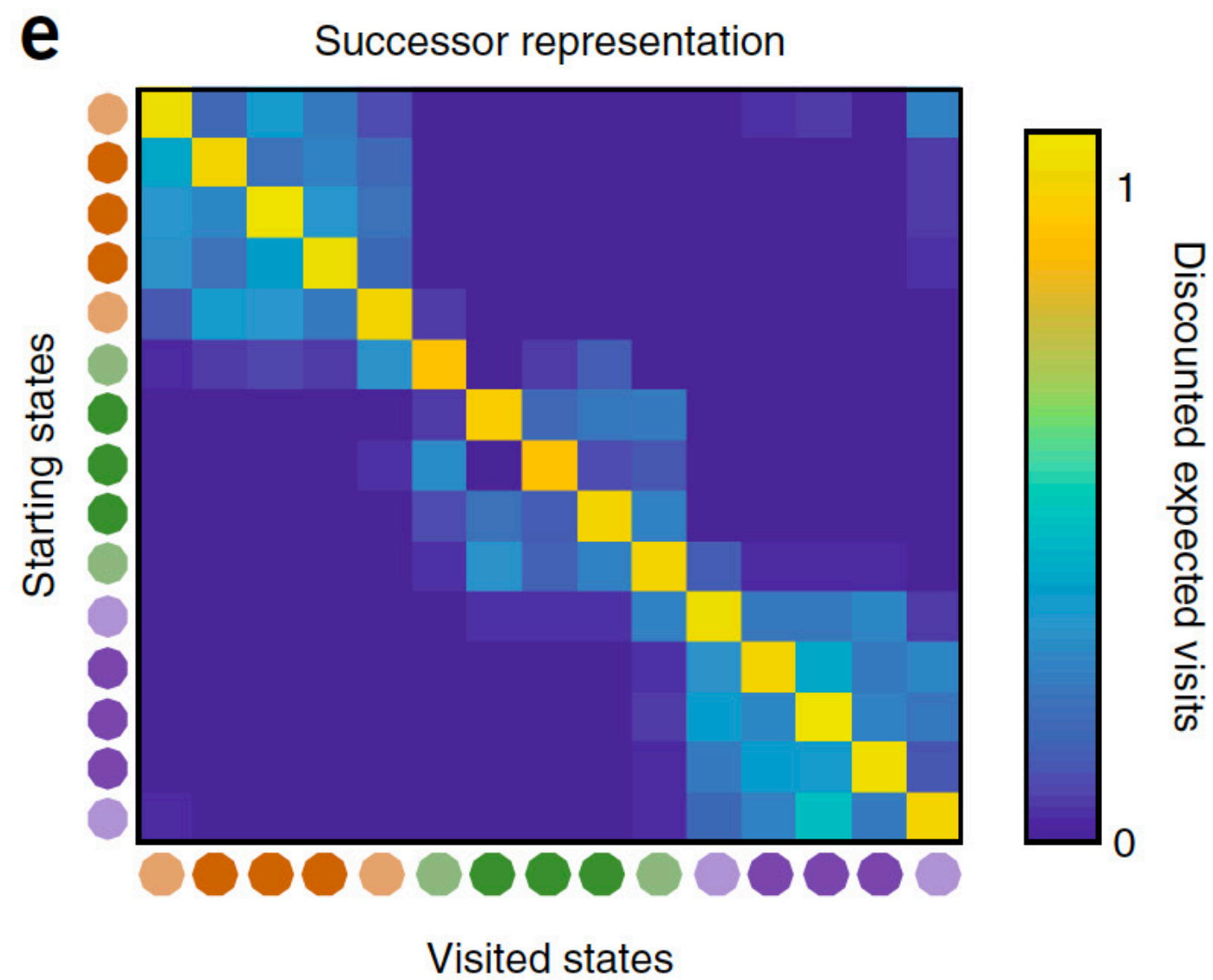
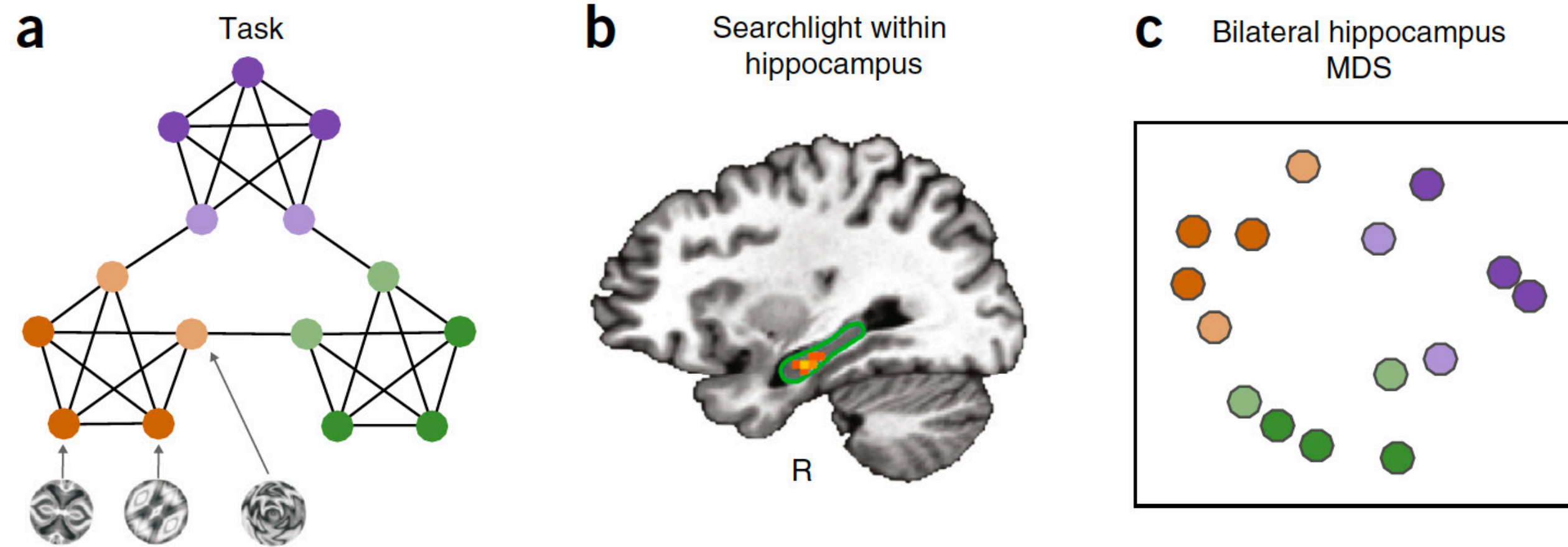
# Accounts for non-spatial brain representations too

Schapiro *et al.* (2015)



# Accounts for non-spatial brain representations too

Schapiro *et al.* (2015)



**We've understood a lot already**

**We've understood a lot already**

**... but the brain uses neurons**



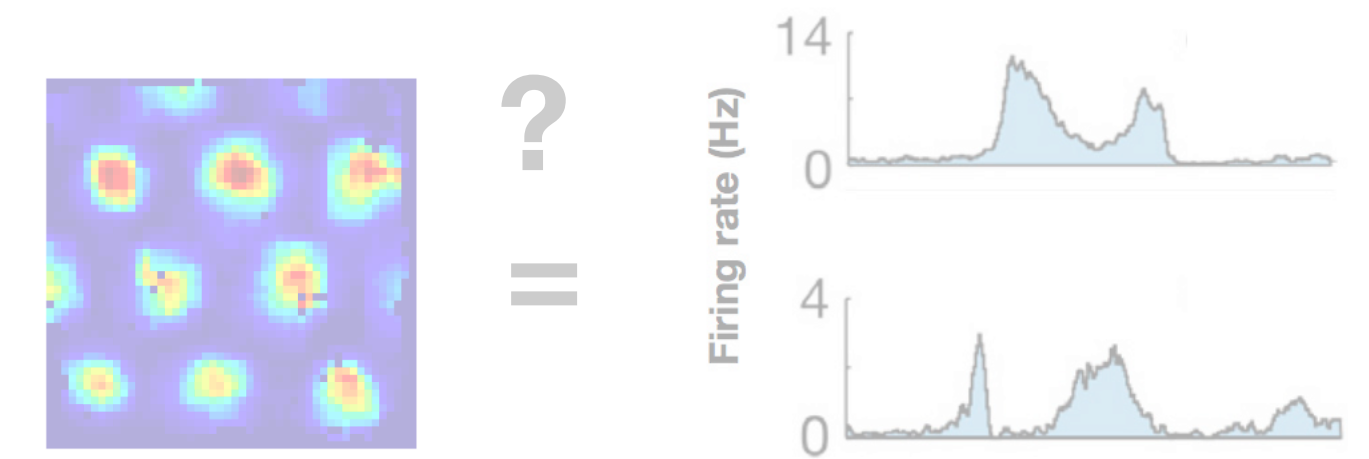
**We've understood a lot already**

**... but the brain uses neurons**

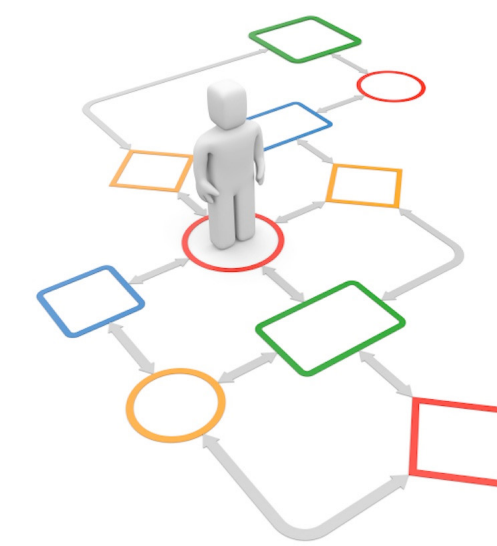
**... and learns from sequential experience**

# Puzzles of cognitive maps in the brain

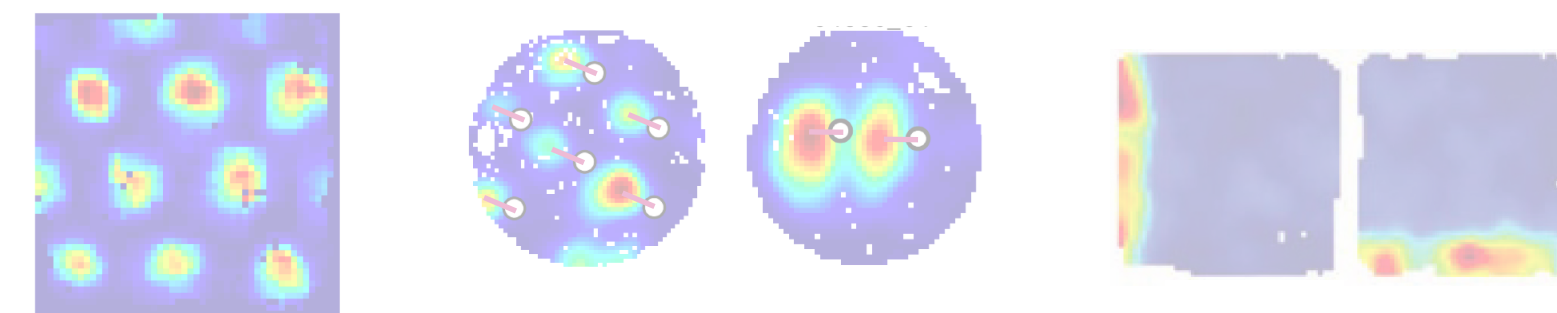
How does the same system do space and non-space?



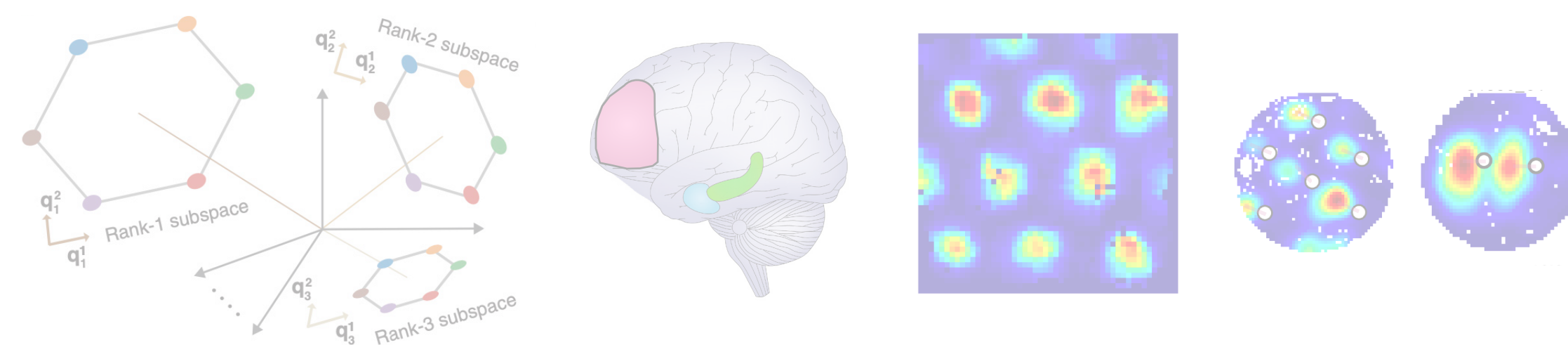
How can brains learn these maps?



Why do the neurons look the ways they do?



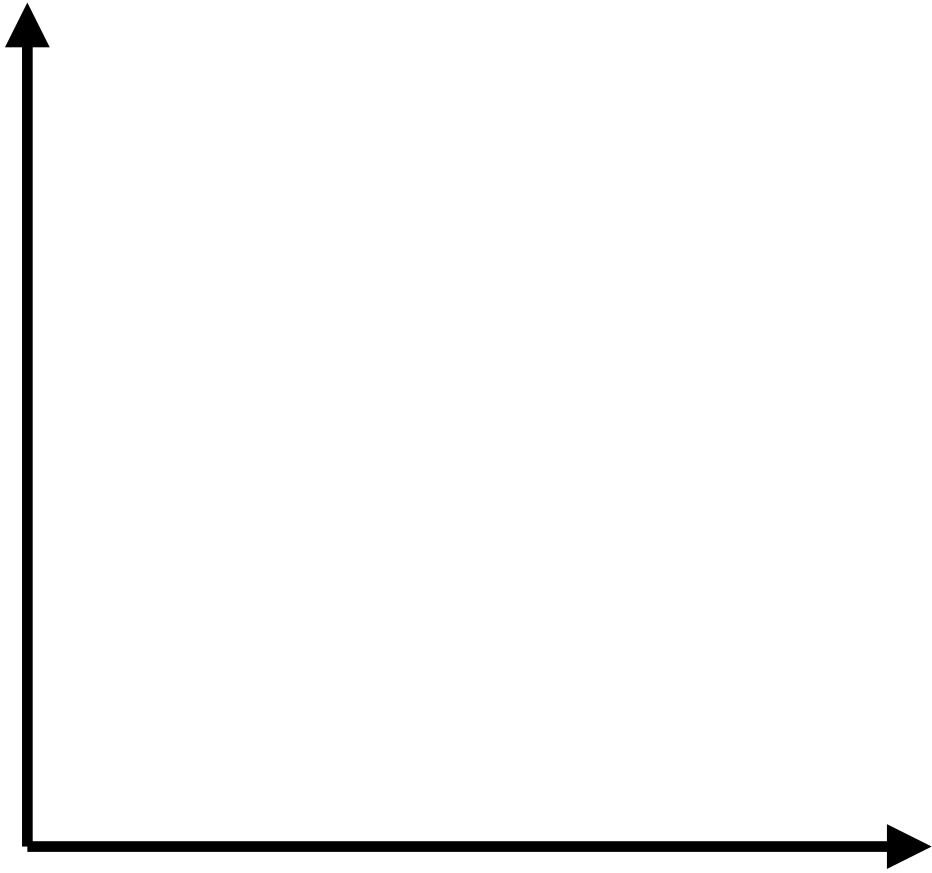
How do different brain regions solve the same problem in different ways?



# **Recurrent neural networks learn from sequences**

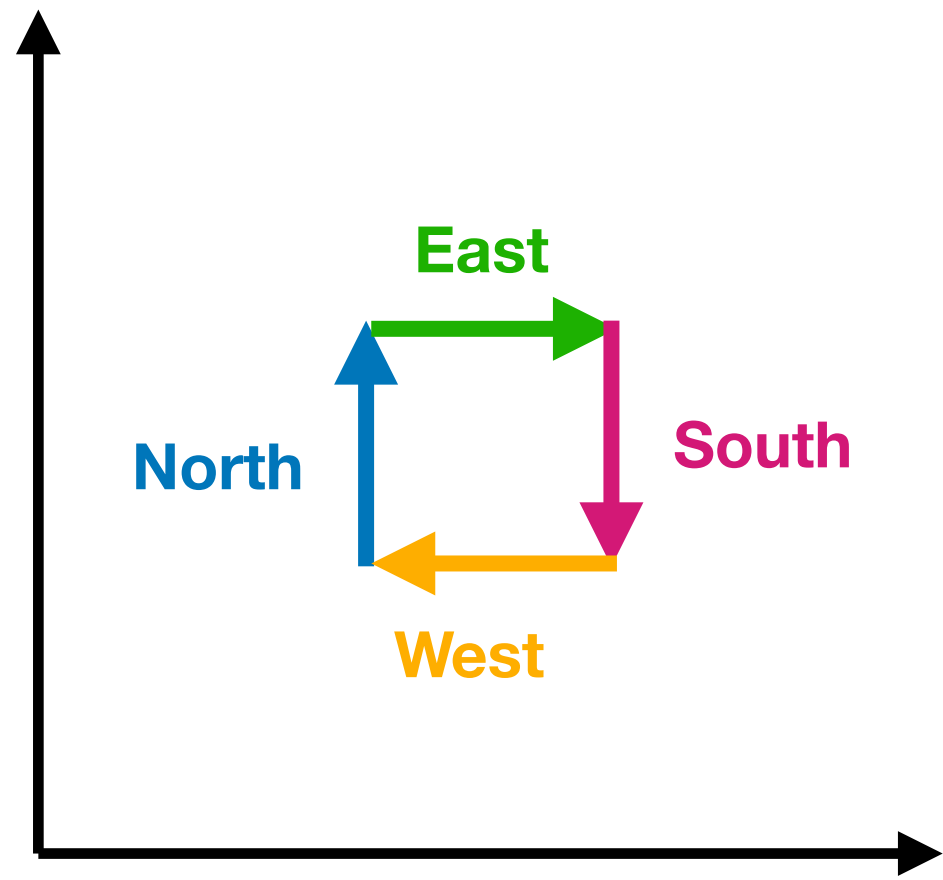
# Recurrent neural networks learn from sequences

Physical space

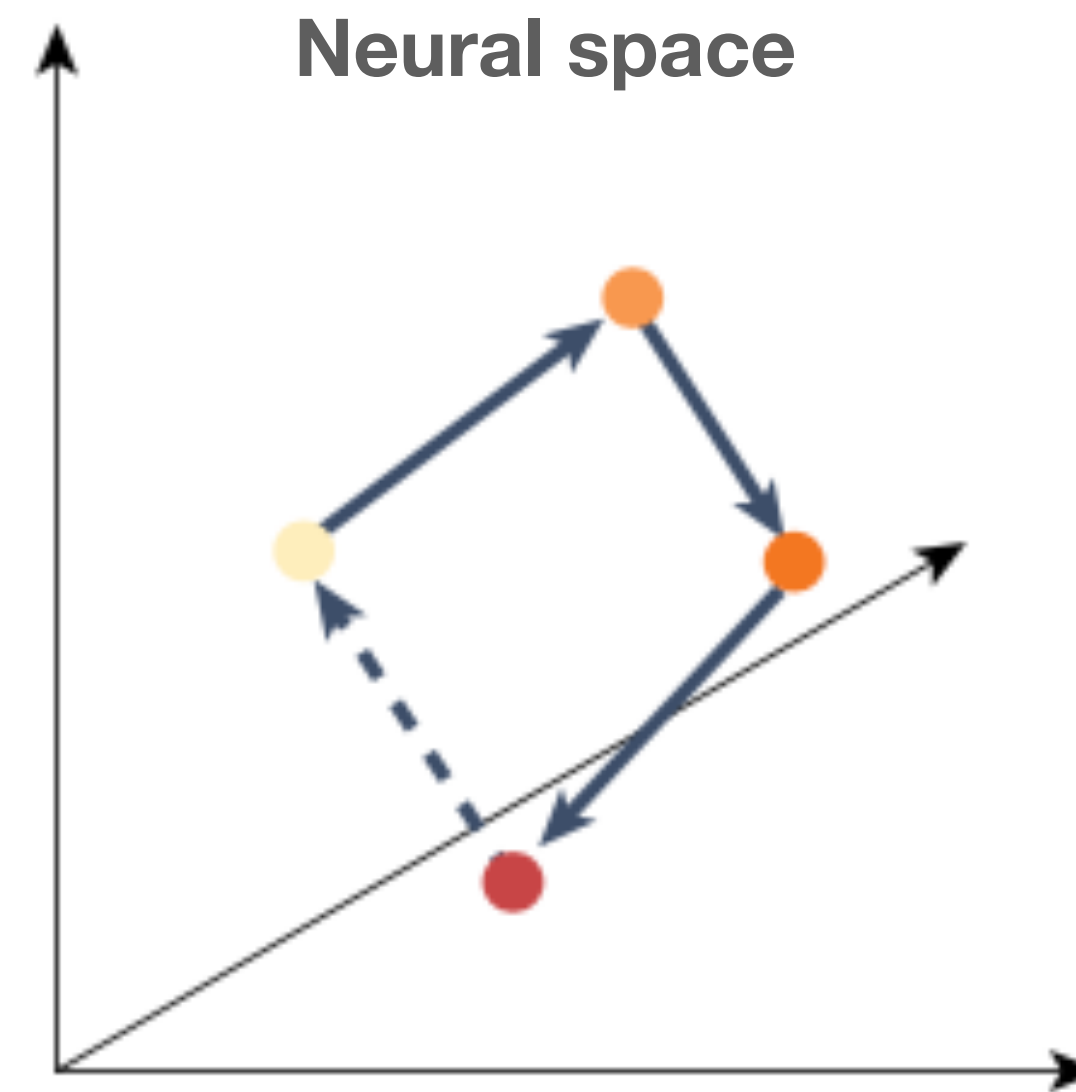
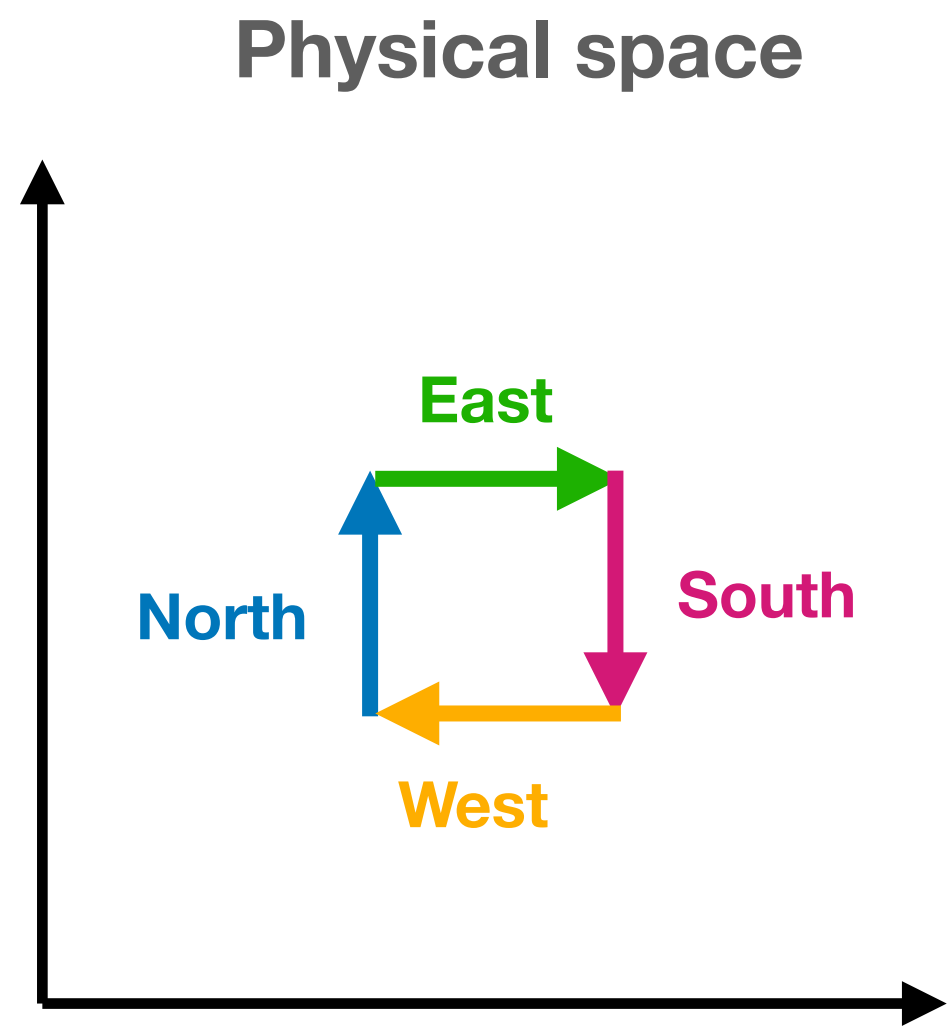


# Recurrent neural networks learn from sequences

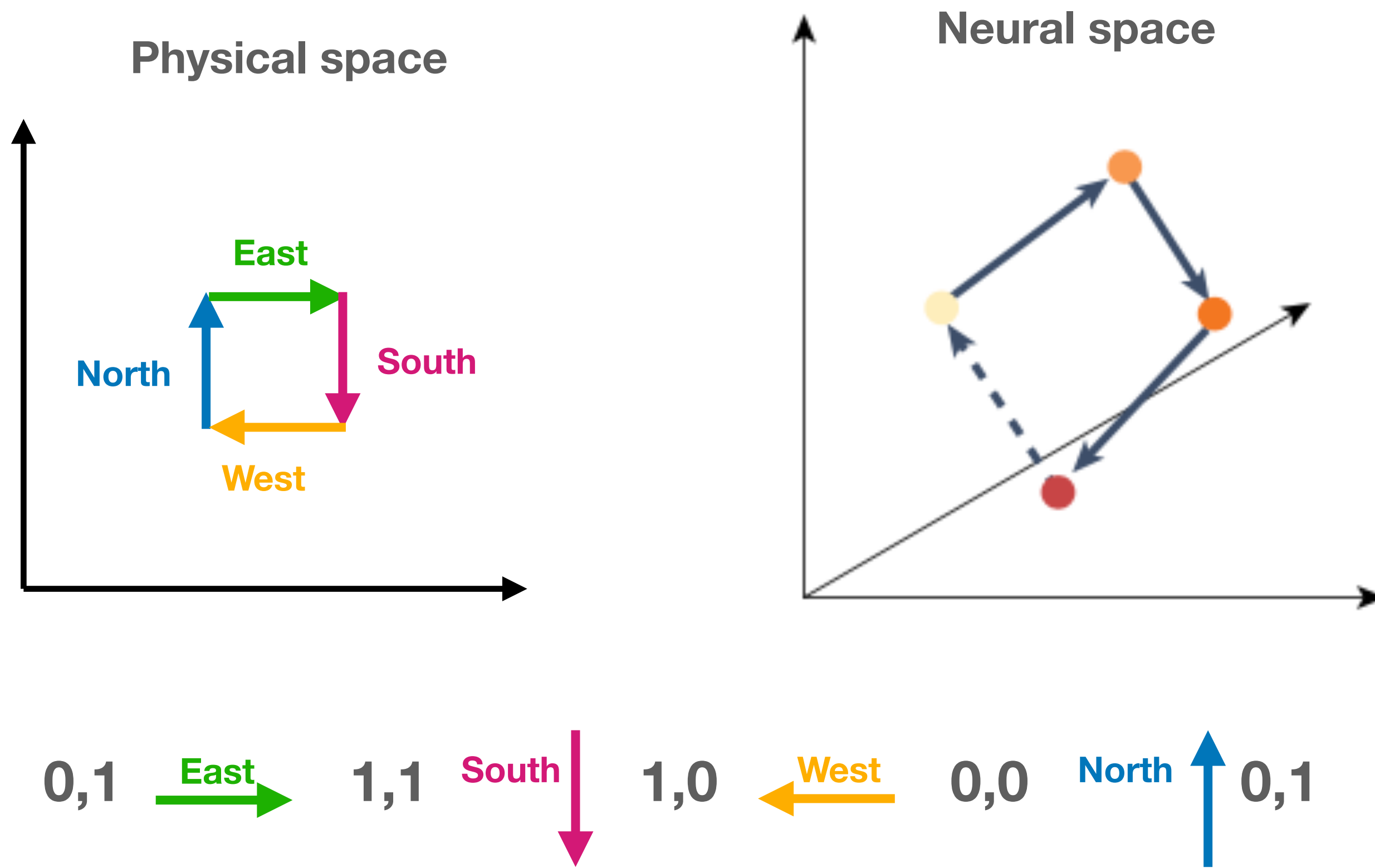
Physical space



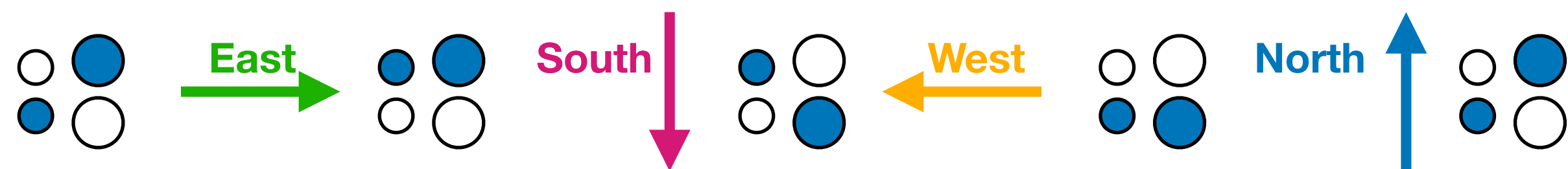
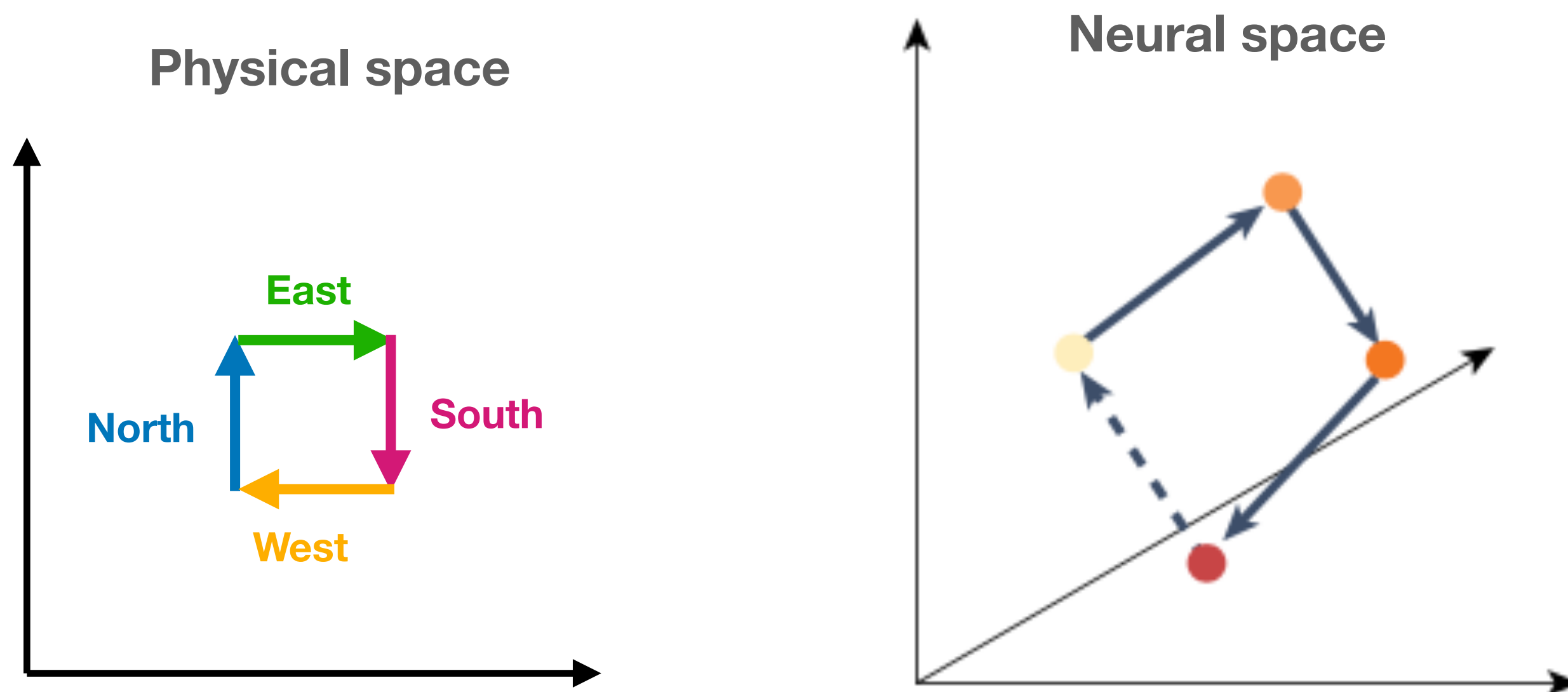
# Recurrent neural networks learn from sequences



# Recurrent neural networks learn from sequences

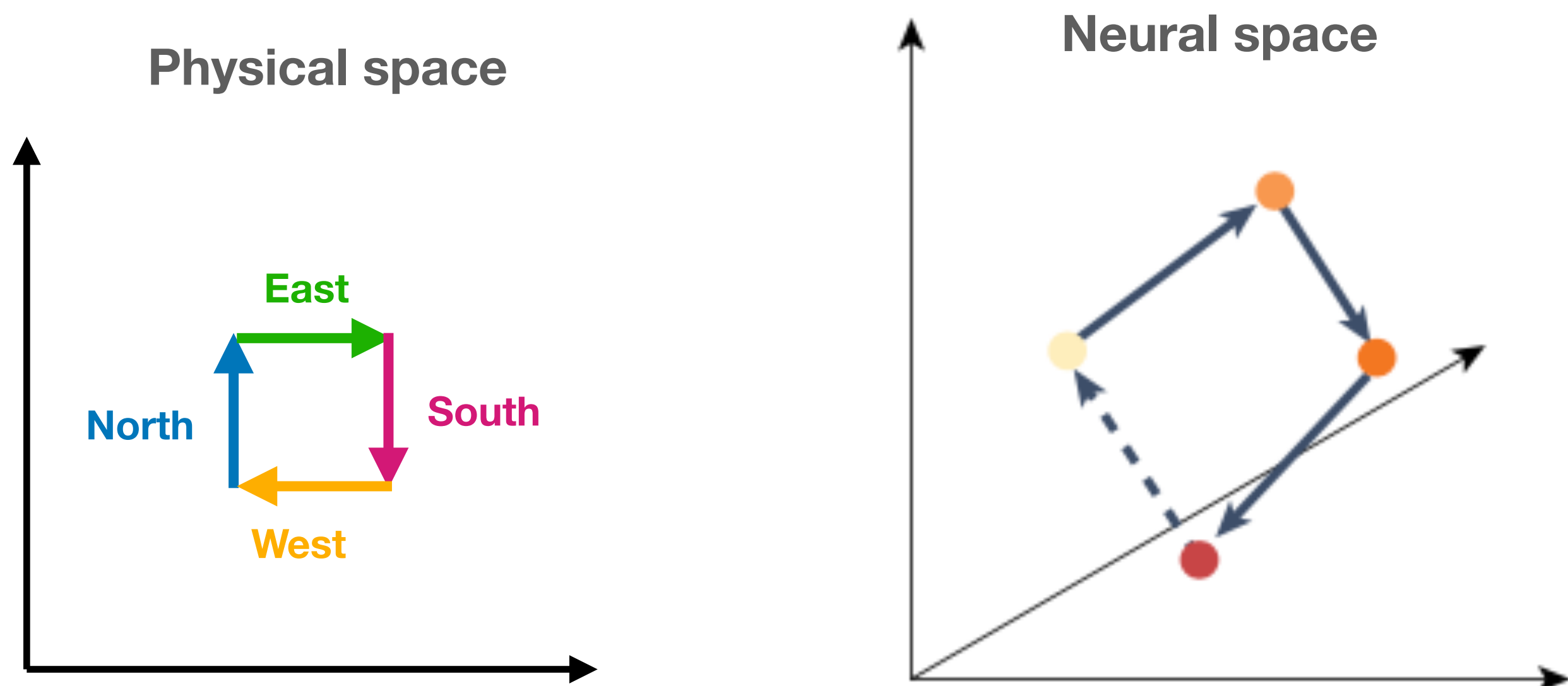


# Recurrent neural networks learn from sequences

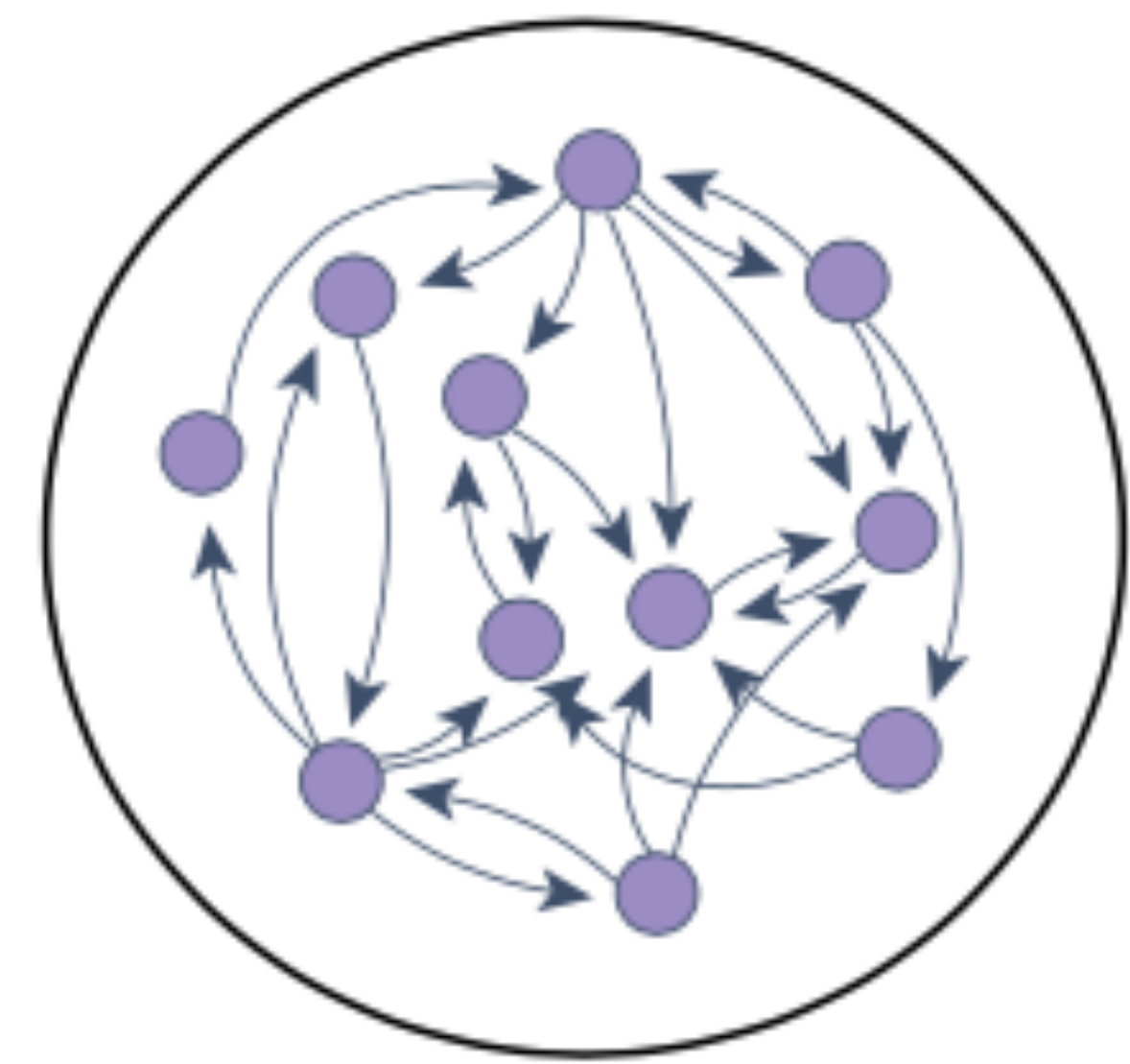
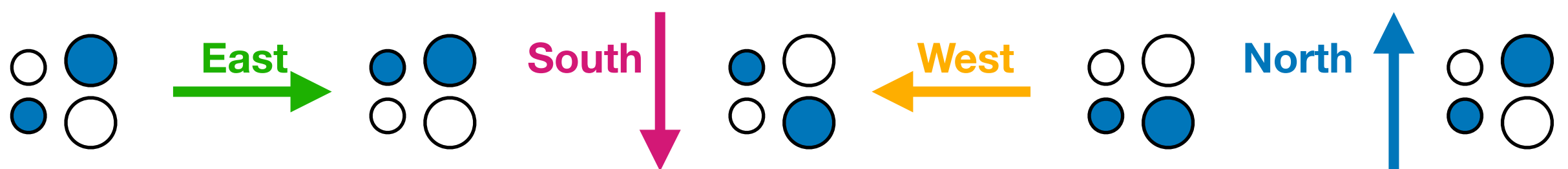




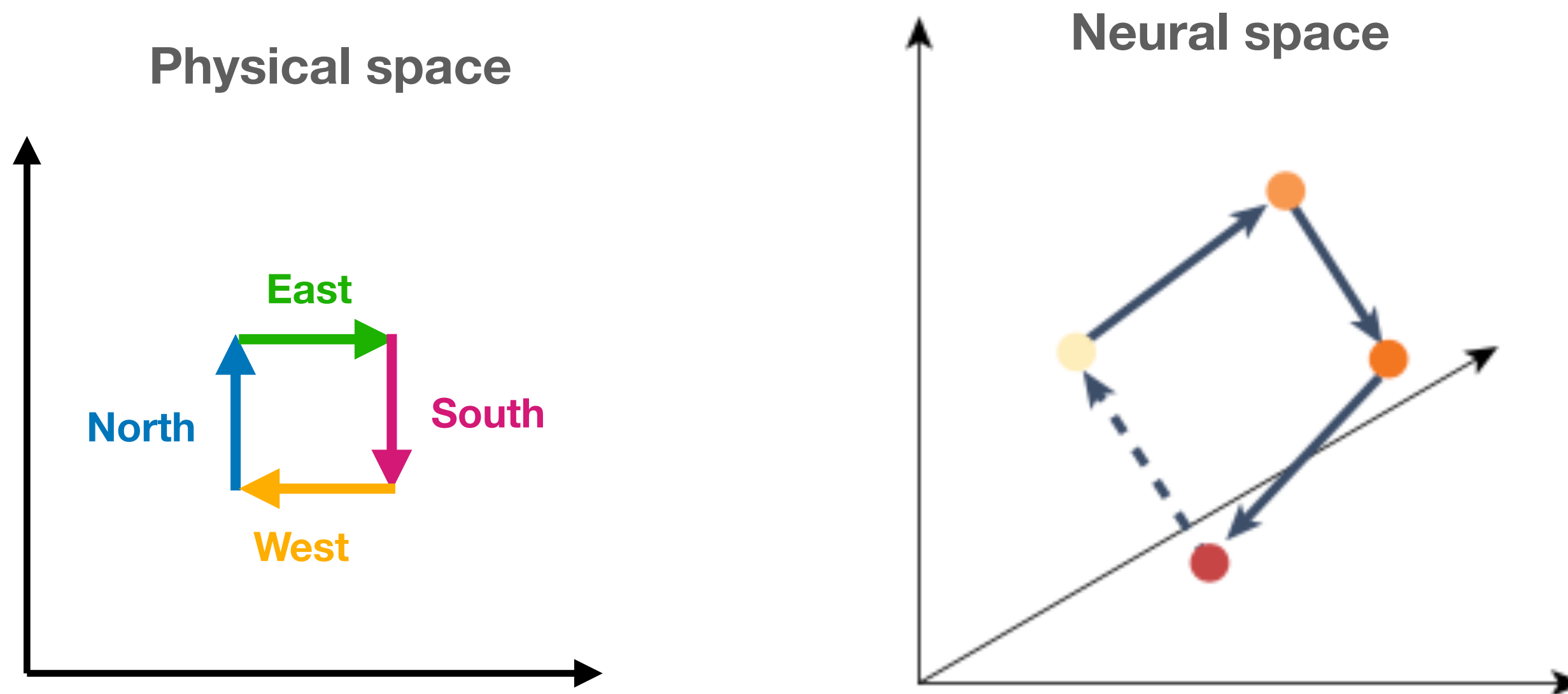
# Recurrent neural networks learn from sequences



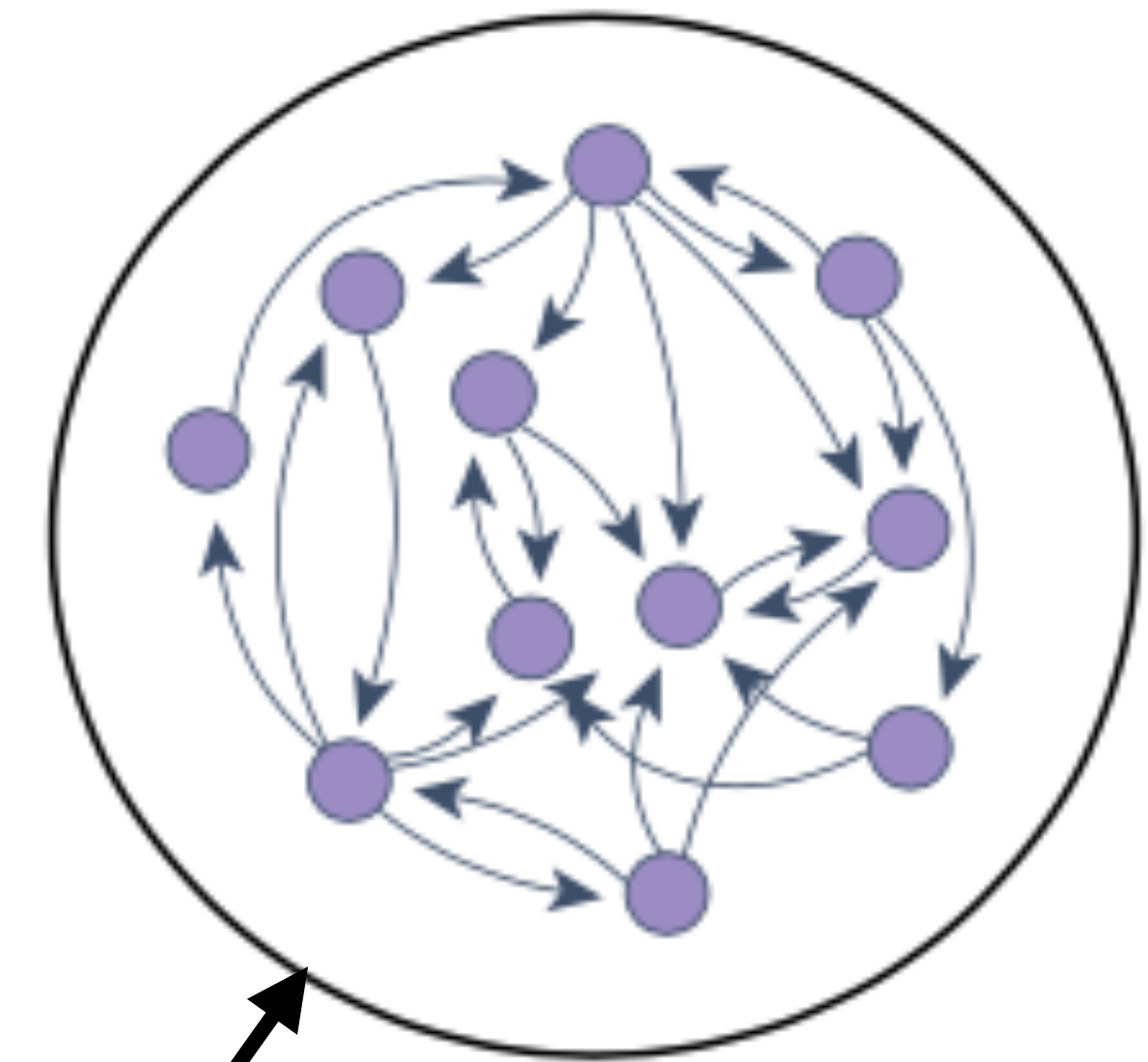
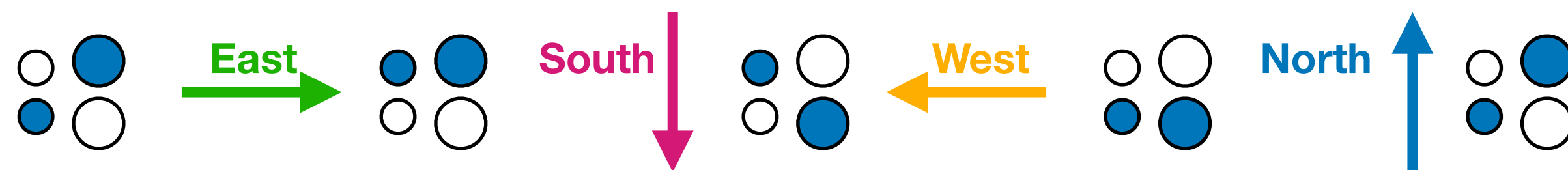
0,1  $\xrightarrow{\text{East}}$  1,1  $\downarrow$  South 1,0  $\xleftarrow{\text{West}}$  0,0  $\uparrow$  North 0,1



# Recurrent neural networks learn from sequences

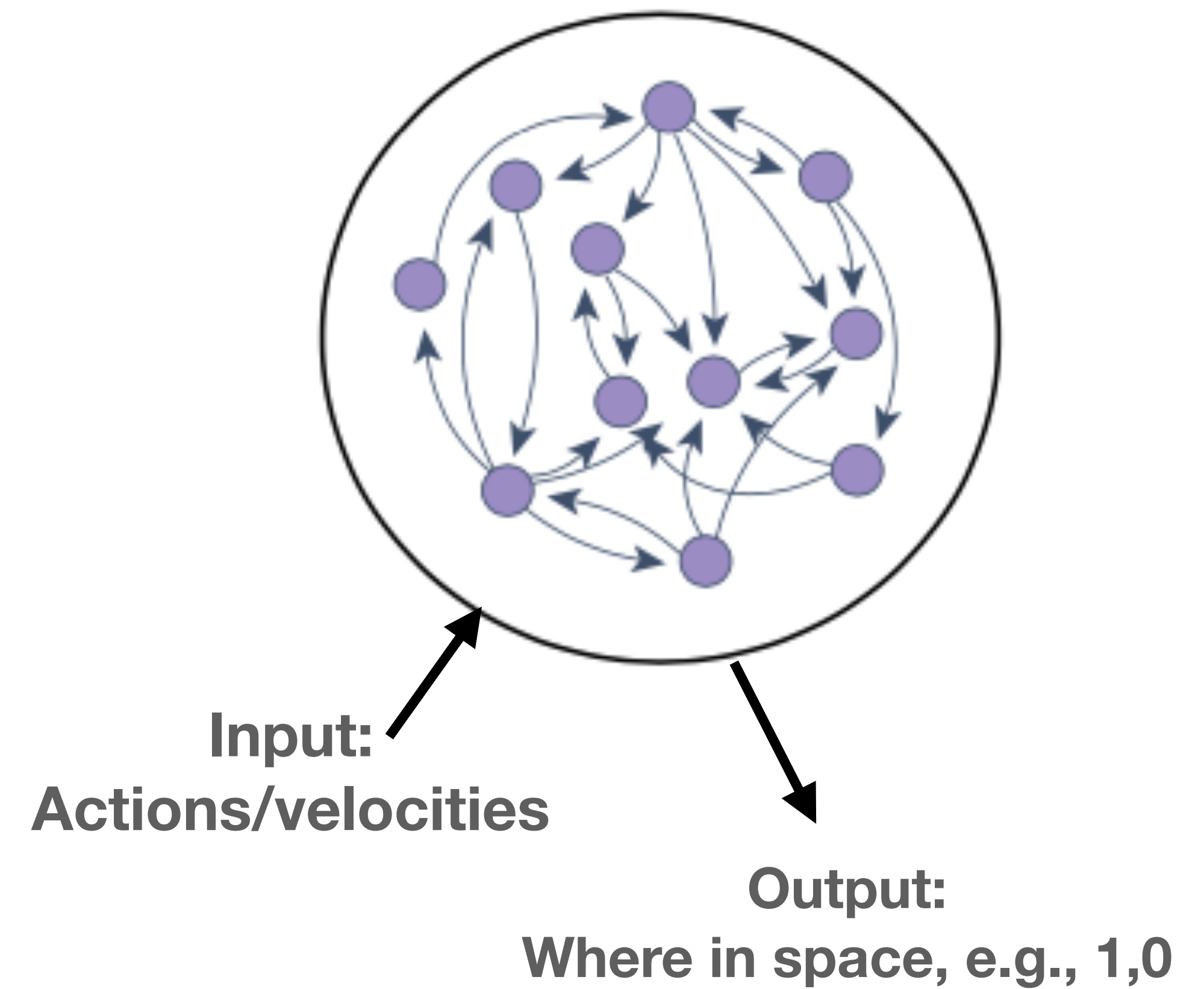
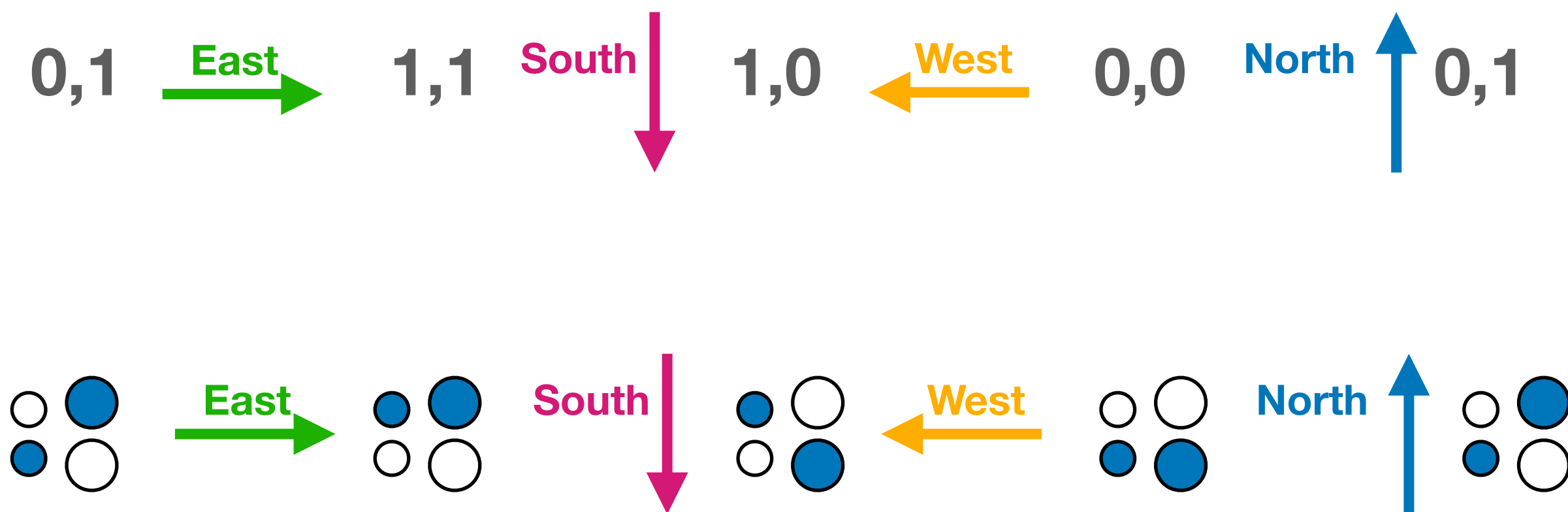
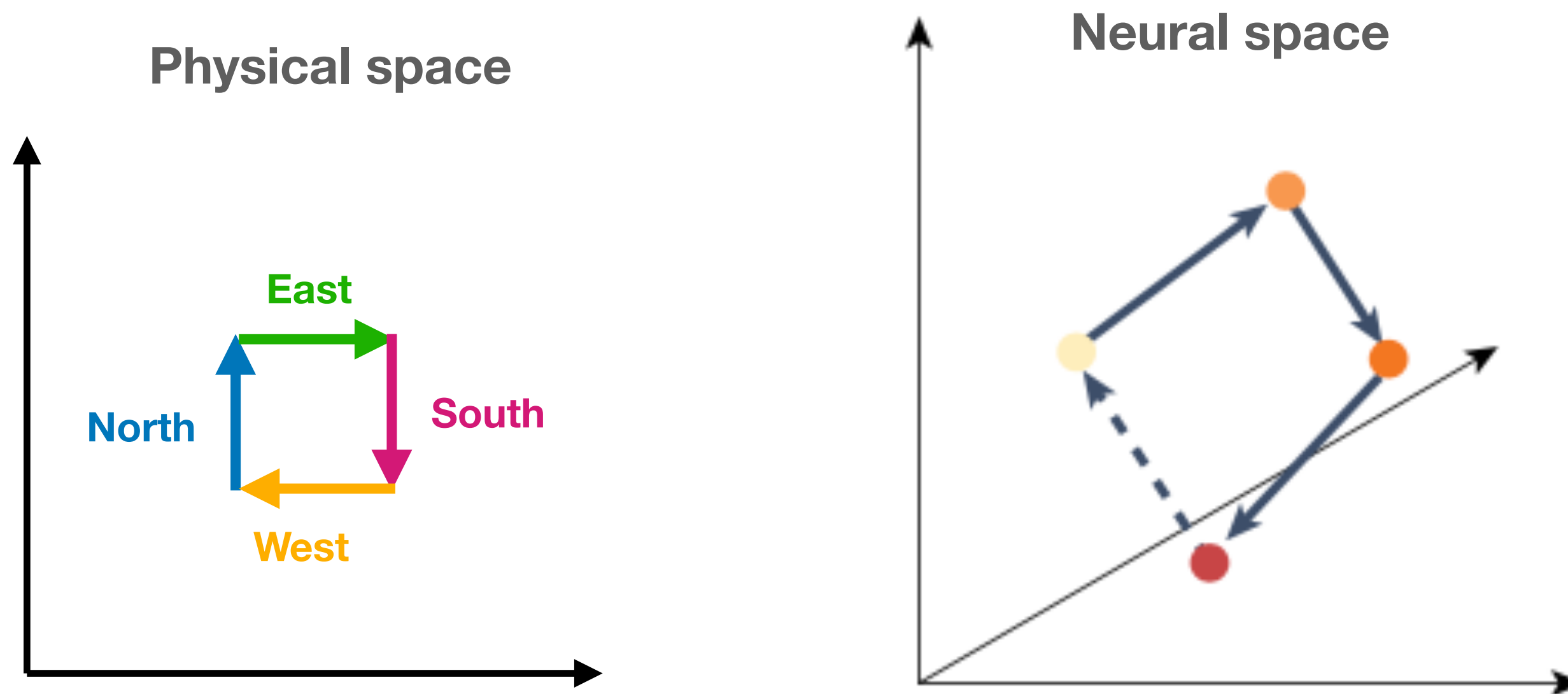


0,1 **East** 1,1 **South** 1,0 **West** 0,0 **North** 0,1

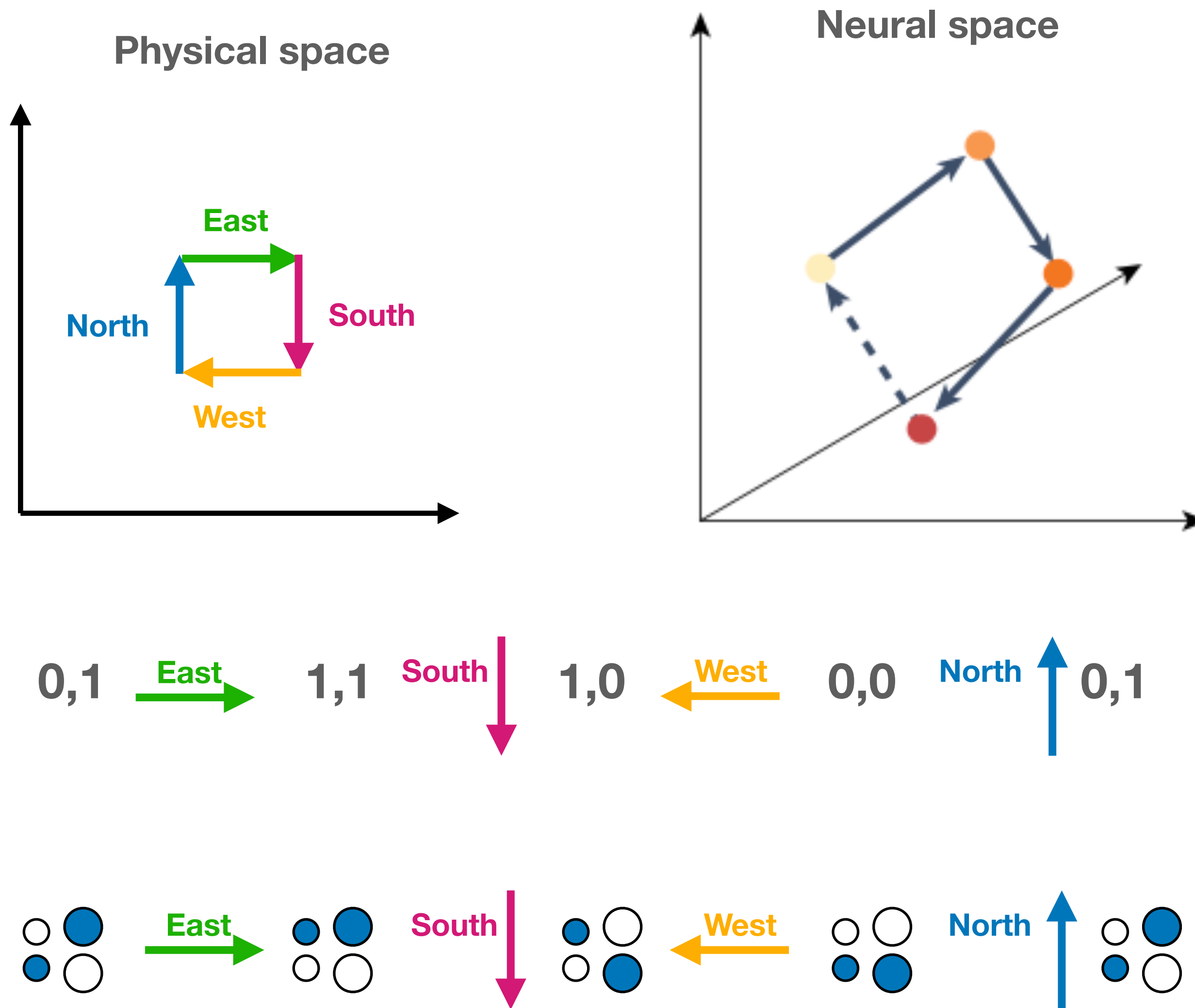


Input:  
Actions/velocities

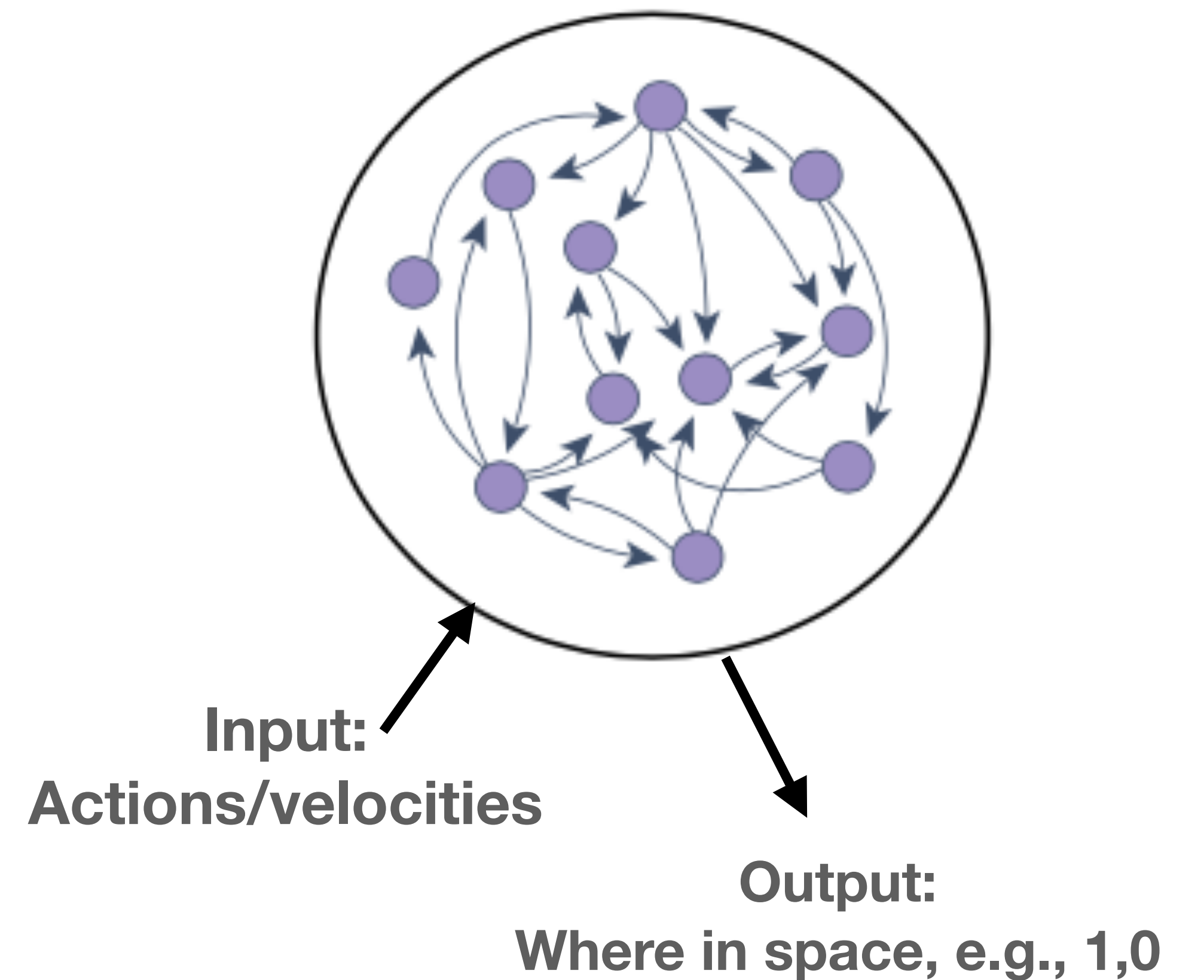
# Recurrent neural networks learn from sequences



# Recurrent neural networks learn from sequences

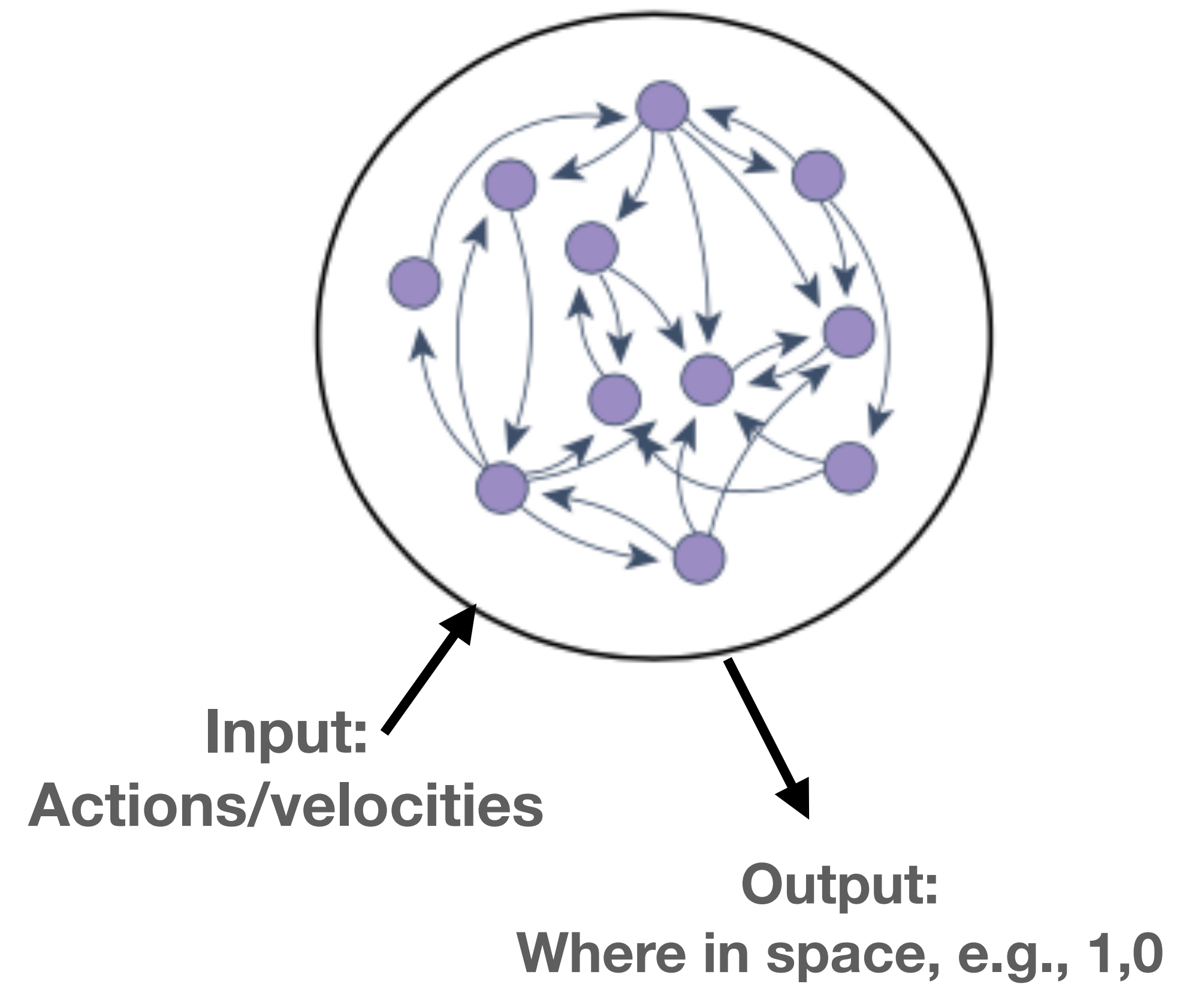


$$\frac{d\bar{s}}{dt} = -\bar{s} + f(W\bar{s} + B\bar{a})$$

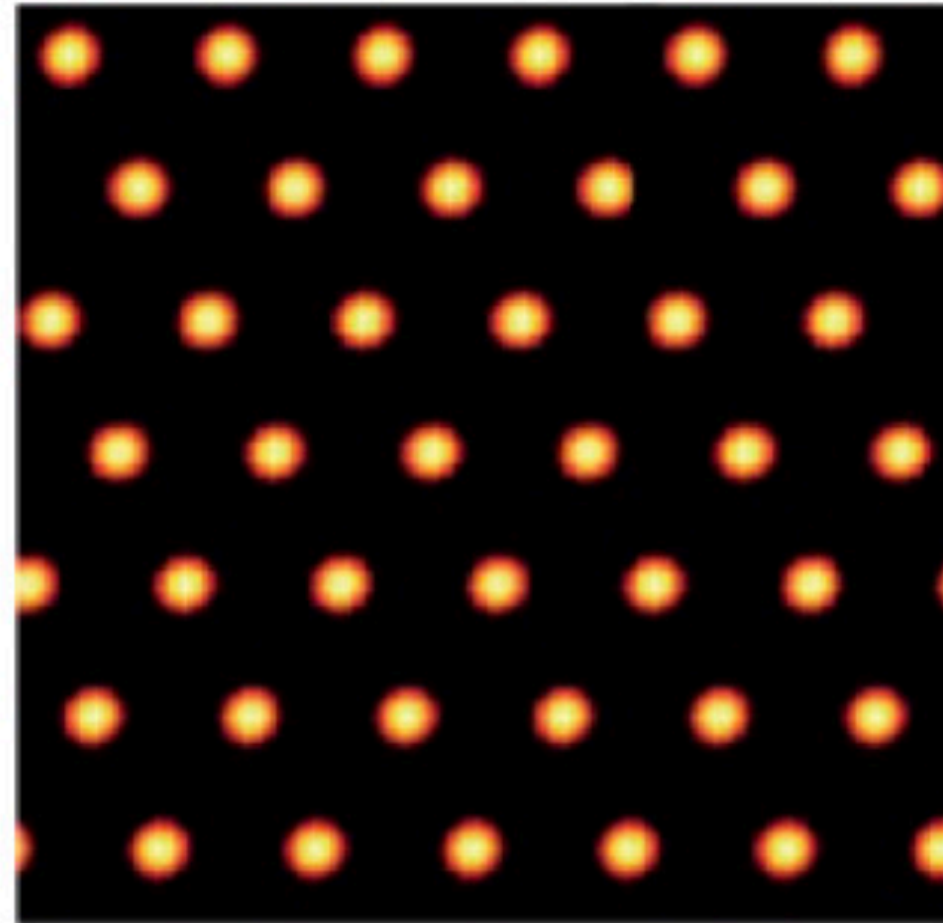


# These networks learn grid cells

$$\frac{d\bar{s}}{dt} = -\bar{s} + f(W\bar{s} + B\bar{a})$$

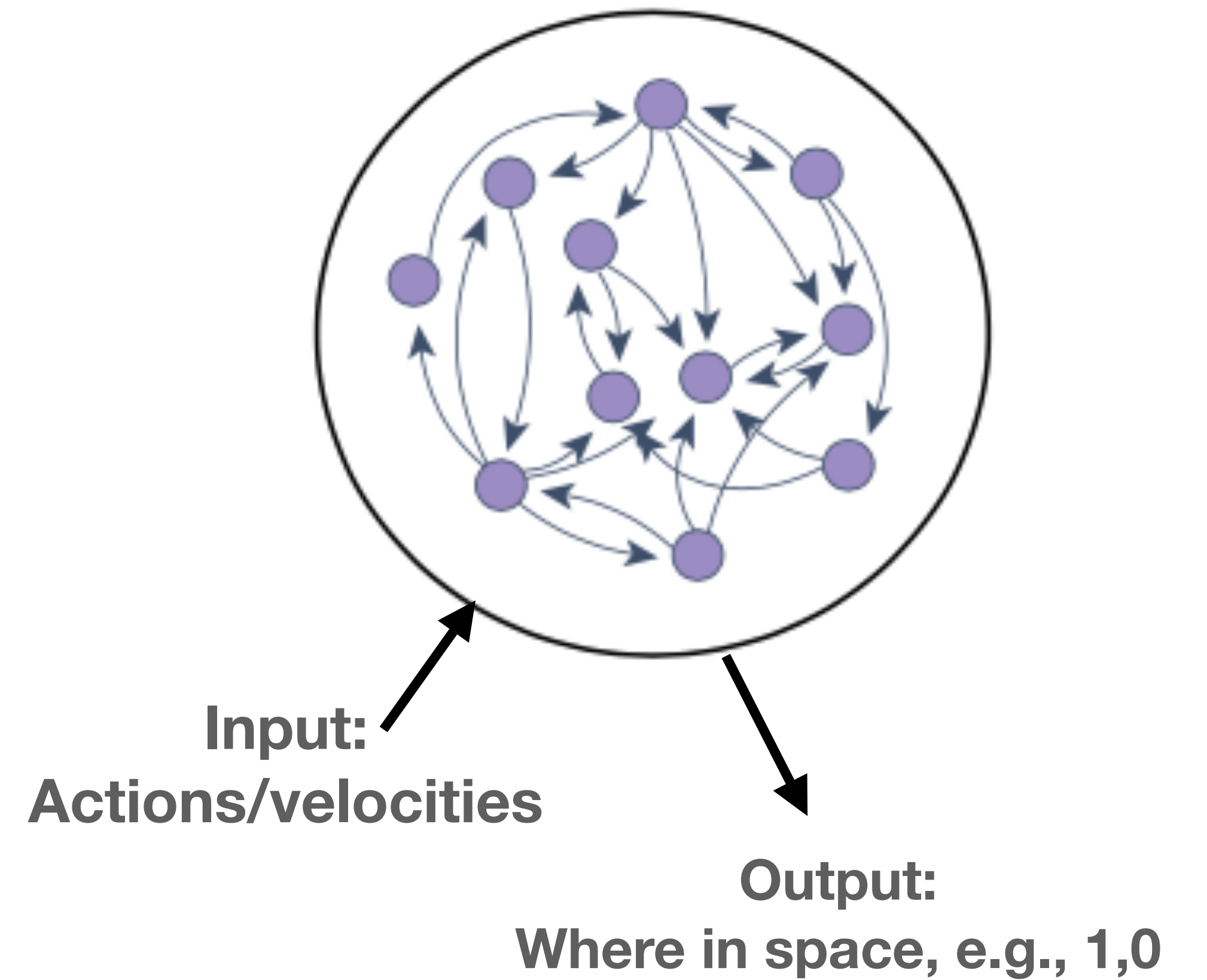


# These networks learn grid cells

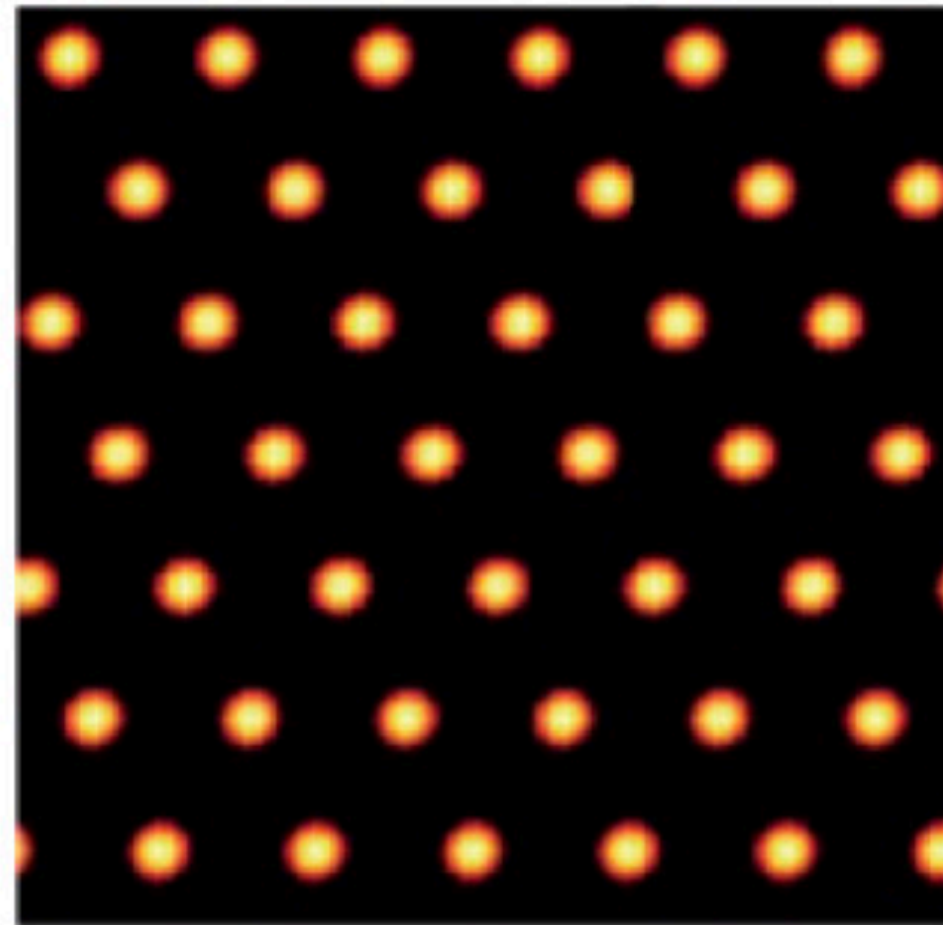


Burak & Feite (2009)

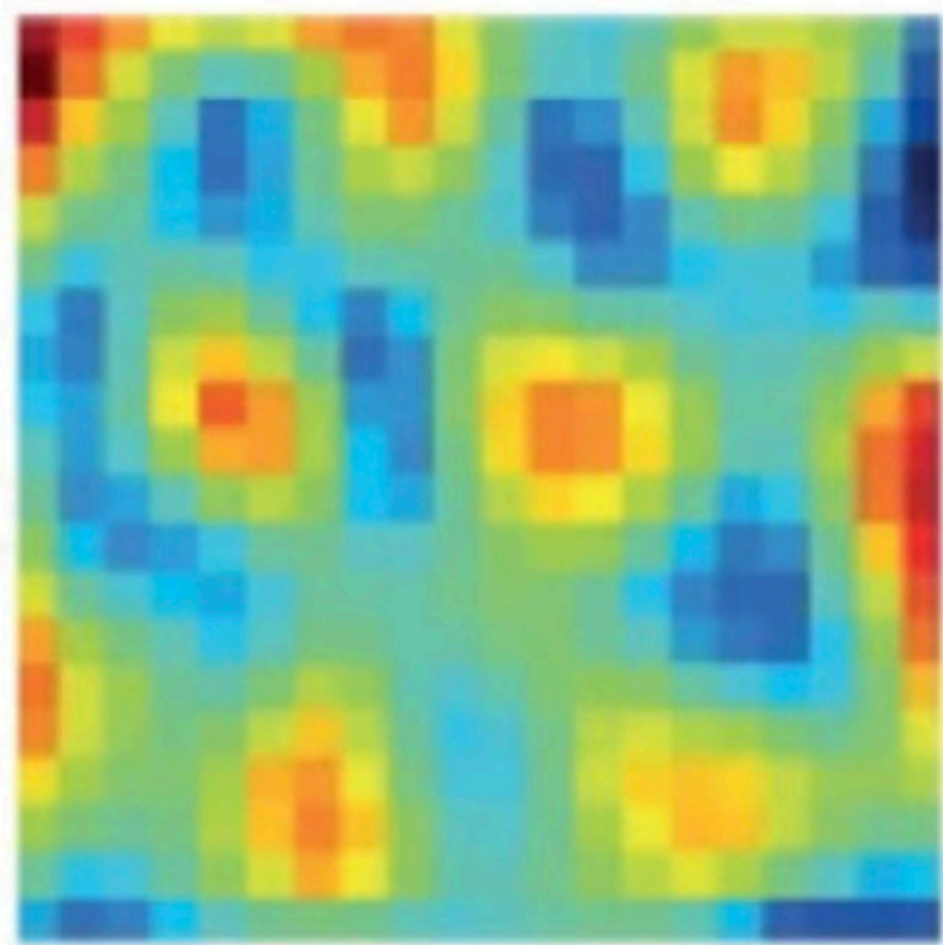
$$\frac{d\bar{s}}{dt} = -\bar{s} + f(W\bar{s} + B\bar{a})$$



# These networks learn grid cells

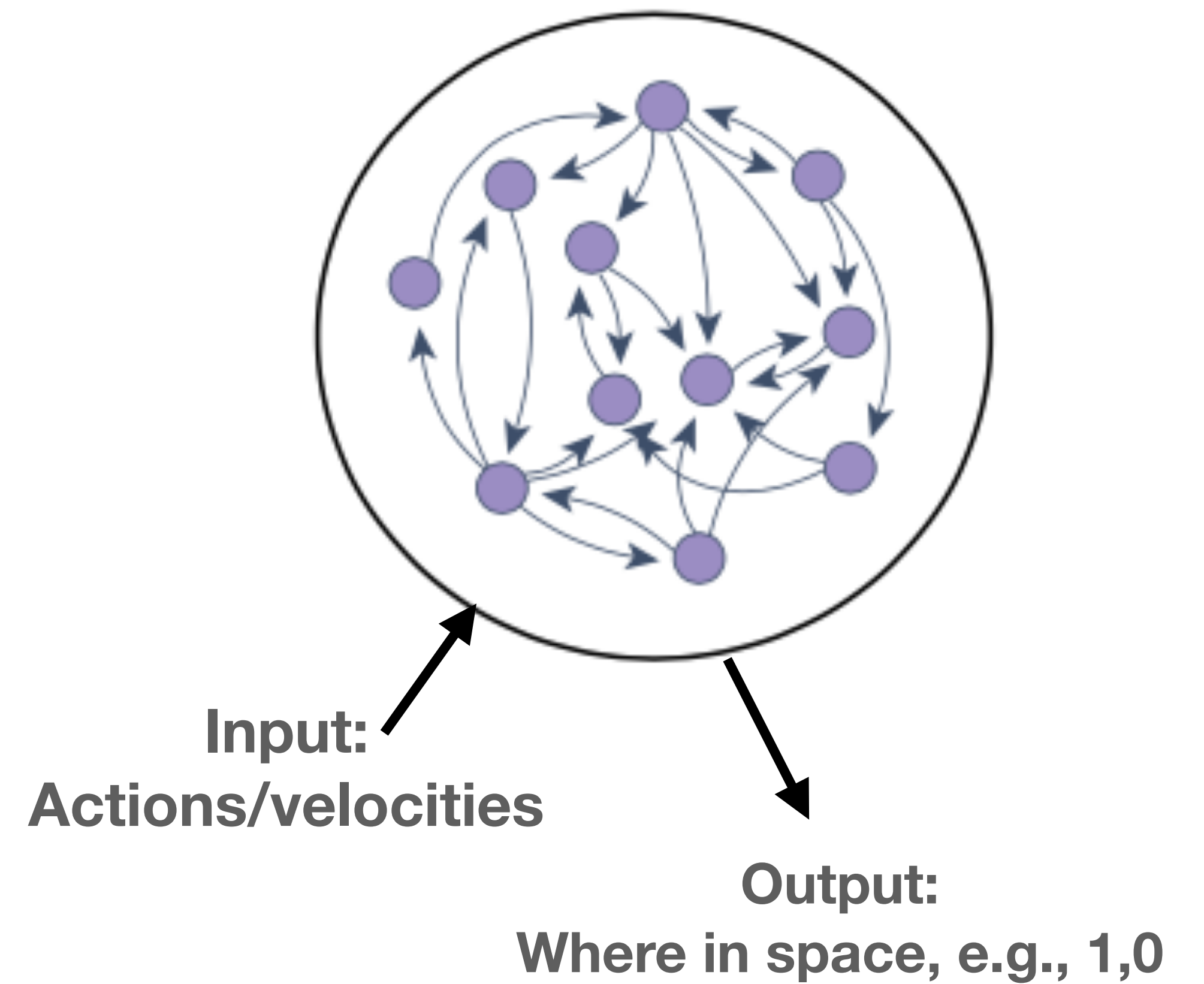


Burak & Feite (2009)

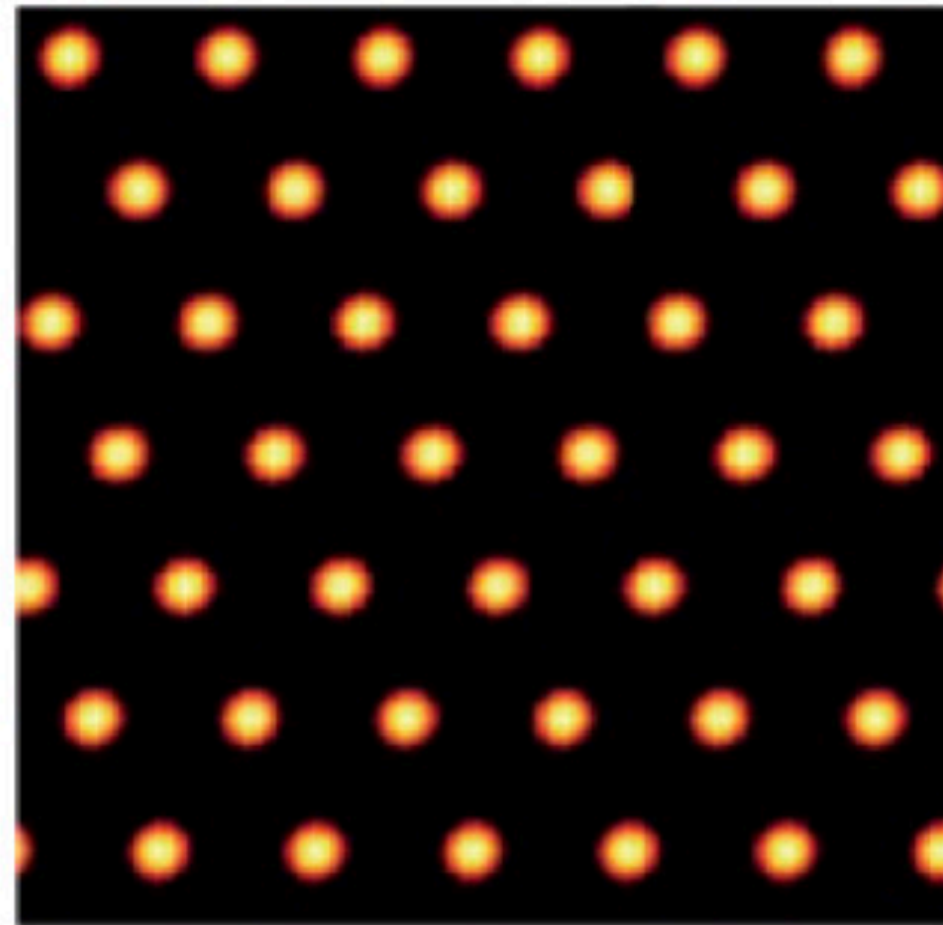


Banino et al., (2018)

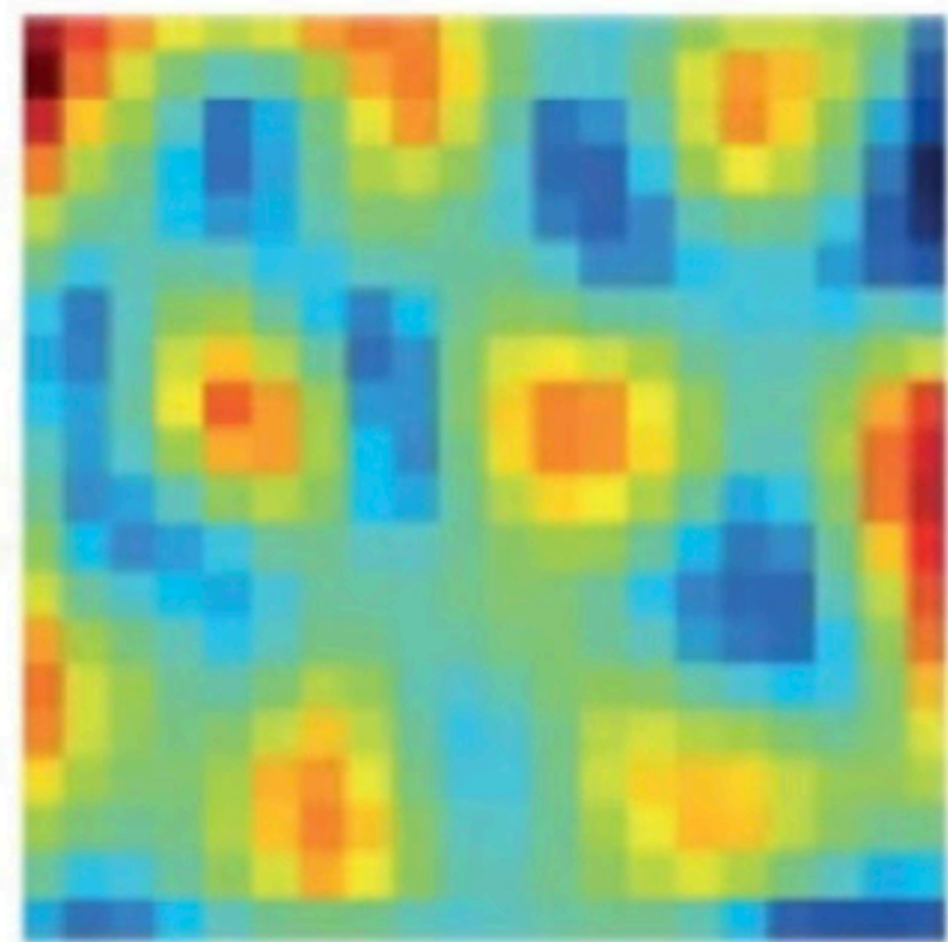
$$\frac{d\bar{s}}{dt} = -\bar{s} + f(W\bar{s} + B\bar{a})$$



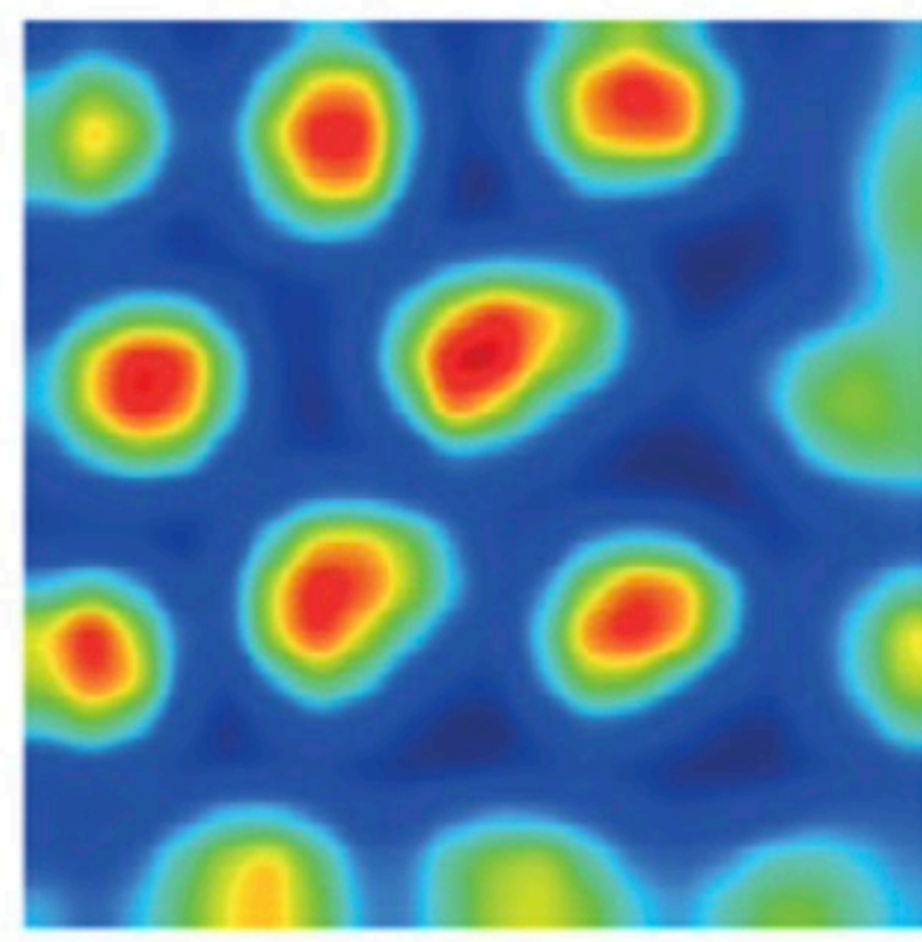
# These networks learn grid cells



Burak & Feite (2009)

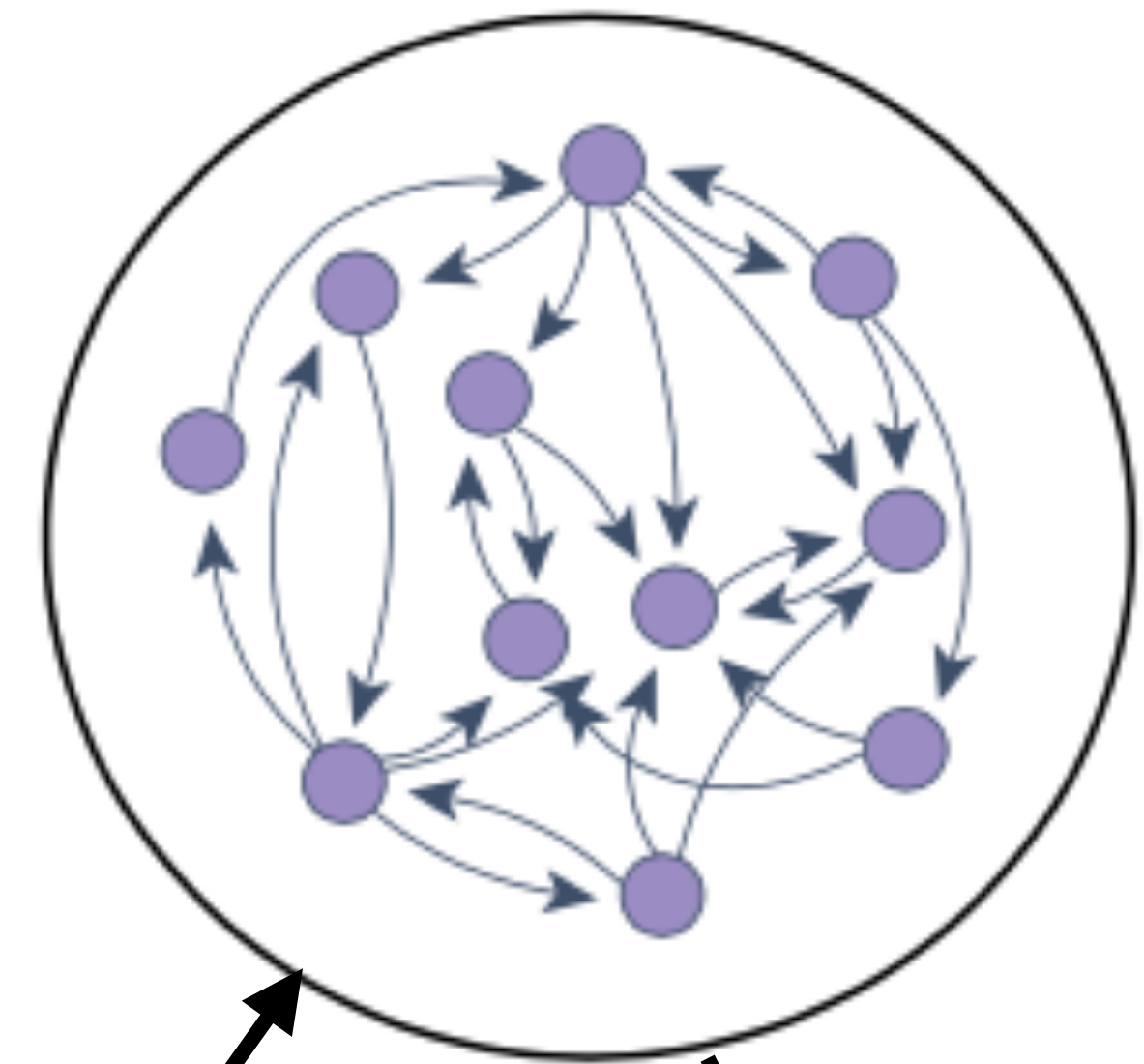


Banino et al., (2018)



Sorscher et al., (2023)

$$\frac{d\bar{s}}{dt} = -\bar{s} + f(W\bar{s} + B\bar{a})$$



**Input:**  
Actions/velocities

**Output:**  
Where in space, e.g., 1,0



**But these sequences are not what the brain sees**

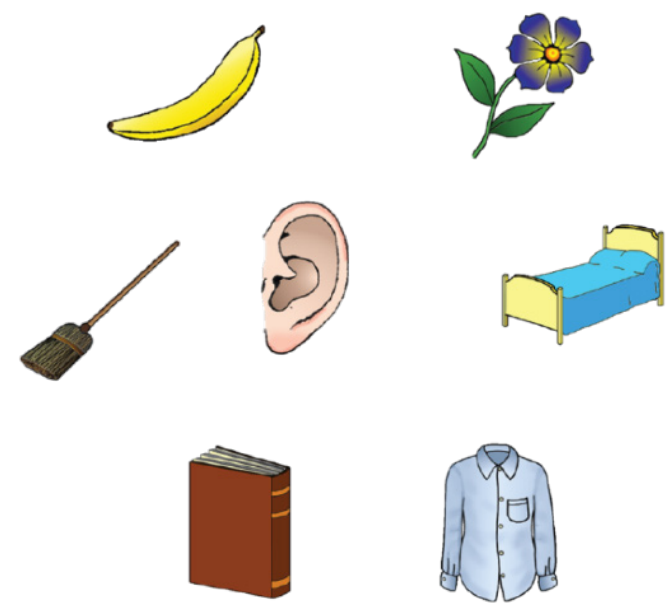


**The brain sees sequences of external observations**

# The brain sees sequences of external observations

Given a stream of sensory stimuli and relations

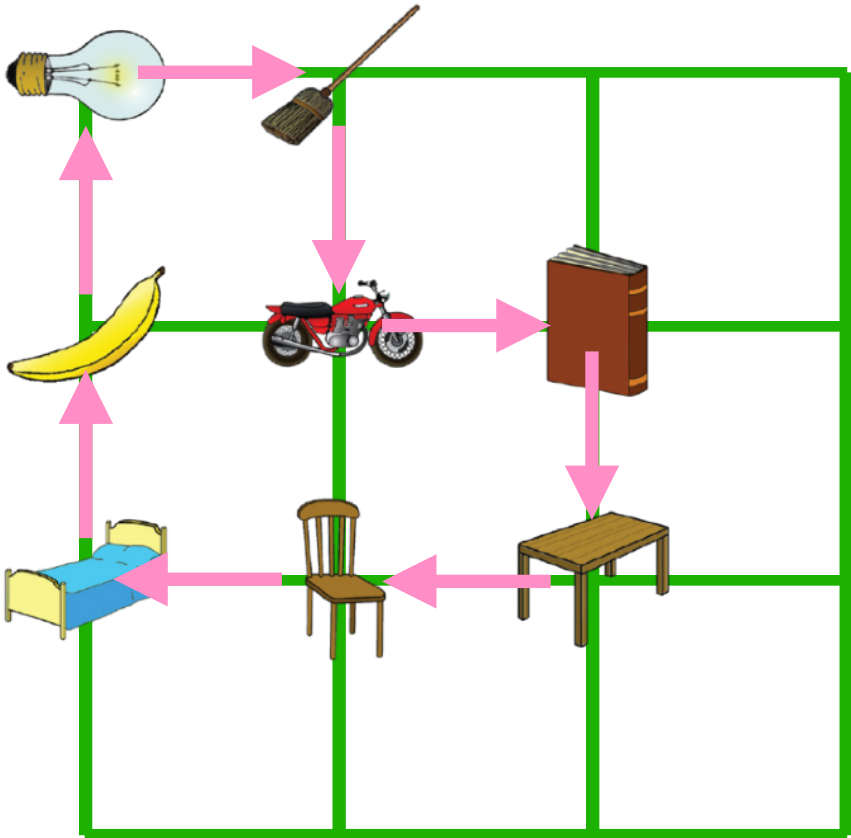
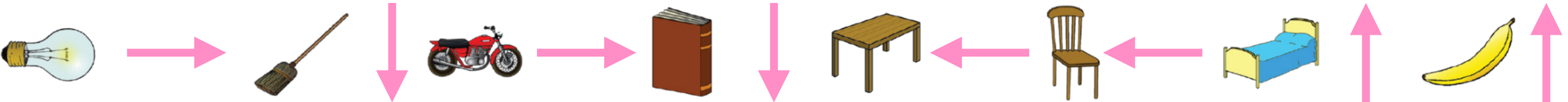
?




# The brain sees sequences of external observations

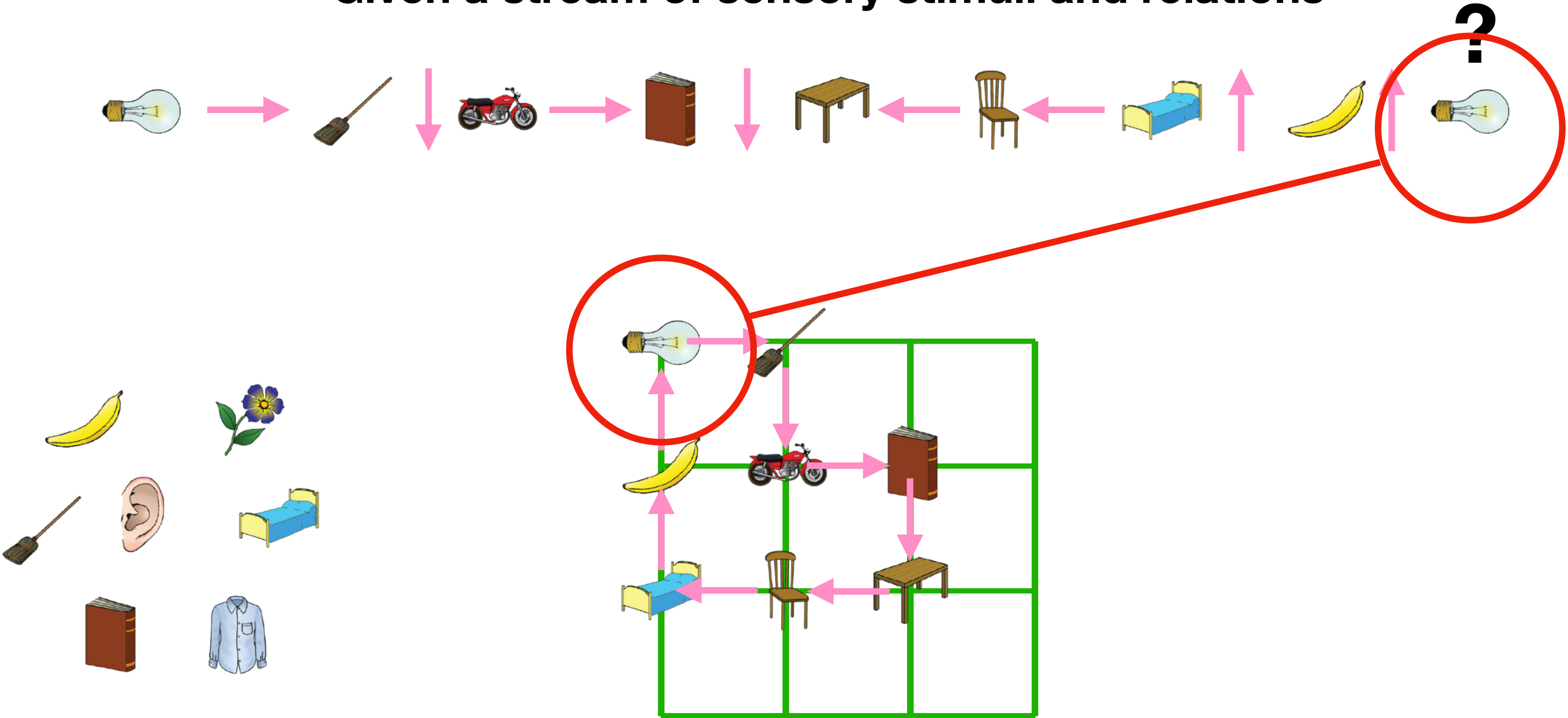
Given a stream of sensory stimuli and relations

?



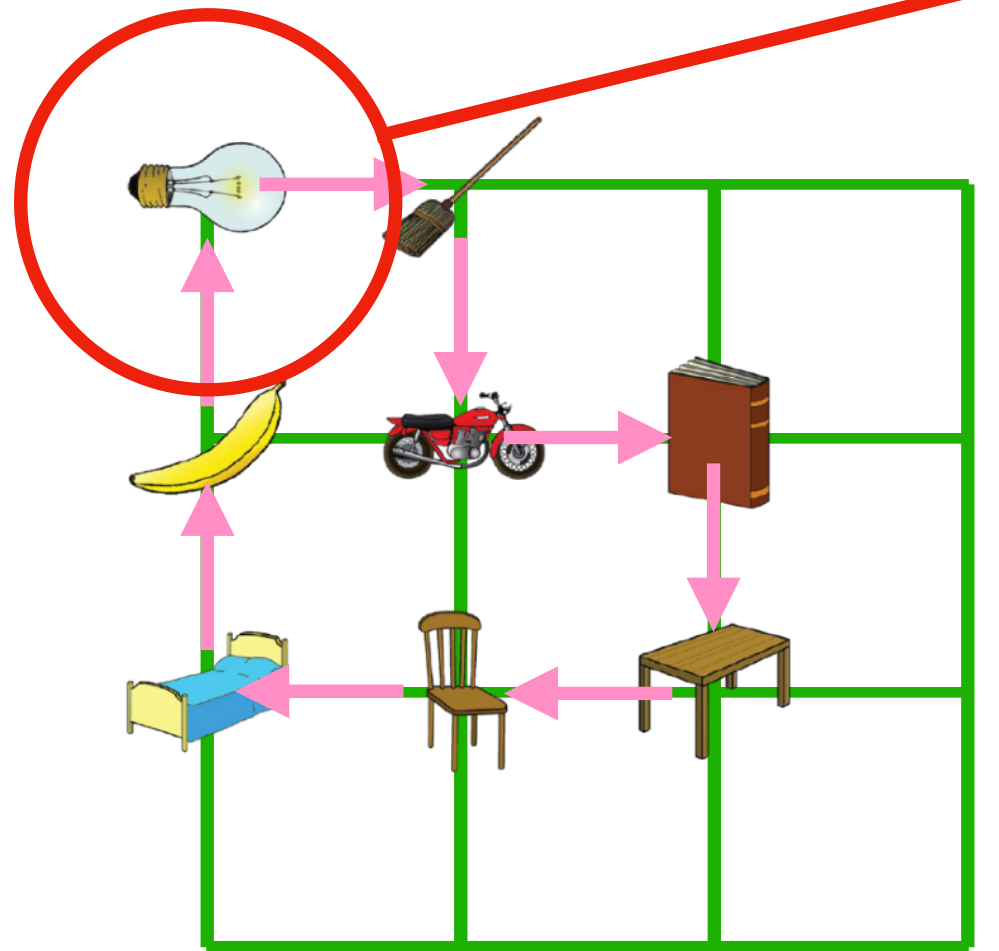
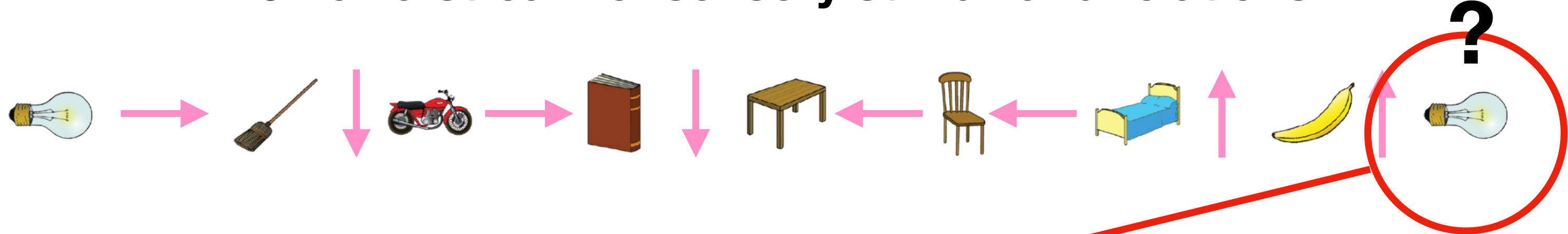
# The brain sees sequences of external observations

Given a stream of sensory stimuli and relations



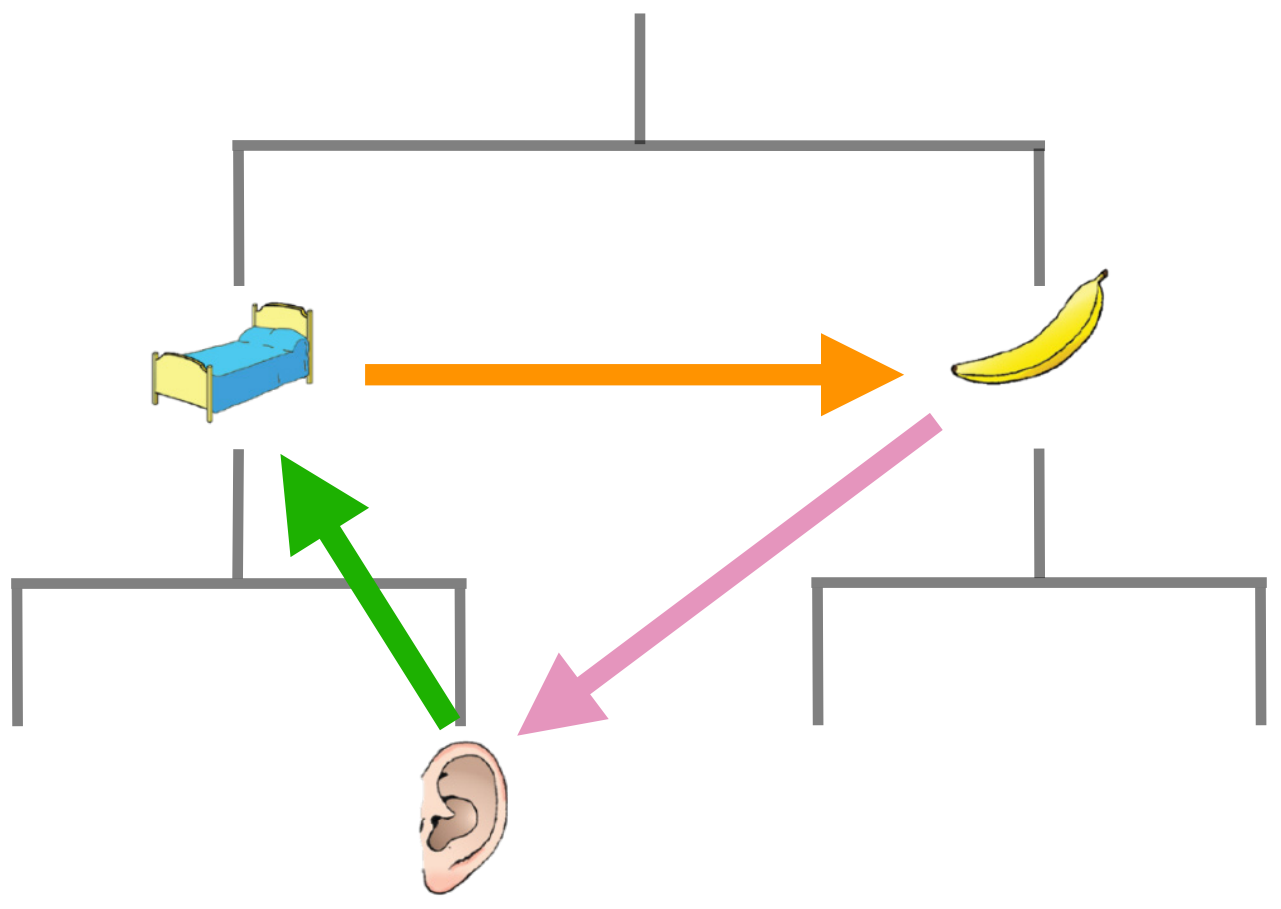
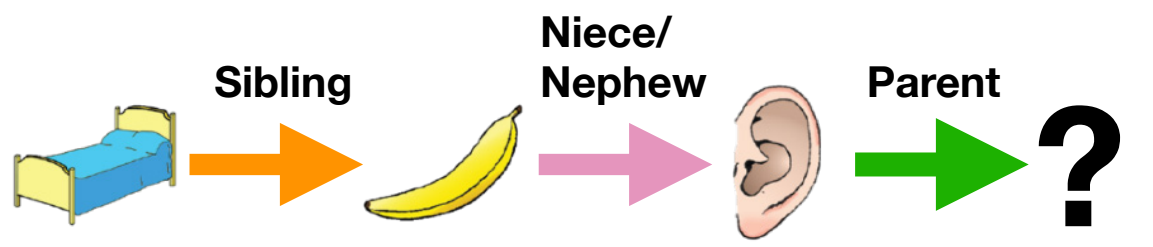
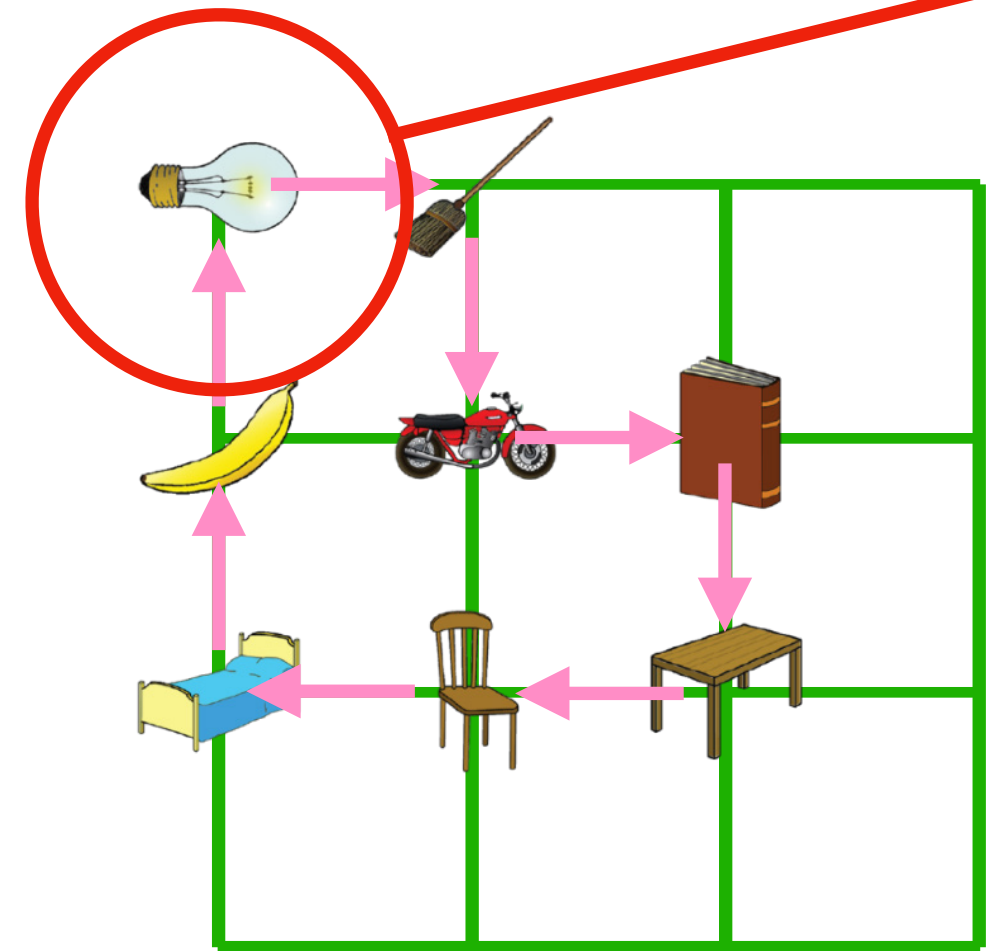
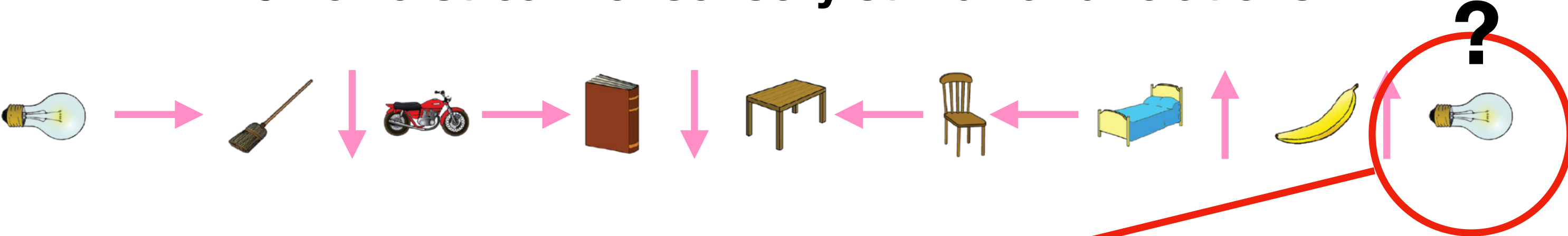
# The brain sees sequences of external observations

Given a stream of sensory stimuli and relations



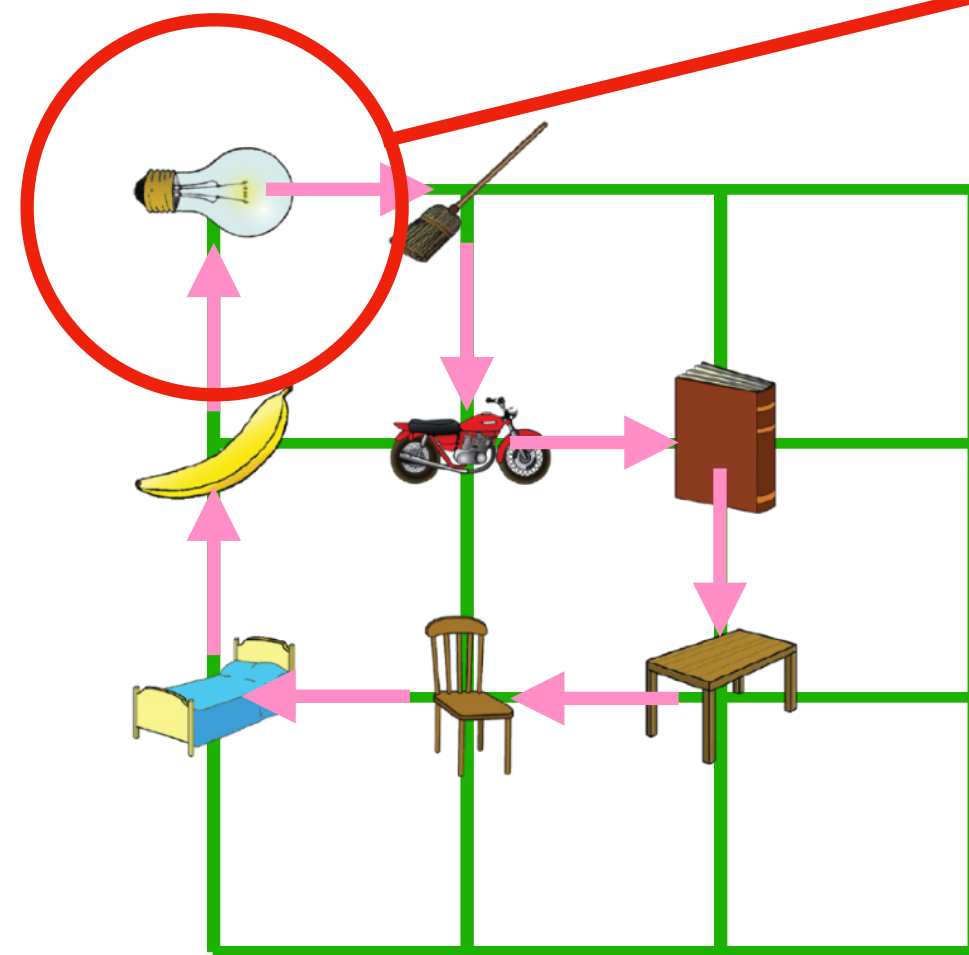
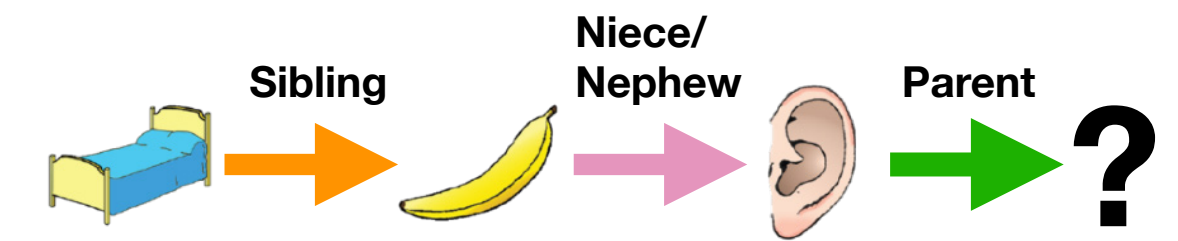
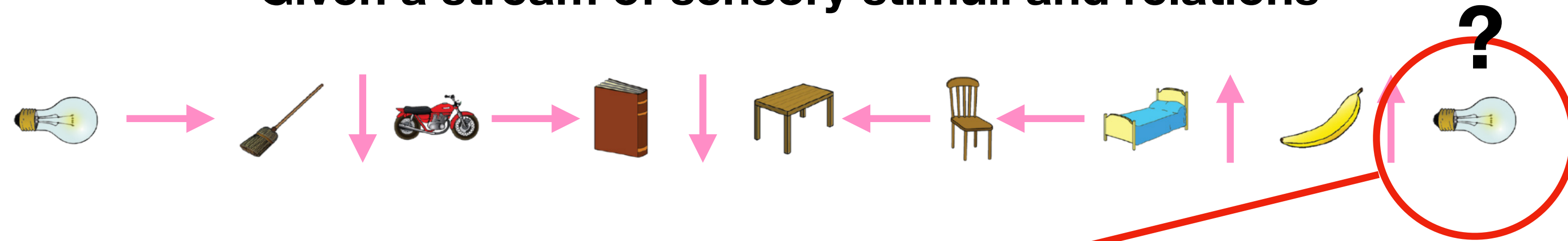
# The brain sees sequences of external observations

Given a stream of sensory stimuli and relations

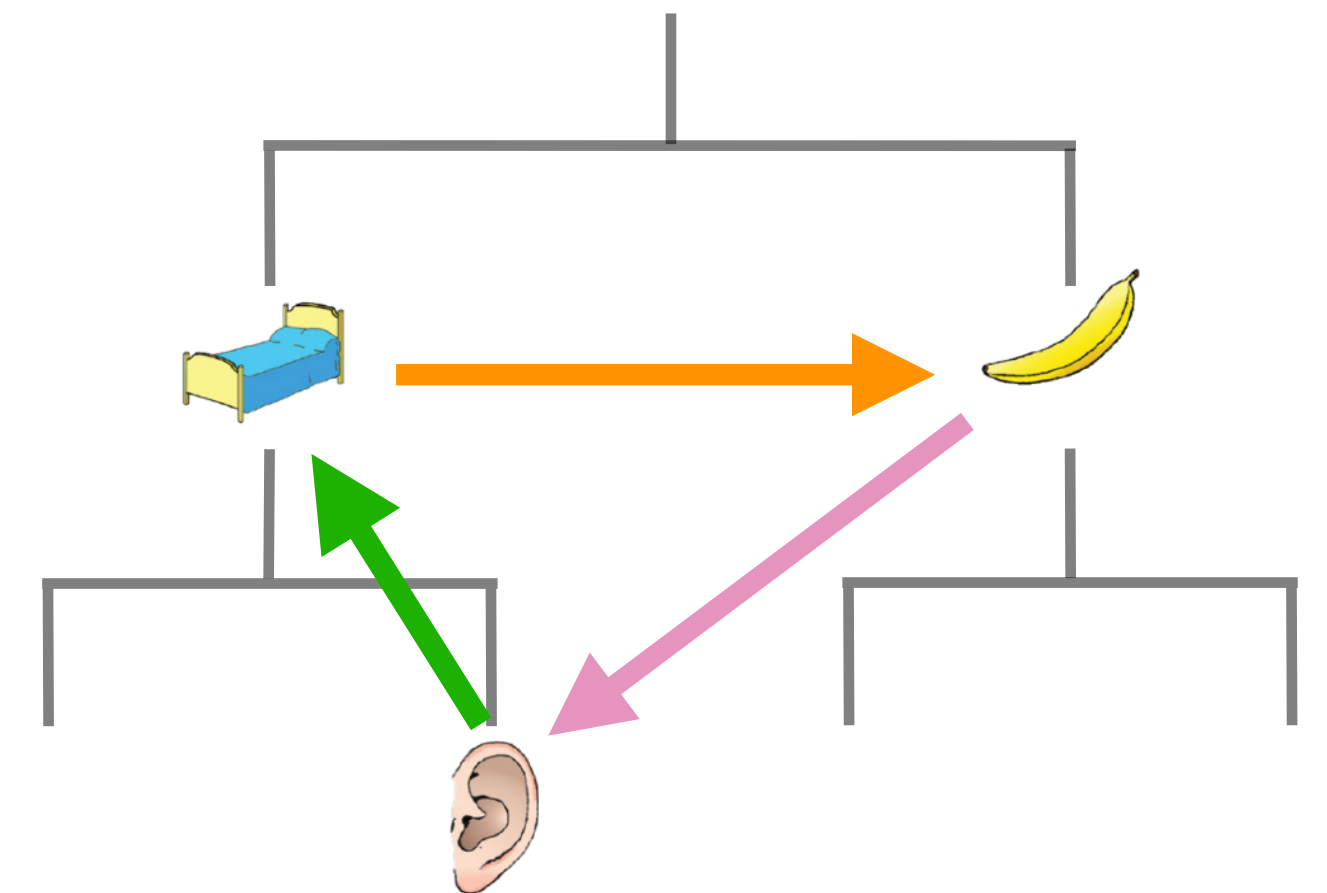


# The brain sees sequences of external observations

Given a stream of sensory stimuli and relations



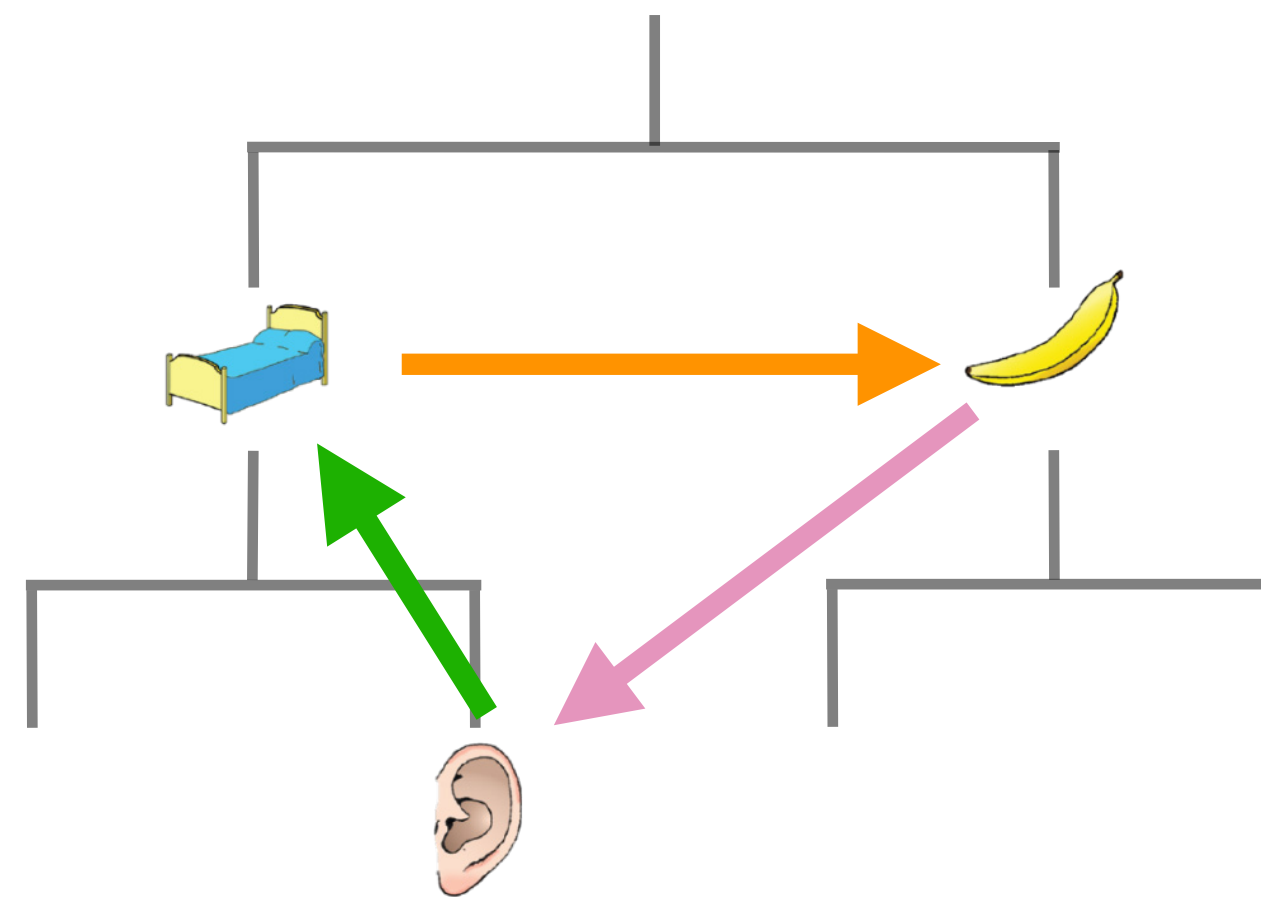
Knowing the map is only way to correctly predict these shortcut inferences zero shot



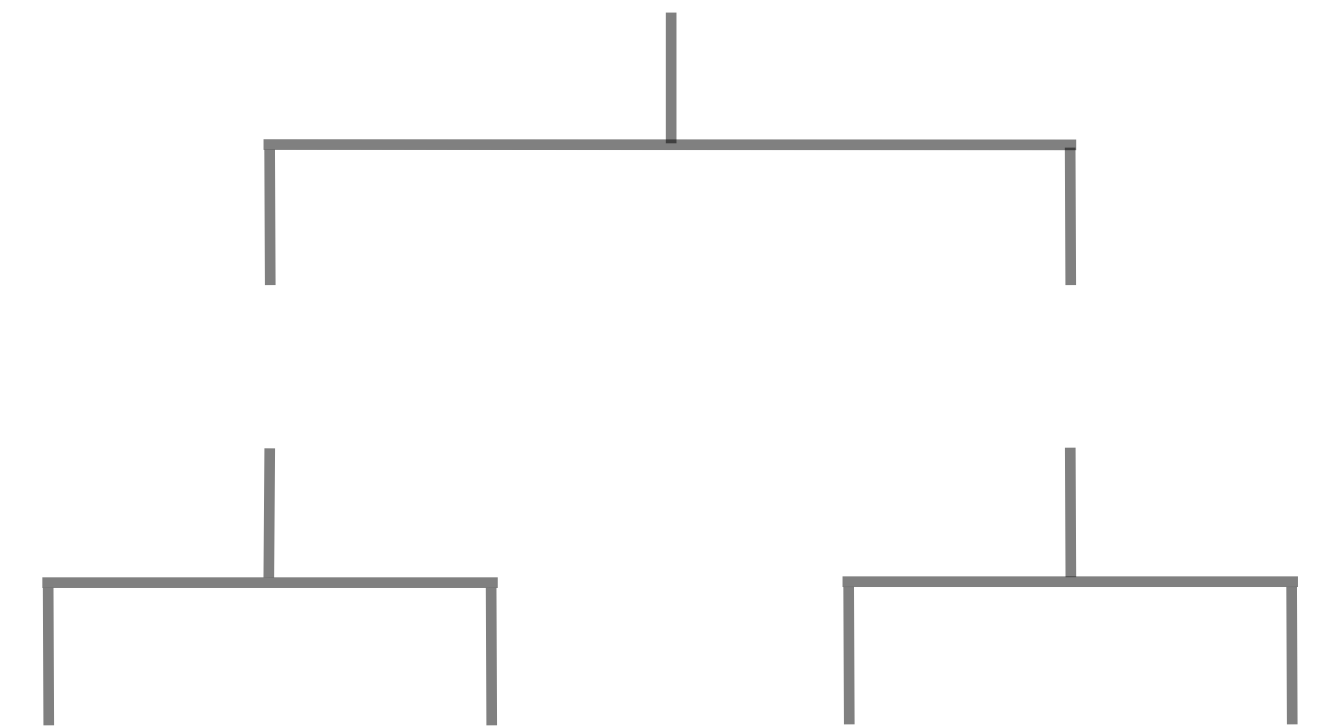
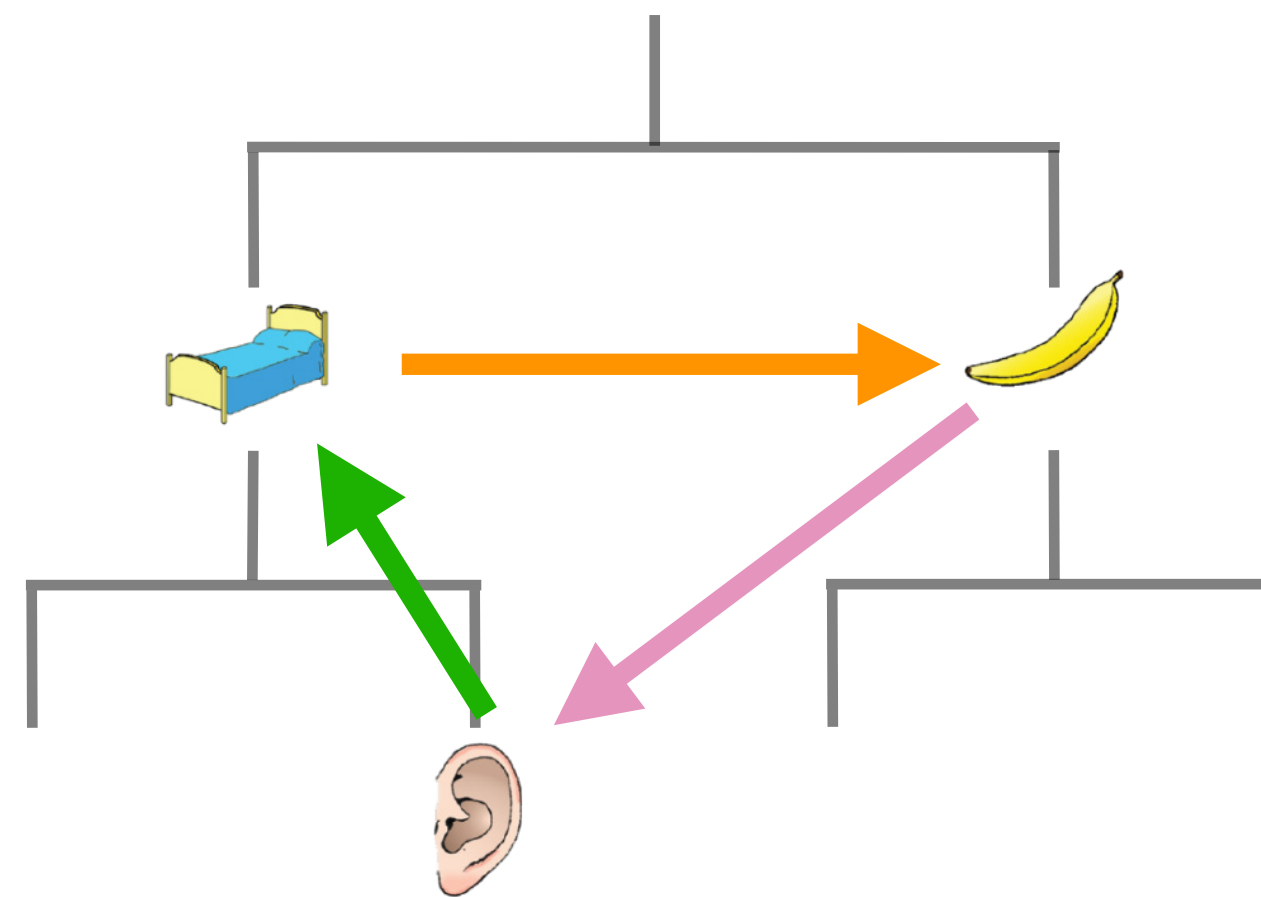
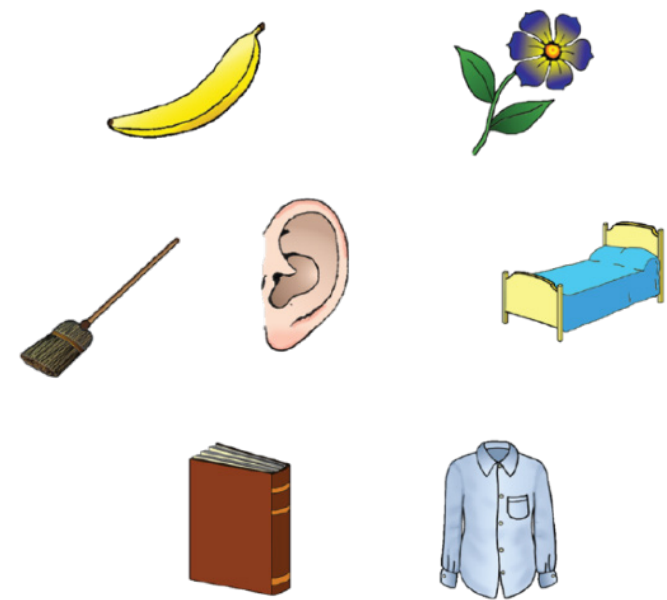


# **Representing abstract location and making conjunctive memories**

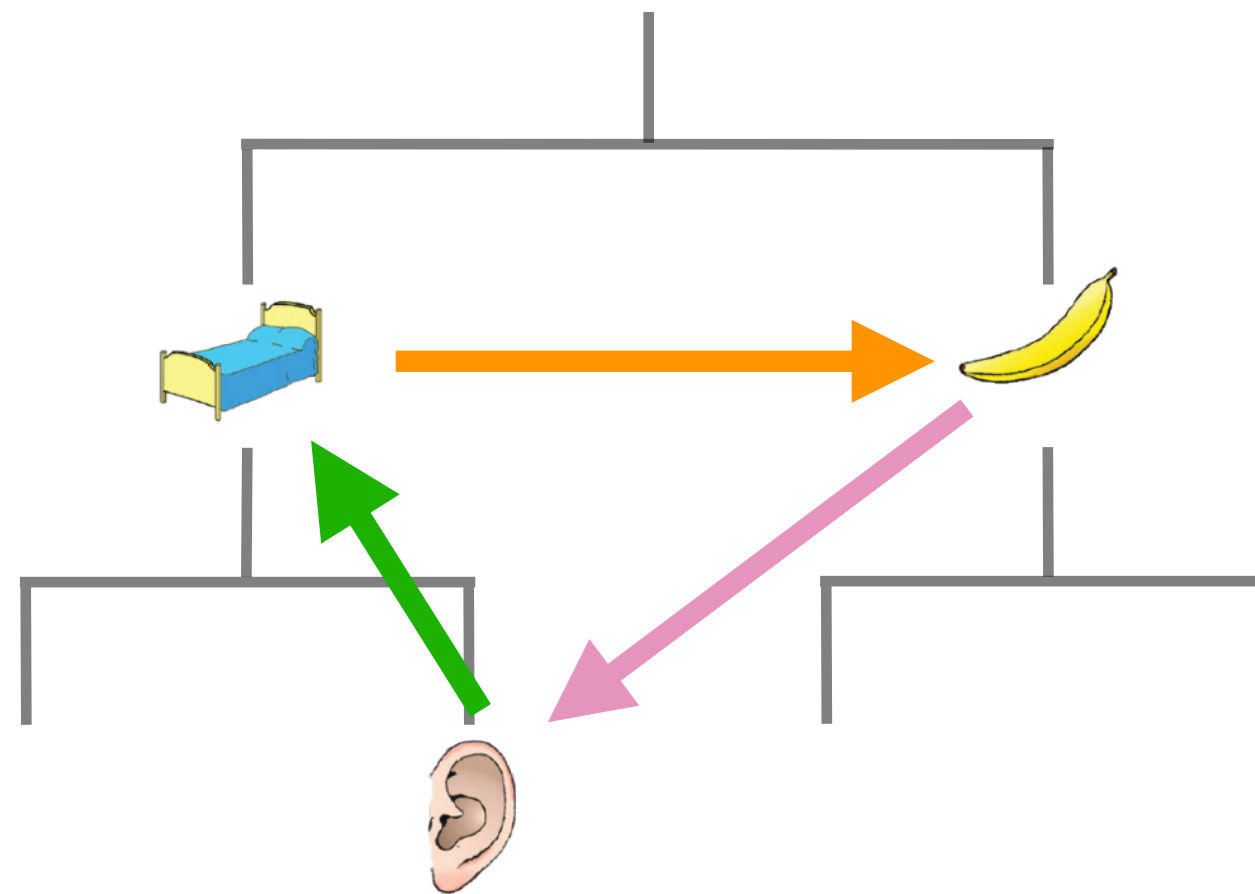
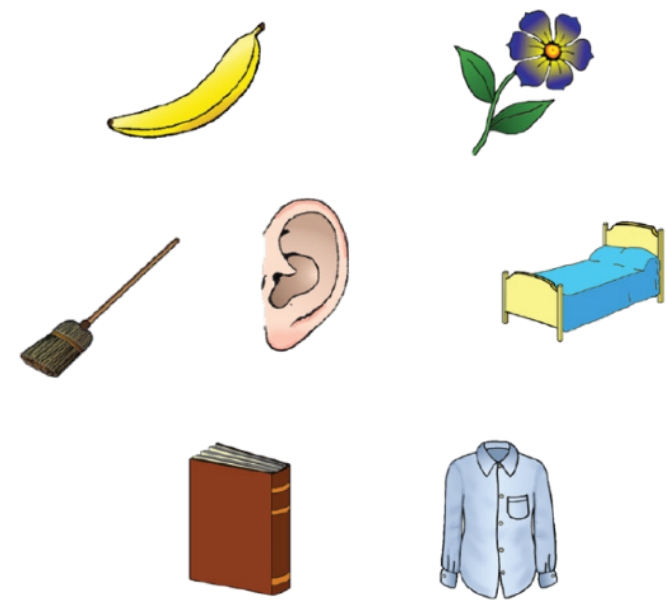
# Representing abstract location and making conjunctive memories



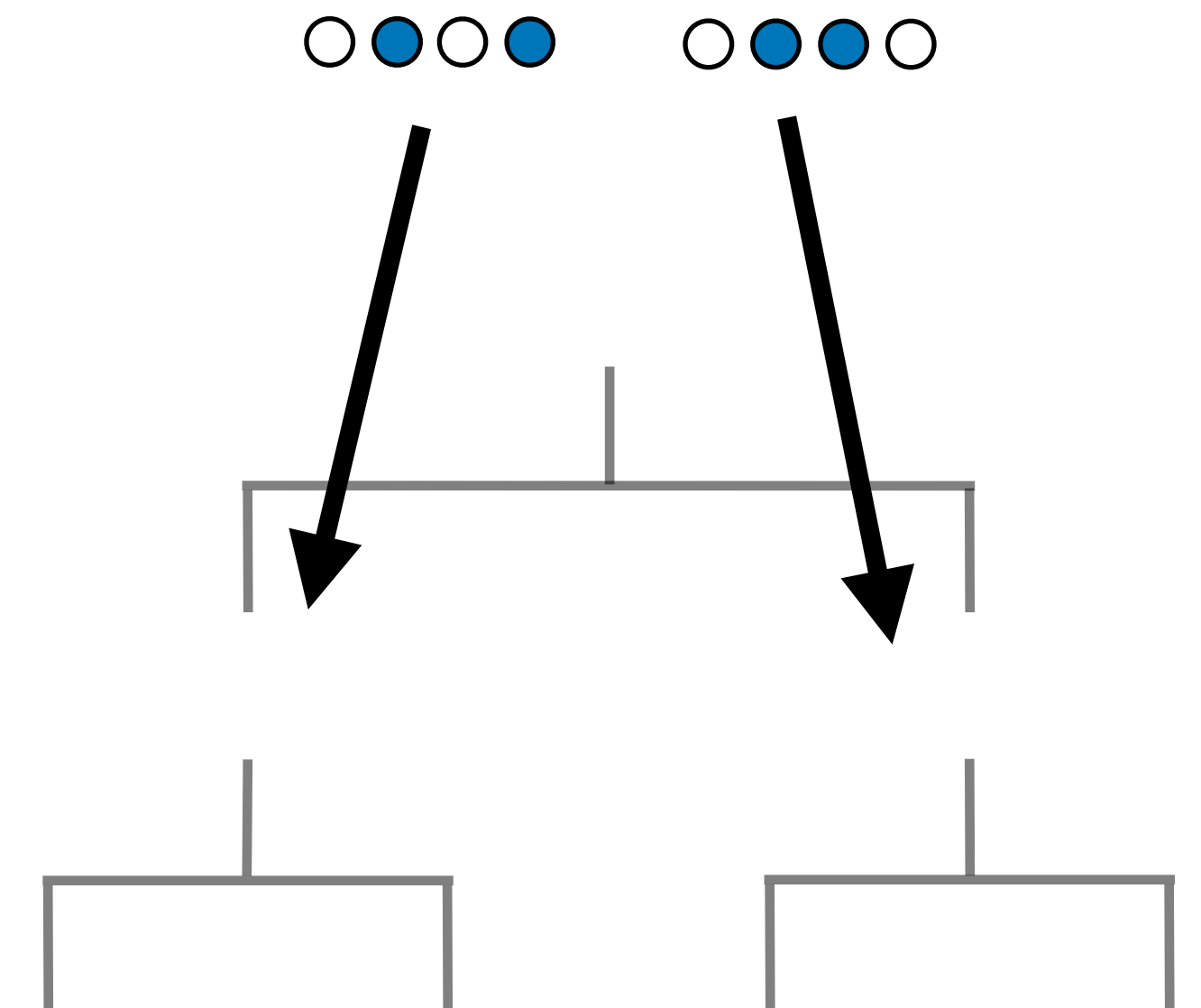
# Representing abstract location and making conjunctive memories



# Representing abstract location and making conjunctive memories



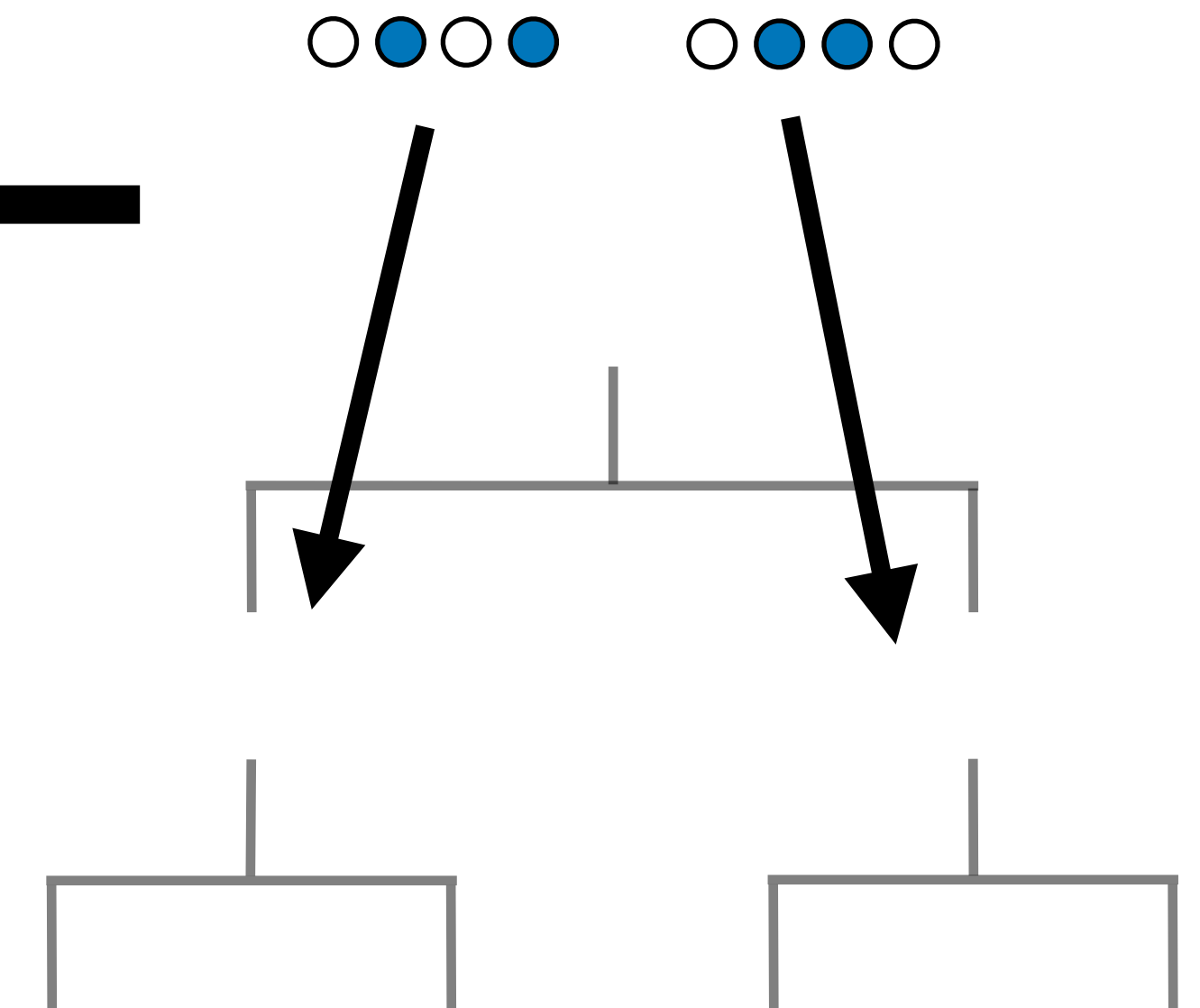
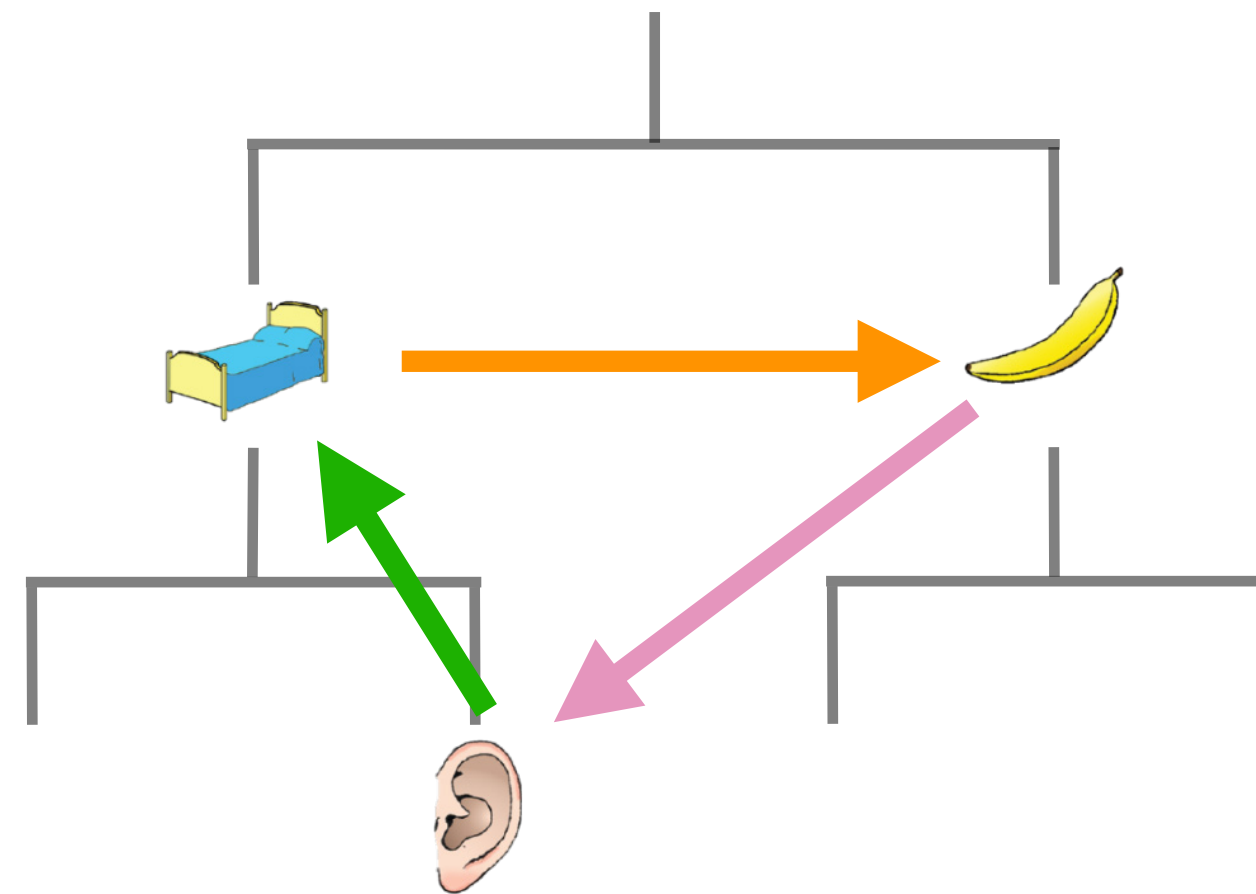
Represent **location** in abstract structure



# Representing abstract location and making conjunctive memories

Compose building blocks and form **memories** to know who is where in each family

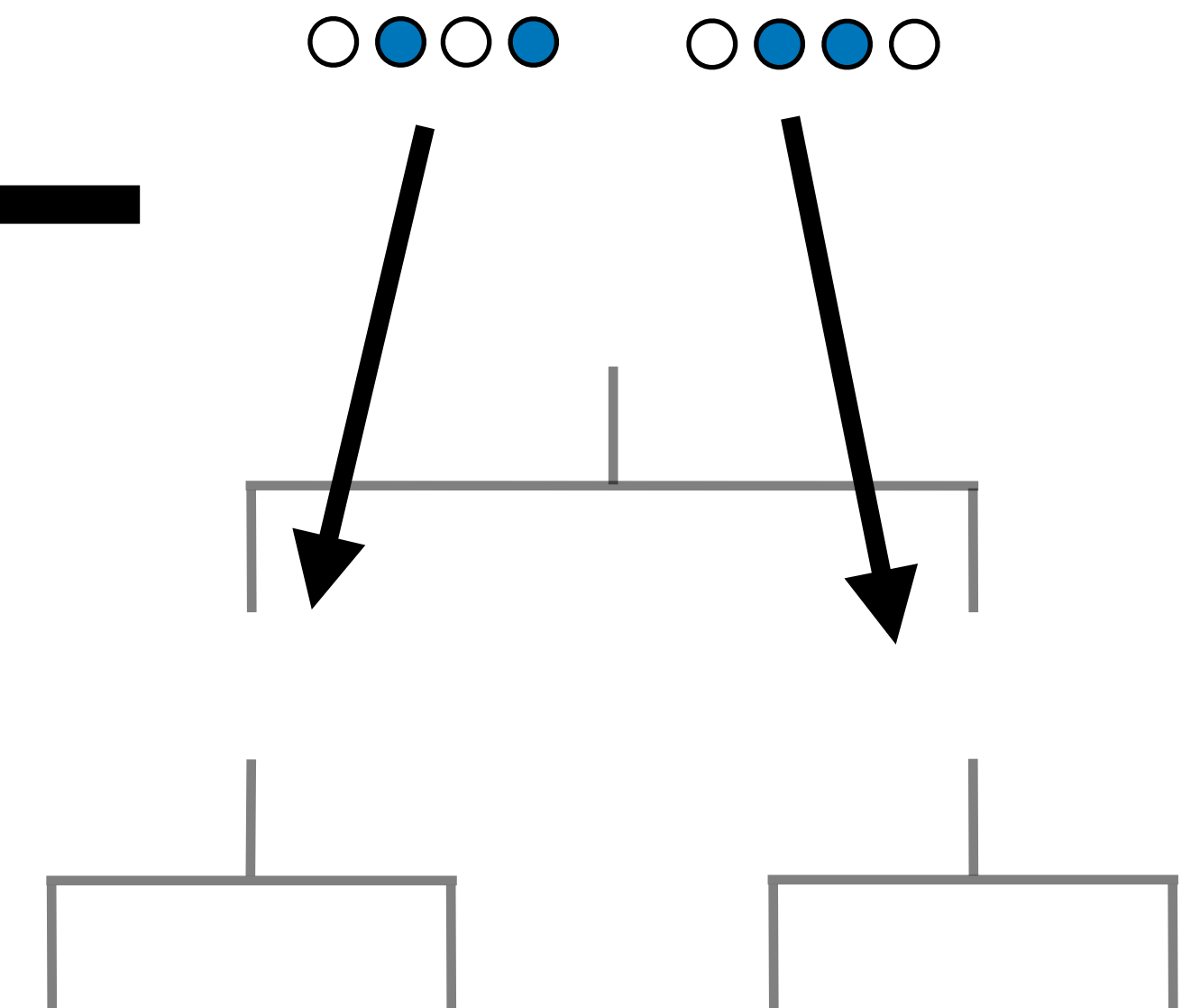
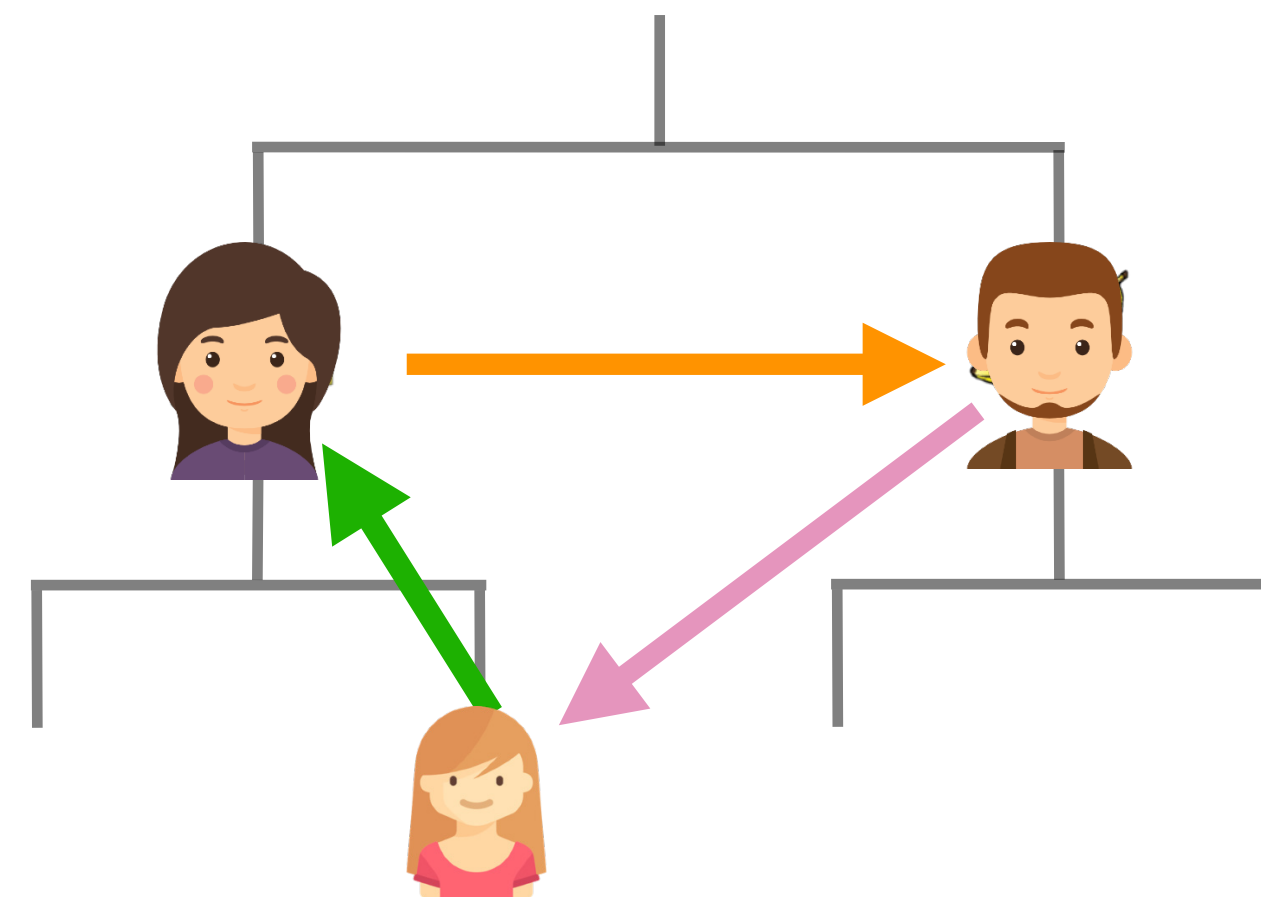
Represent **location** in abstract structure



# Representing abstract location and making conjunctive memories

Compose building blocks and form **memories** to know who is where in each family

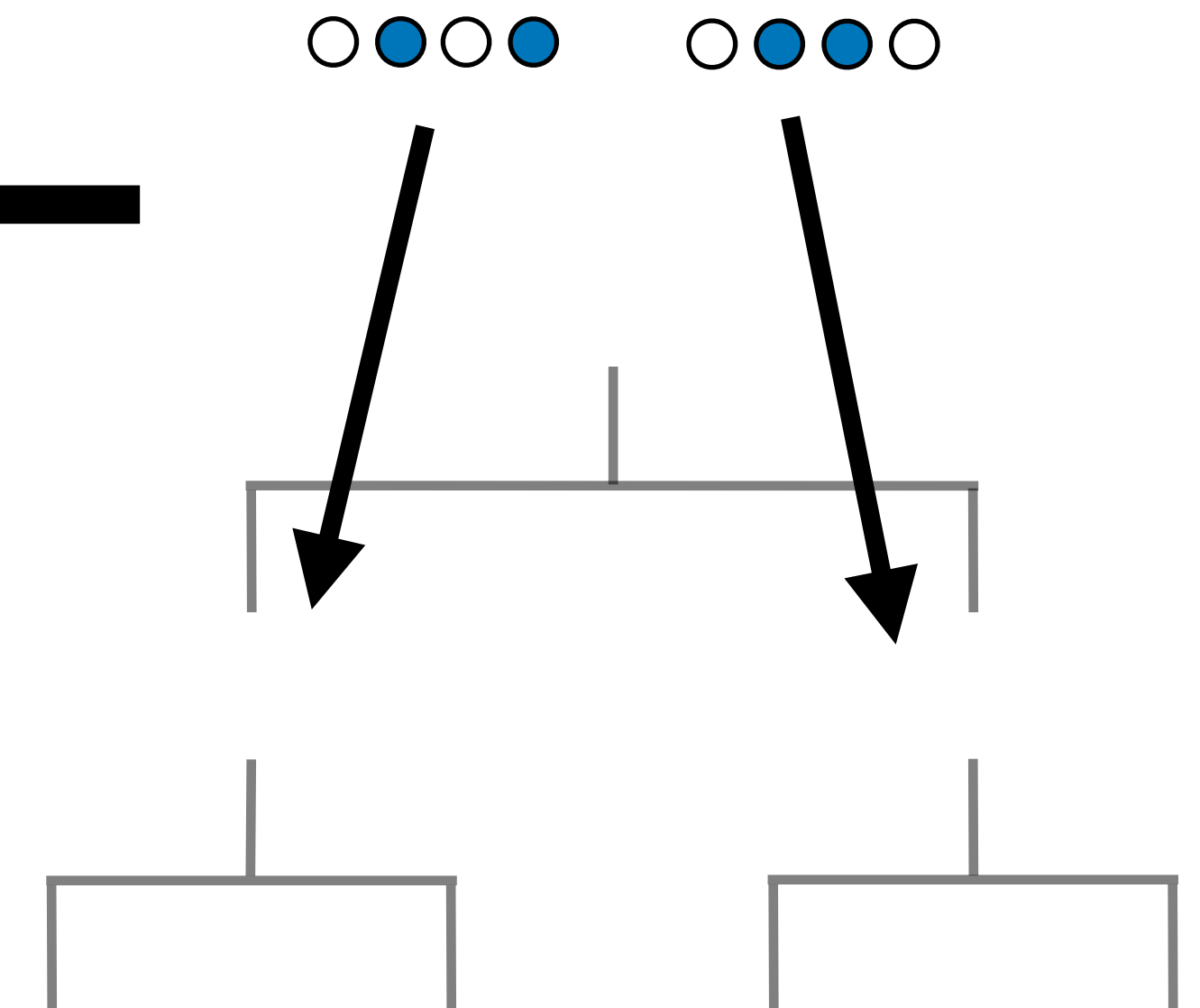
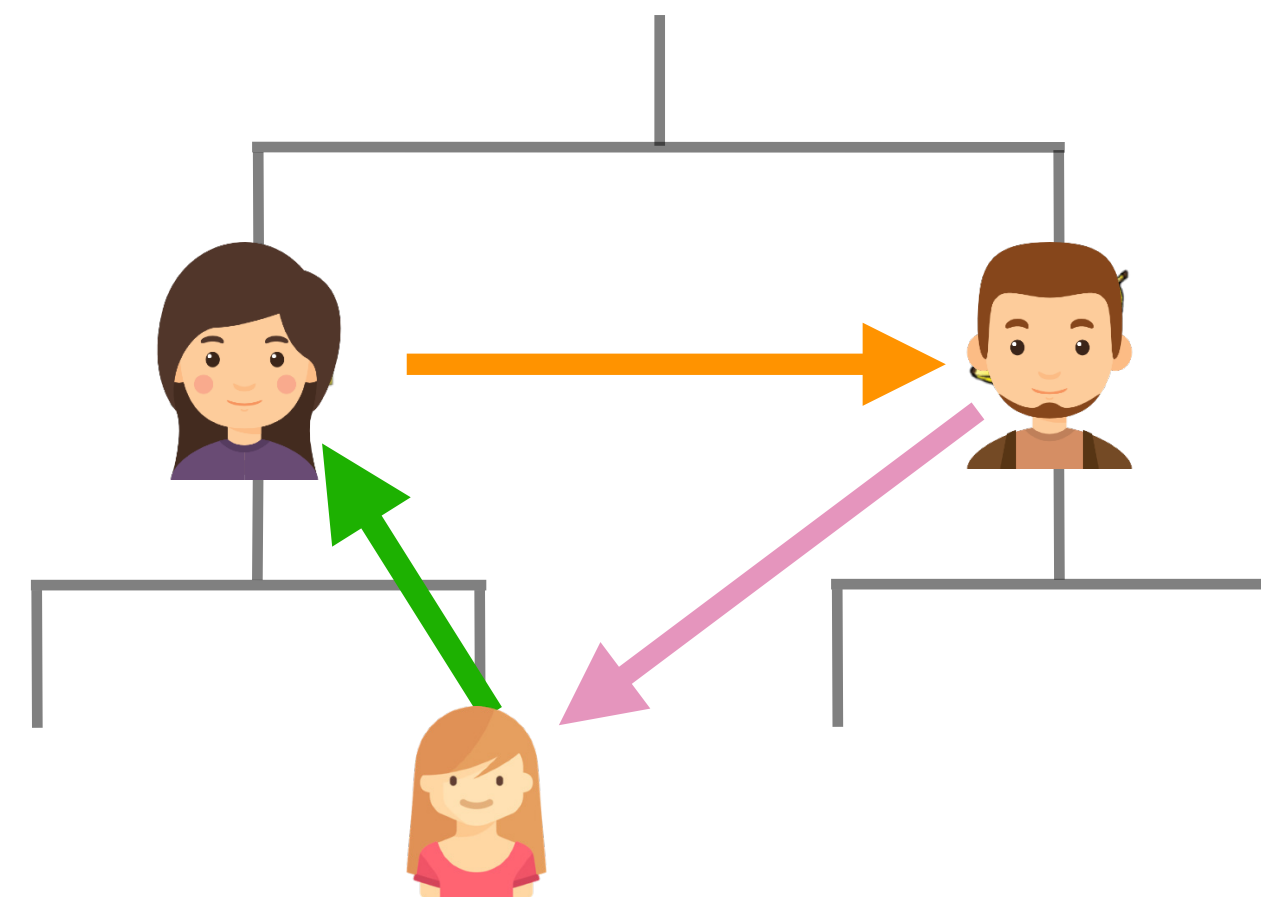
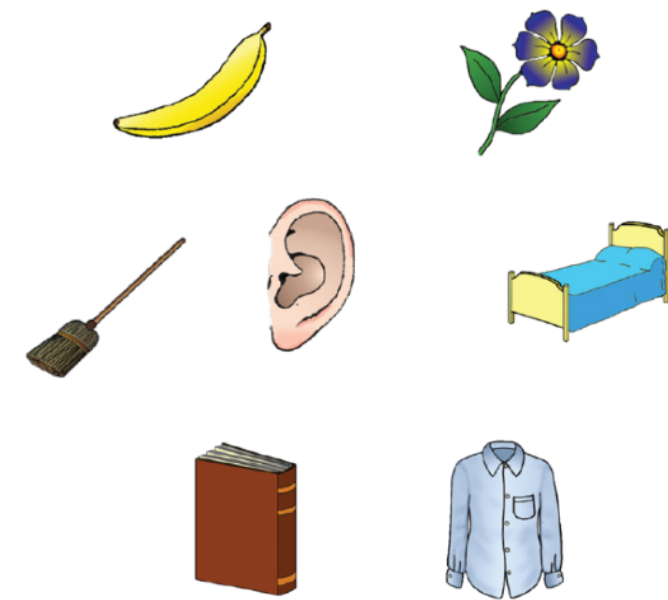
Represent **location** in abstract structure



# Representing abstract location and making conjunctive memories

Compose building blocks and form **memories** to know who is where in each family

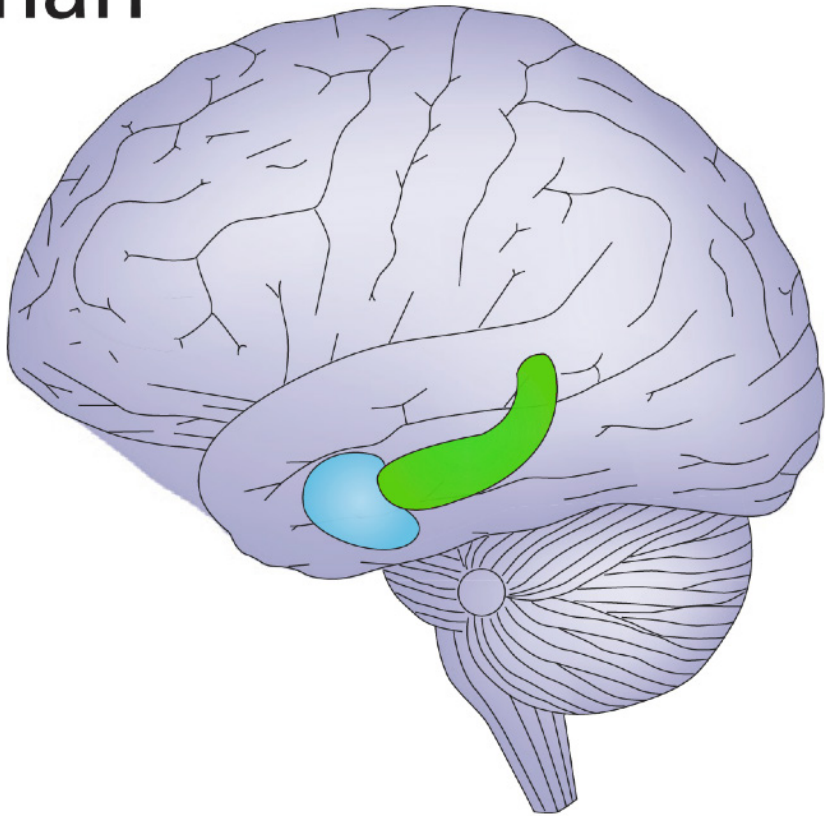
Represent **location** in abstract structure



How do we do this and what on earth will the representations look like?

# TEM learns abstract location representations by prediction errors

Human

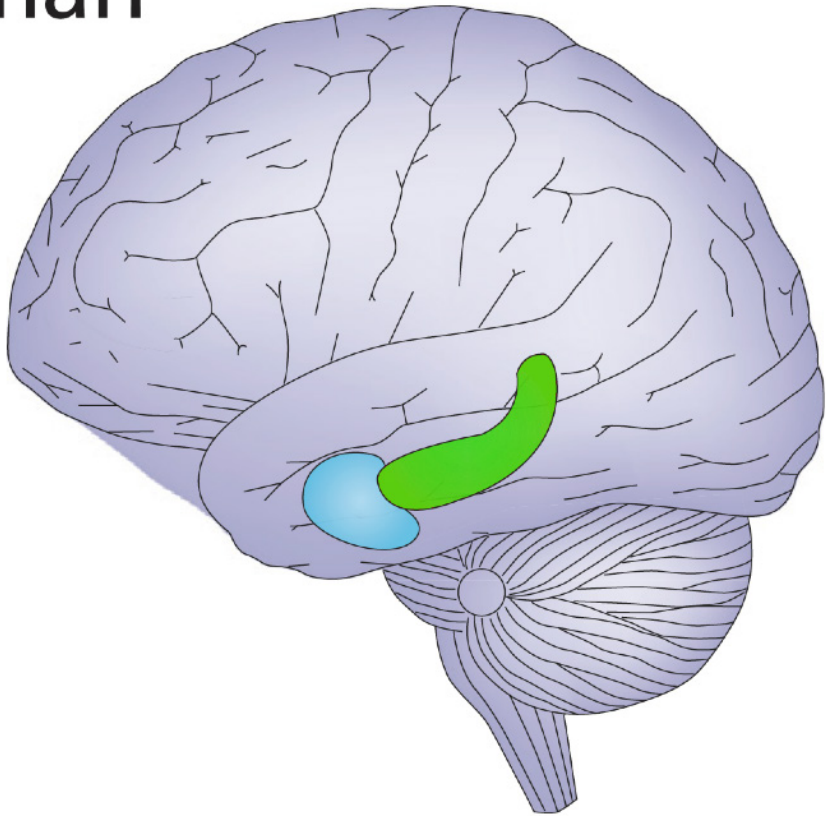


● Hippocampus ● EC

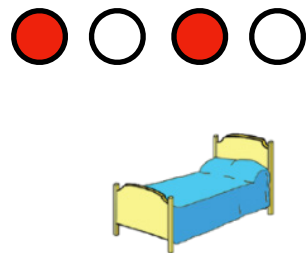


# TEM learns abstract location representations by prediction errors

Human

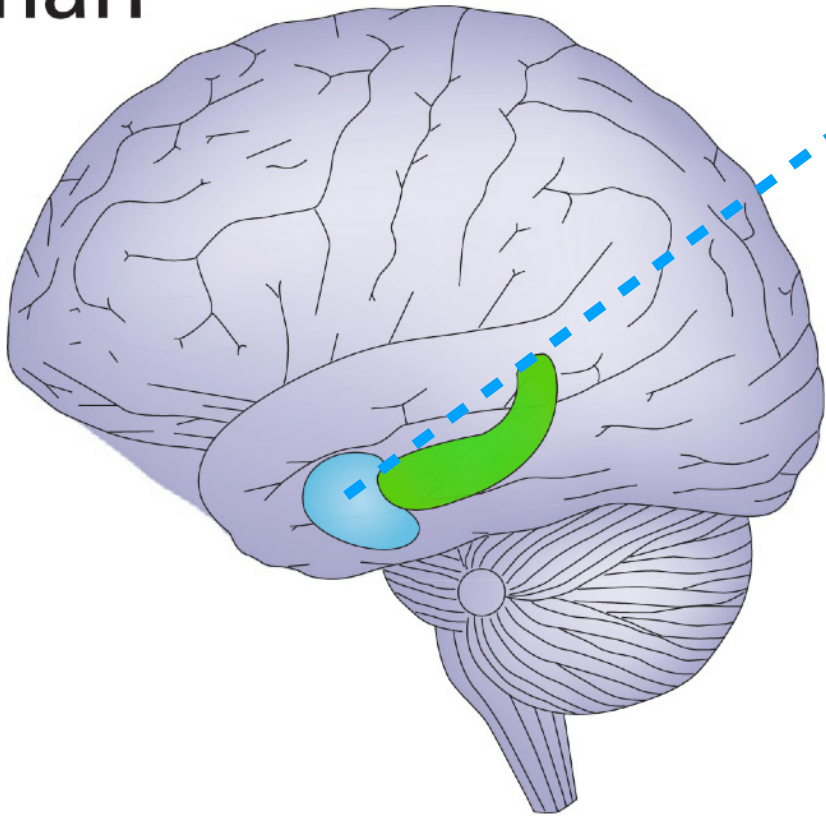


● Hippocampus ● EC



# TEM learns abstract location representations by prediction errors

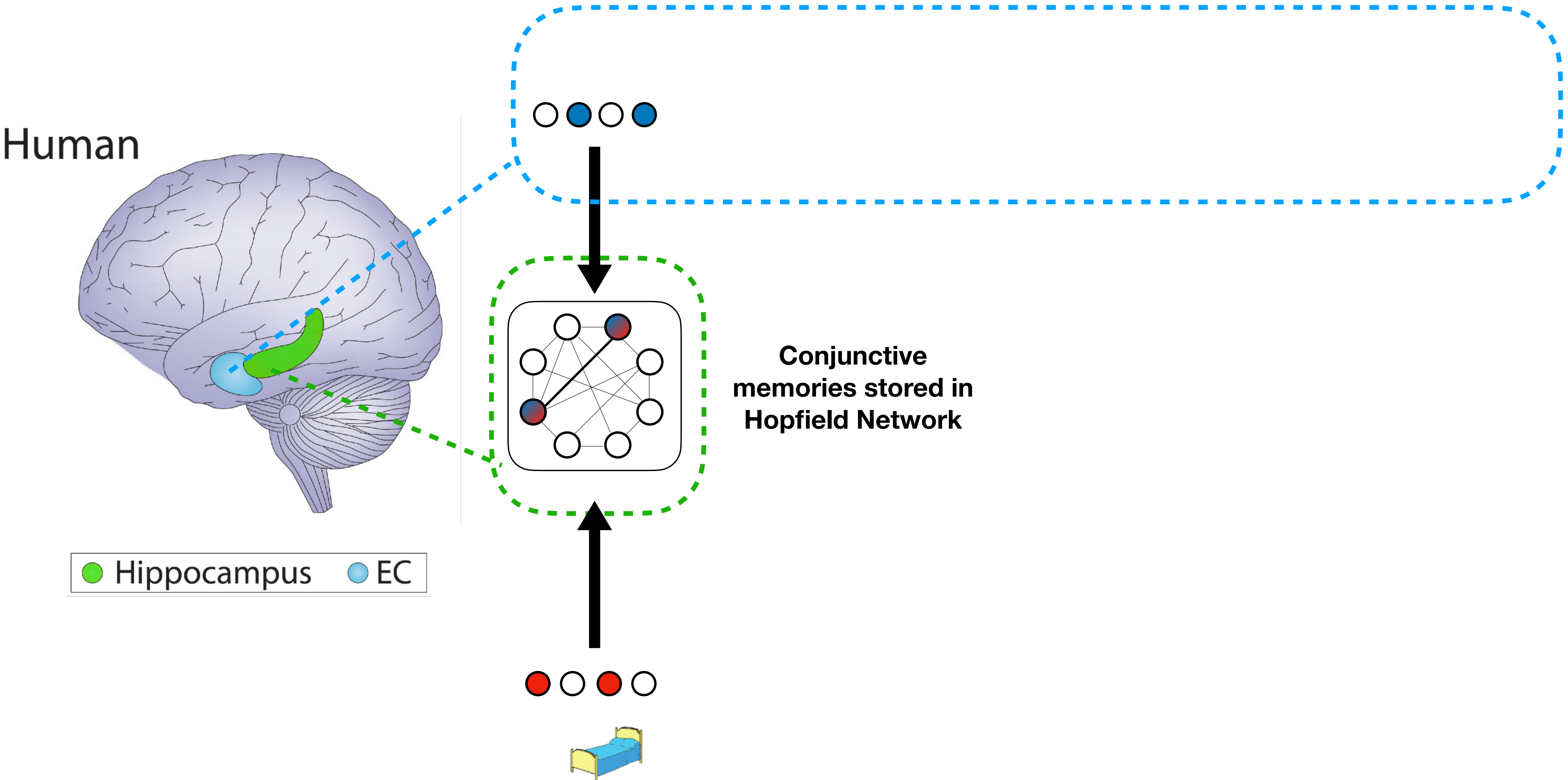
Human



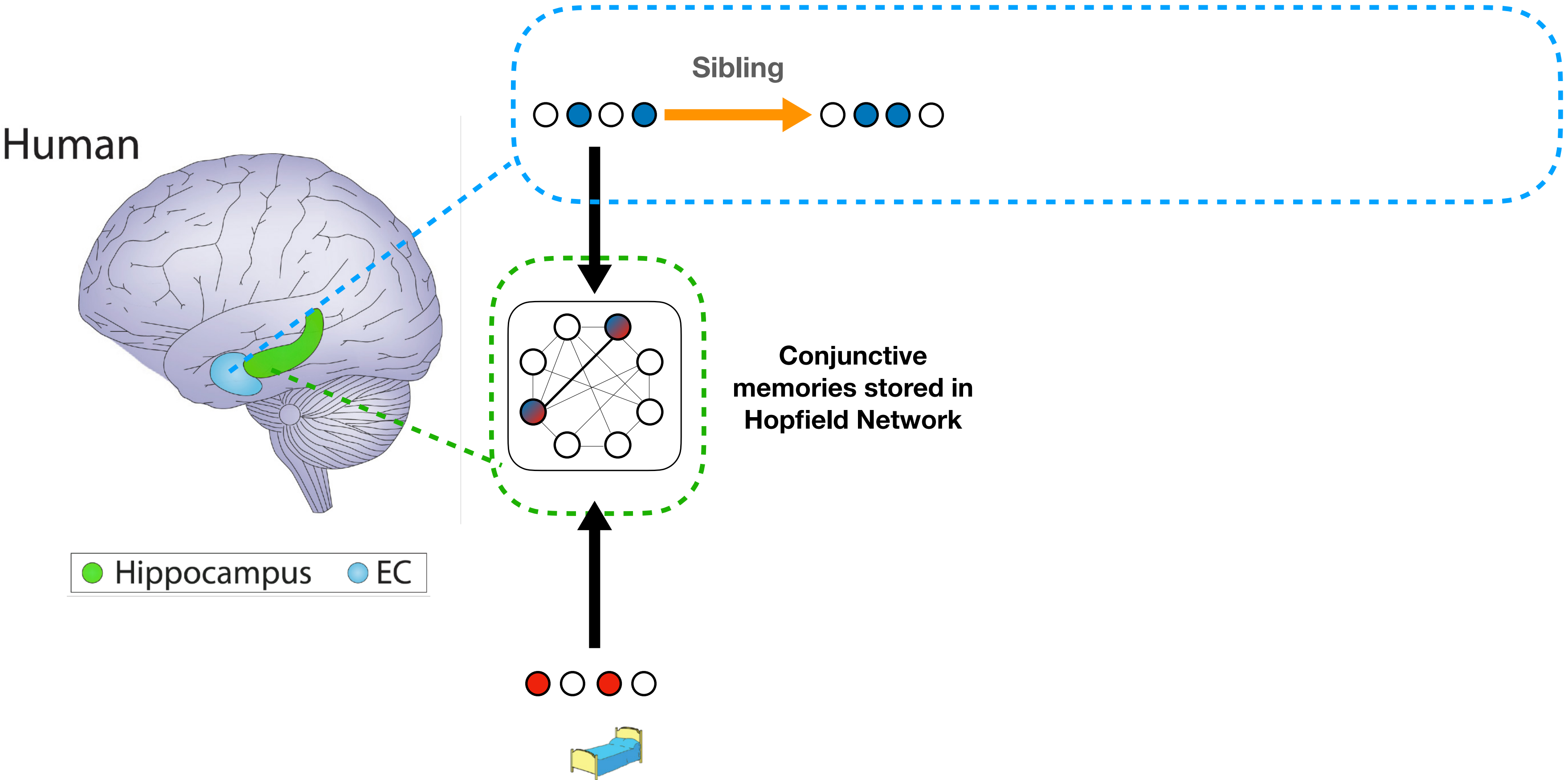
● Hippocampus ● EC



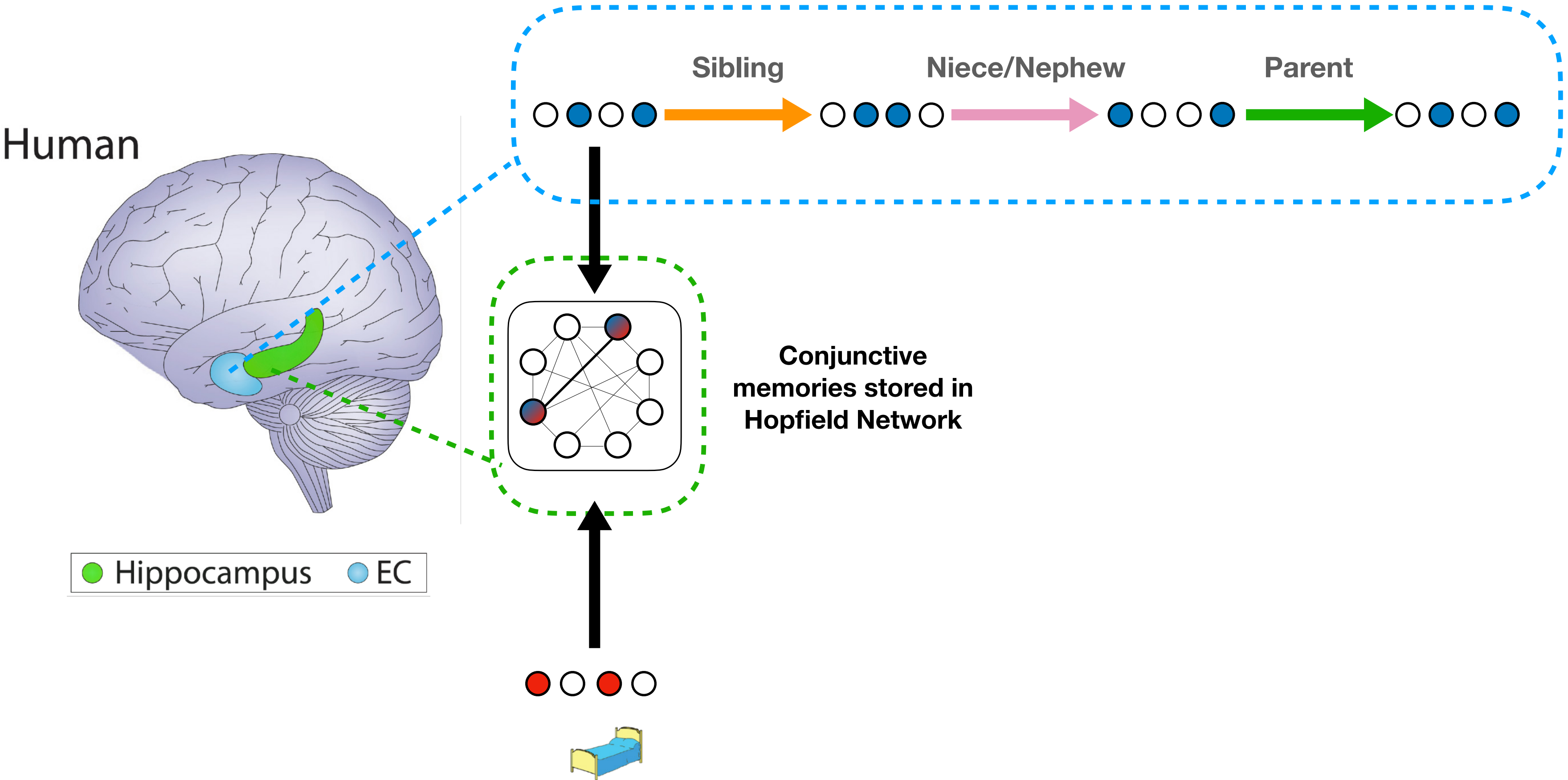
# TEM learns abstract location representations by prediction errors



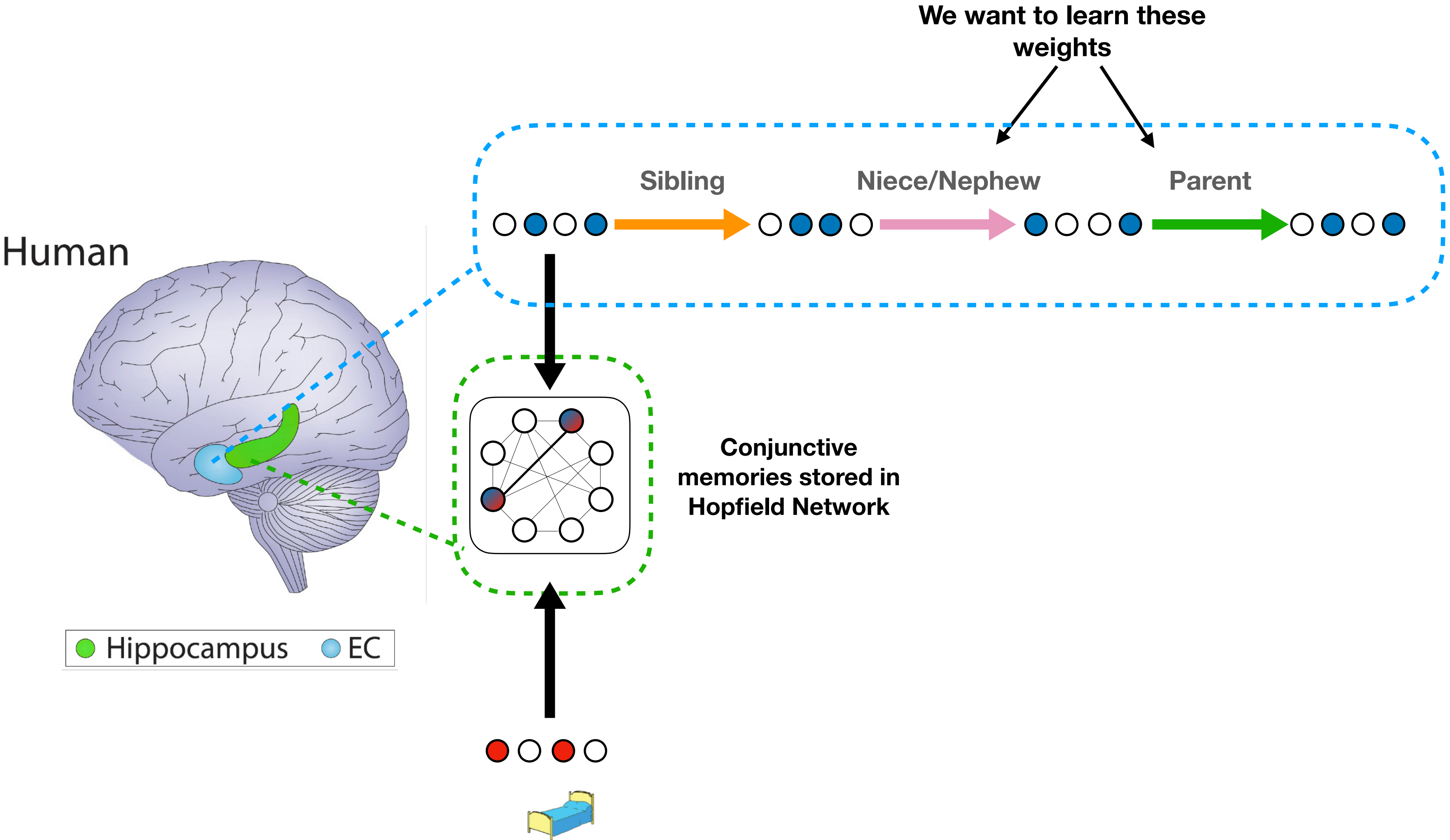
# TEM learns abstract location representations by prediction errors



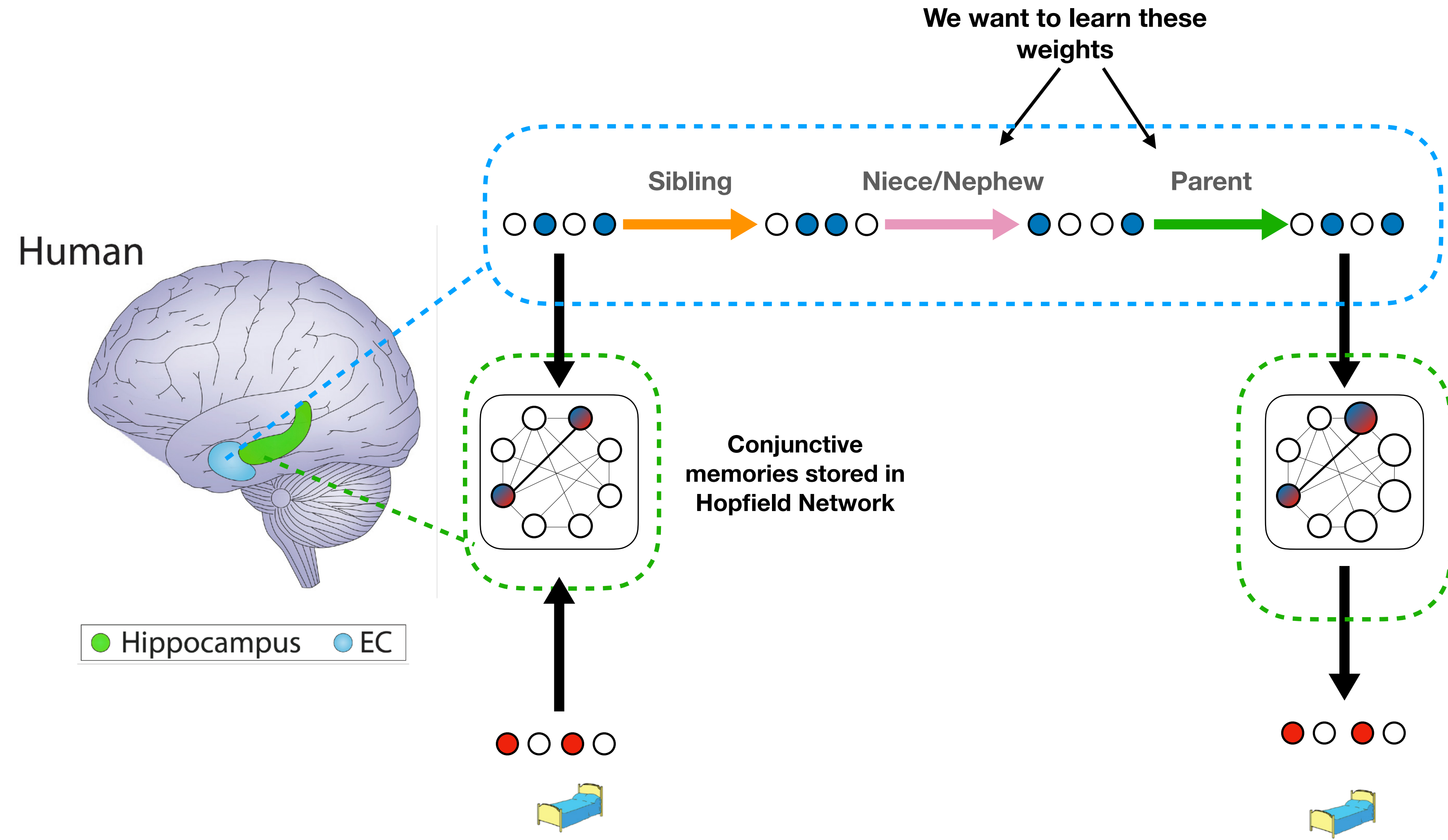
# TEM learns abstract location representations by prediction errors



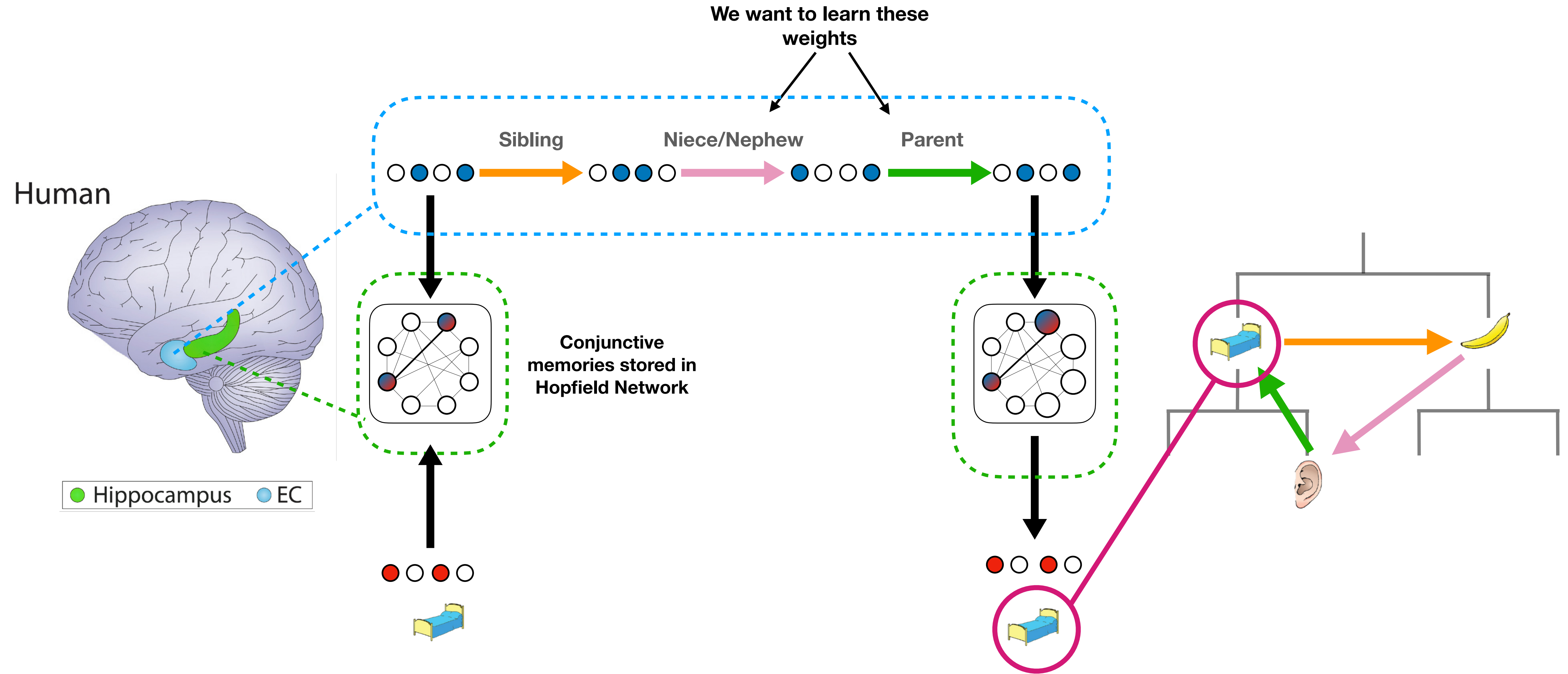
# TEM learns abstract location representations by prediction errors



# TEM learns abstract location representations by prediction errors

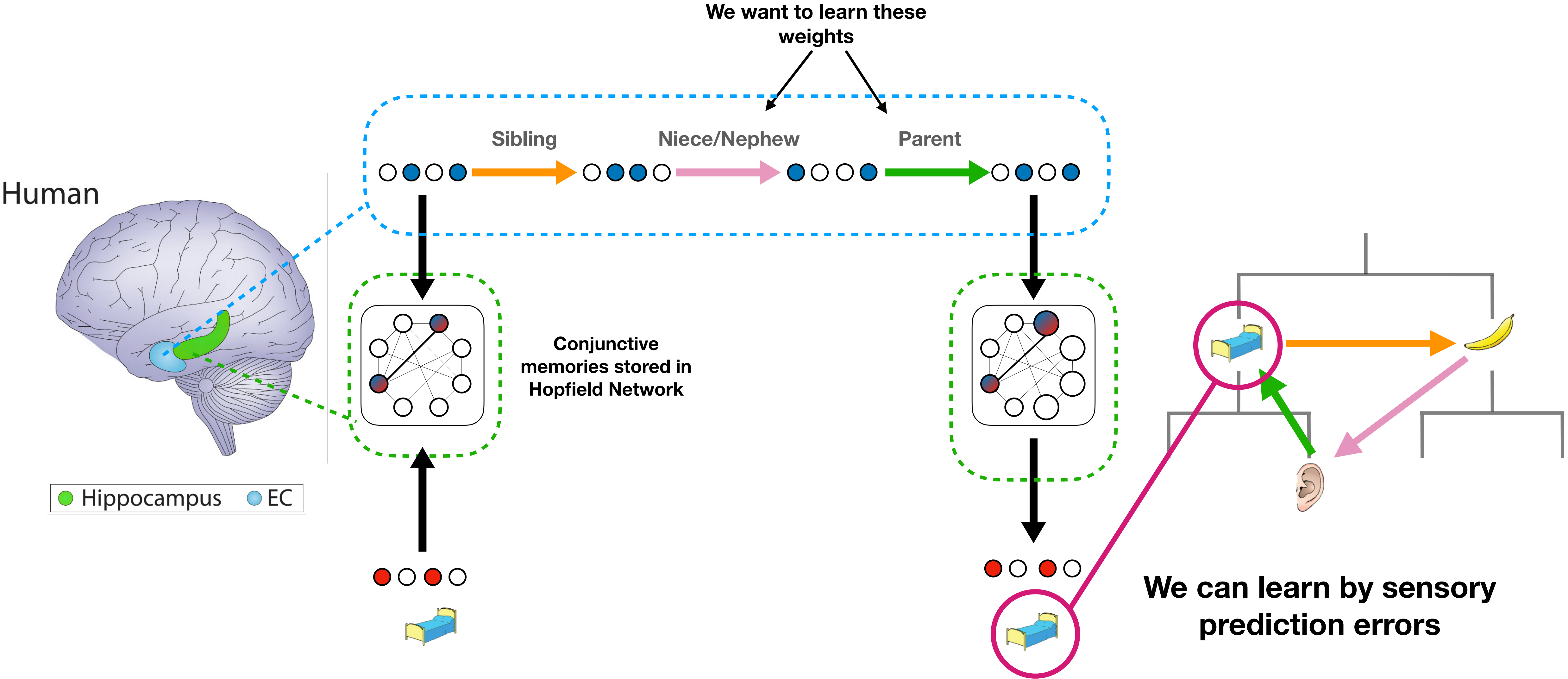


# TEM learns abstract location representations by prediction errors

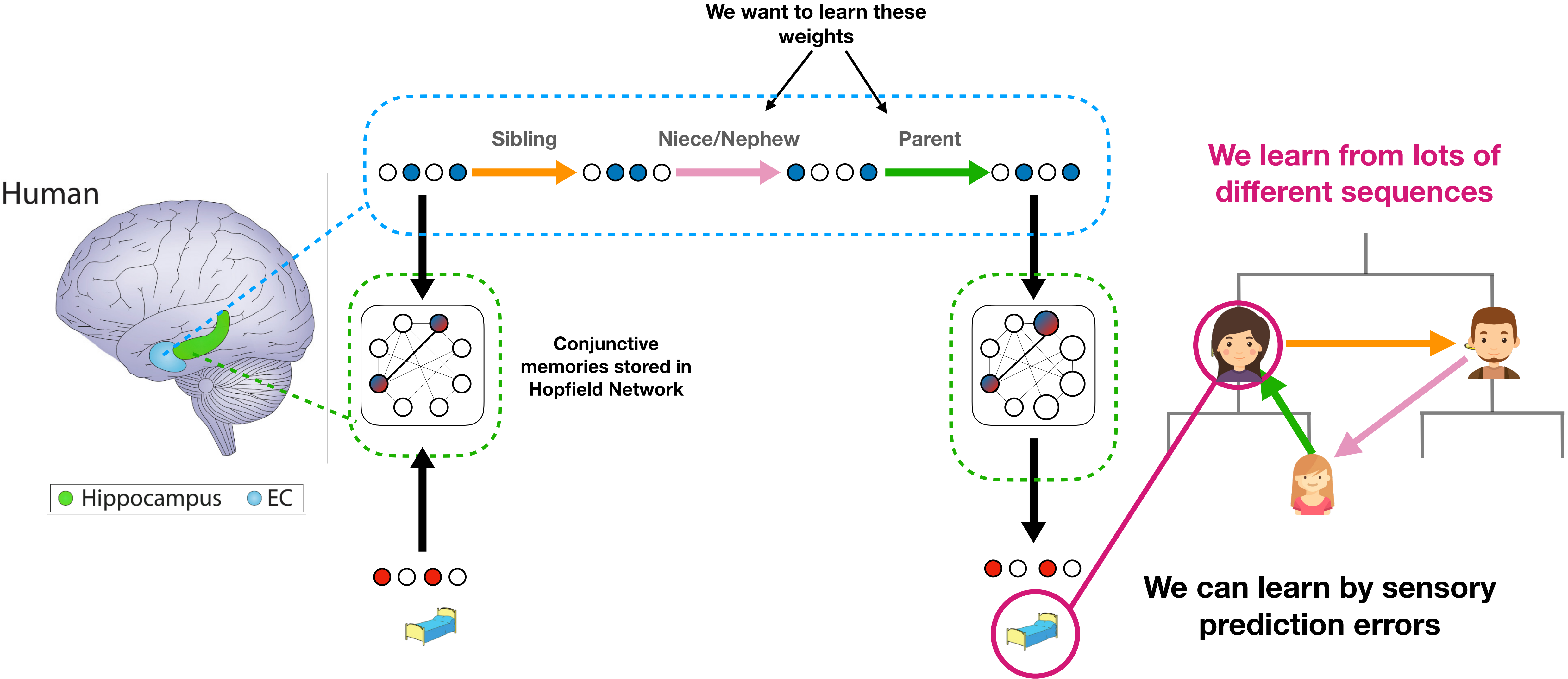




# TEM learns abstract location representations by prediction errors



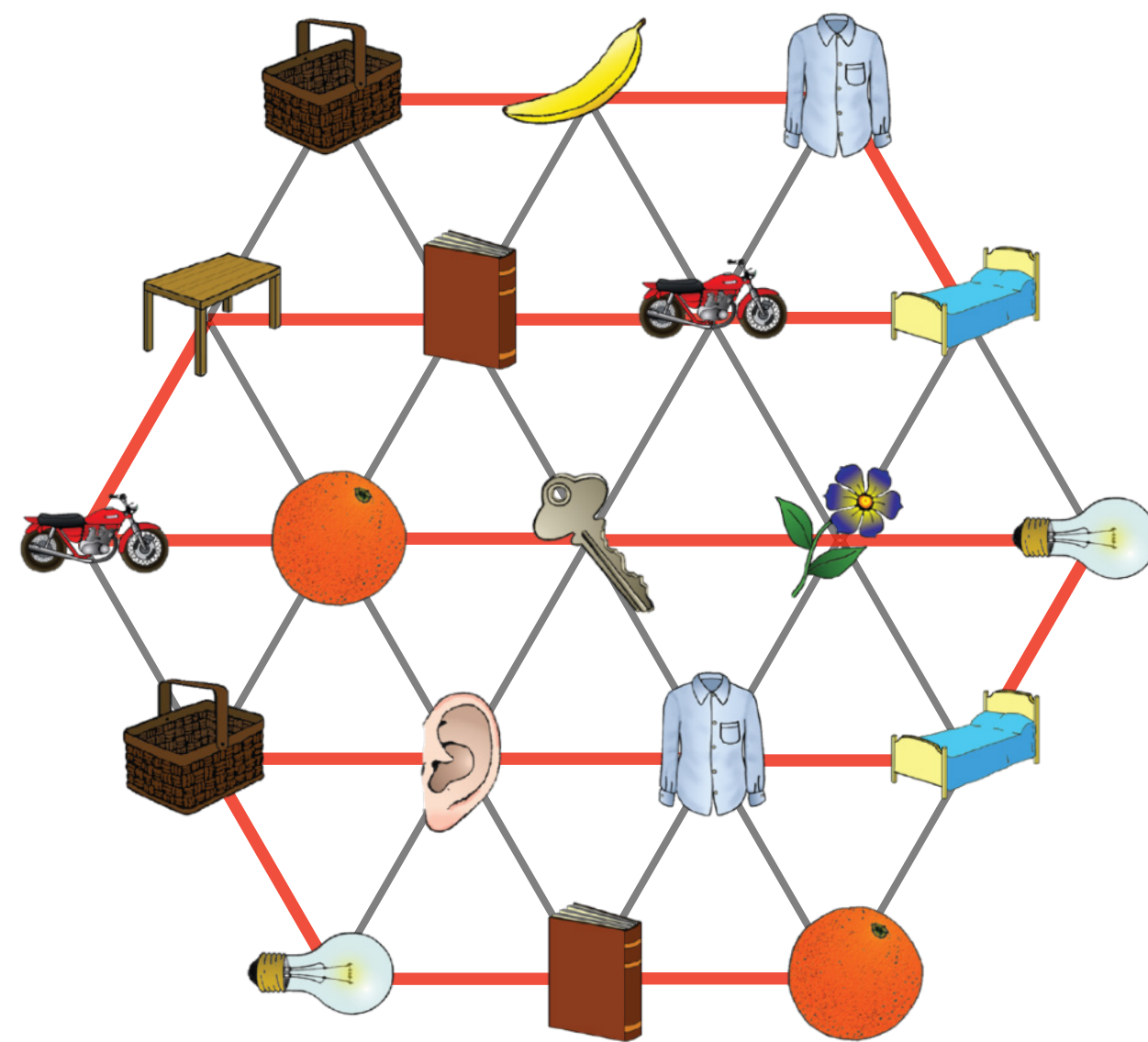
# TEM learns abstract location representations by prediction errors



**Does it work?**

# TEM infers relationships when it sees a new observation

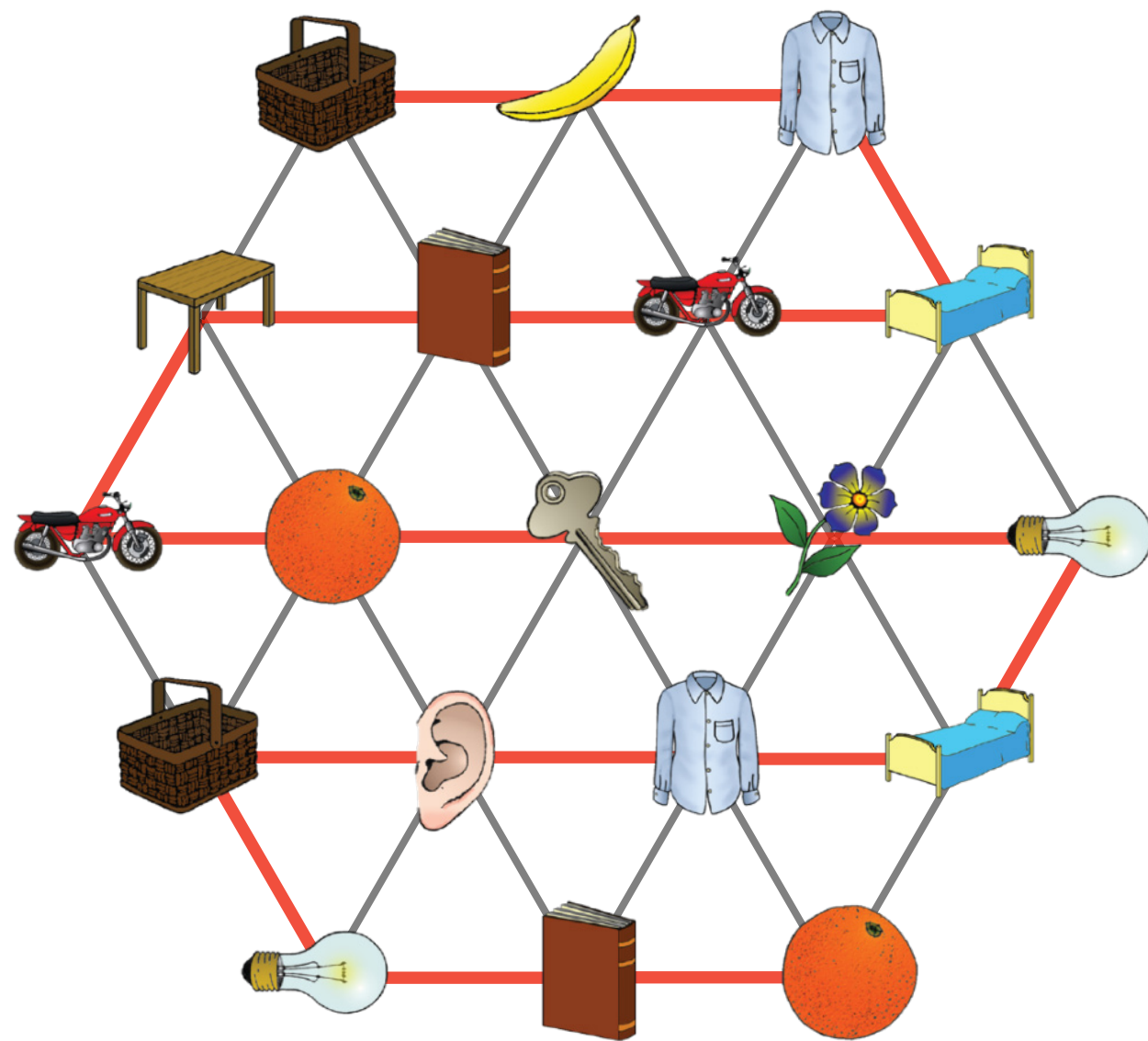
## Testing on a novel spatial environment



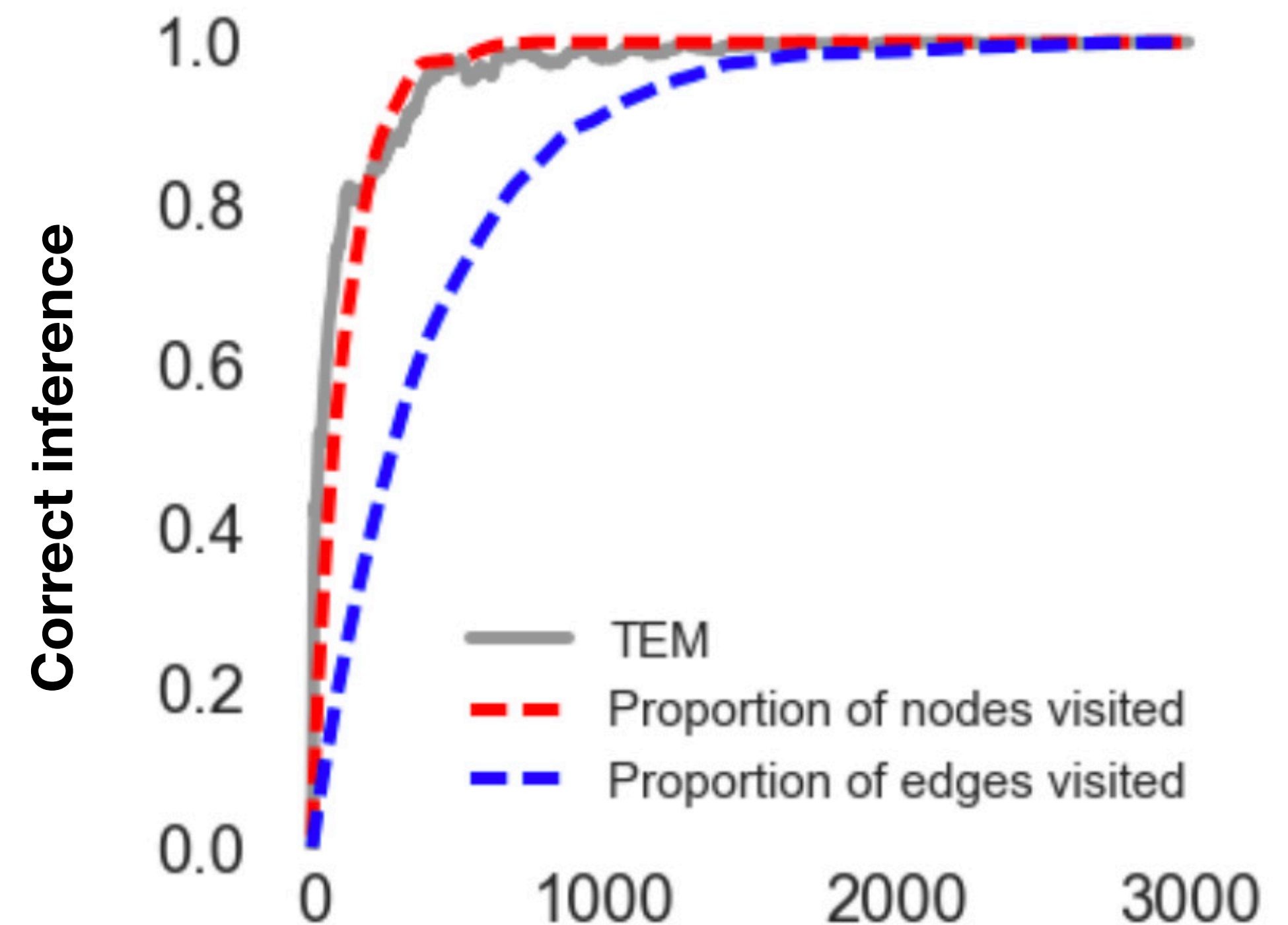
Performance in new worlds should depend on number of **nodes** it has seen, not number of **edges**

# TEM infers relationships when it sees a new observation

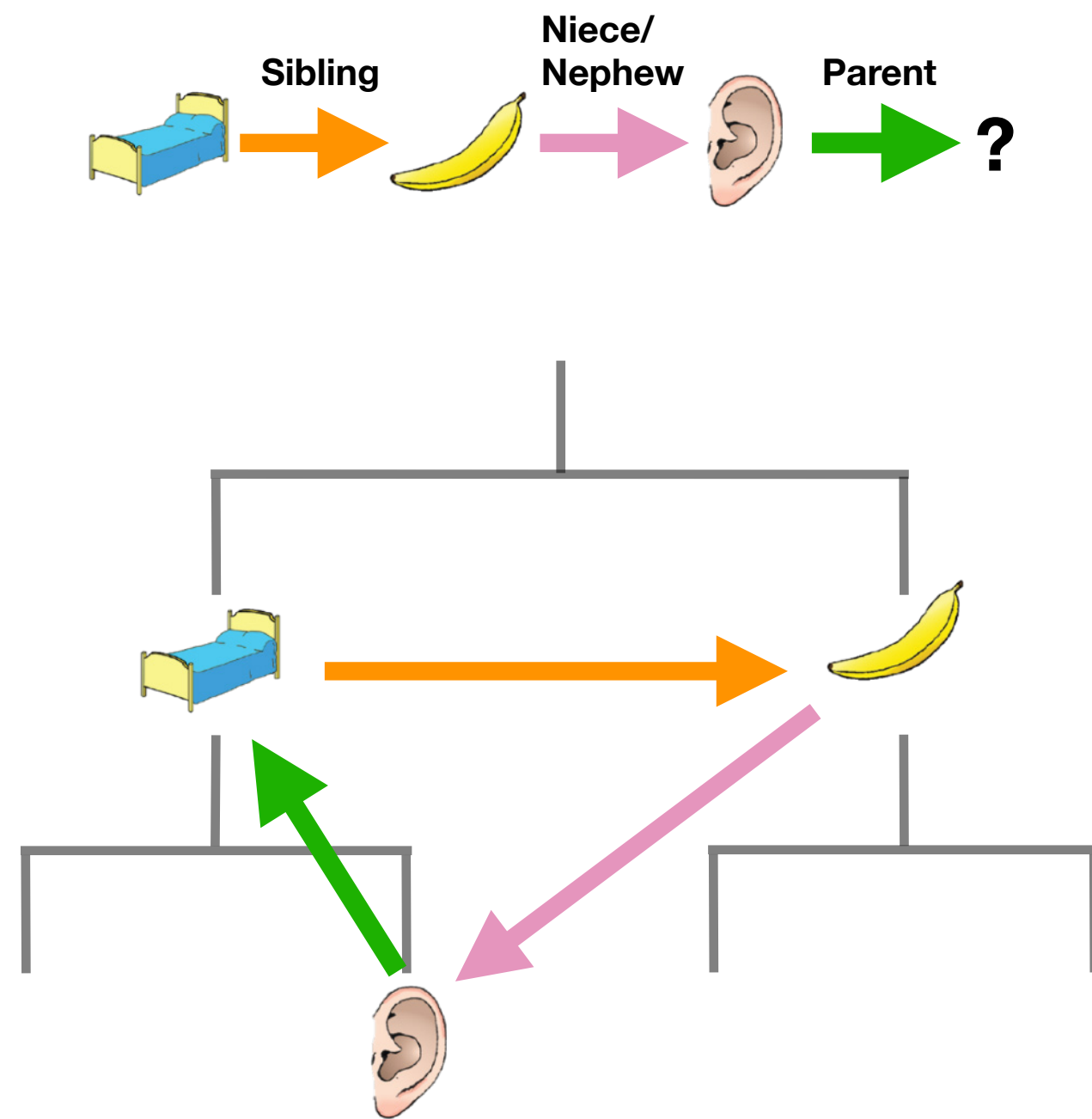
## Testing on a novel spatial environment



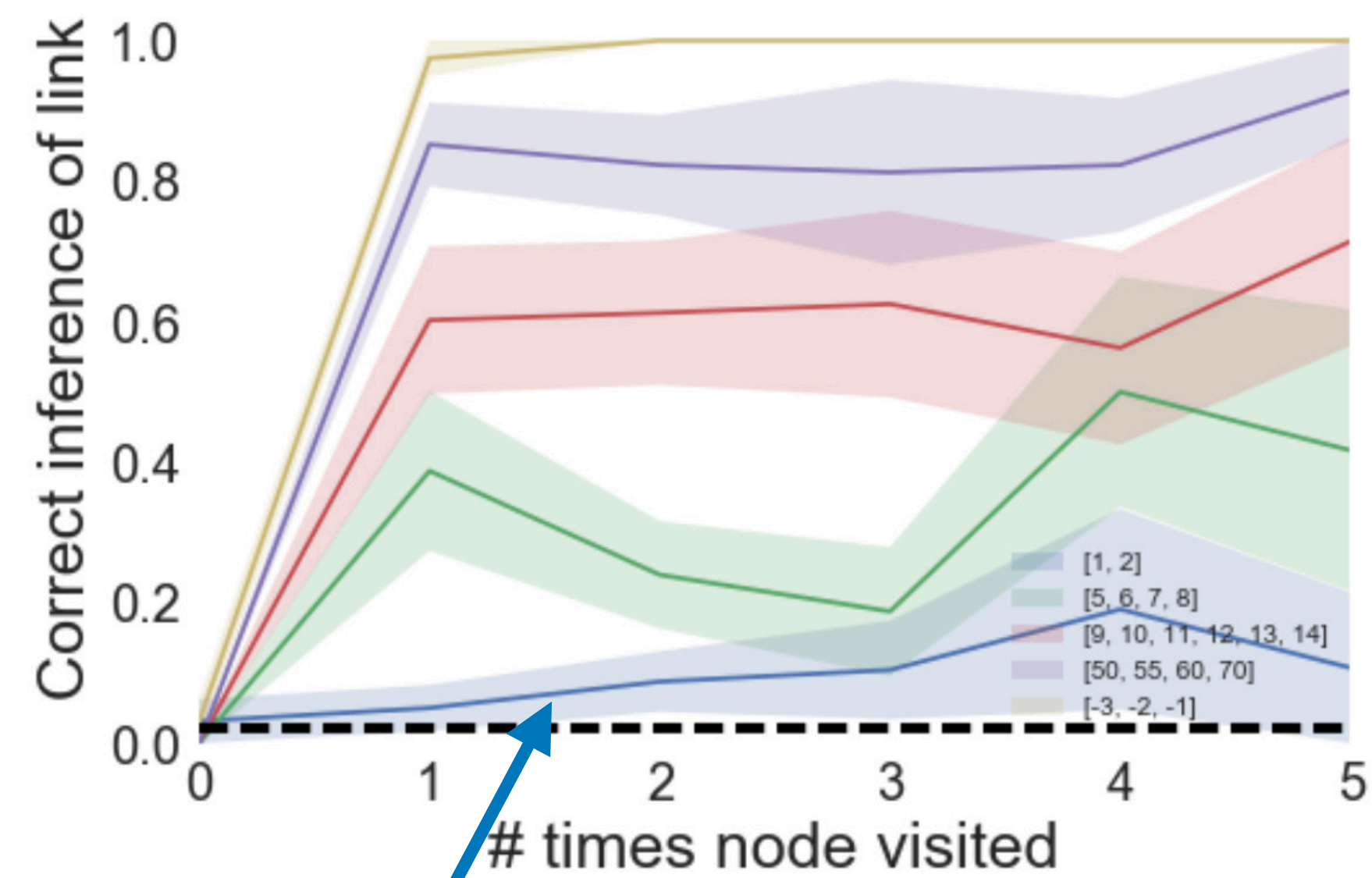
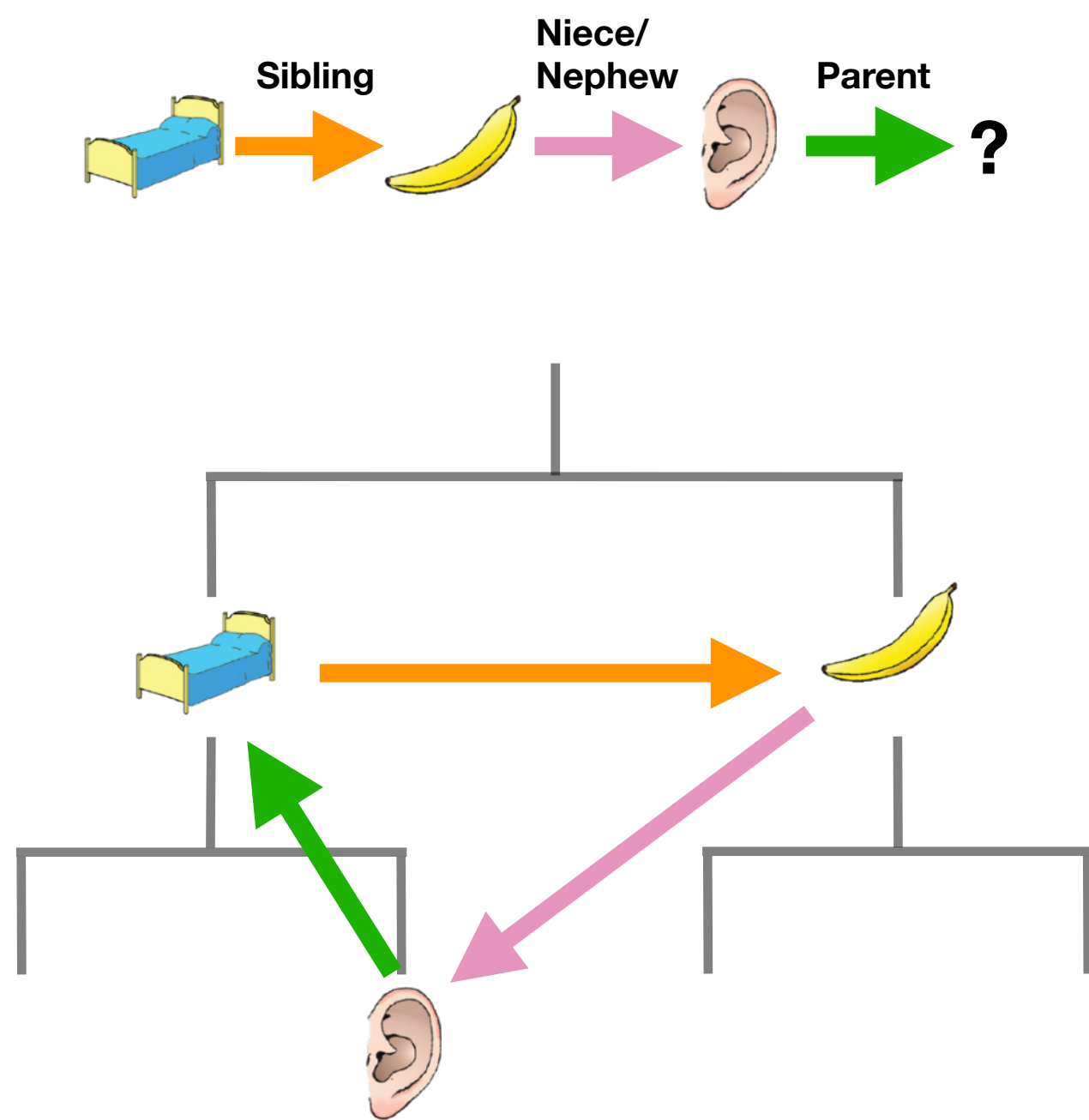
Performance in new worlds should depend on number of **nodes** it has seen, not number of **edges**



# TEM makes inferences because it has learned in similar structures

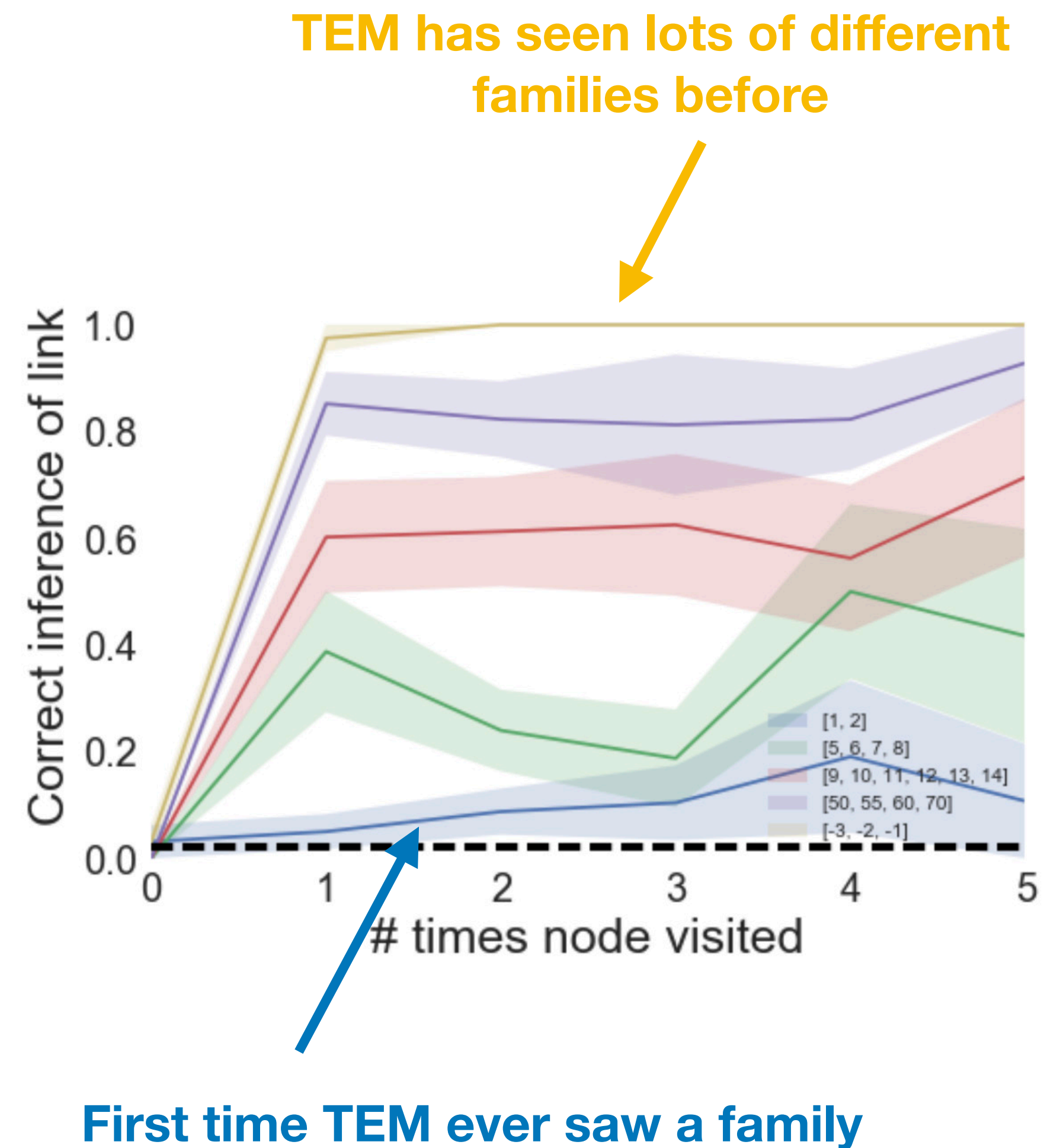
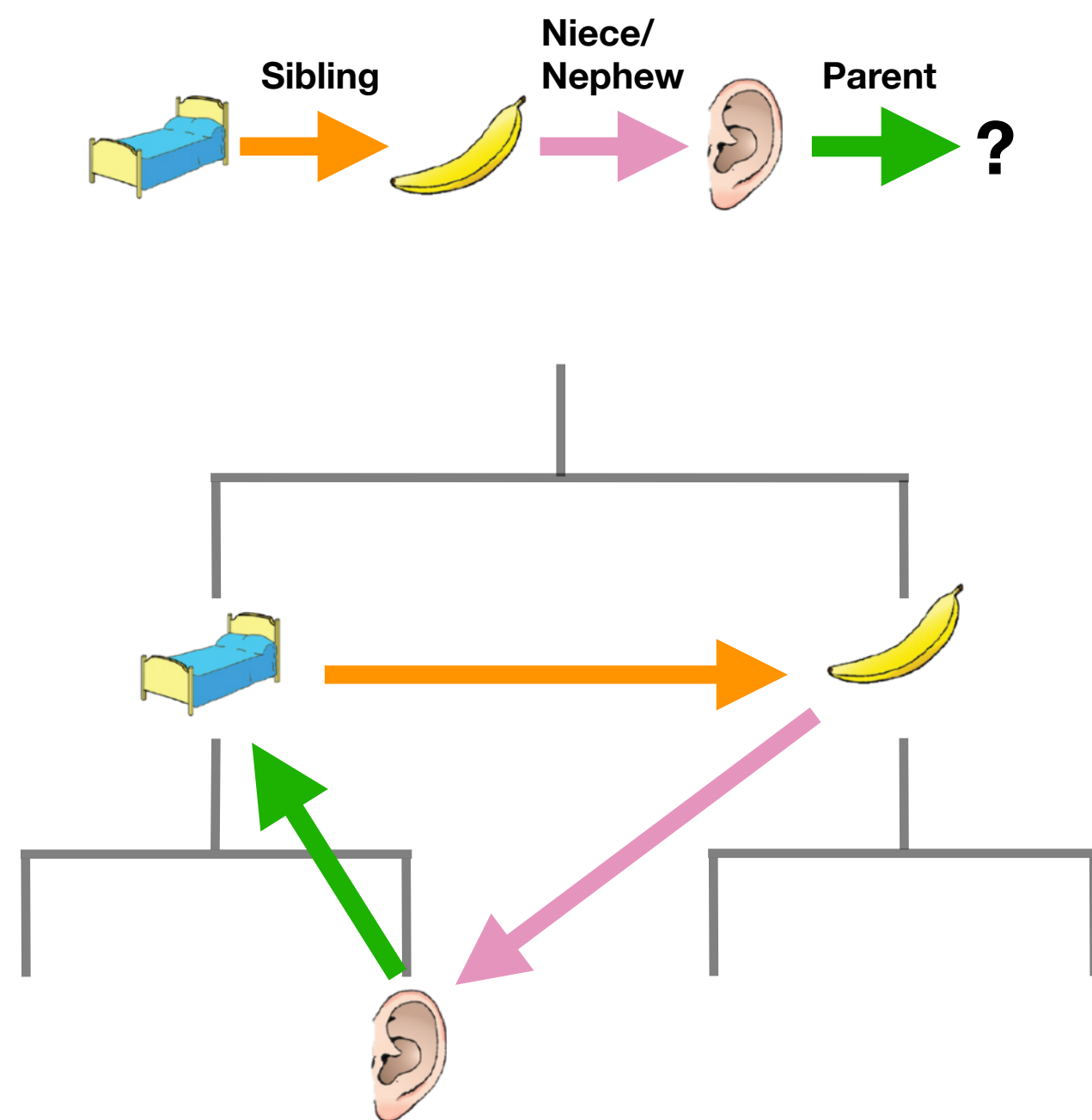


# TEM makes inferences because it has learned in similar structures



First time TEM ever saw a family

# TEM makes inferences because it has learned in similar structures

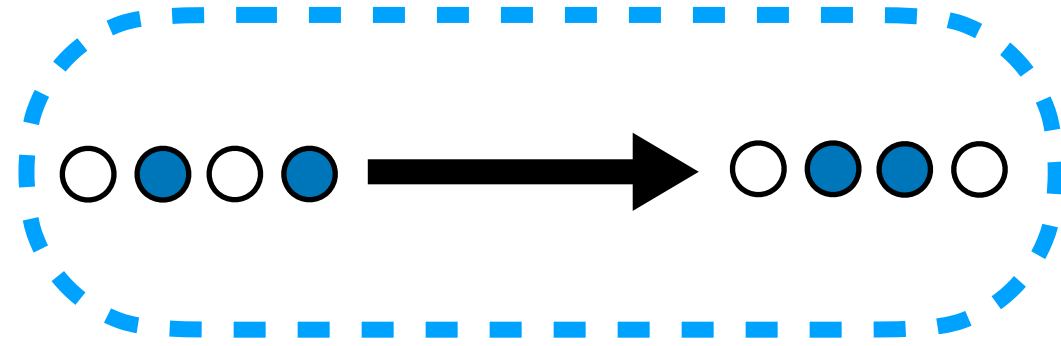




# Learned **building block** representations

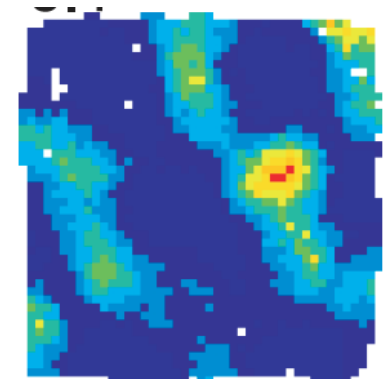
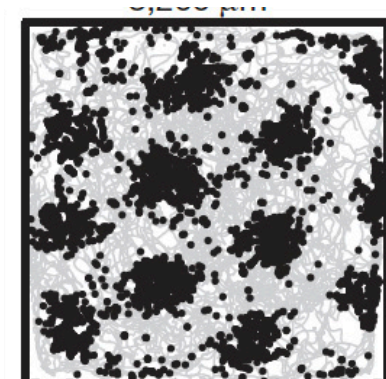
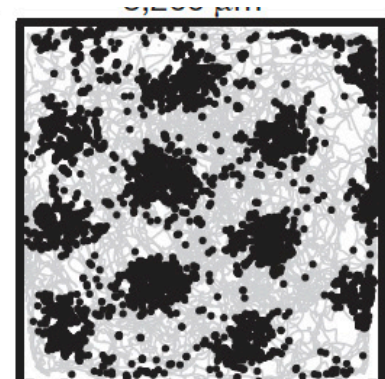
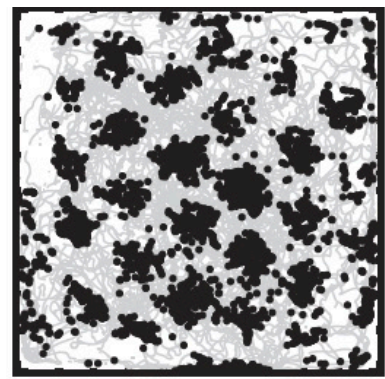
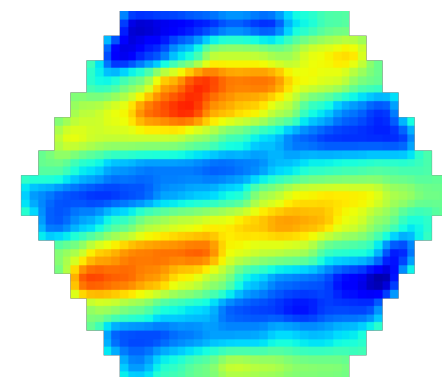
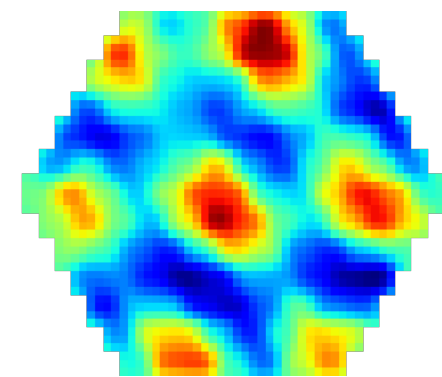
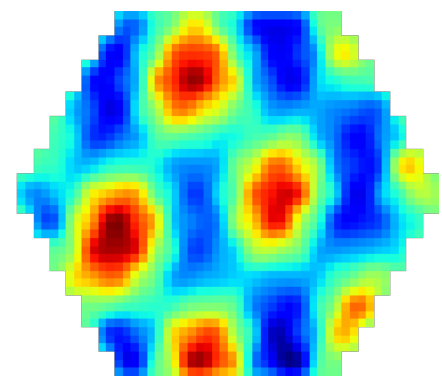
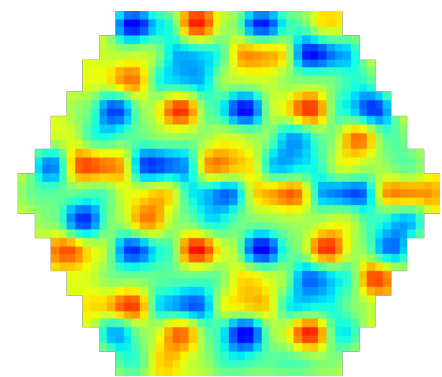
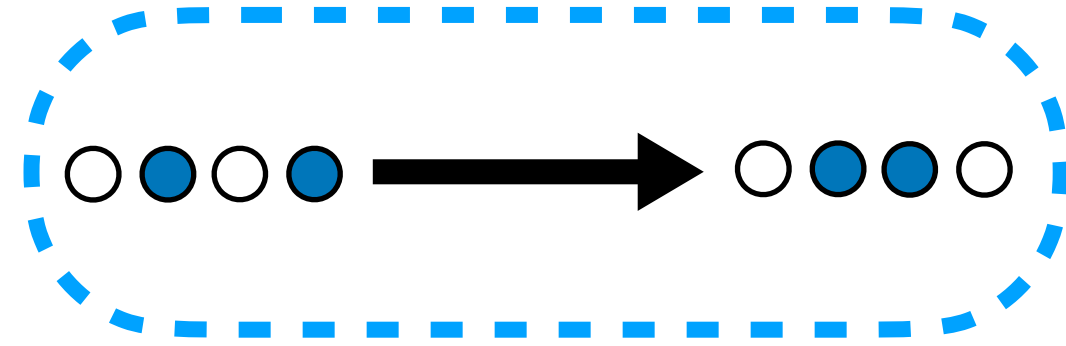
# Learned **building block** representations

Path  
integrating  
RNN



# Learned building block representations

Path  
integrating  
RNN

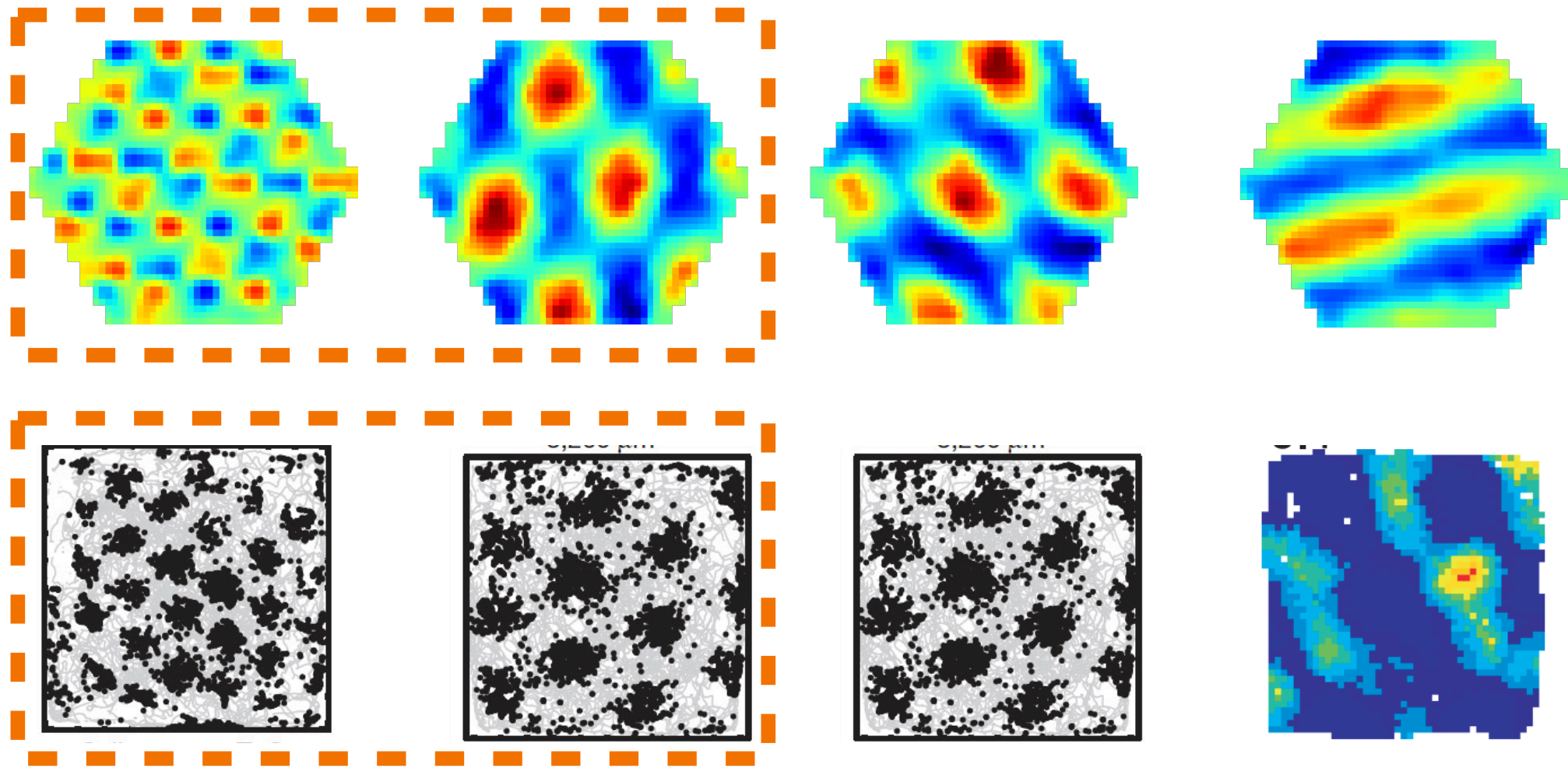
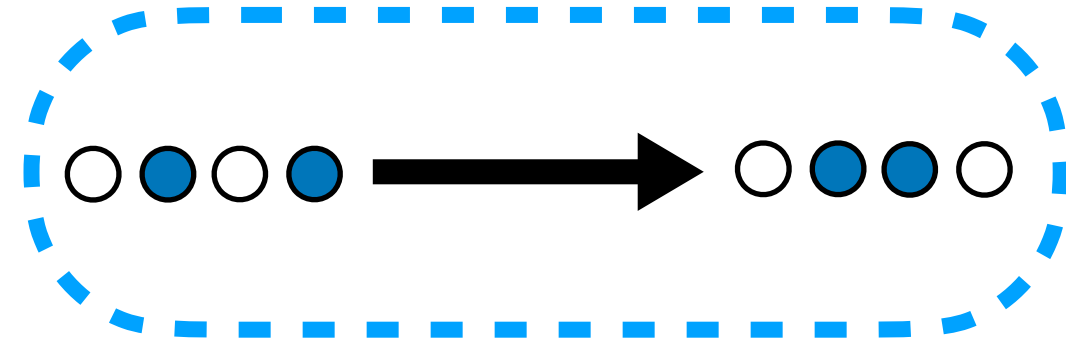


Grid cells  
Hafting et al., 2005

Band cells  
Krupic et al., 2012

# Learned building block representations

Path  
integrating  
RNN

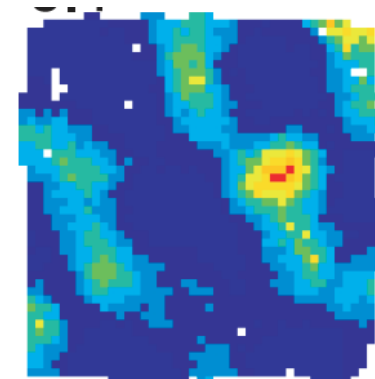
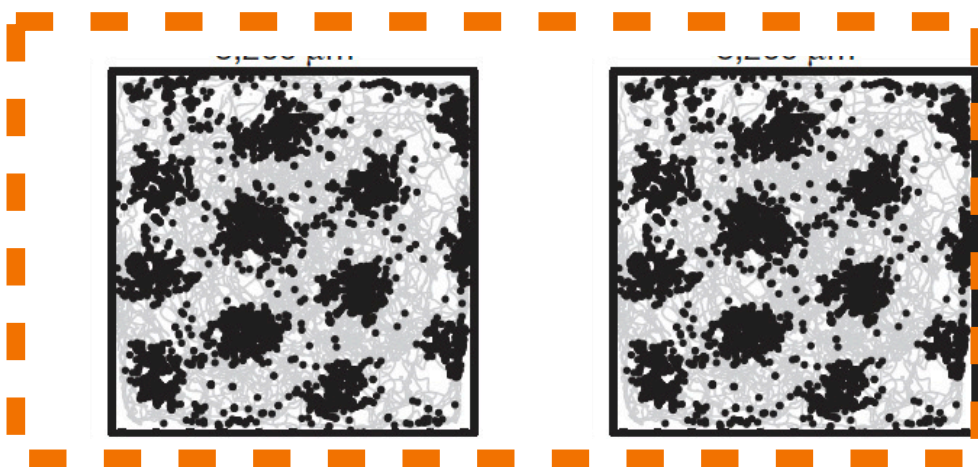
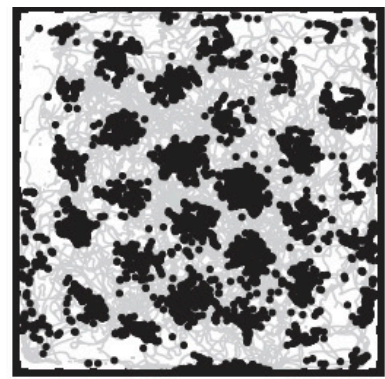
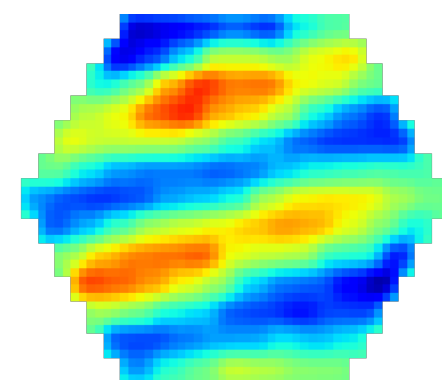
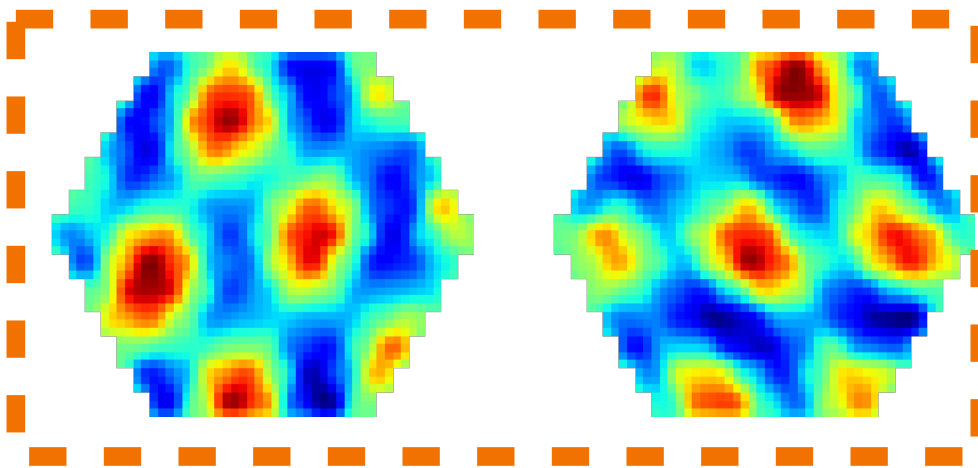
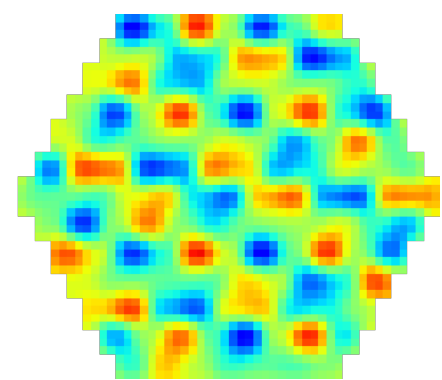
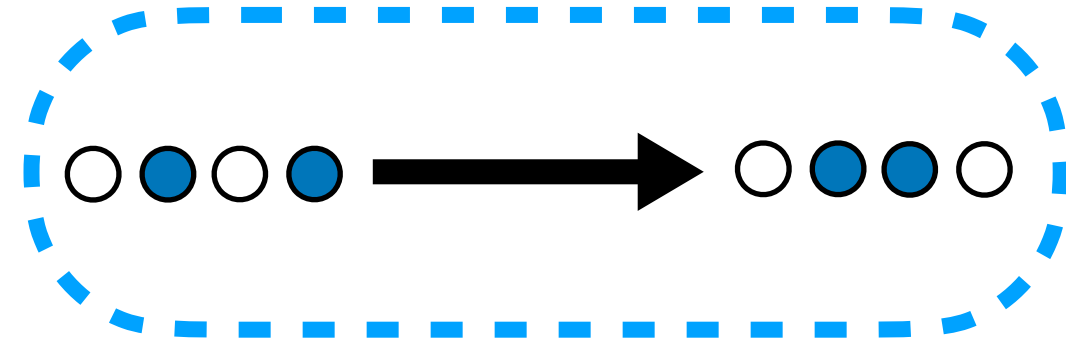


Grid cells  
Hafting et al., 2005

Band cells  
Krupic et al., 2012

# Learned building block representations

Path  
integrating  
RNN

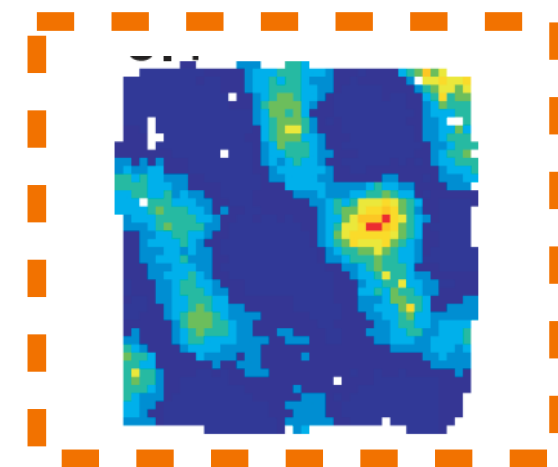
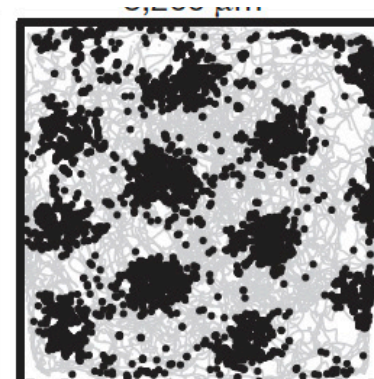
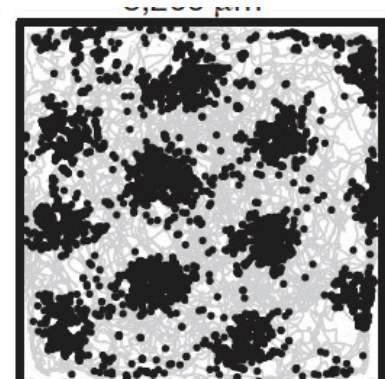
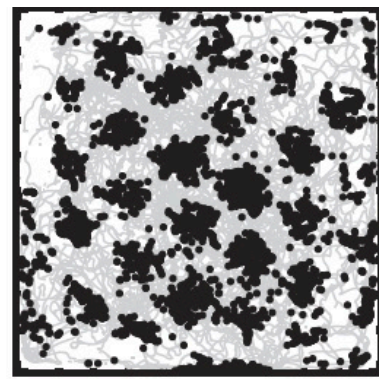
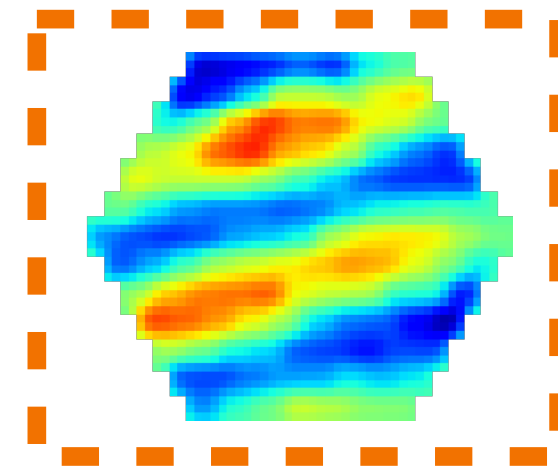
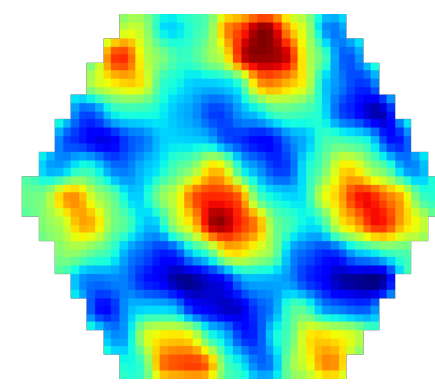
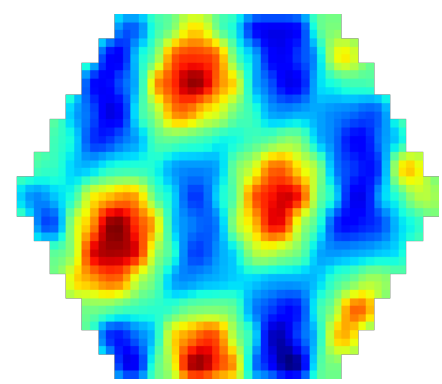
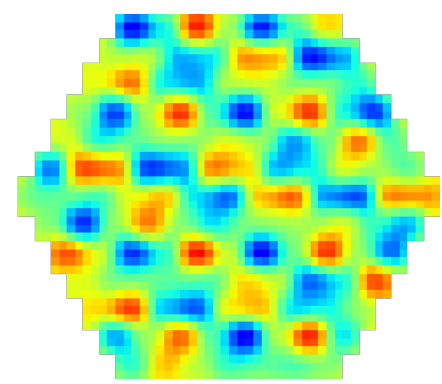
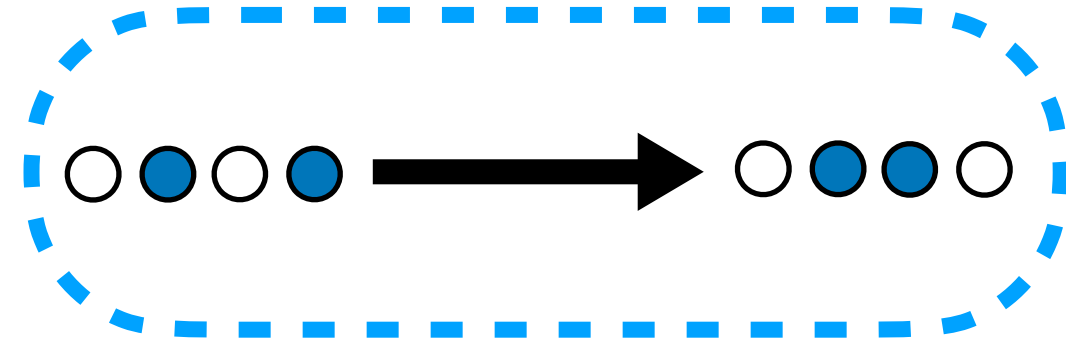


Grid cells  
Hafting et al., 2005

Band cells  
Krupic et al., 2012

# Learned building block representations

Path  
integrating  
RNN

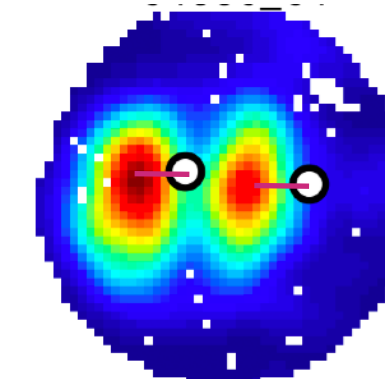
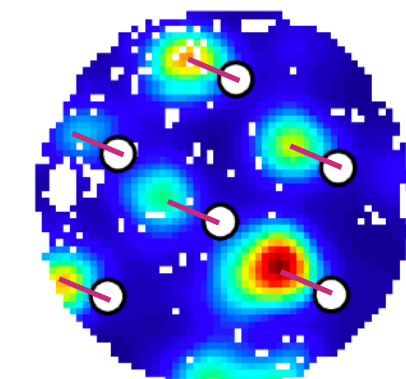
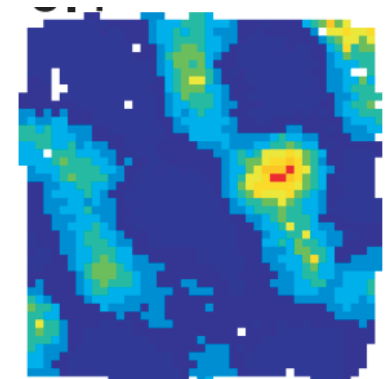
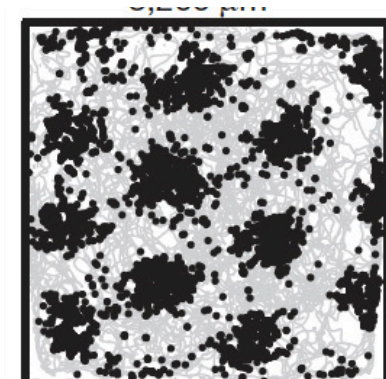
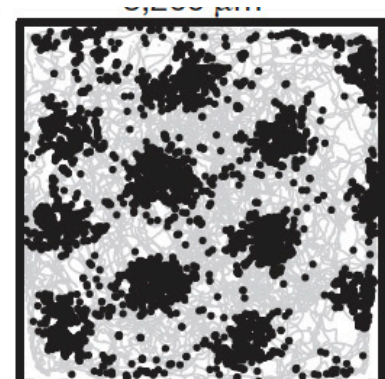
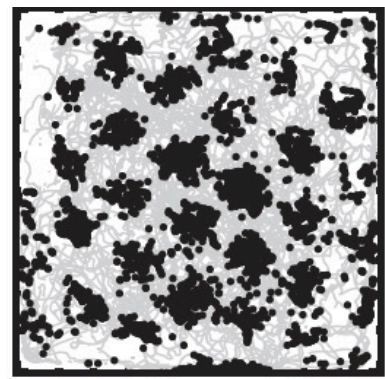
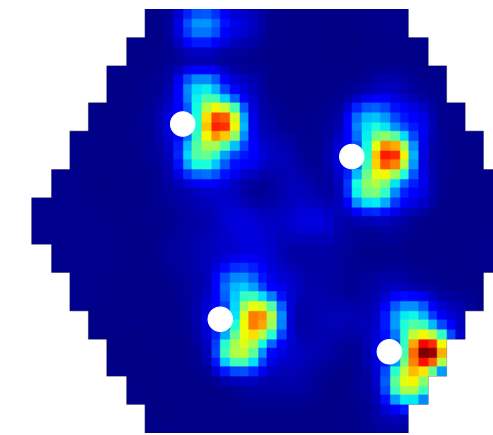
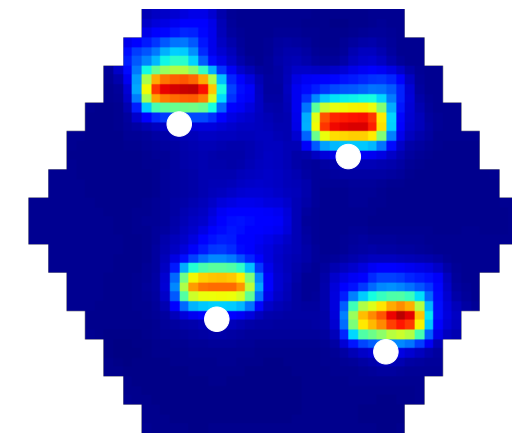
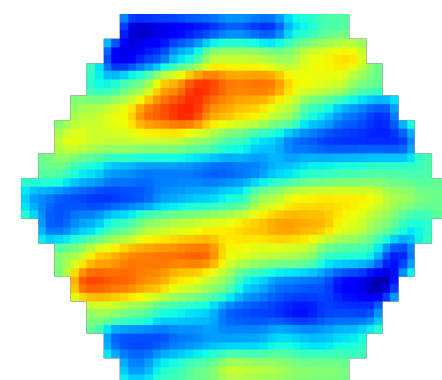
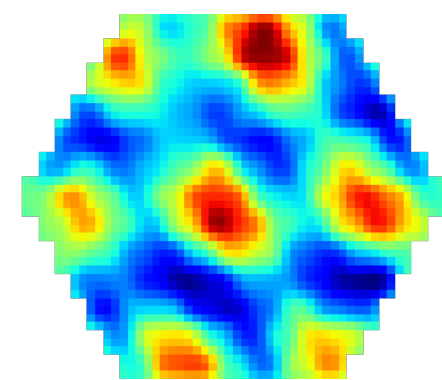
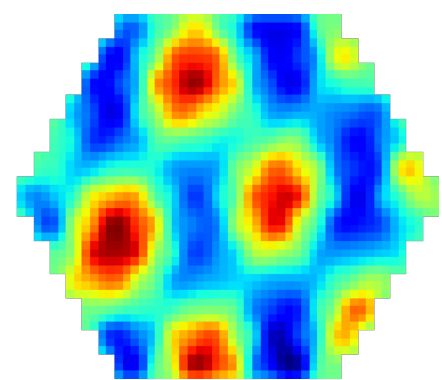
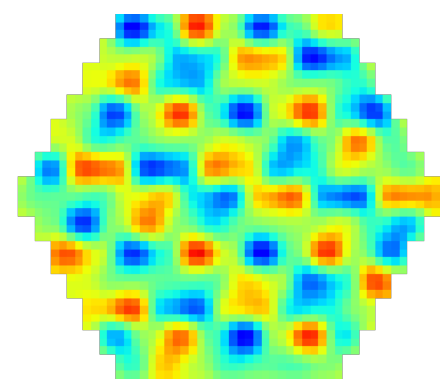
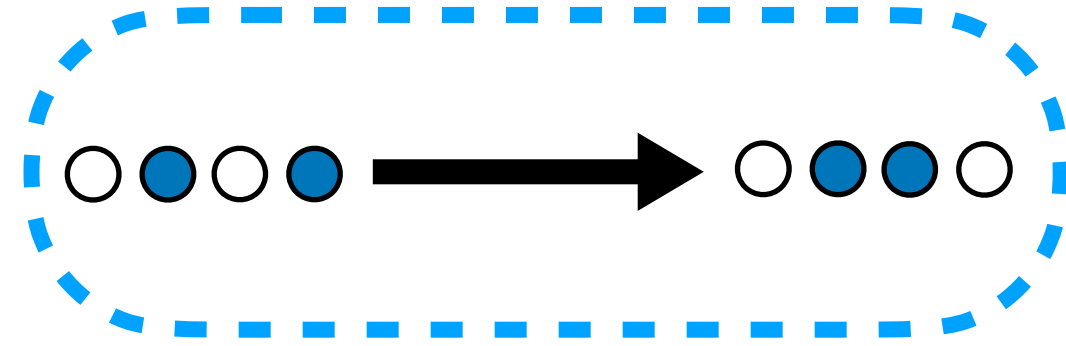


Grid cells  
Hafting et al., 2005

Band cells  
Krupic et al., 2012

# Learned building block representations

Path  
integrating  
RNN



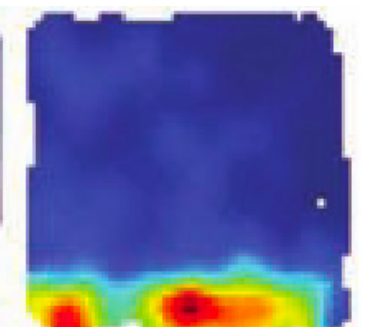
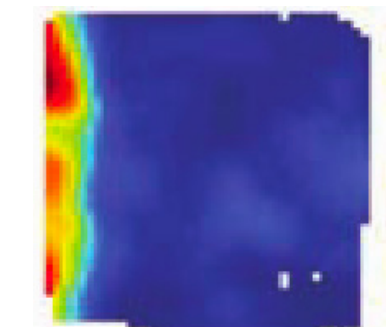
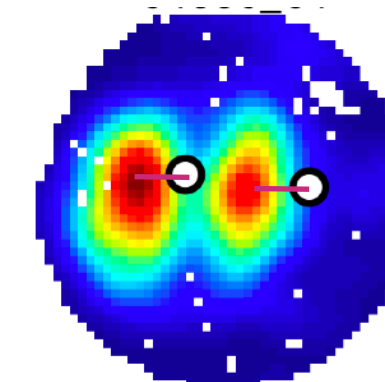
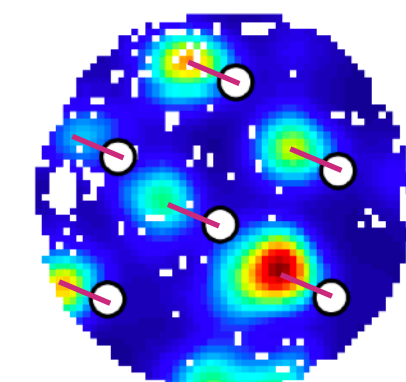
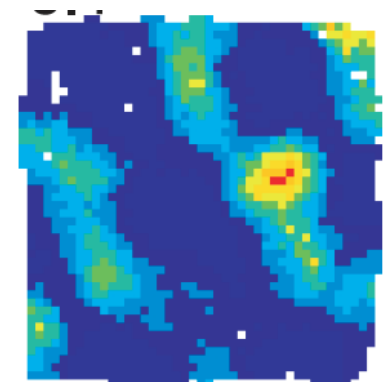
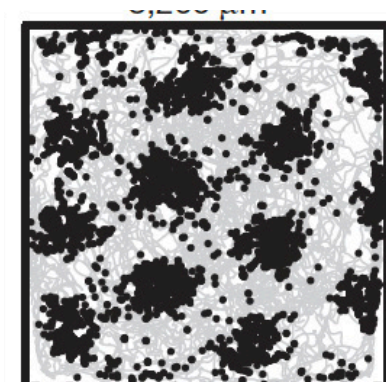
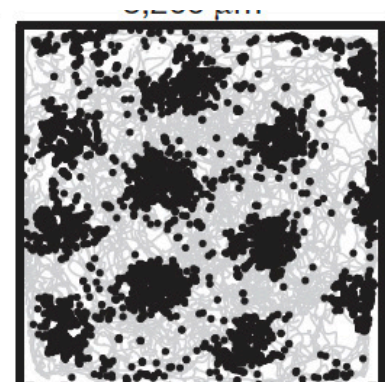
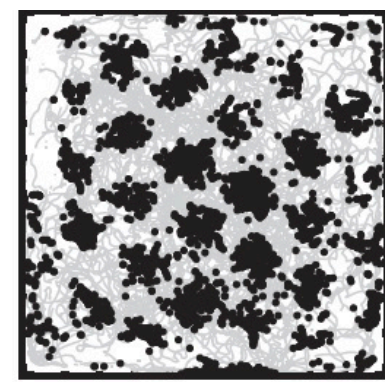
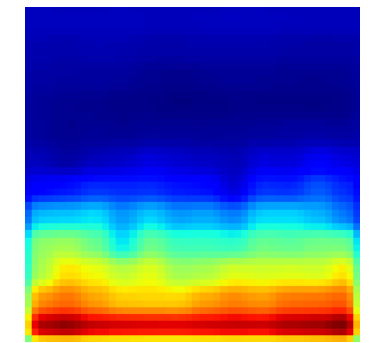
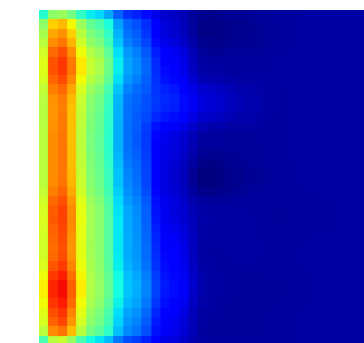
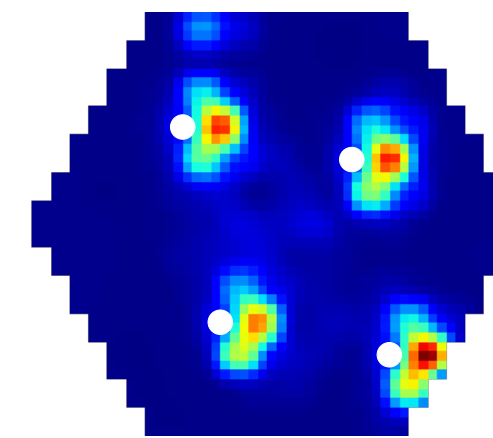
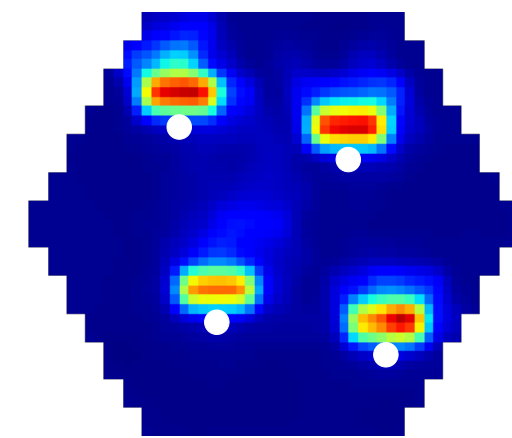
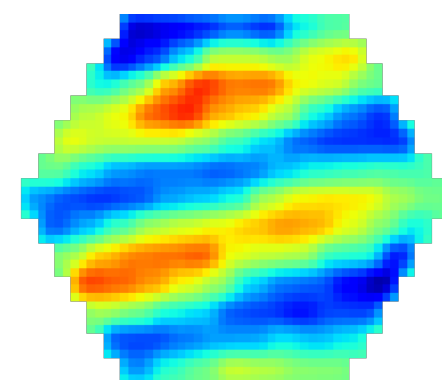
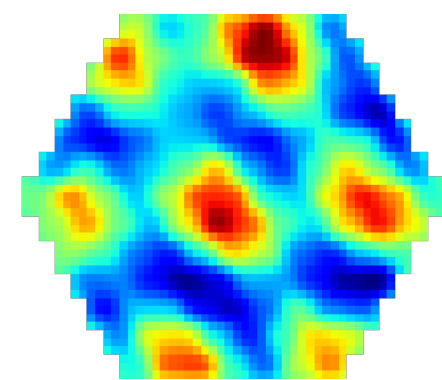
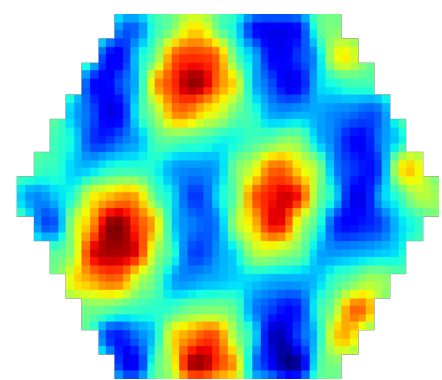
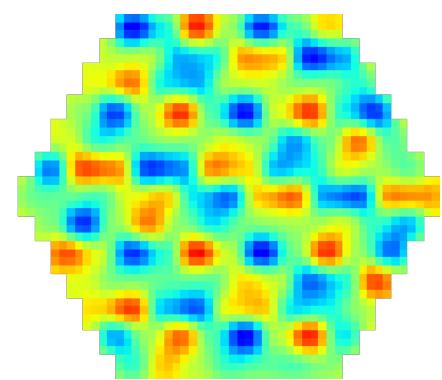
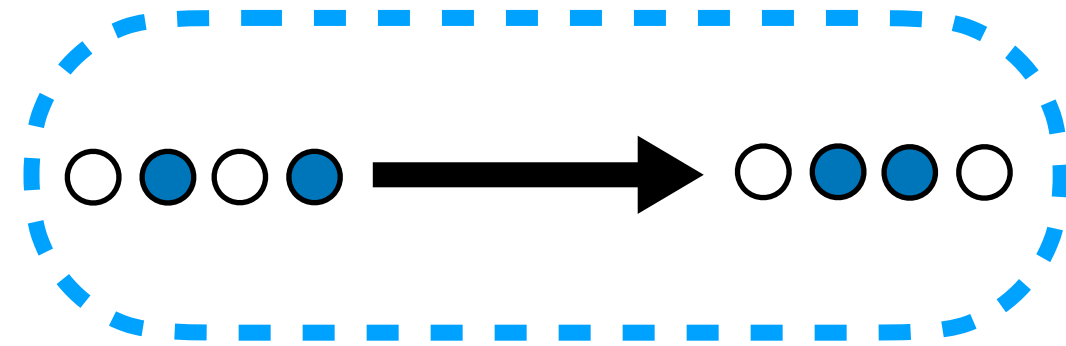
Grid cells  
Haftting et al., 2005

Band cells  
Krupic et al., 2012

Object vector cells  
Hoydal et al., 2018

# Learned building block representations

Path  
integrating  
RNN



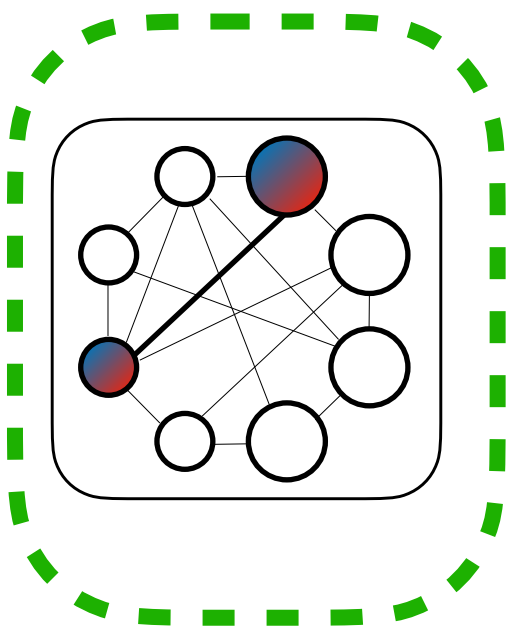
Grid cells  
Haftting et al., 2005

Band cells  
Krupic et al., 2012

Object vector cells  
Hoydal et al., 2018

Border cells  
Solstad et al., 2008

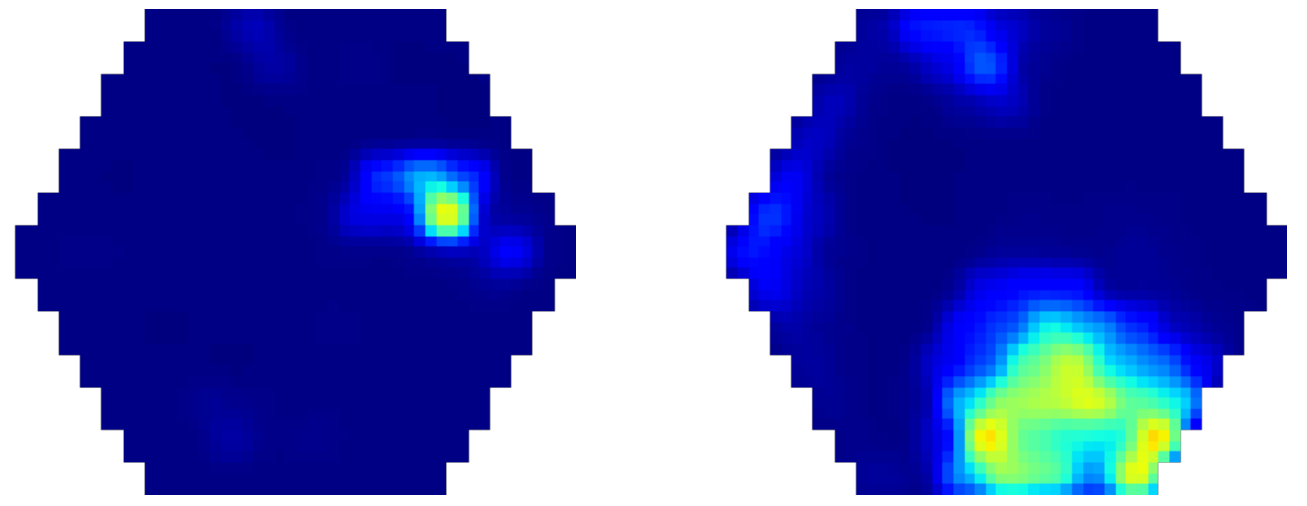




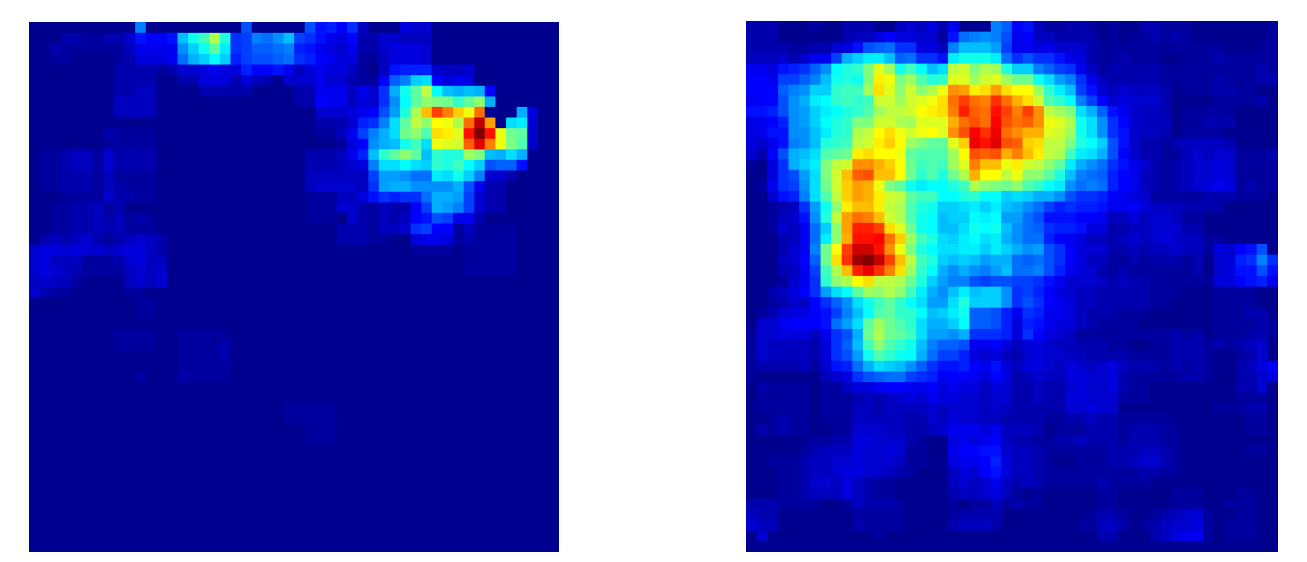
Hopfield Network

# Learned **Hippocampal** representations

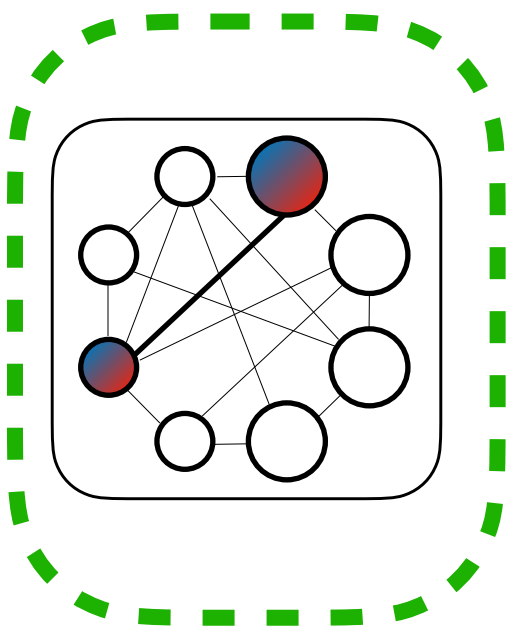
Example TEM cells



Example real cells



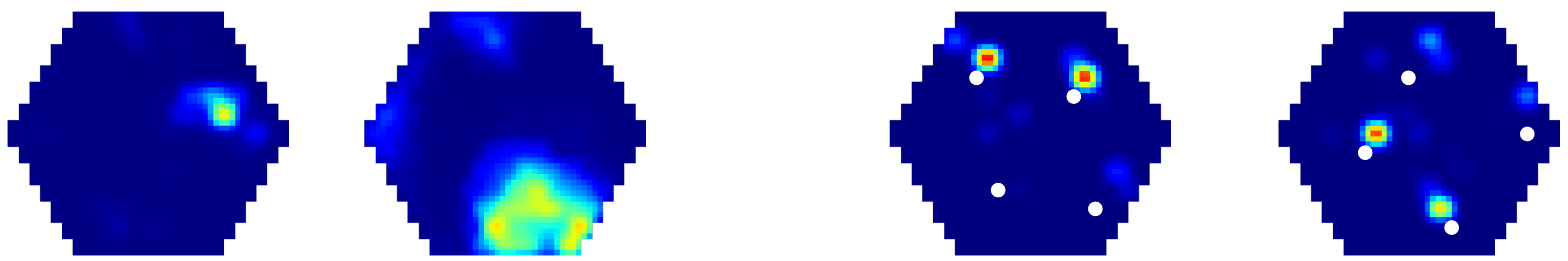
Place cells  
O'Keefe & Dostrovsky, 1971



Hopfield Network

# Learned **Hippocampal** representations

Example TEM cells



Example real cells



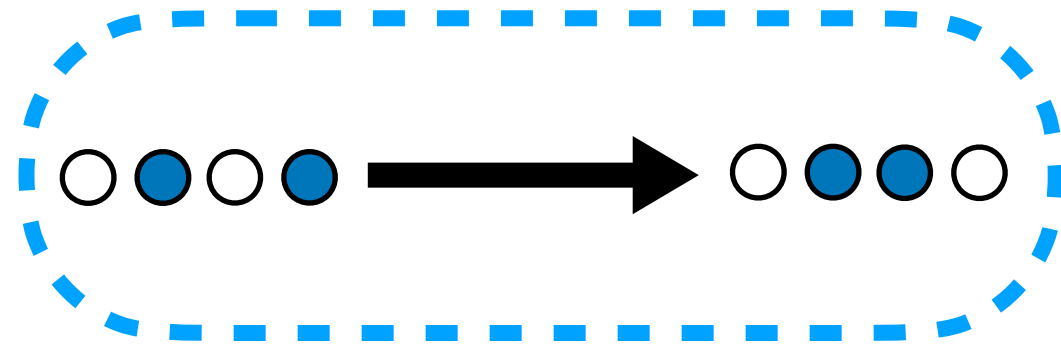
Place cells  
O'Keefe & Dostrovsky, 1971

Landmark cells  
Deshmukh & Knierim, 2013

**How do these representations generalise?**

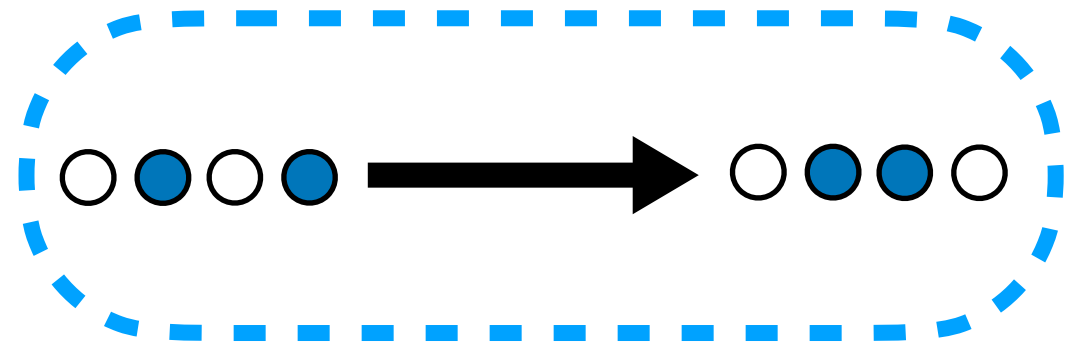
# Building block remapping

Path  
integrating  
RNN



# Building block remapping

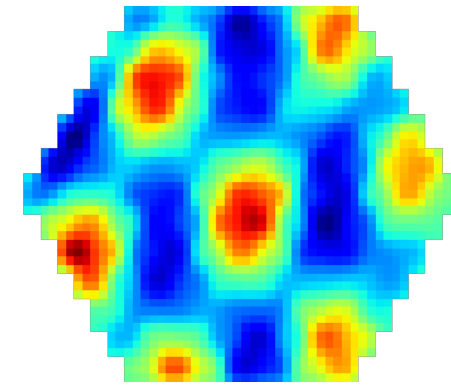
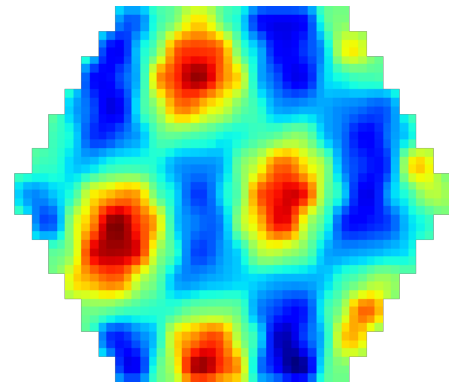
Path  
integrating  
RNN



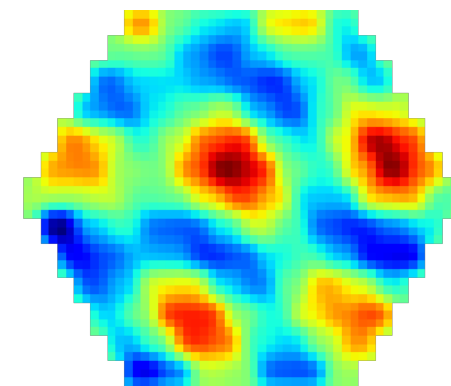
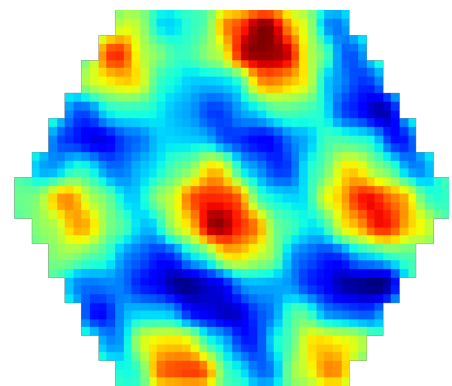
World 1

World 2

Cell 1

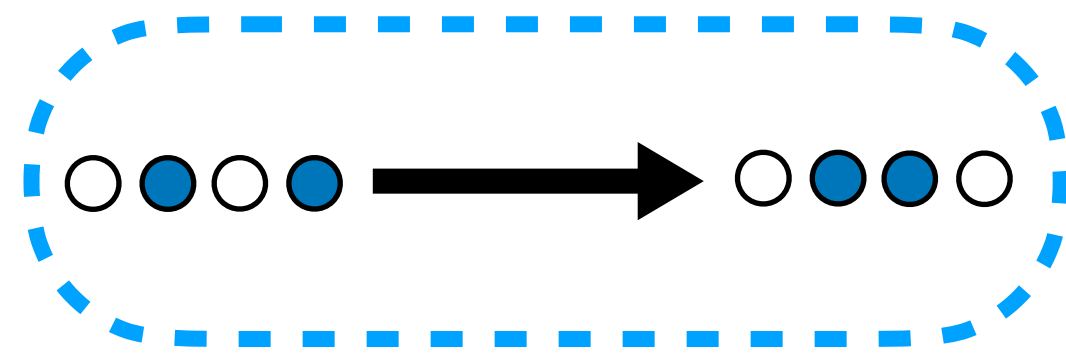


Cell 2



# Building block remapping

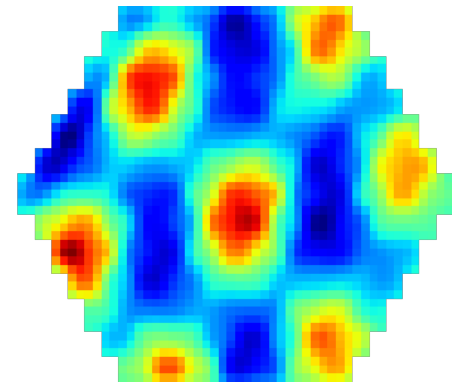
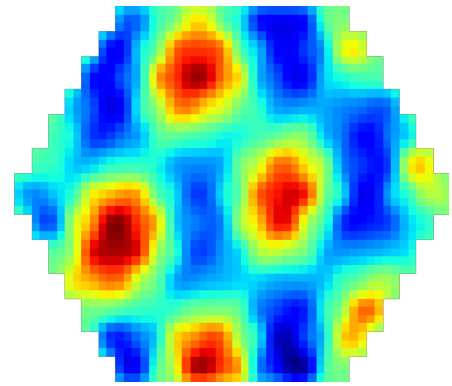
Path  
integrating  
RNN



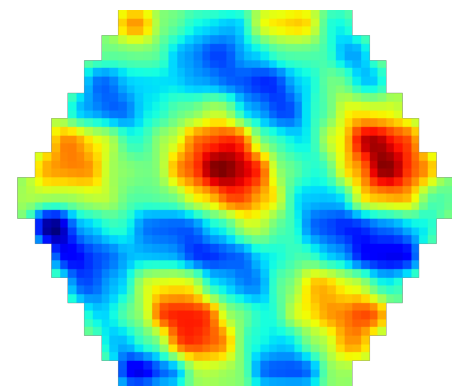
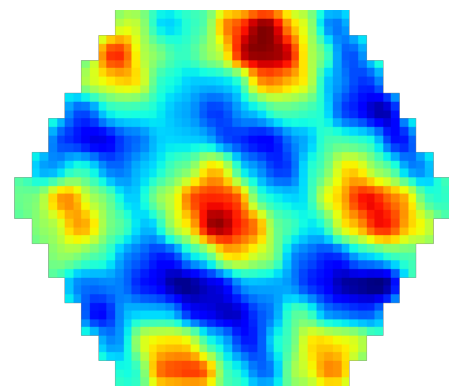
World 1

World 2

Cell 1

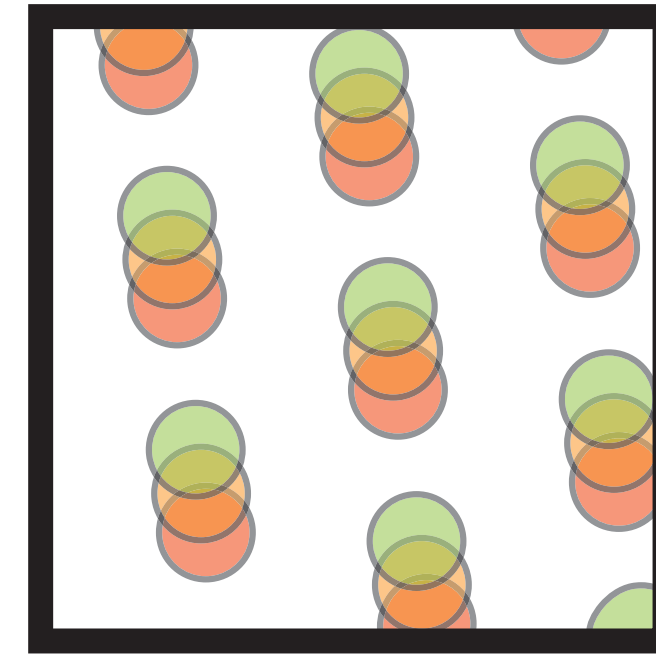
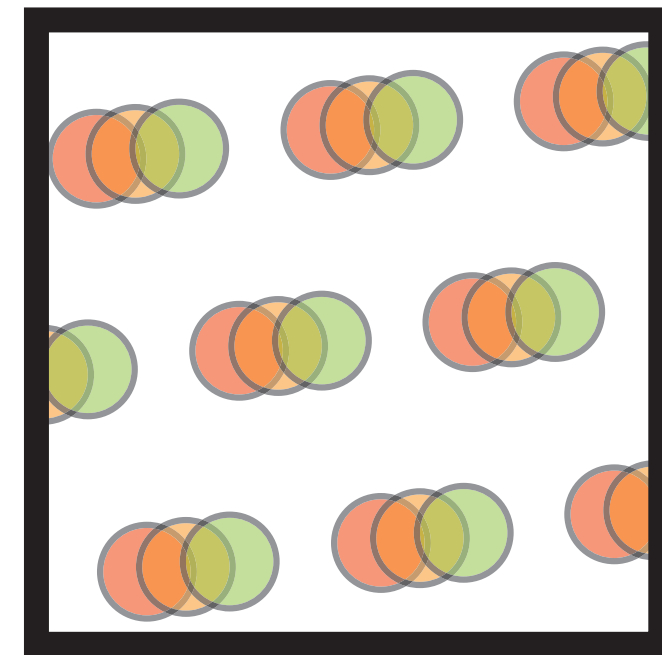


Cell 2



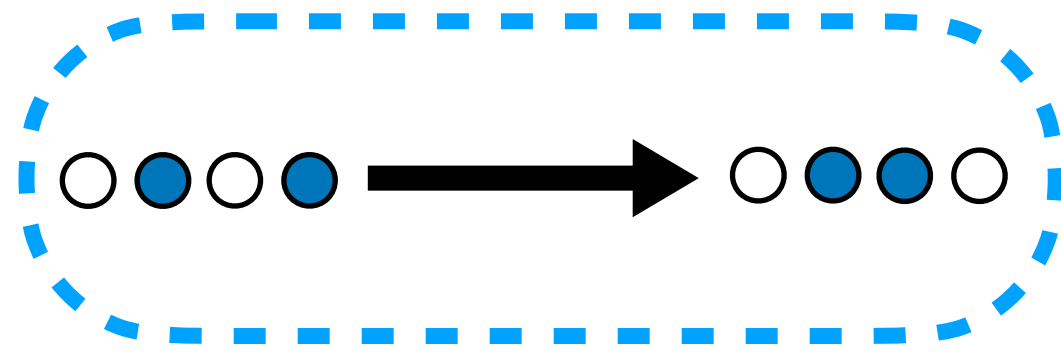
World 1

World 2



# Building block remapping

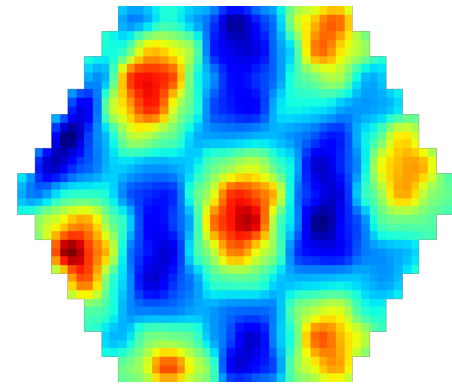
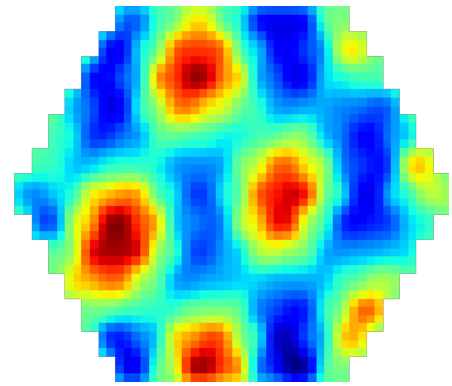
Path  
integrating  
RNN



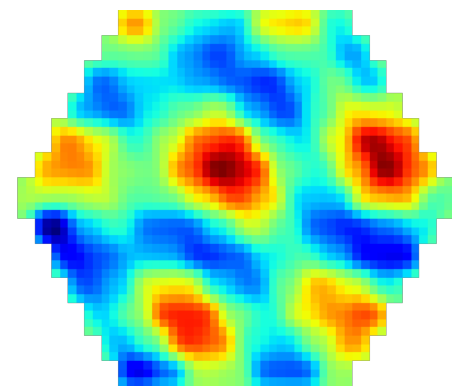
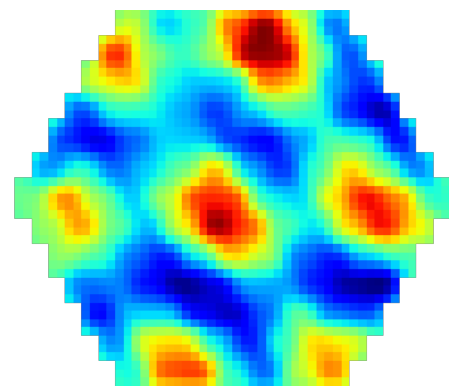
World 1

World 2

Cell 1

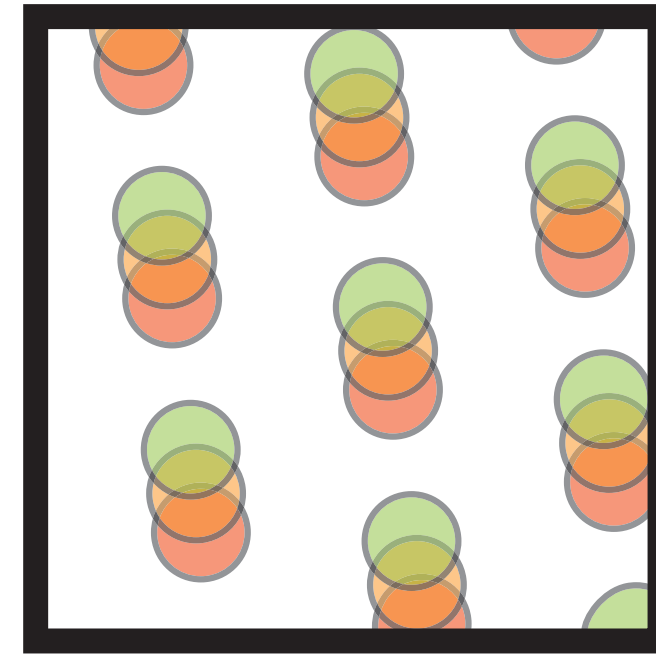
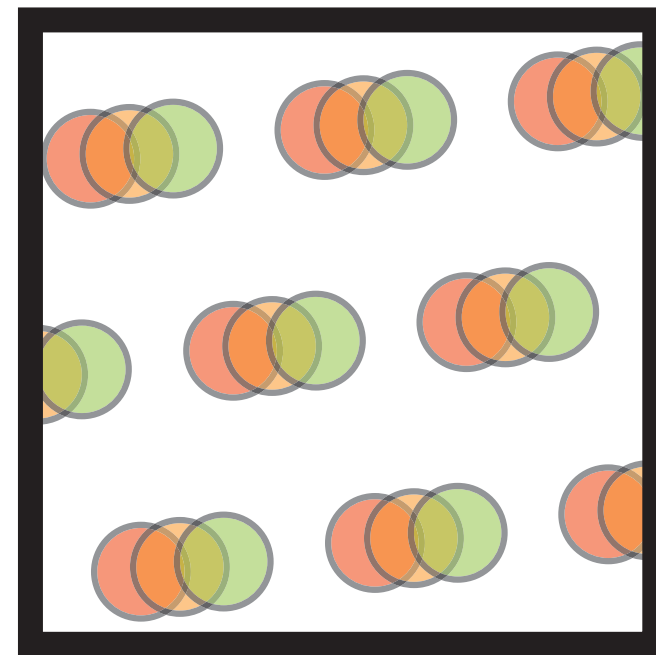


Cell 2



World 1

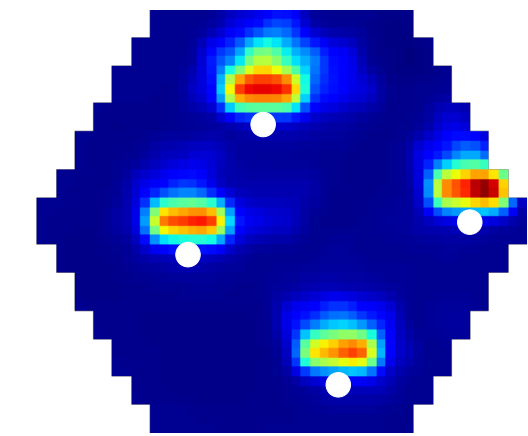
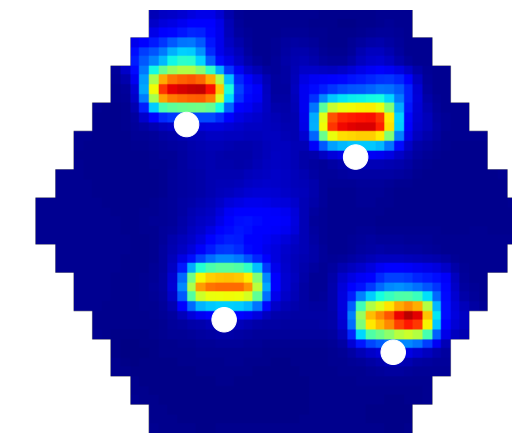
World 2



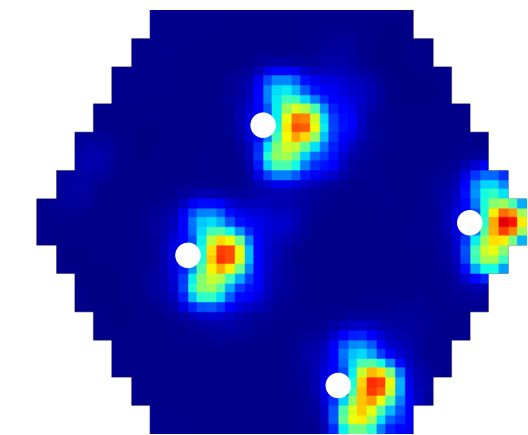
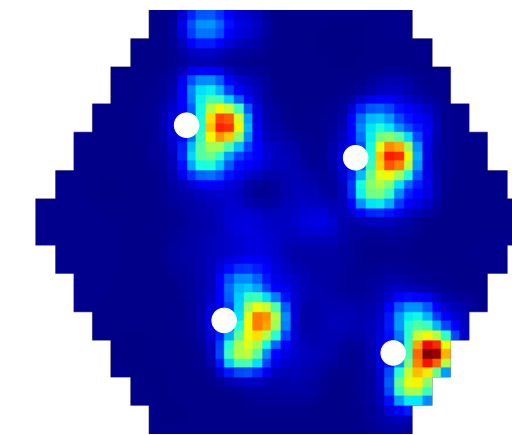
World 1

World 2

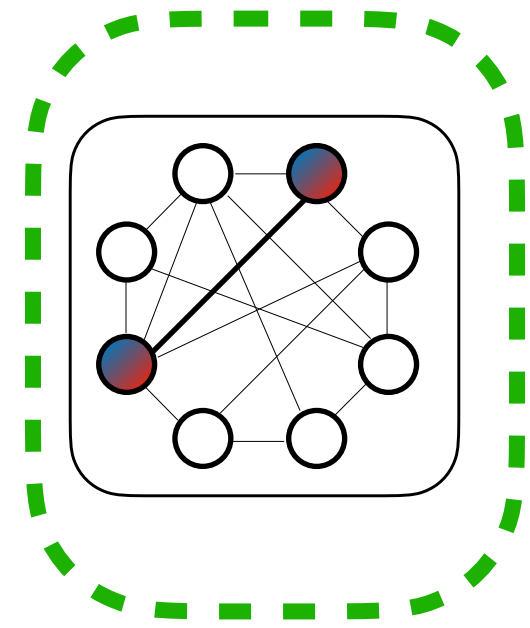
Cell 1



Cell 2

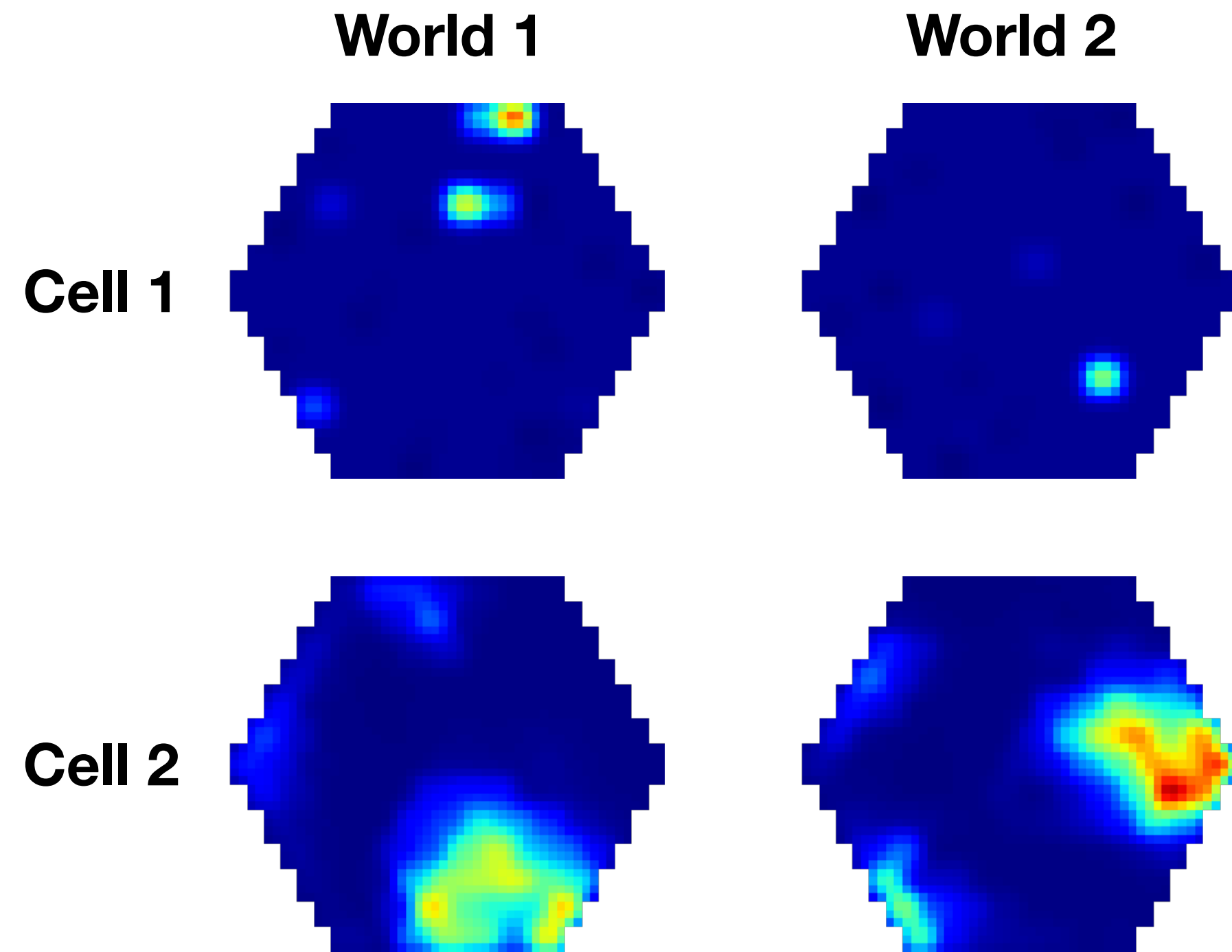


# Hippocampal remapping



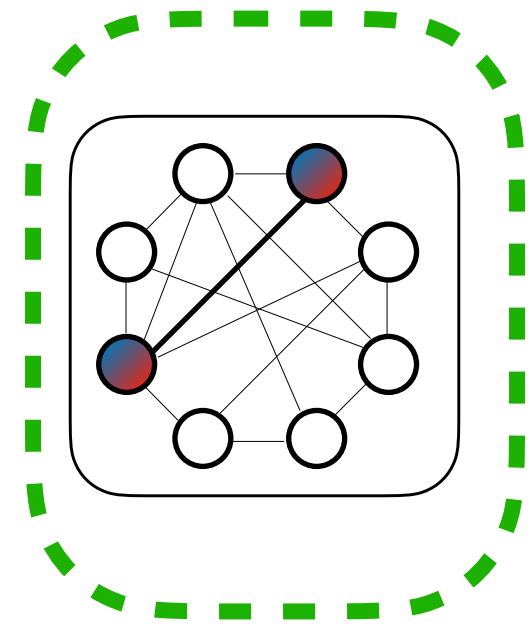
Hopfield Network

Model cells





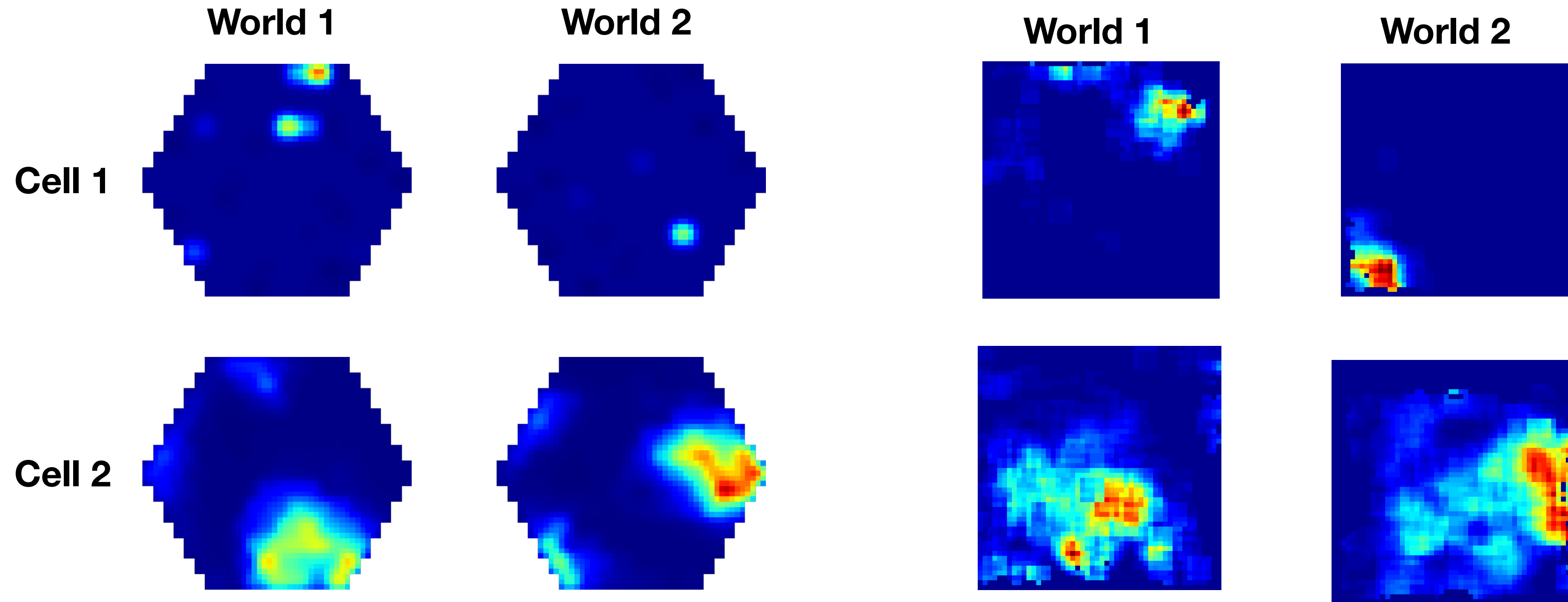
# Hippocampal remapping



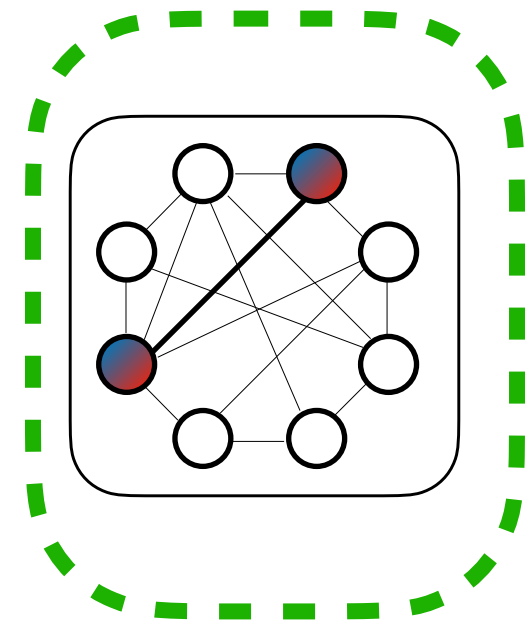
Hopfield Network

Model cells

Real cells



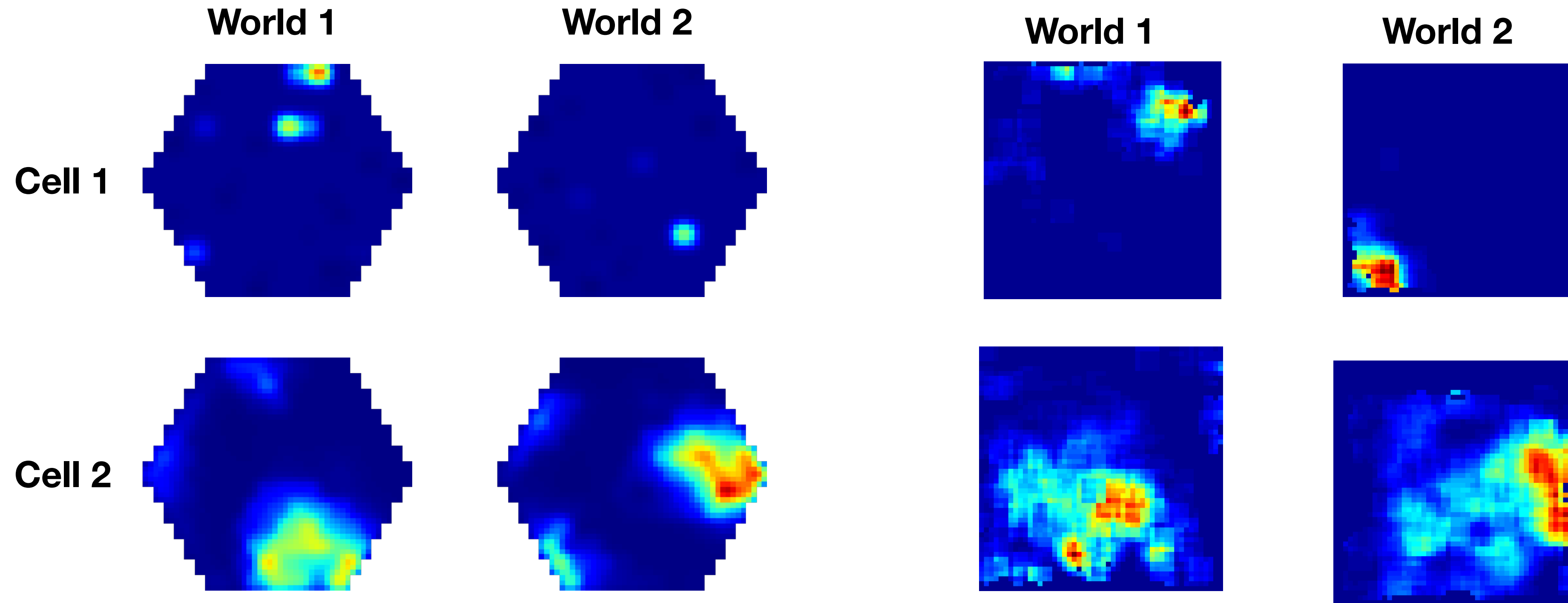
# Hippocampal remapping



Hopfield Network

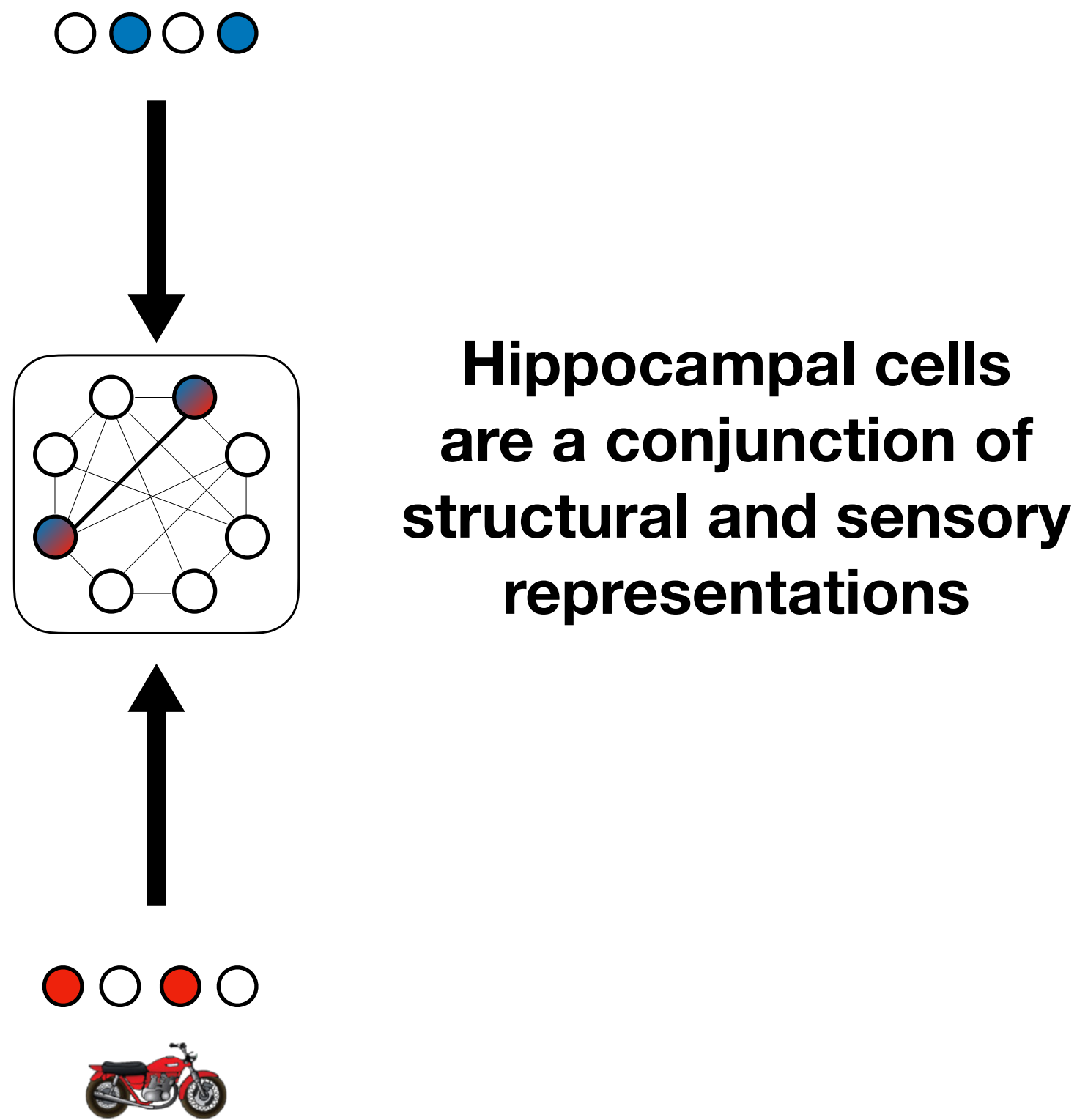
Model cells

Real cells

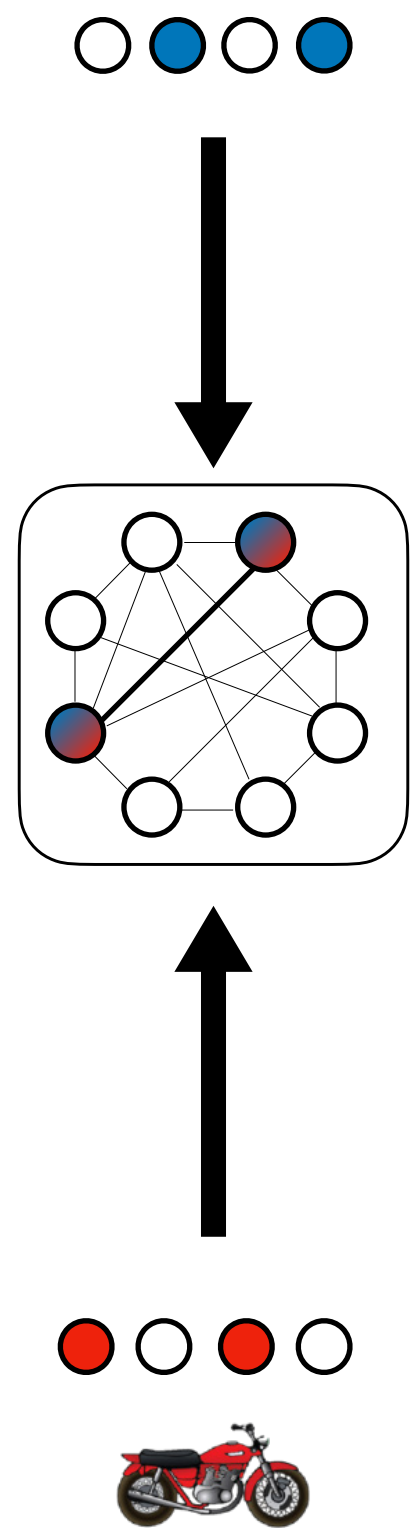


Looks random but is it really?

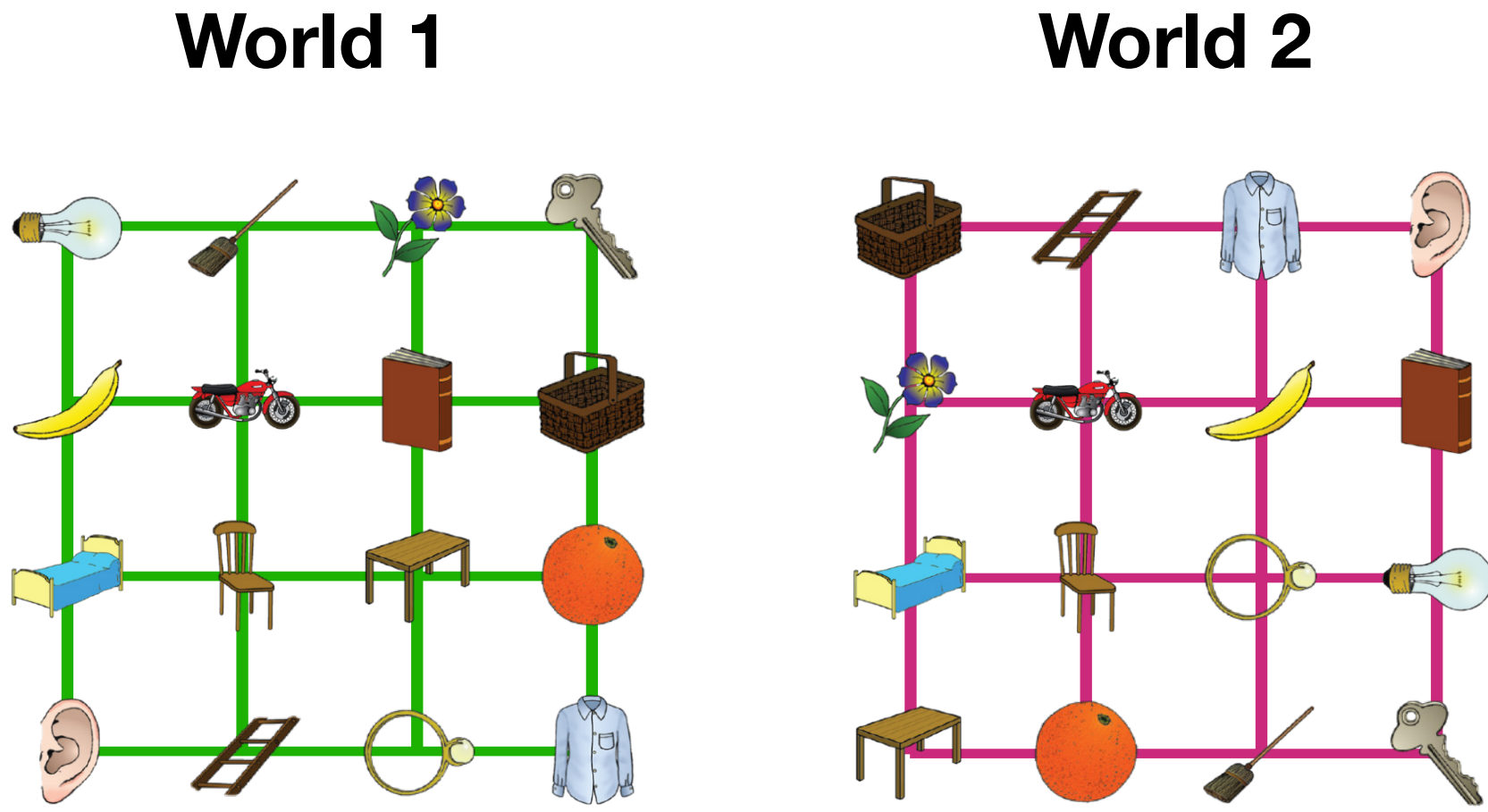
# Remember that place cells are compositions of building blocks



# Remember that place cells are compositions of building blocks

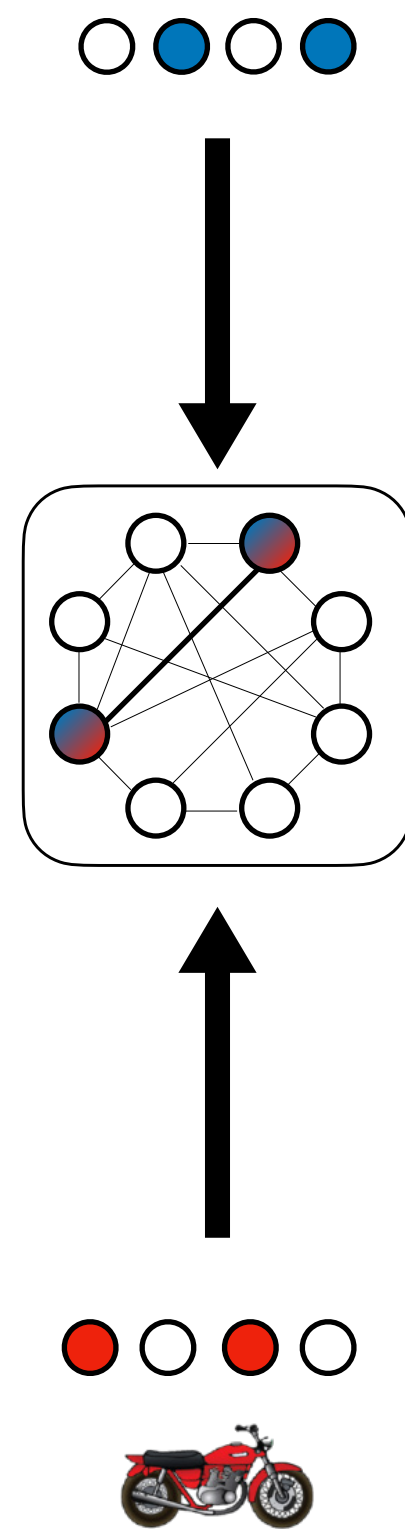


Hippocampal cells are a conjunction of structural and sensory representations

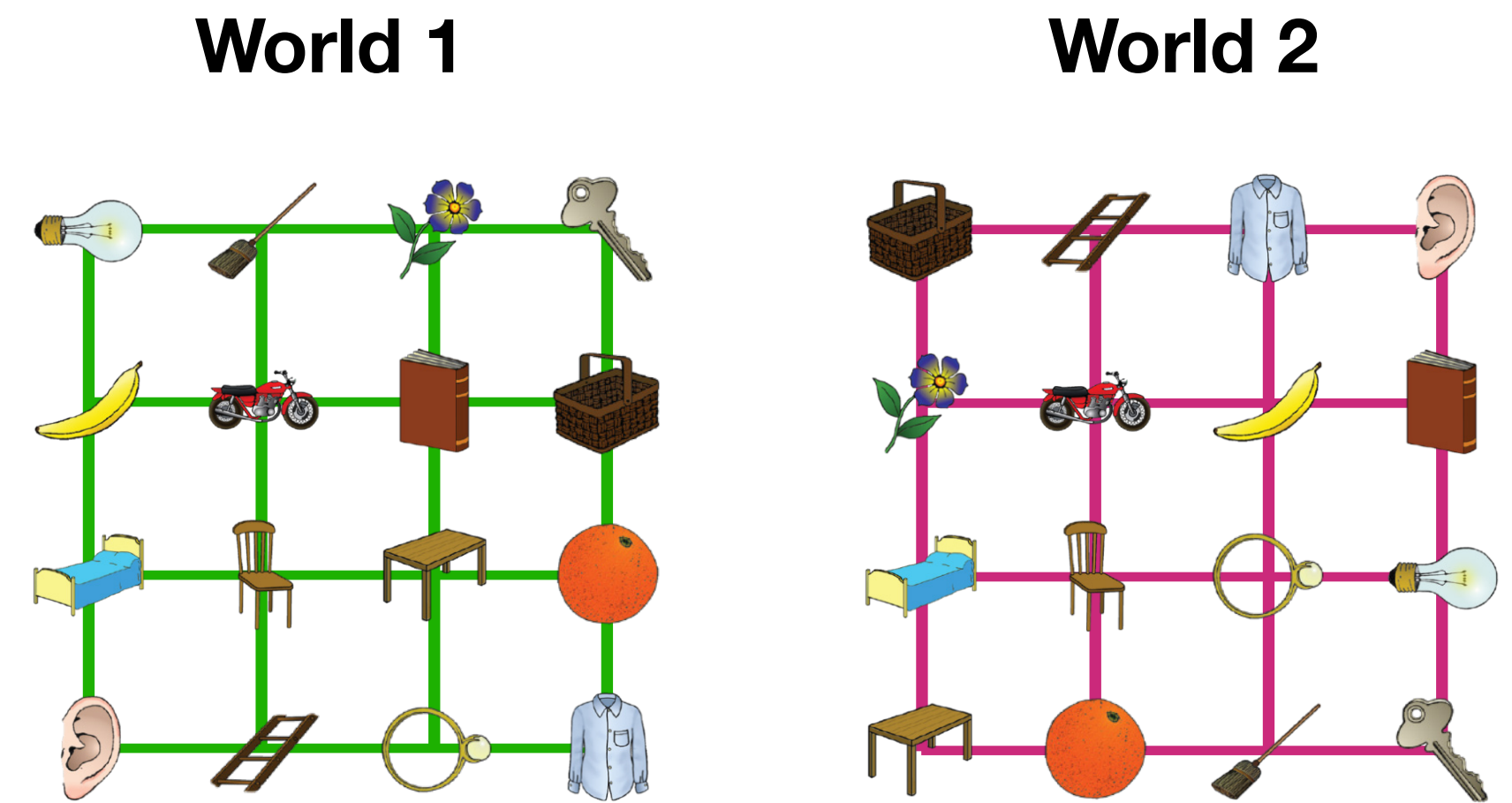


Building blocks are just in different combinations in different worlds

# Remember that place cells are compositions of building blocks



Hippocampal cells are a conjunction of structural and sensory representations

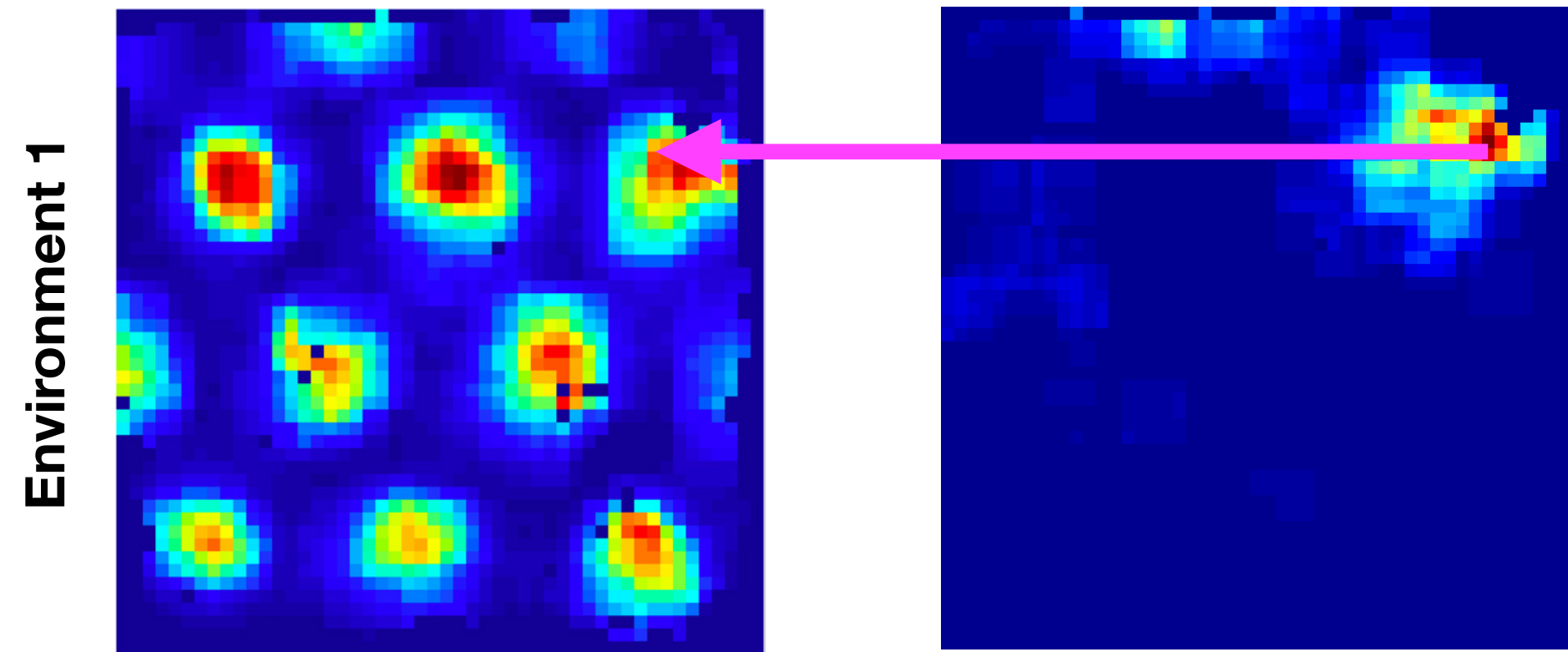


Building blocks are just in different combinations in different worlds

**Hippocampal cells should remap consistent with their building block inputs**

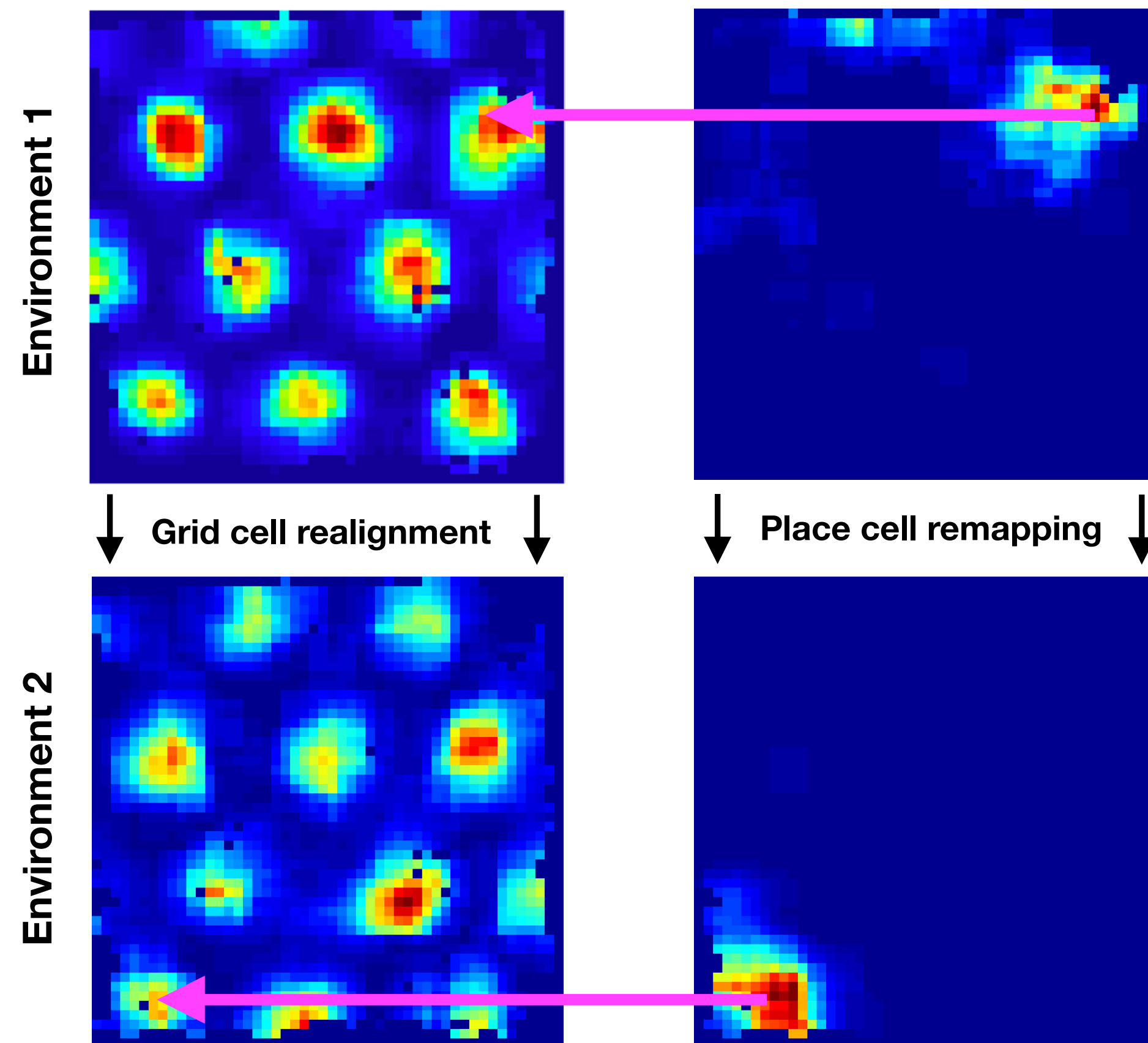
# Prediction: Grid cells control place cell remapping

Place cells locations are constrained by grid cells



# Prediction: Grid cells control place cell remapping

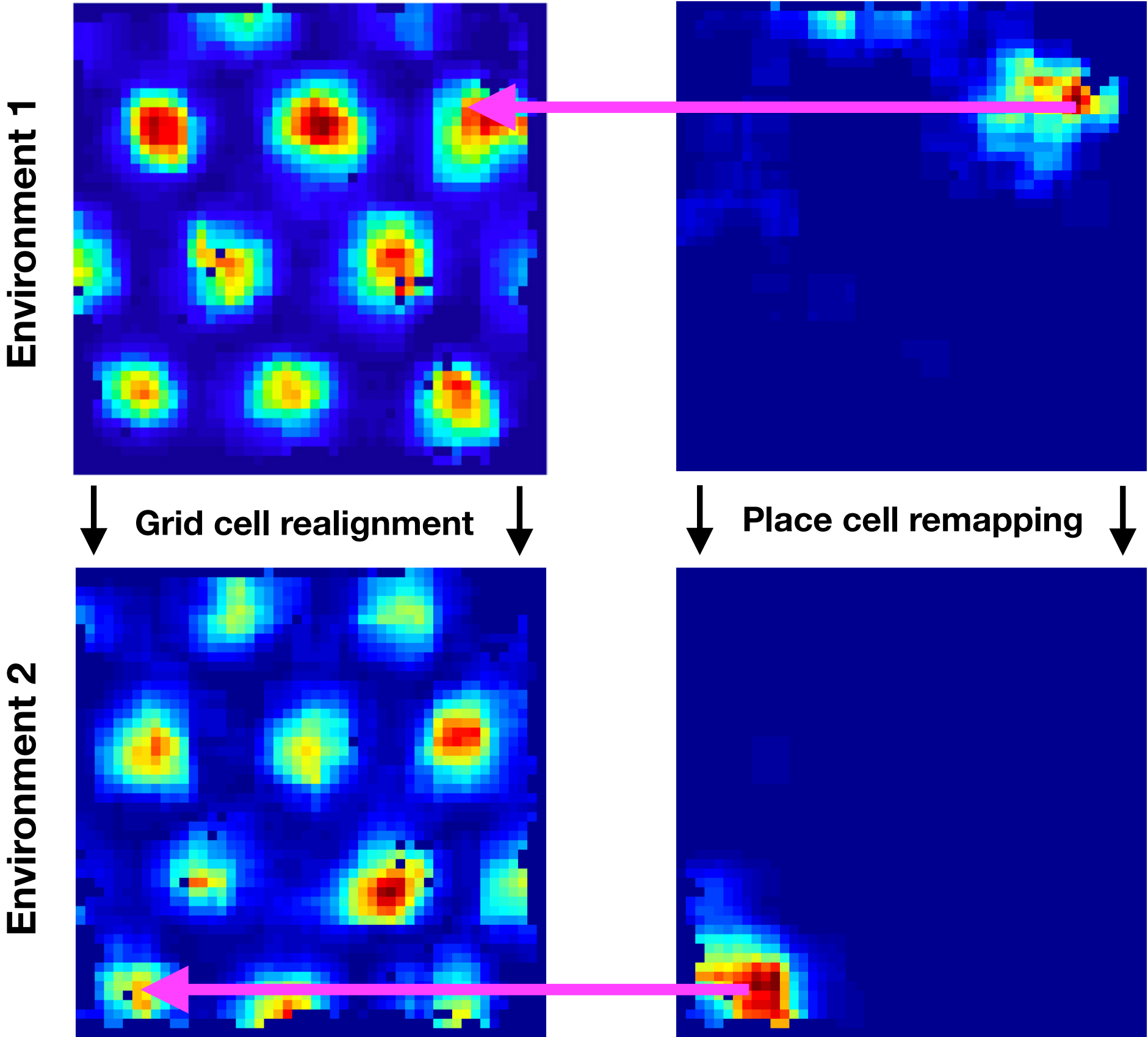
Place cells locations are constrained by grid cells



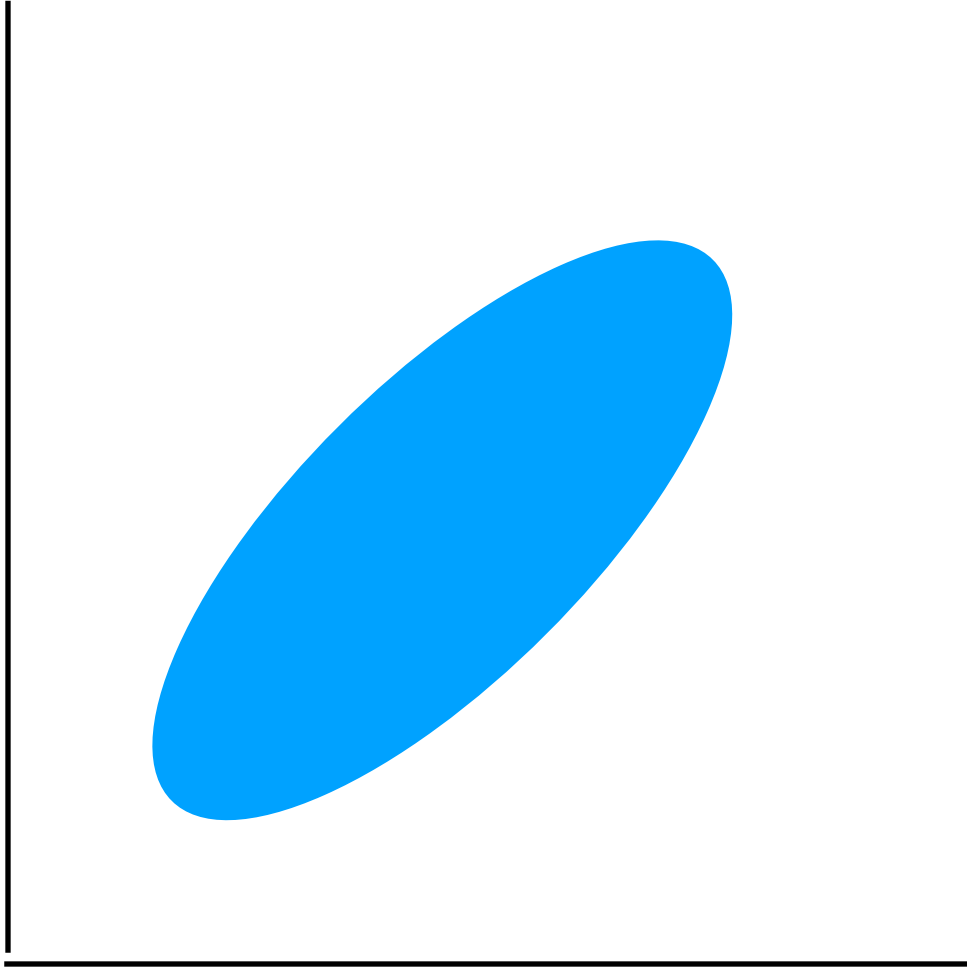
So they should remap to locations consistent with that grid cell

# Prediction: Grid cells control place cell remapping

Place cells locations are constrained by grid cells



grid cell firing at place field env 2

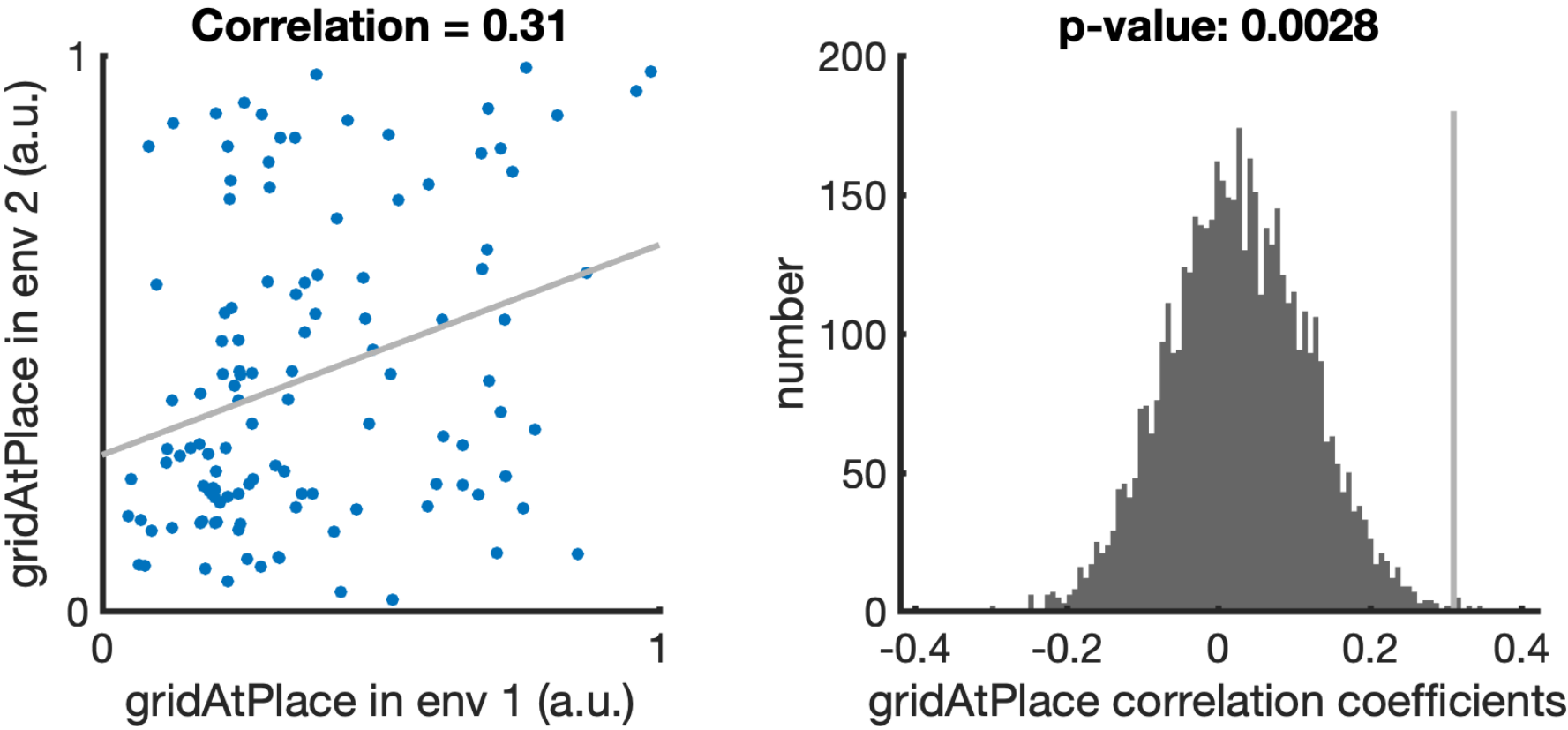


grid cell firing at place field env 1

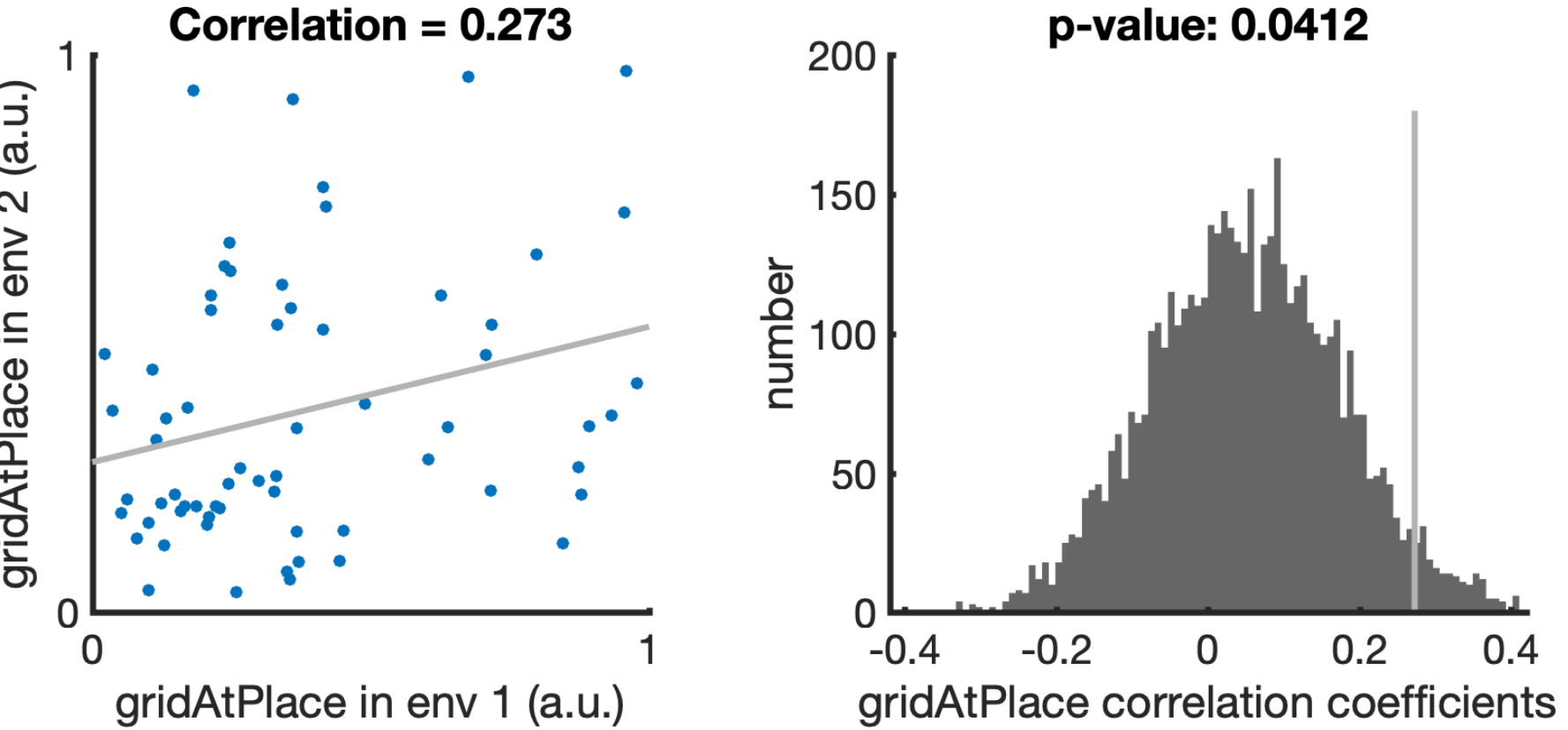
So they should remap to locations consistent with that grid cell



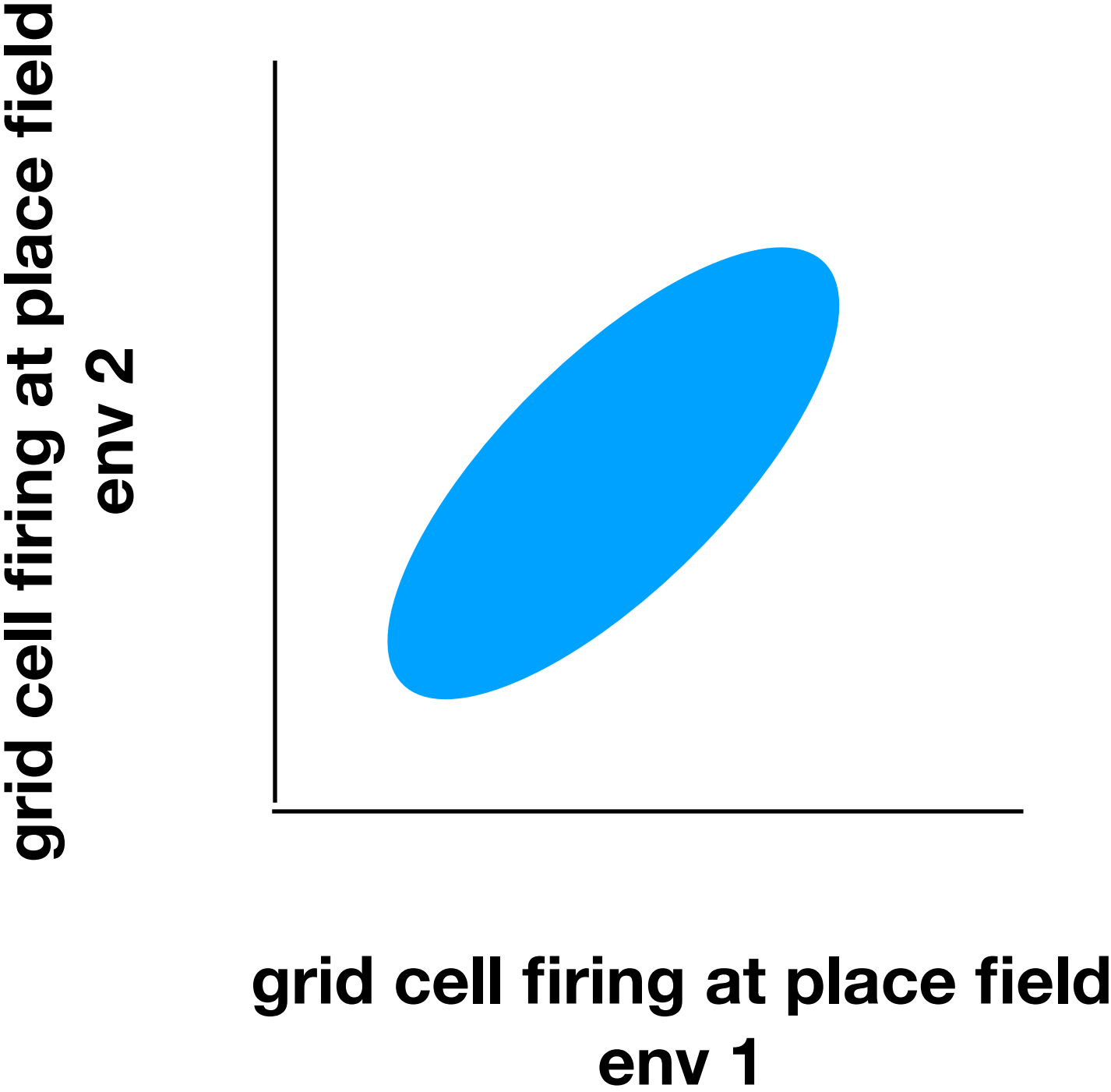
# Prediction: Grid cells control place cell remapping



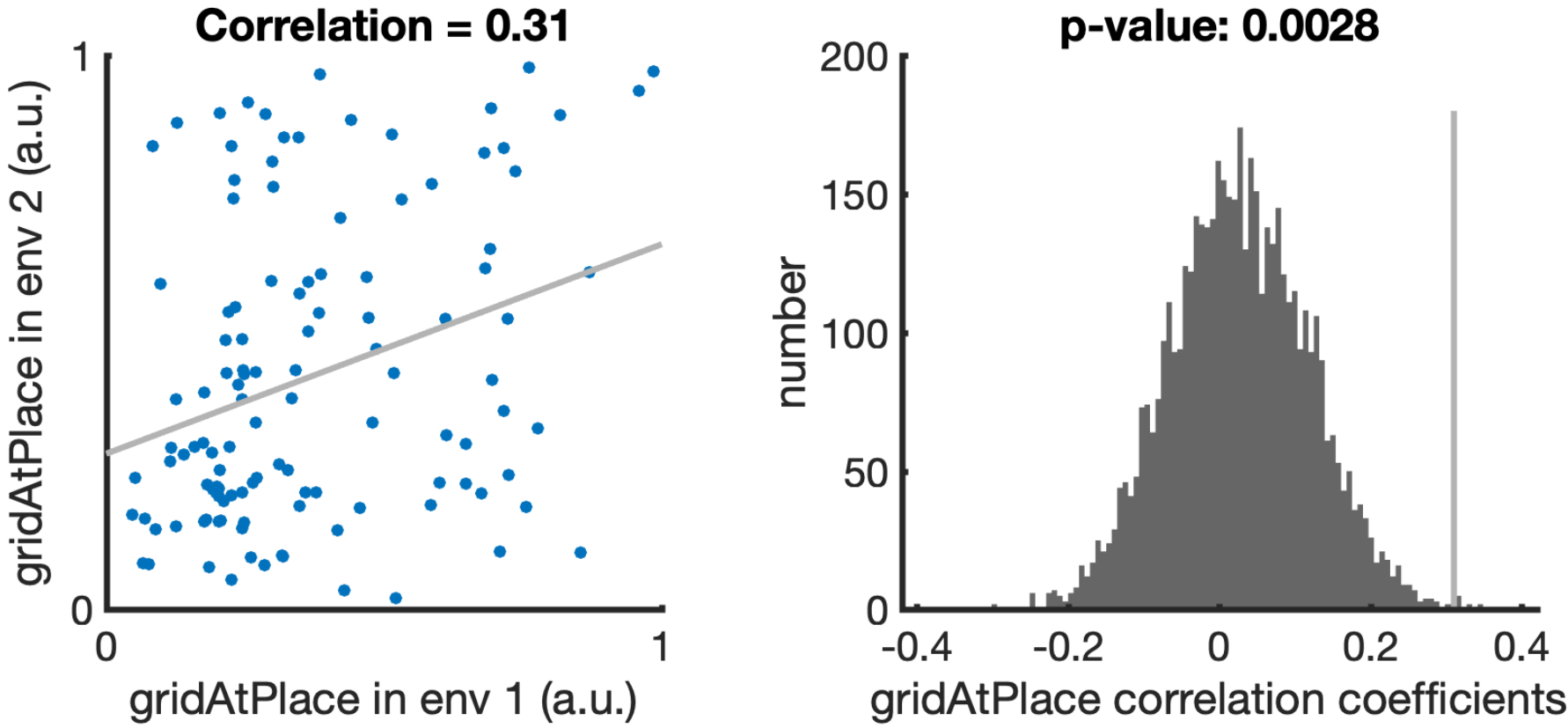
Data from Barry et al Current Biol. 2007



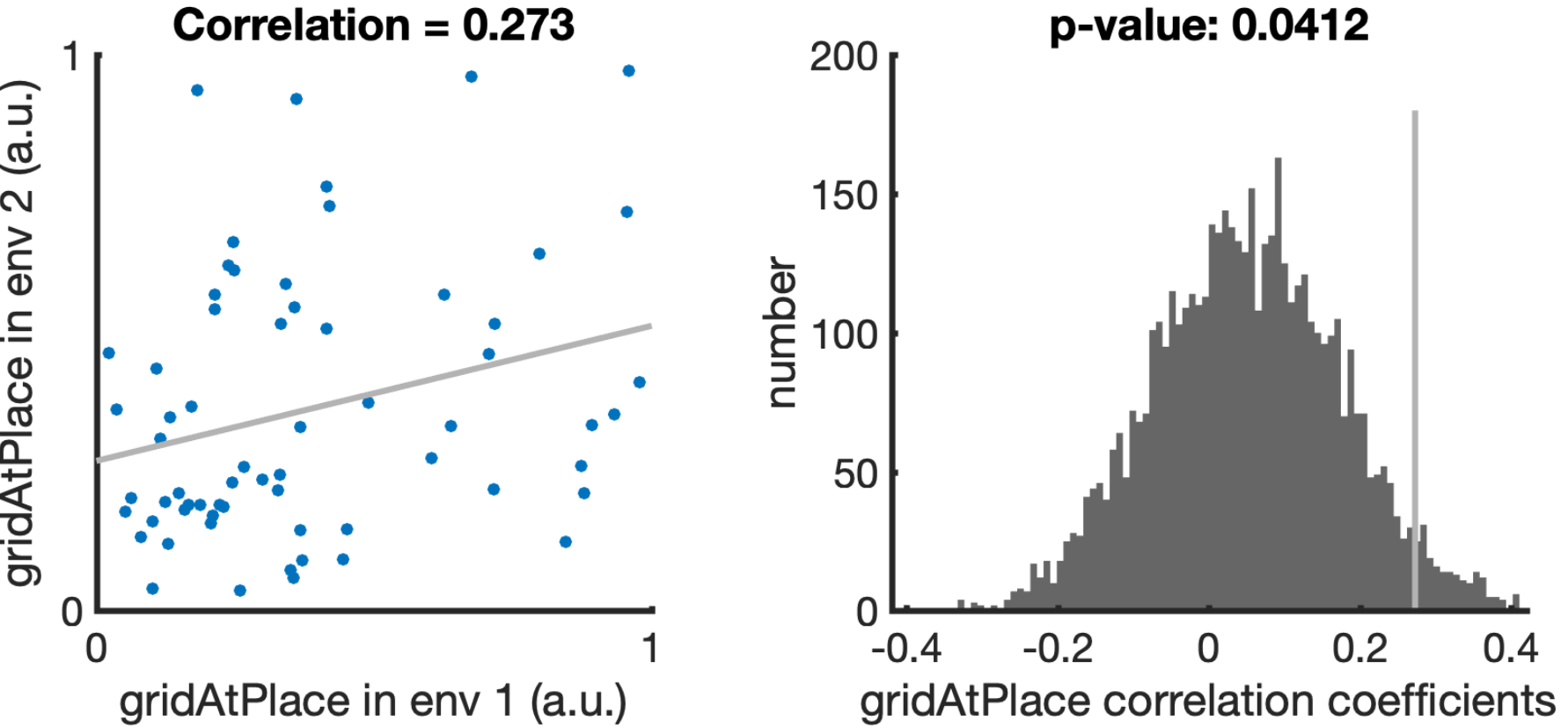
Data from Chen et al eLife 2018



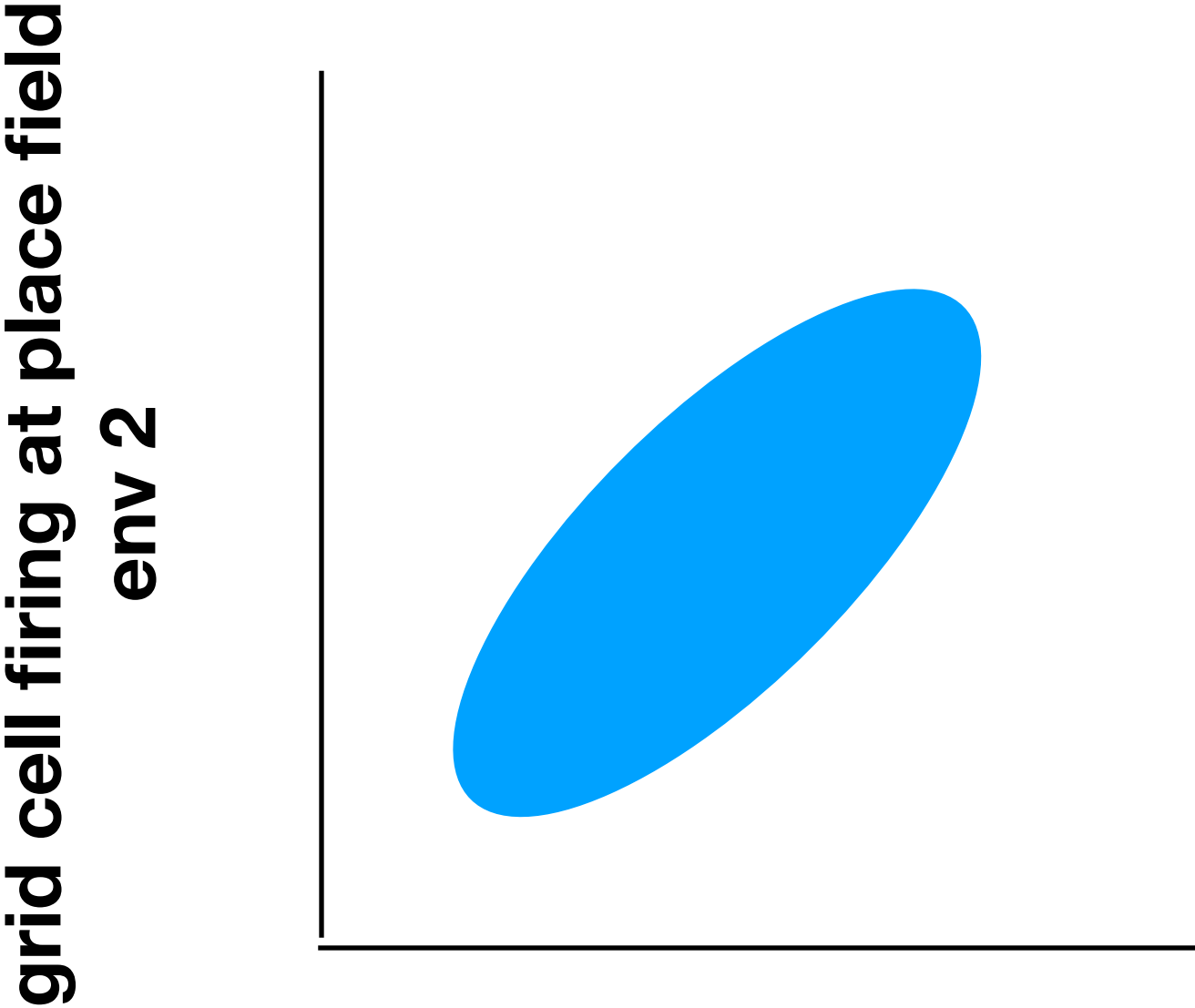
# Prediction: Grid cells control place cell remapping



Data from Barry et al Current Biol. 2007



Data from Chen et al eLife 2018



grid cell firing at place field env 2

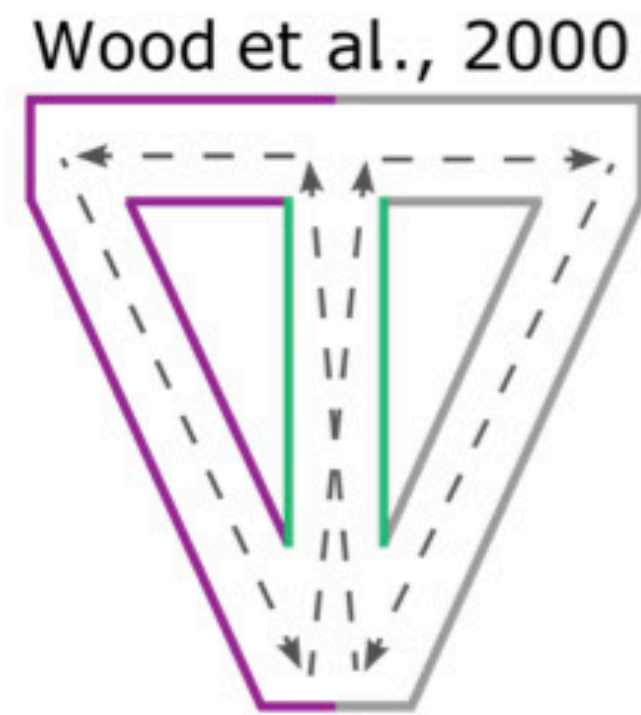
grid cell firing at place field env 1

**Place cells don't remap randomly as previously thought!**

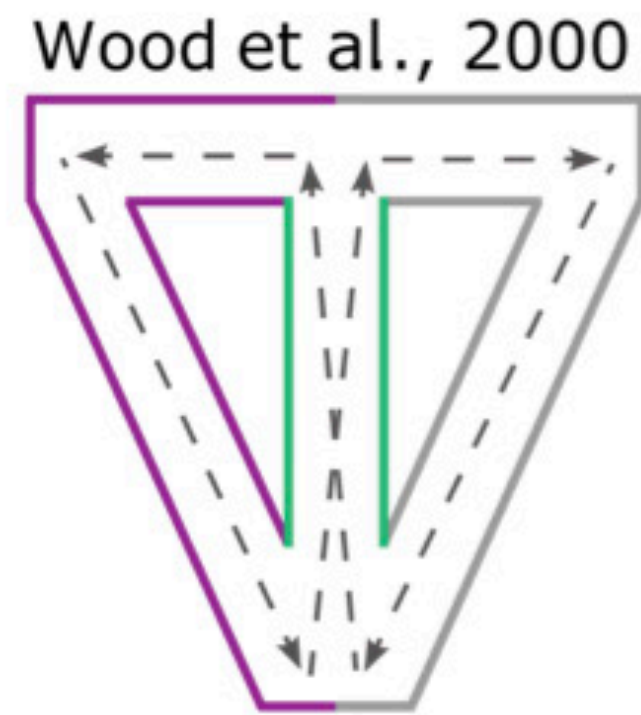
**Representing more complex tasks**

**Lots of other hippocampal representations accounted for**

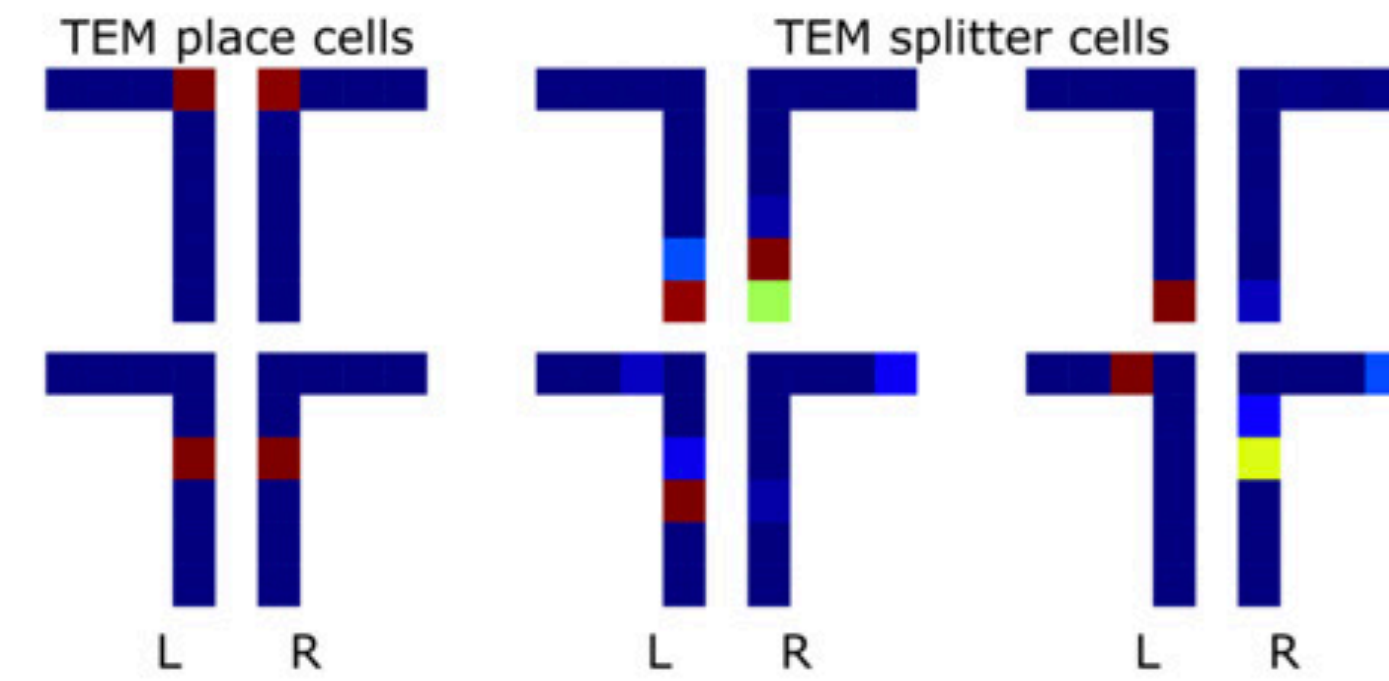
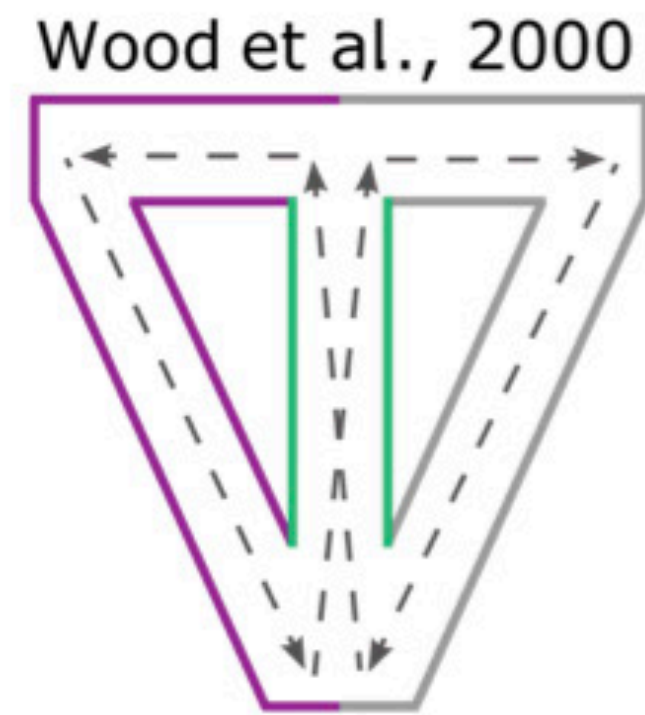
# Lots of other hippocampal representations accounted for



# Lots of other hippocampal representations accounted for

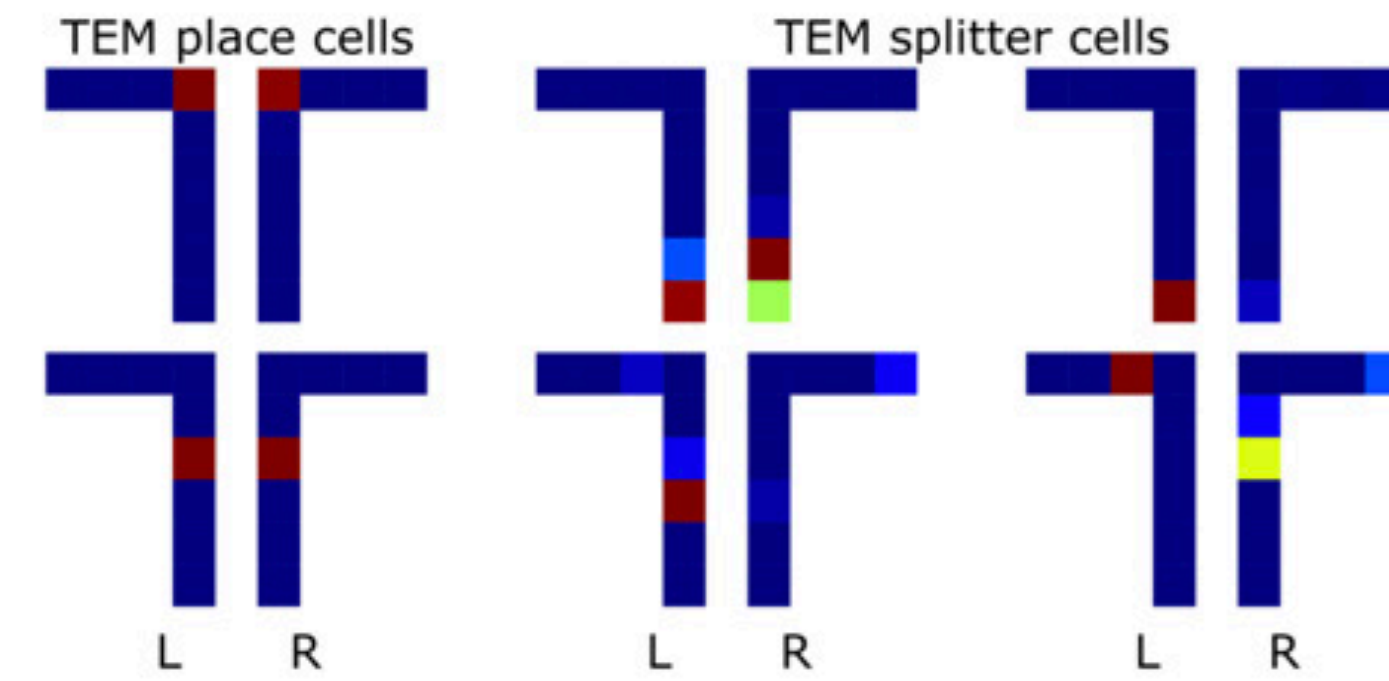
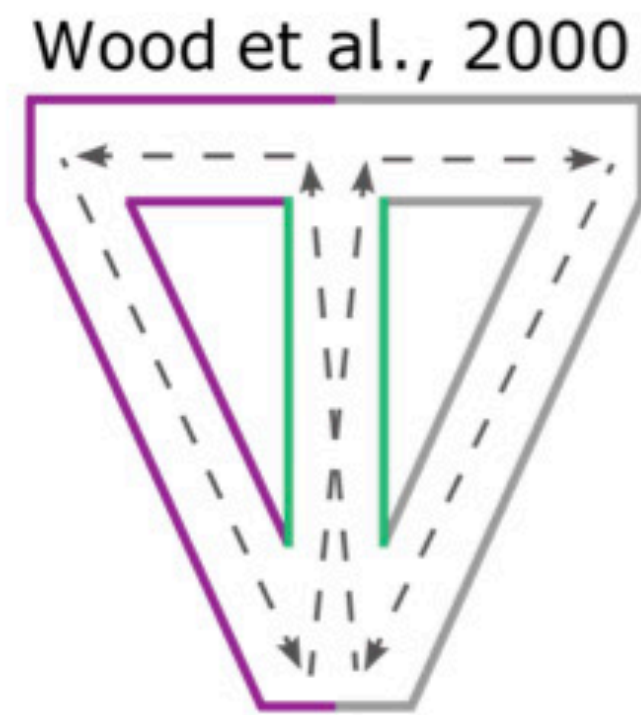


# Lots of other hippocampal representations accounted for



# Lots of other hippocampal representations accounted for

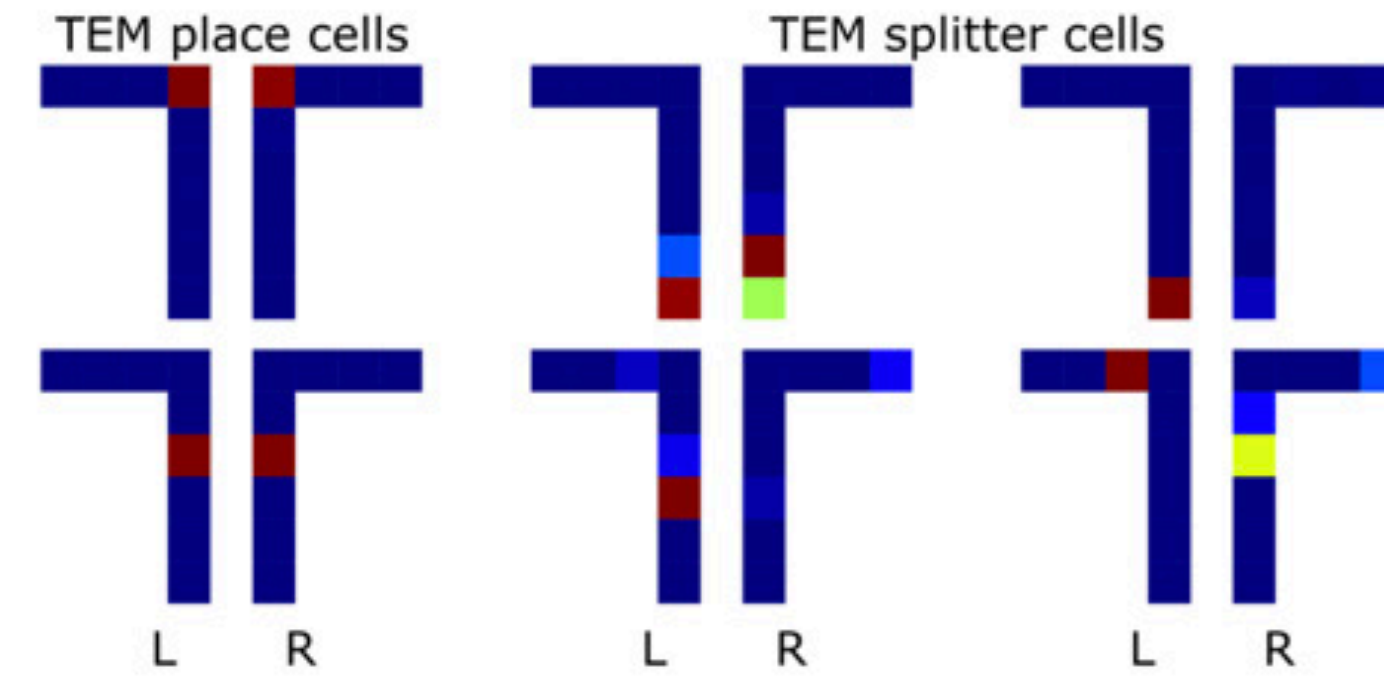
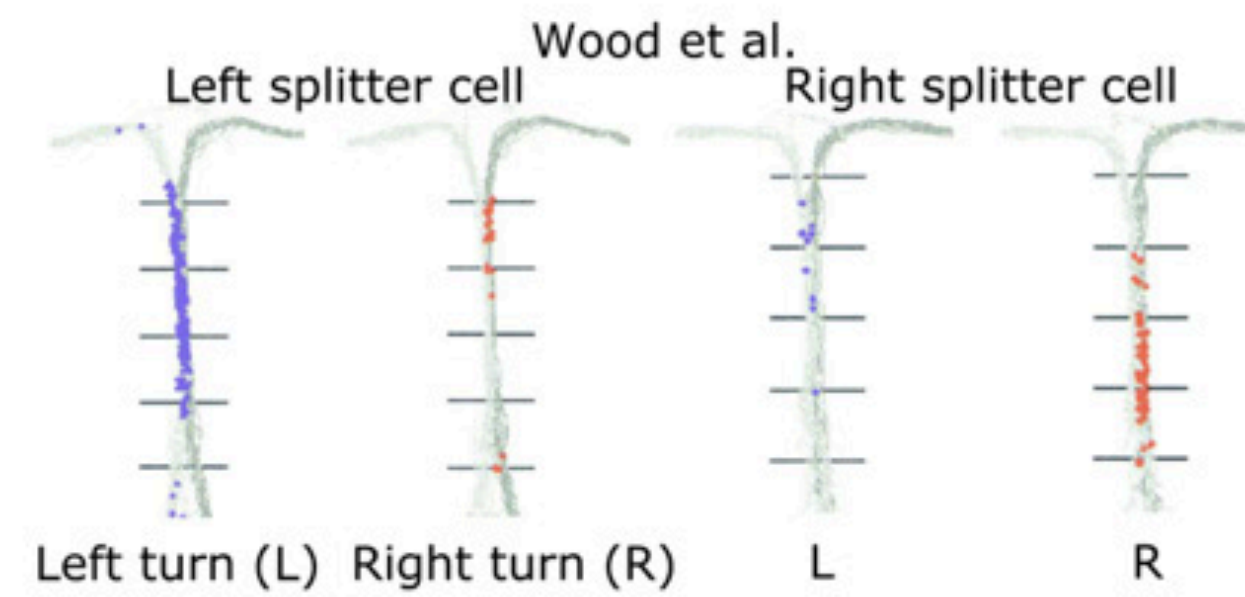
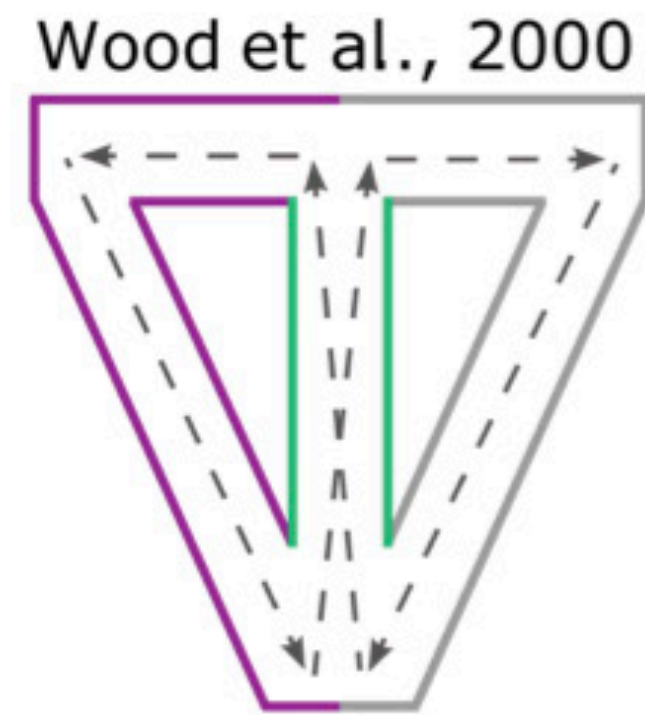
Simultaneous representation of spatial and task location





# Lots of other hippocampal representations accounted for

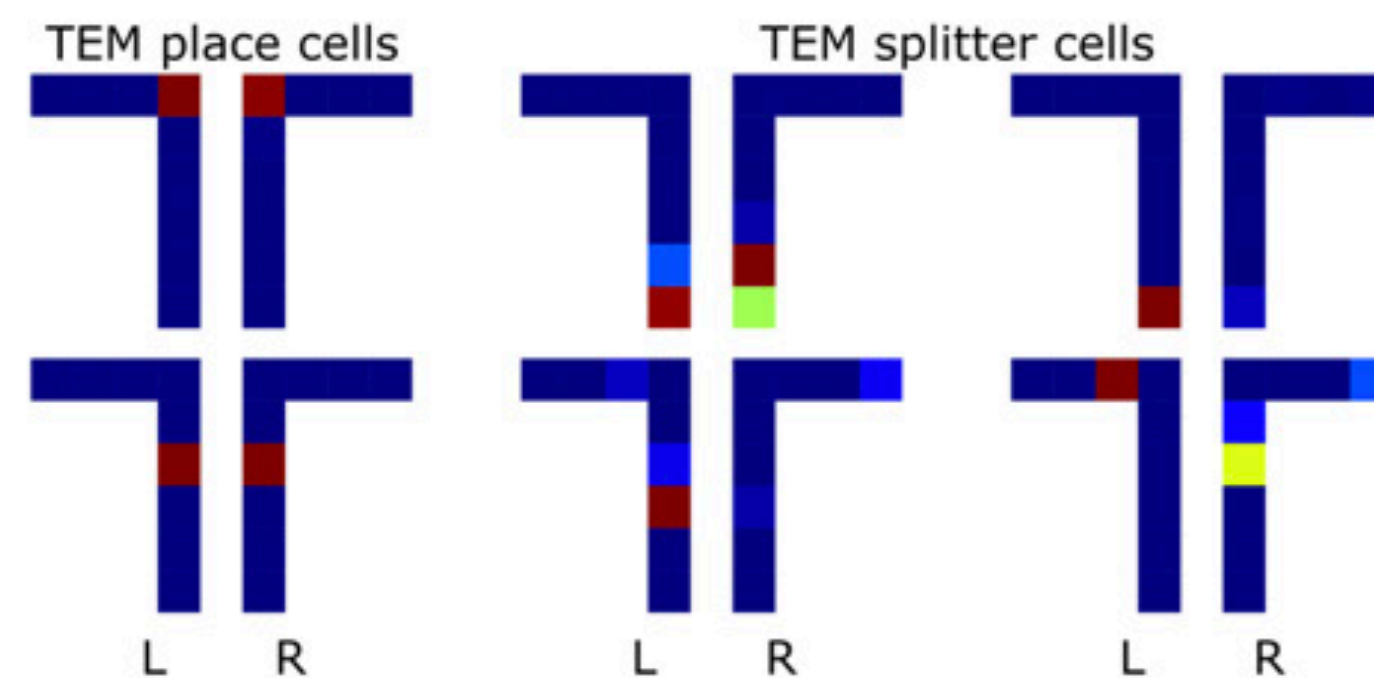
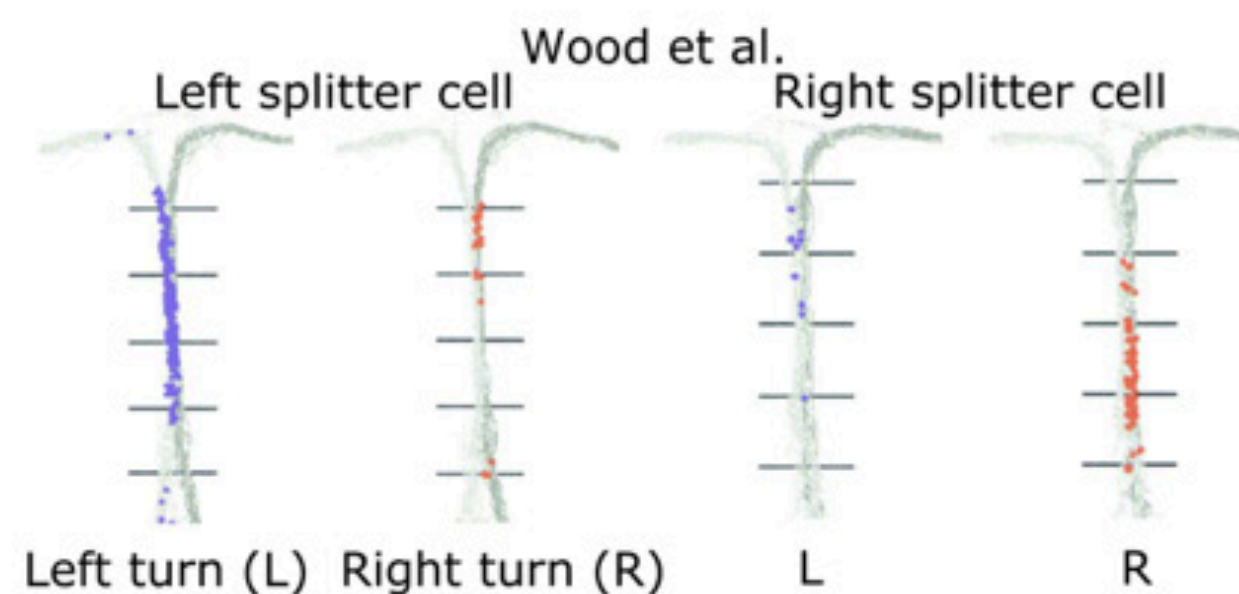
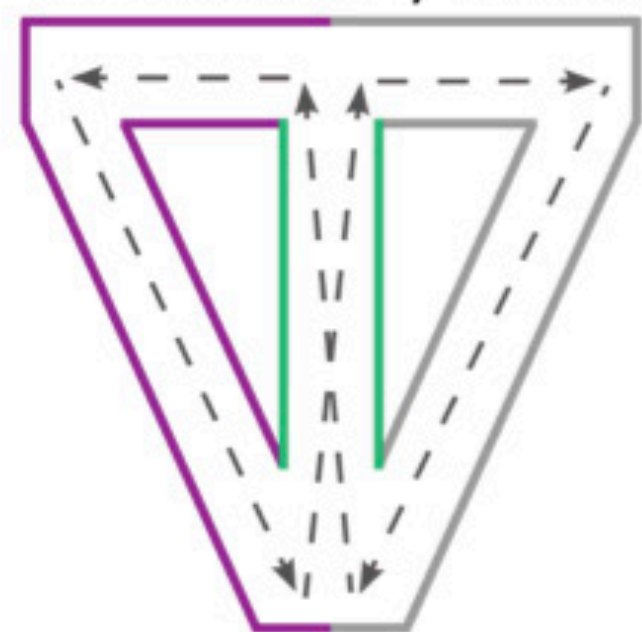
Simultaneous representation of spatial and task location



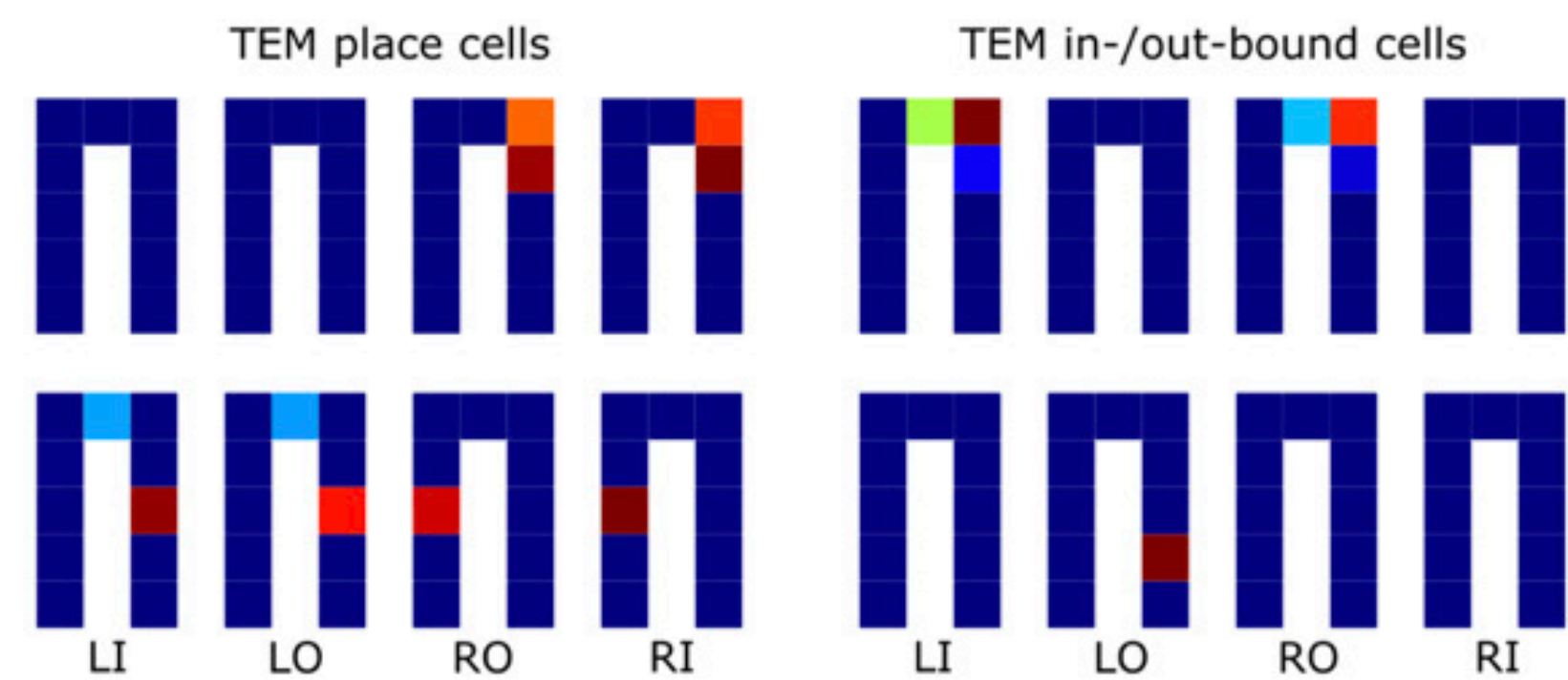
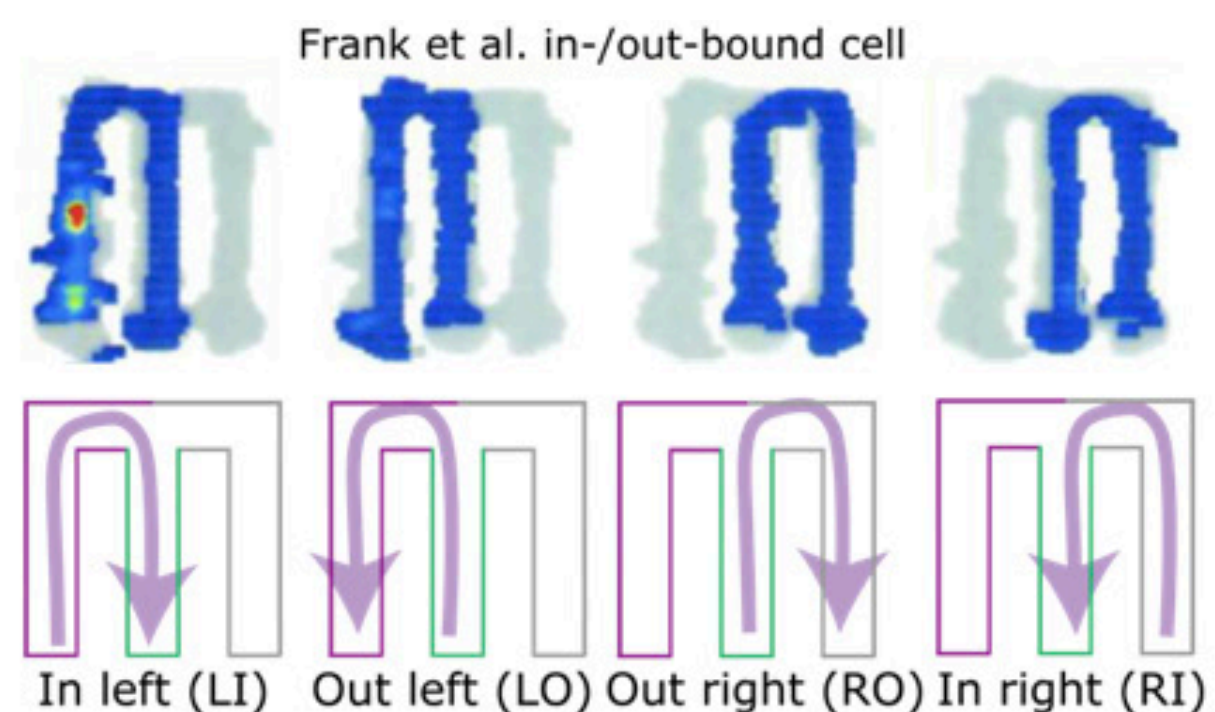
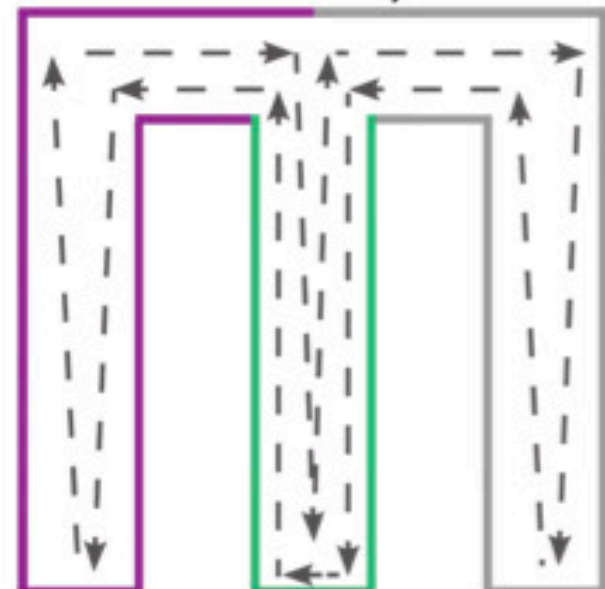
# Lots of other hippocampal representations accounted for

## Simultaneous representation of spatial and task location

Wood et al., 2000



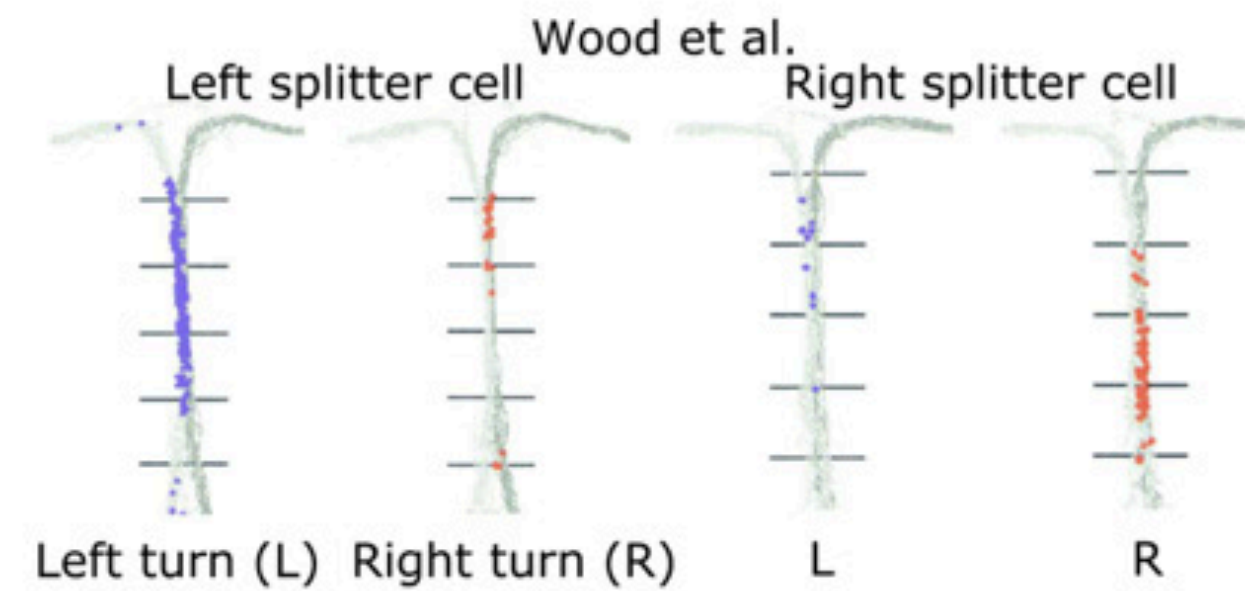
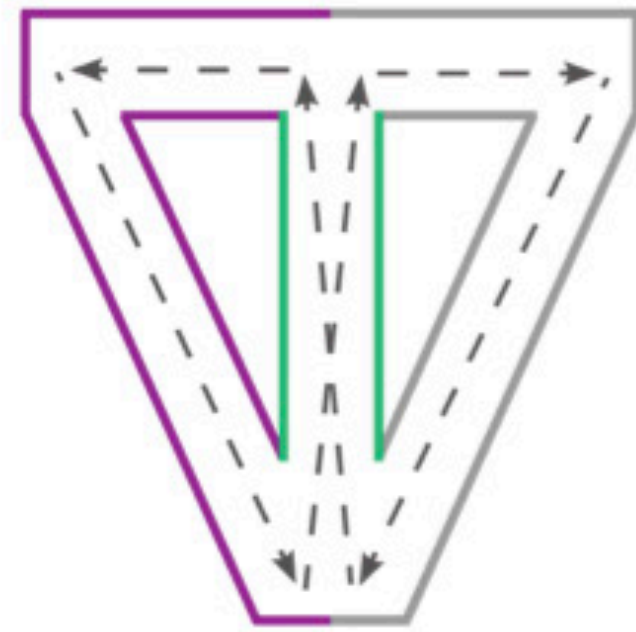
Frank et al., 2000



# Lots of other hippocampal representations accounted for

Simultaneous representation of spatial and task location

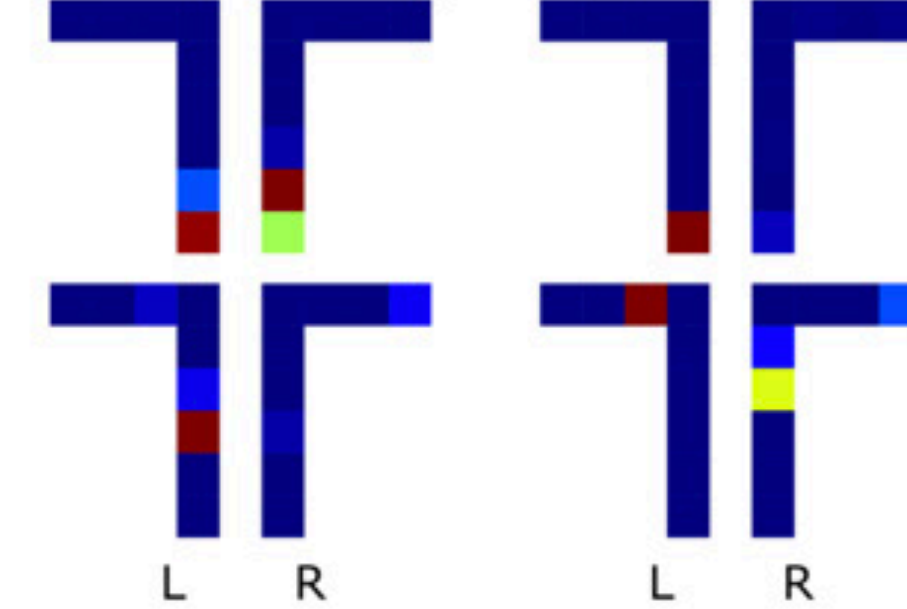
Wood et al., 2000



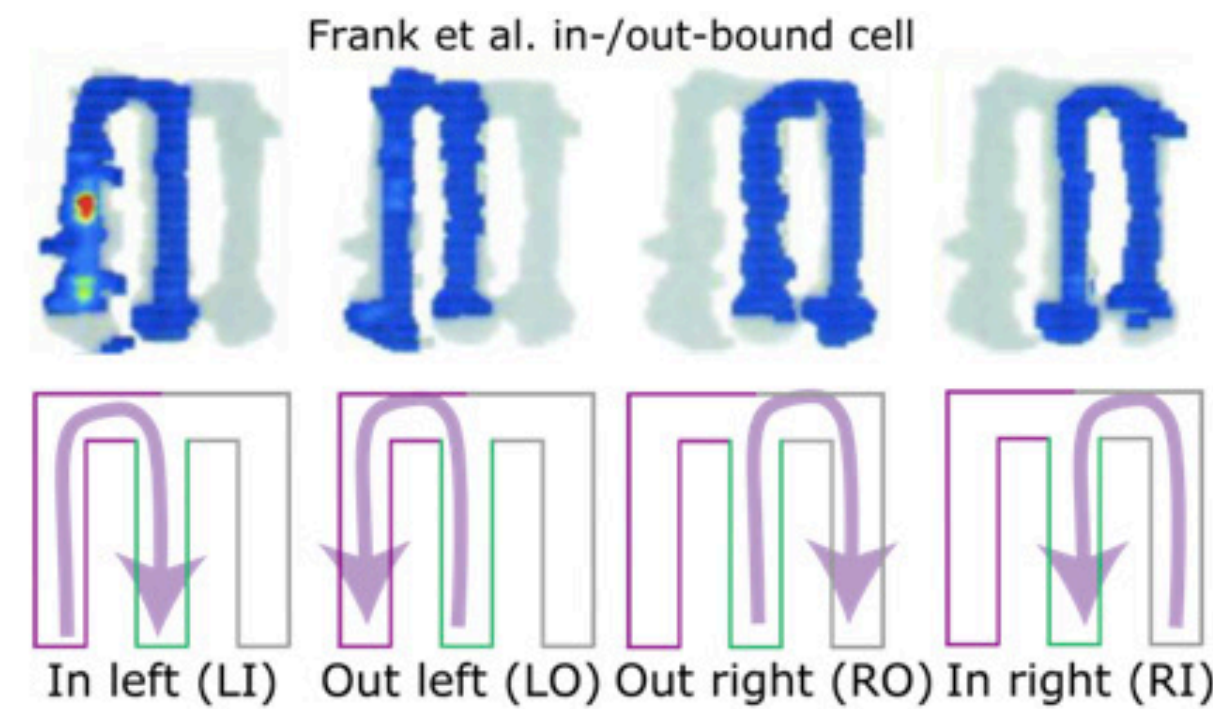
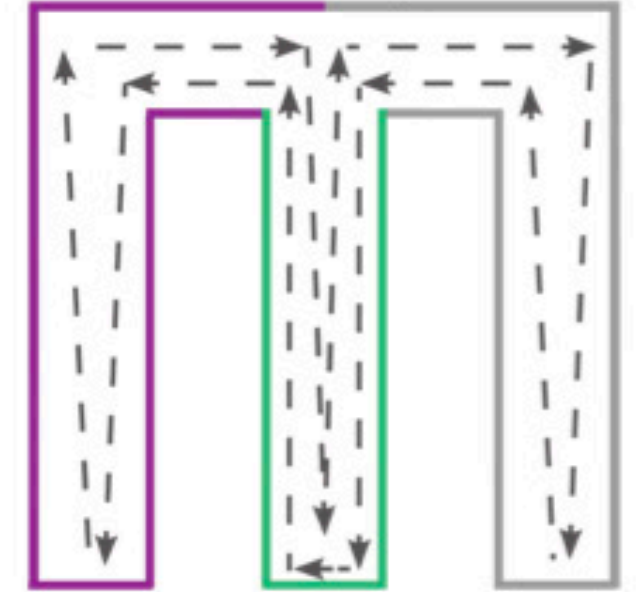
TEM place cells



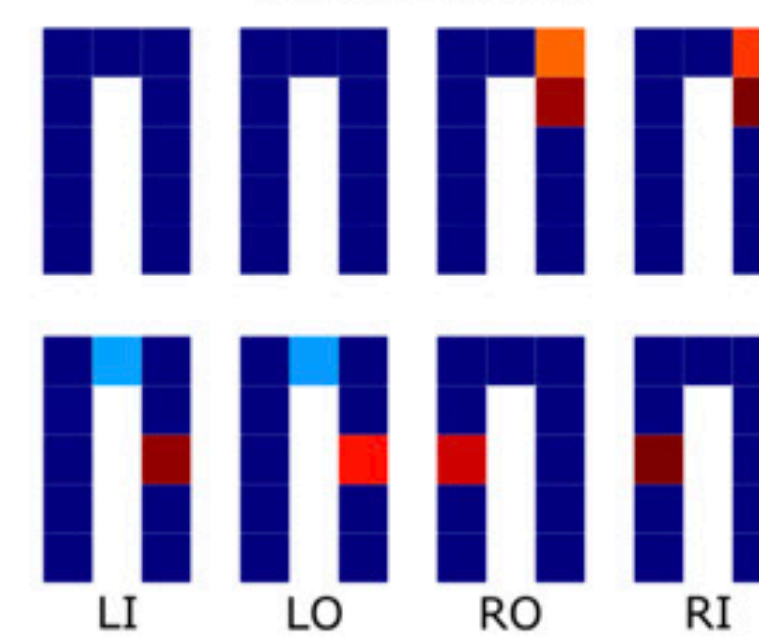
TEM splitter cells



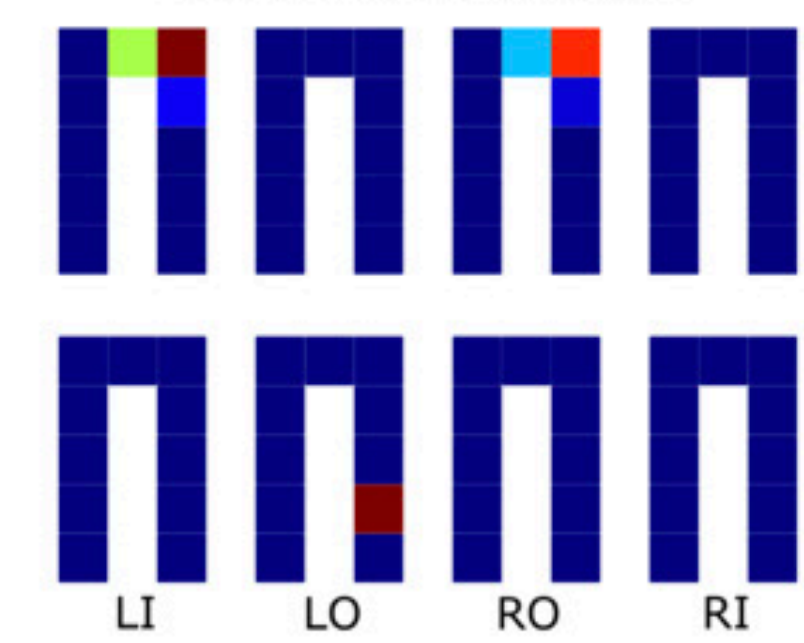
Frank et al., 2000



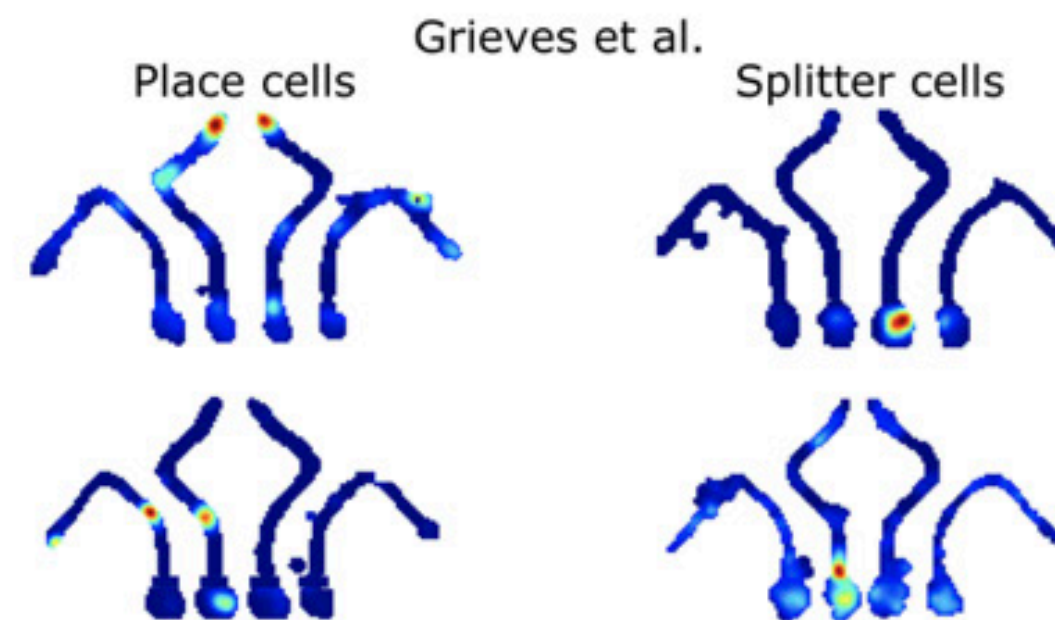
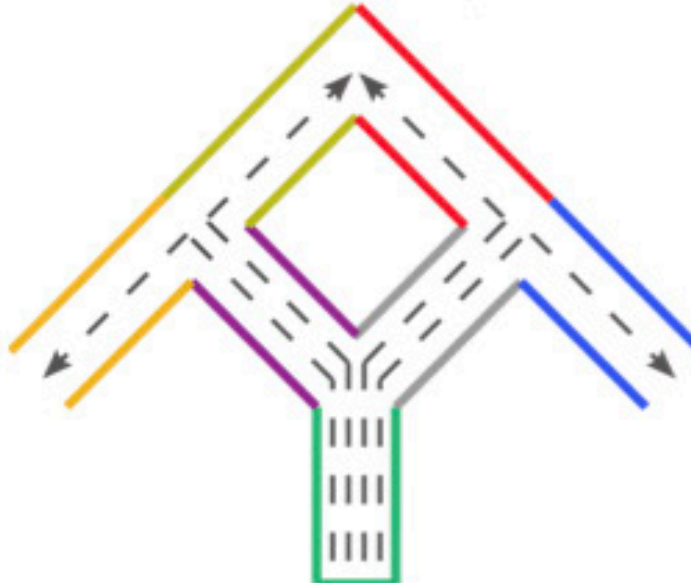
TEM place cells



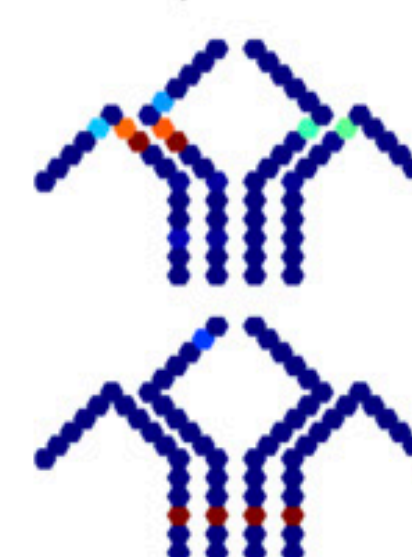
TEM in-/out-bound cells



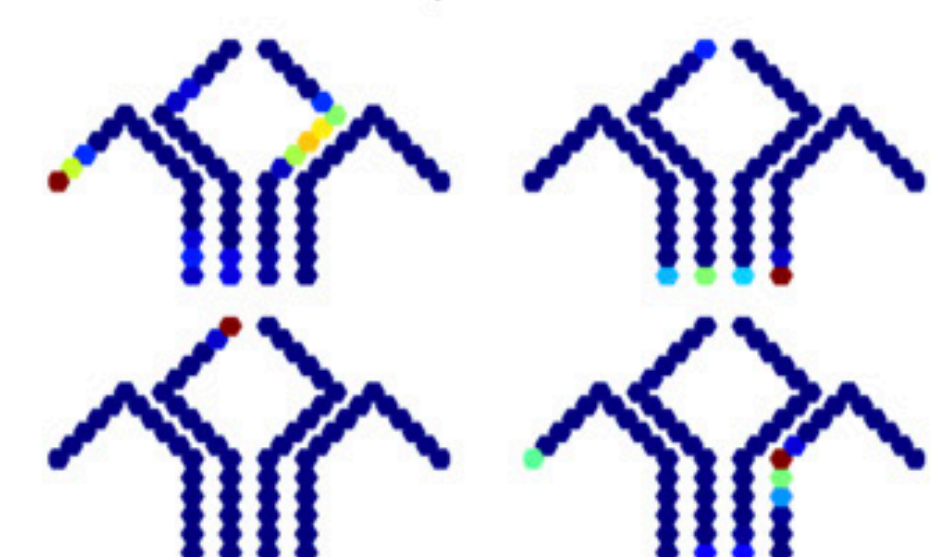
Grieves et al., 2016



TEM place cells



TEM splitter cells

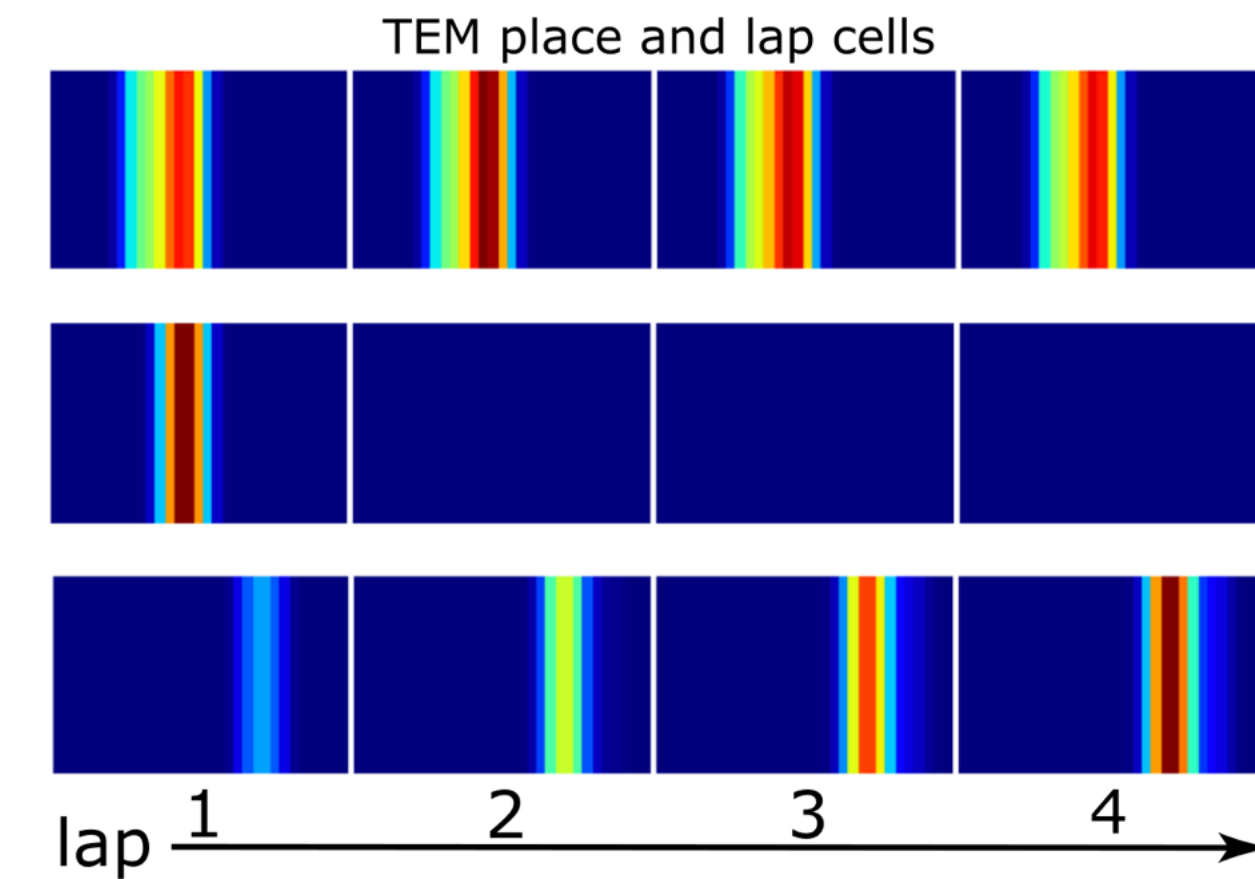
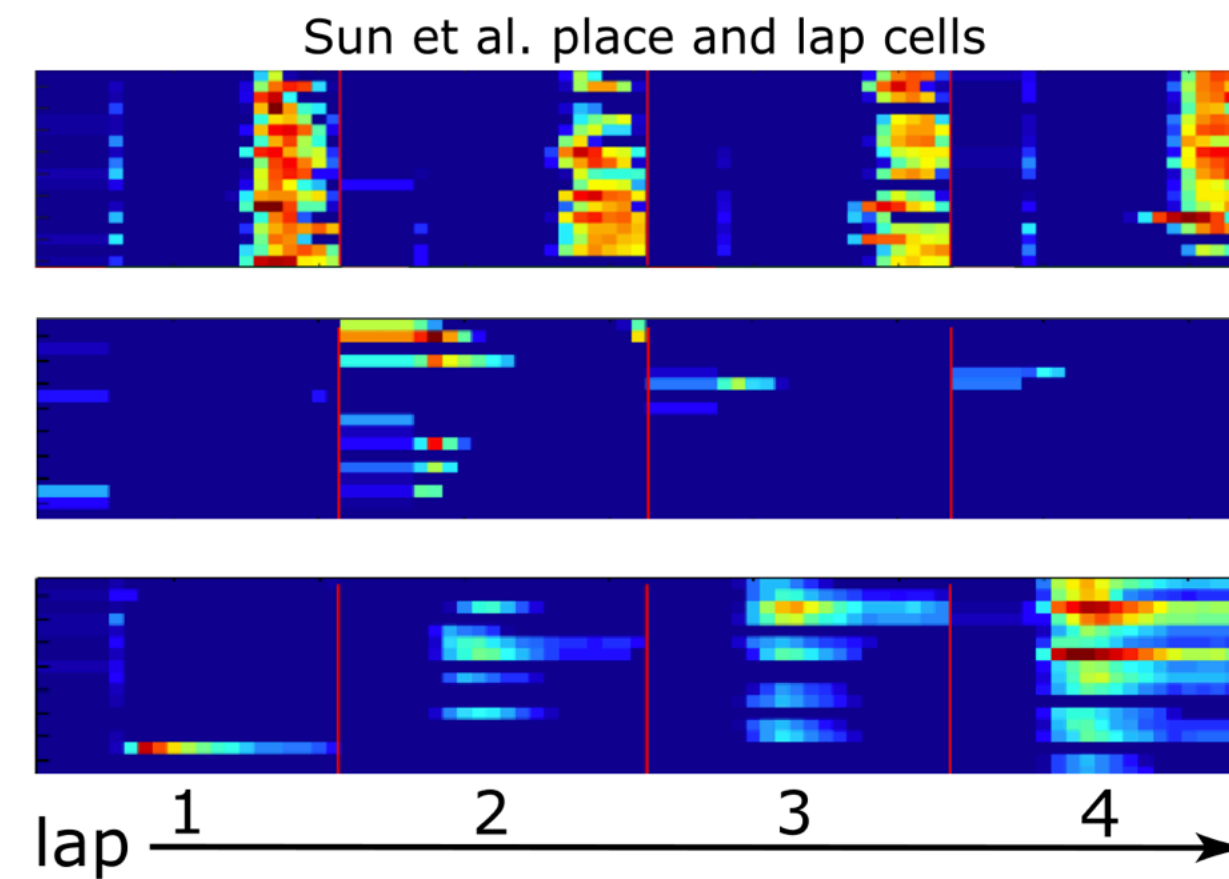
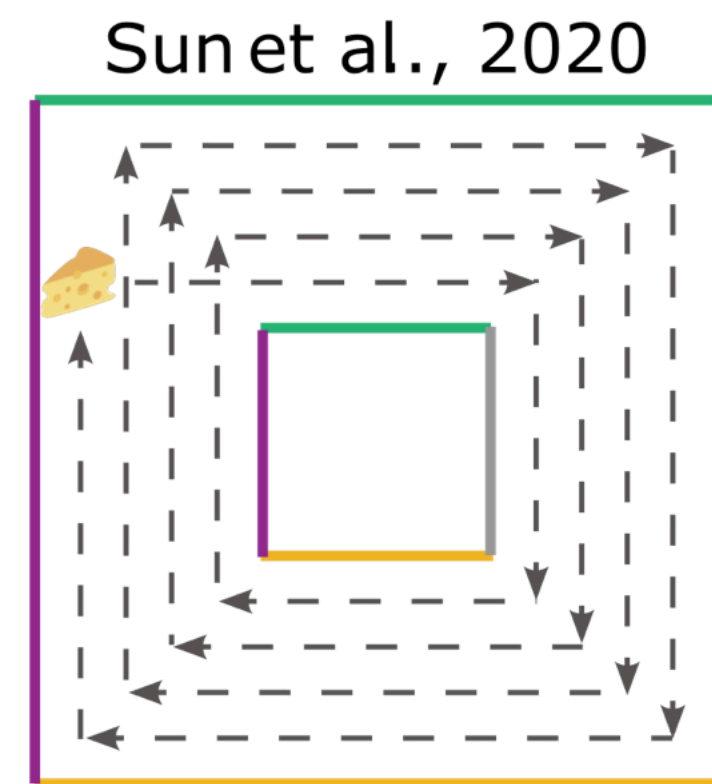


# Lots of other hippocampal representations accounted for

Simultaneous representation of spatial and task location

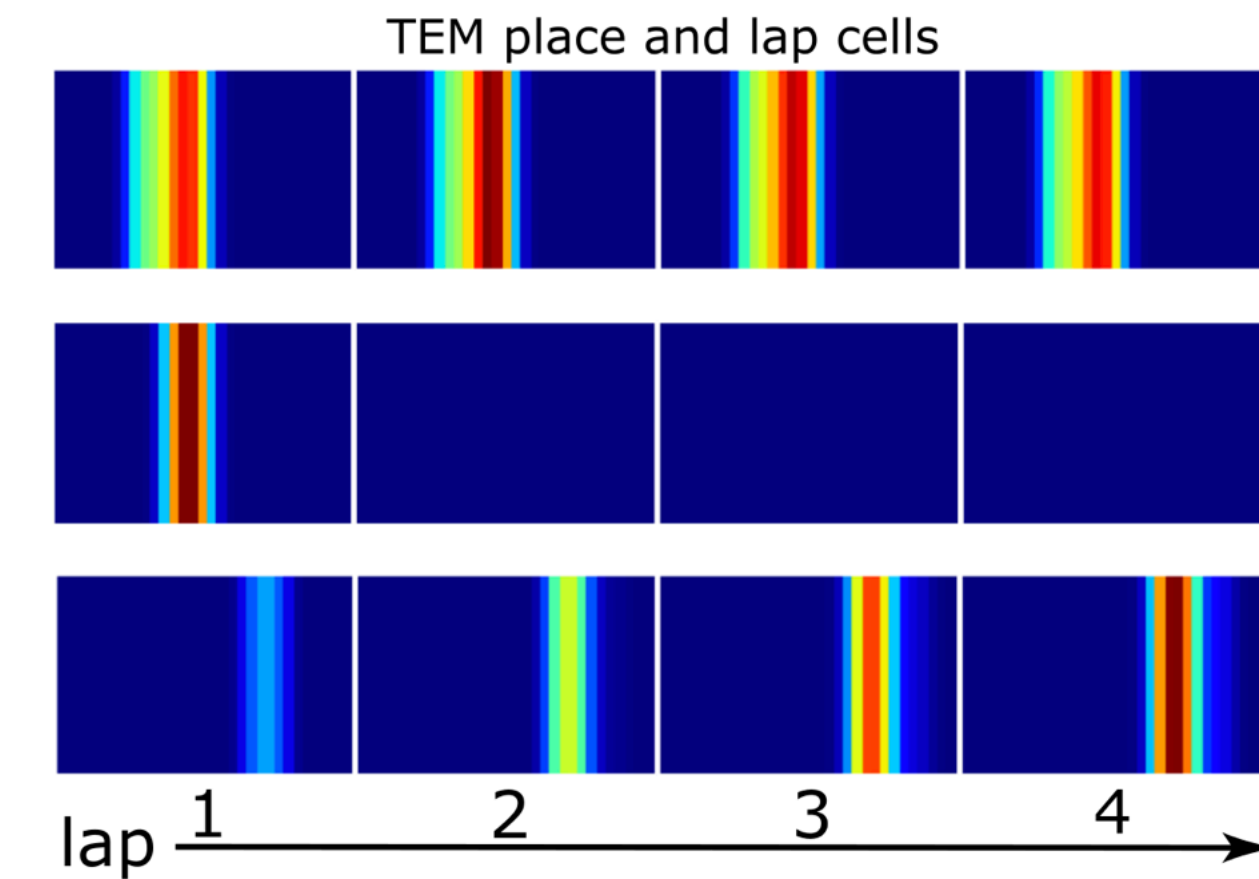
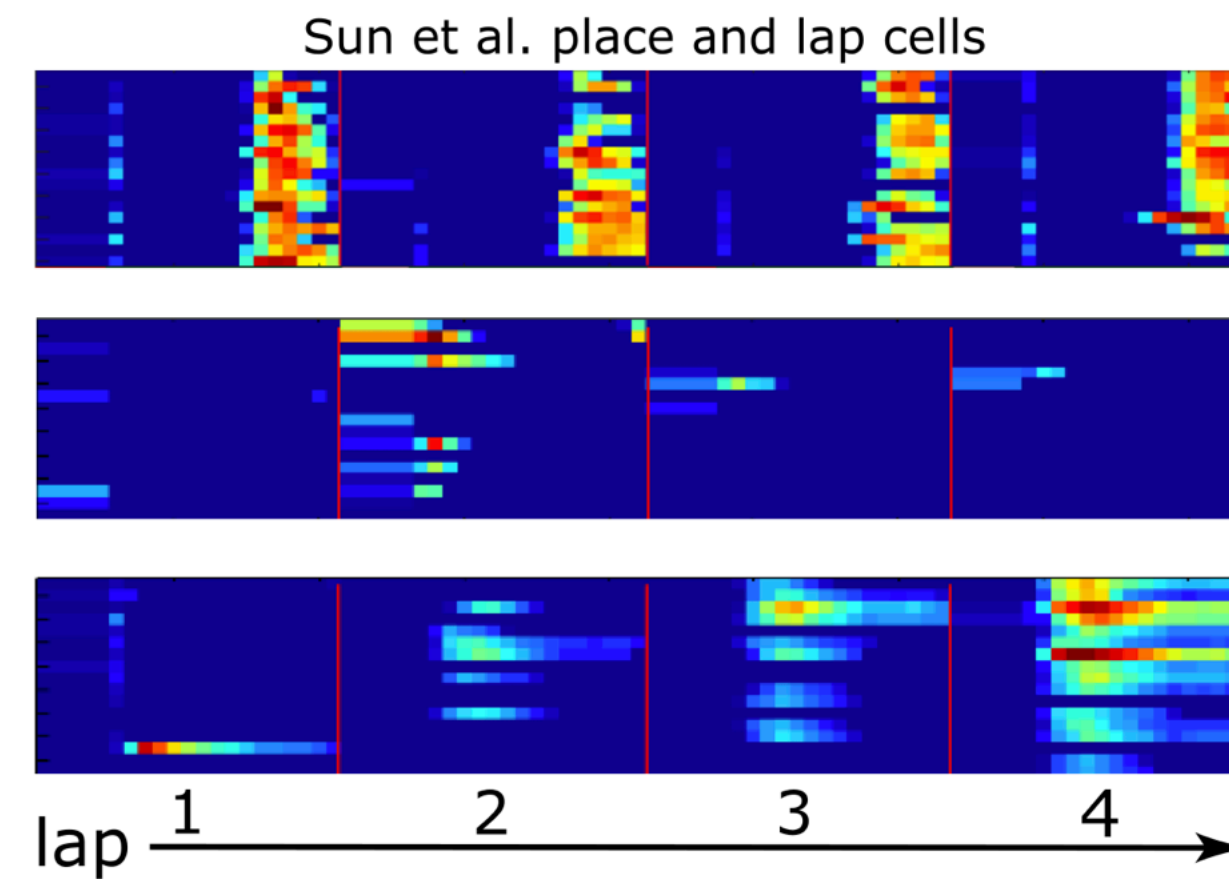
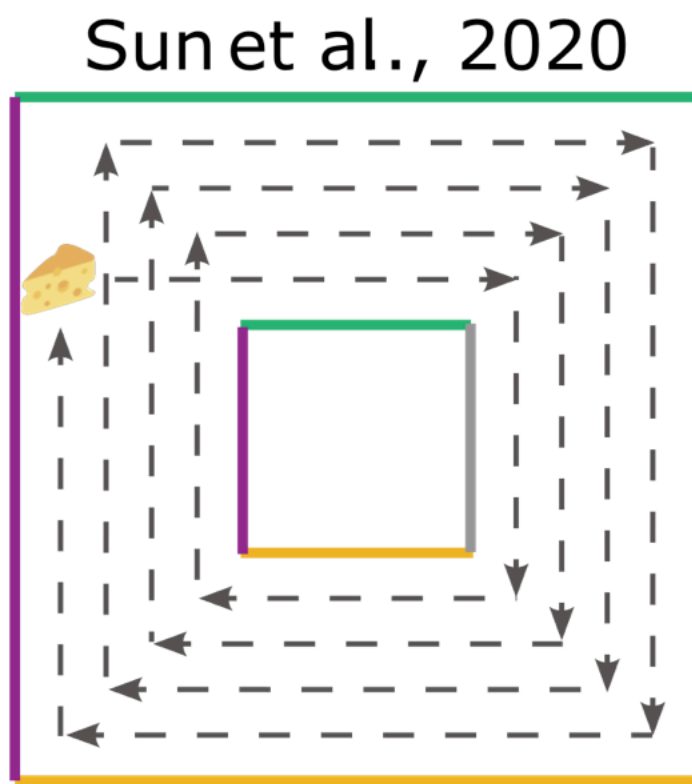
# Lots of other hippocampal representations accounted for

## Simultaneous representation of spatial and task location

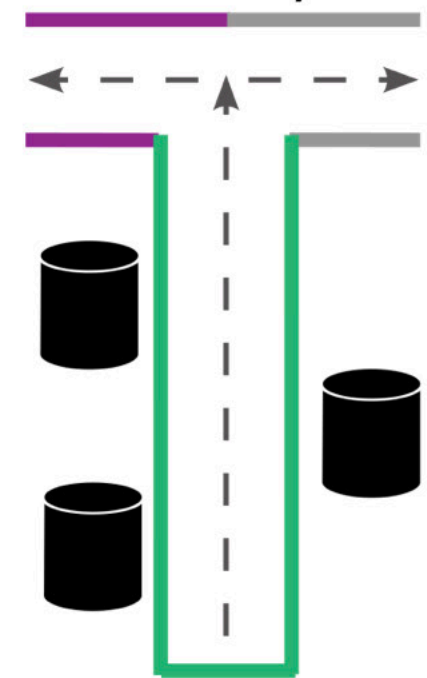


# Lots of other hippocampal representations accounted for

## Simultaneous representation of spatial and task location

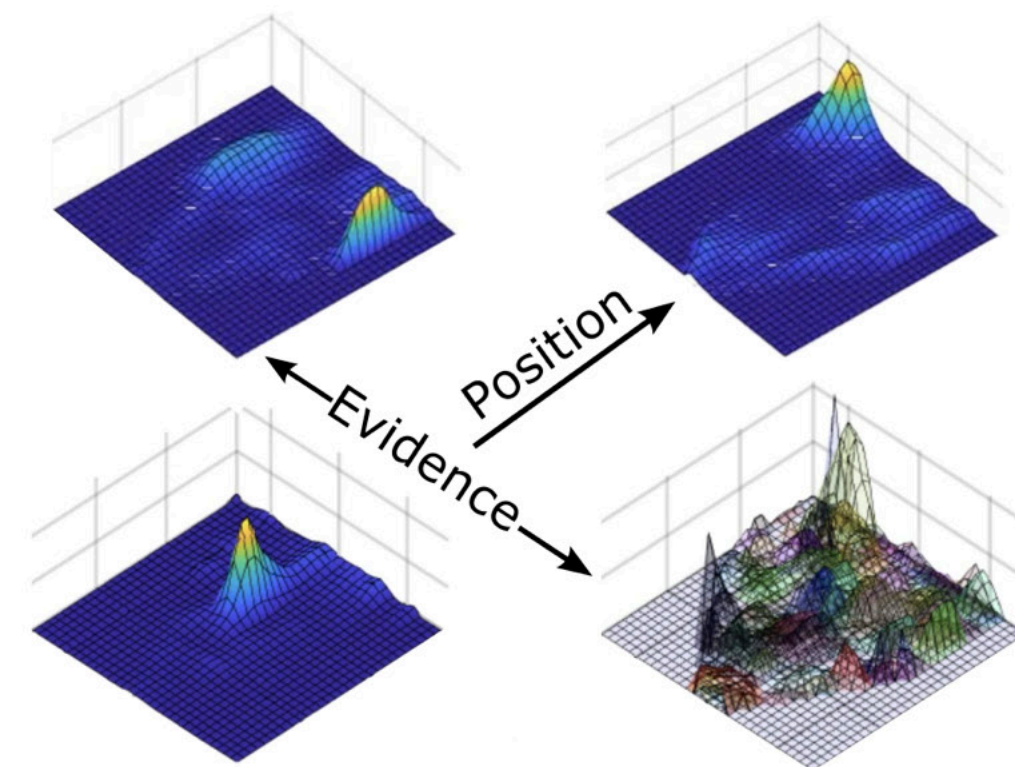


Nieh et al., 2021



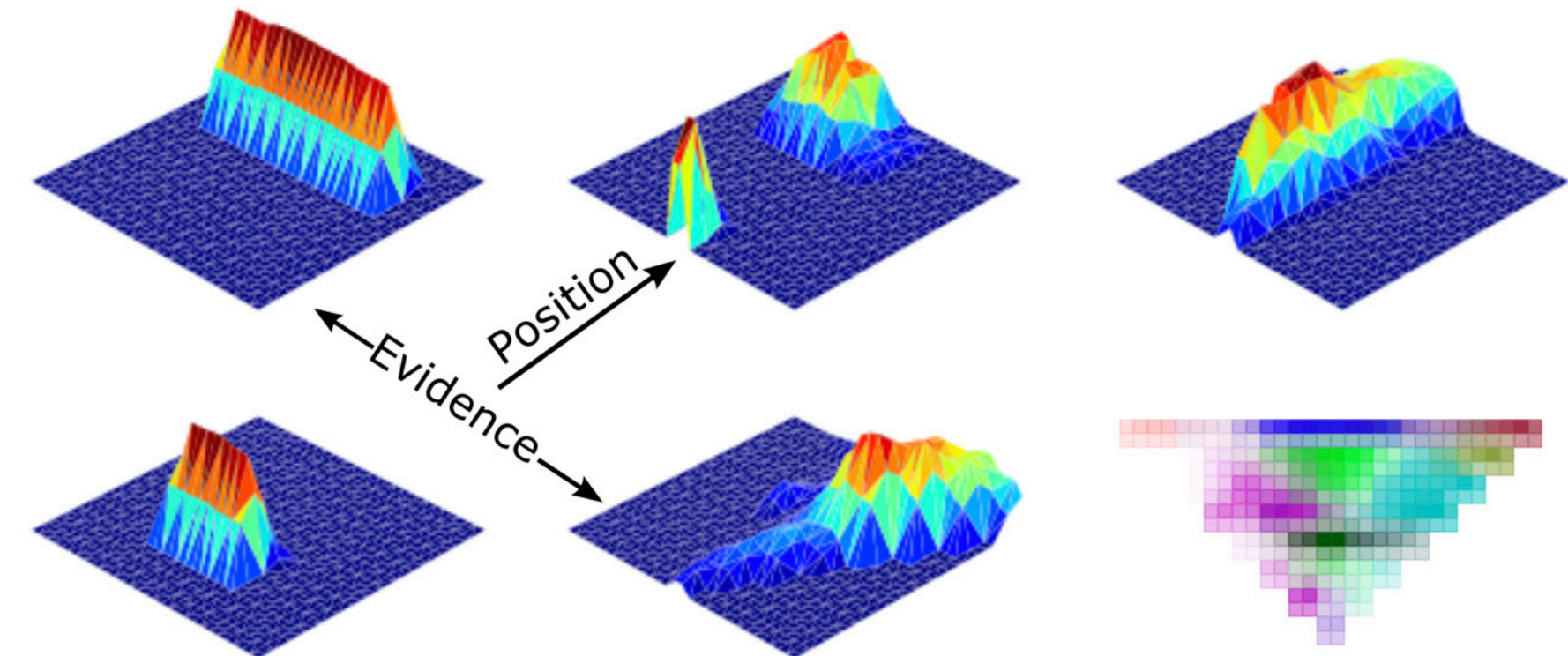
Rodent chooses side with more cues

Nieh et al. abstract position/evidence cells



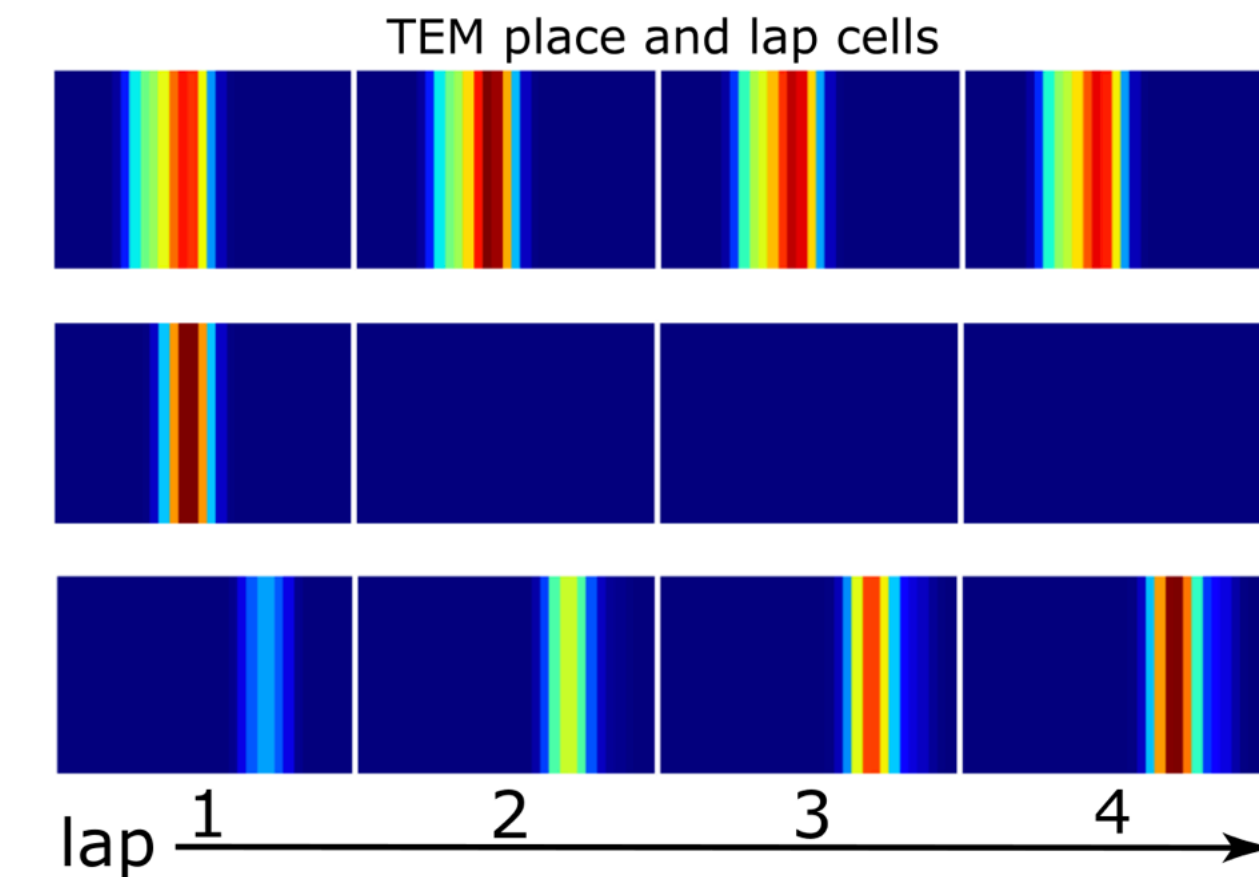
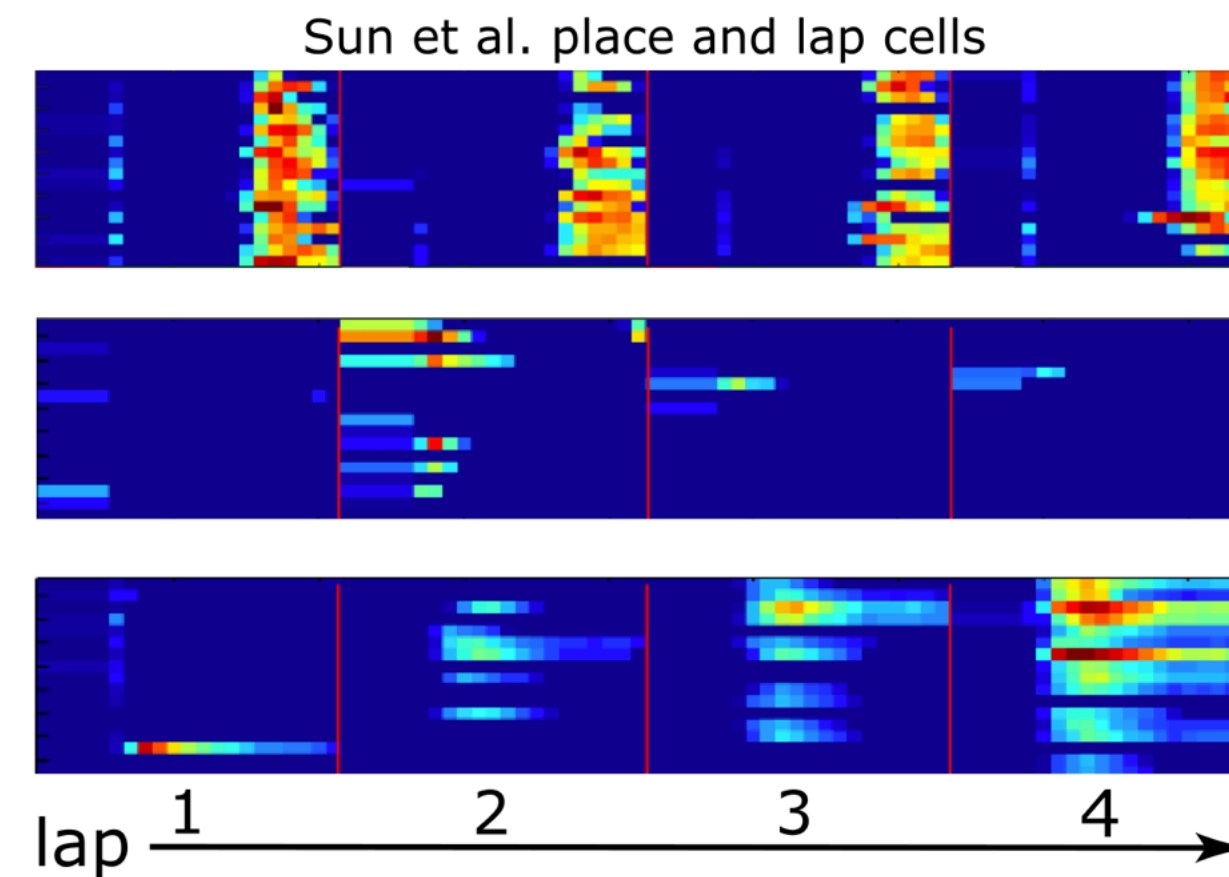
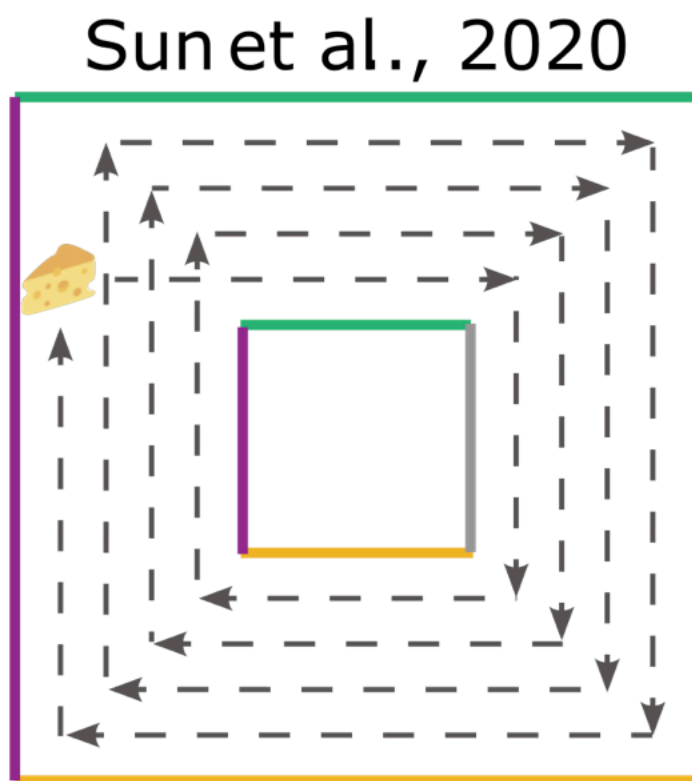
TEM place cells

TEM abstract position/evidence cells

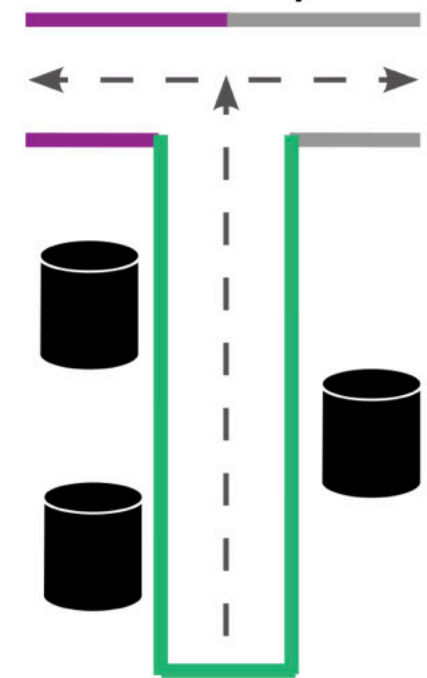


# Lots of other hippocampal representations accounted for

## Simultaneous representation of spatial and task location

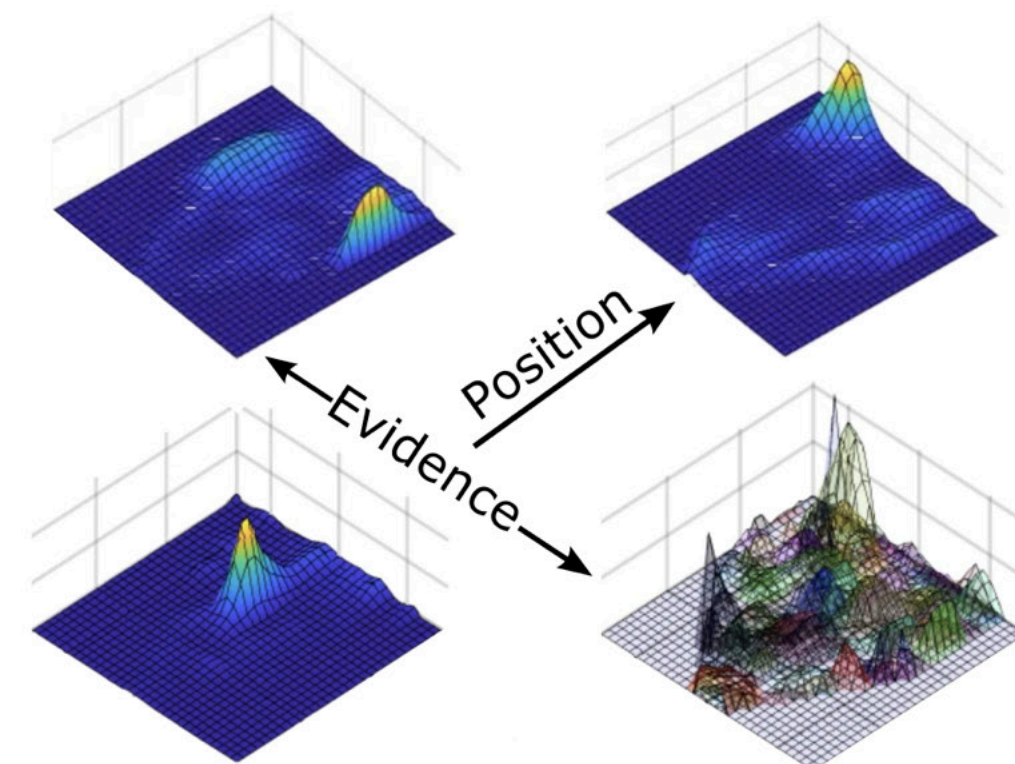


Nieh et al., 2021



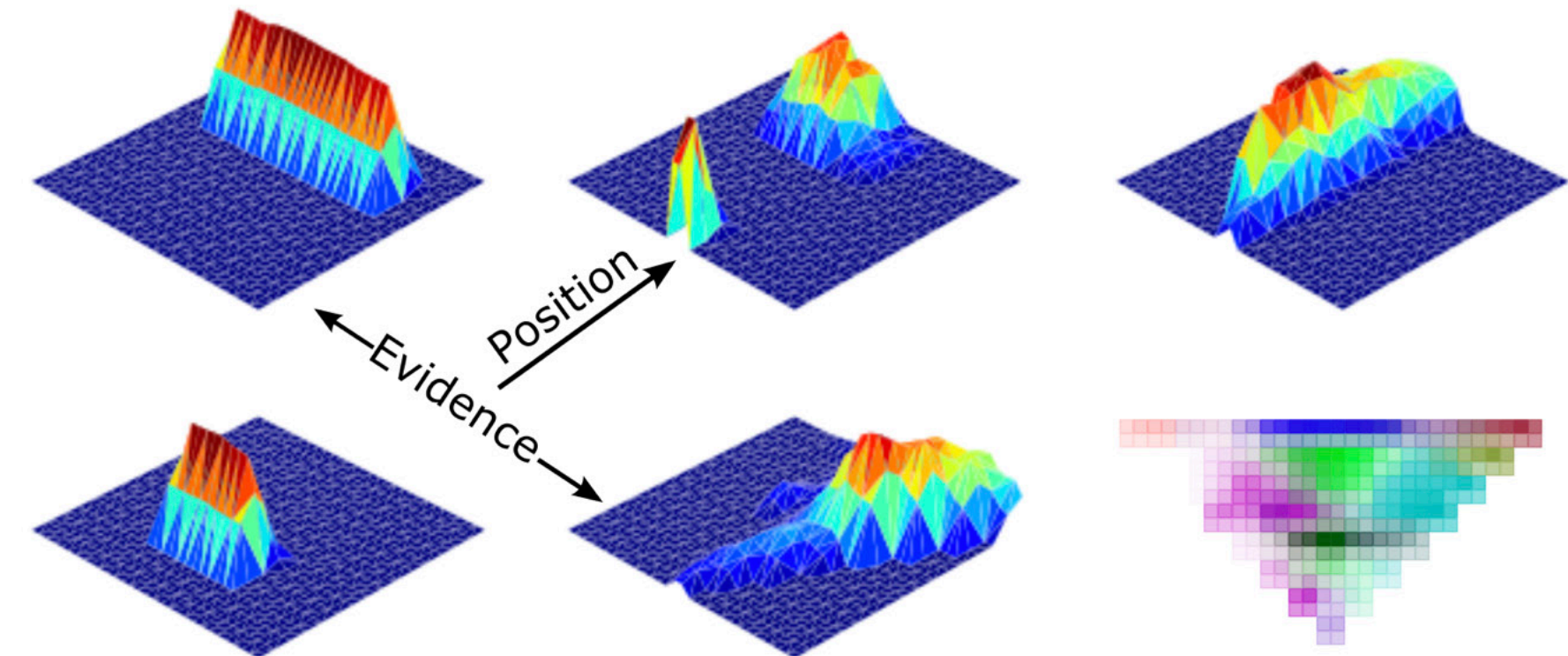
Rodent chooses side with more cues

Nieh et al. abstract position/evidence cells



TEM place cells

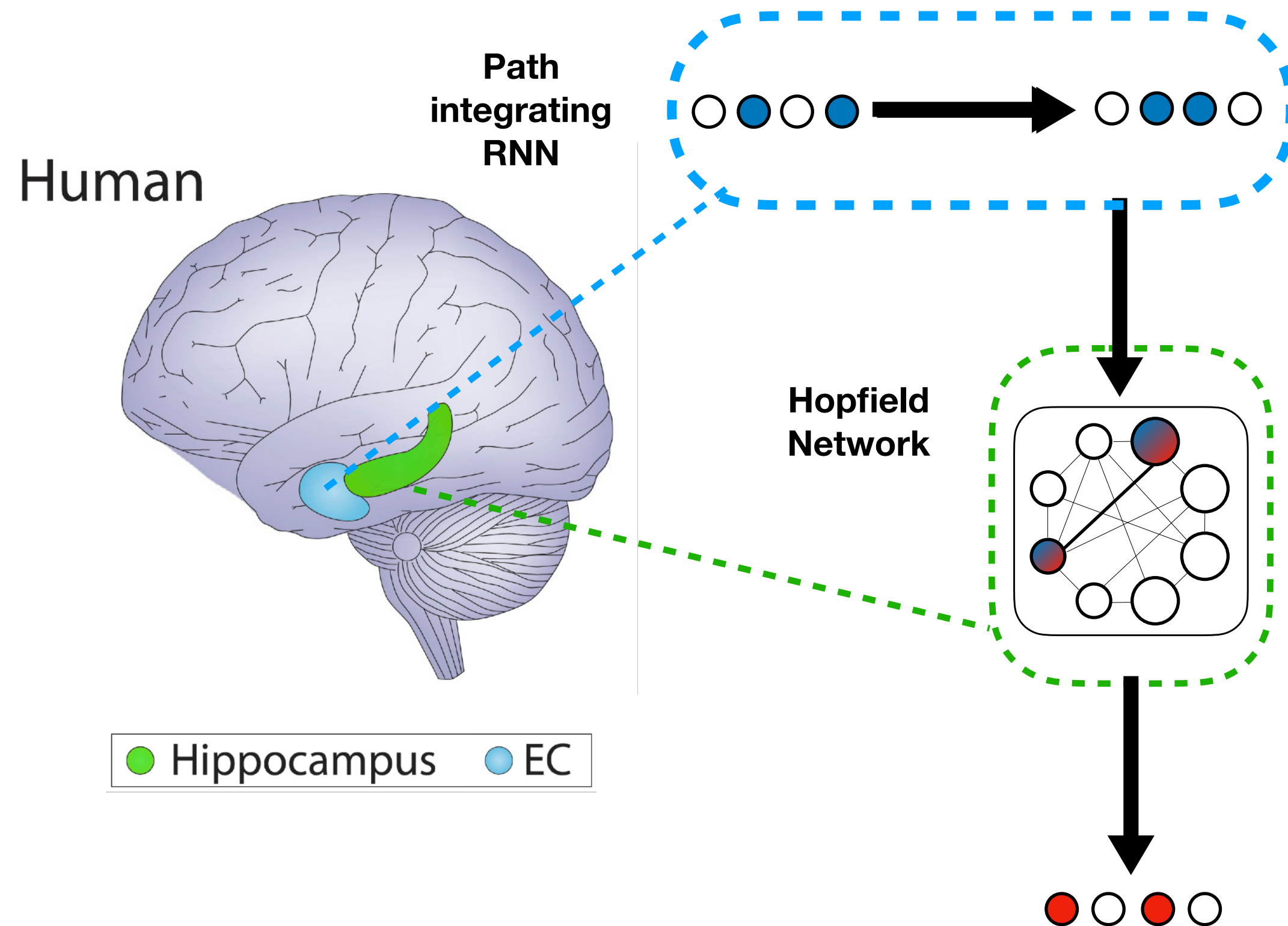
TEM abstract position/evidence cells



**We've explained lots of cells, spatial and non-spatial, by the same principles!**

# Side note: TEM is mathematically related to transformers

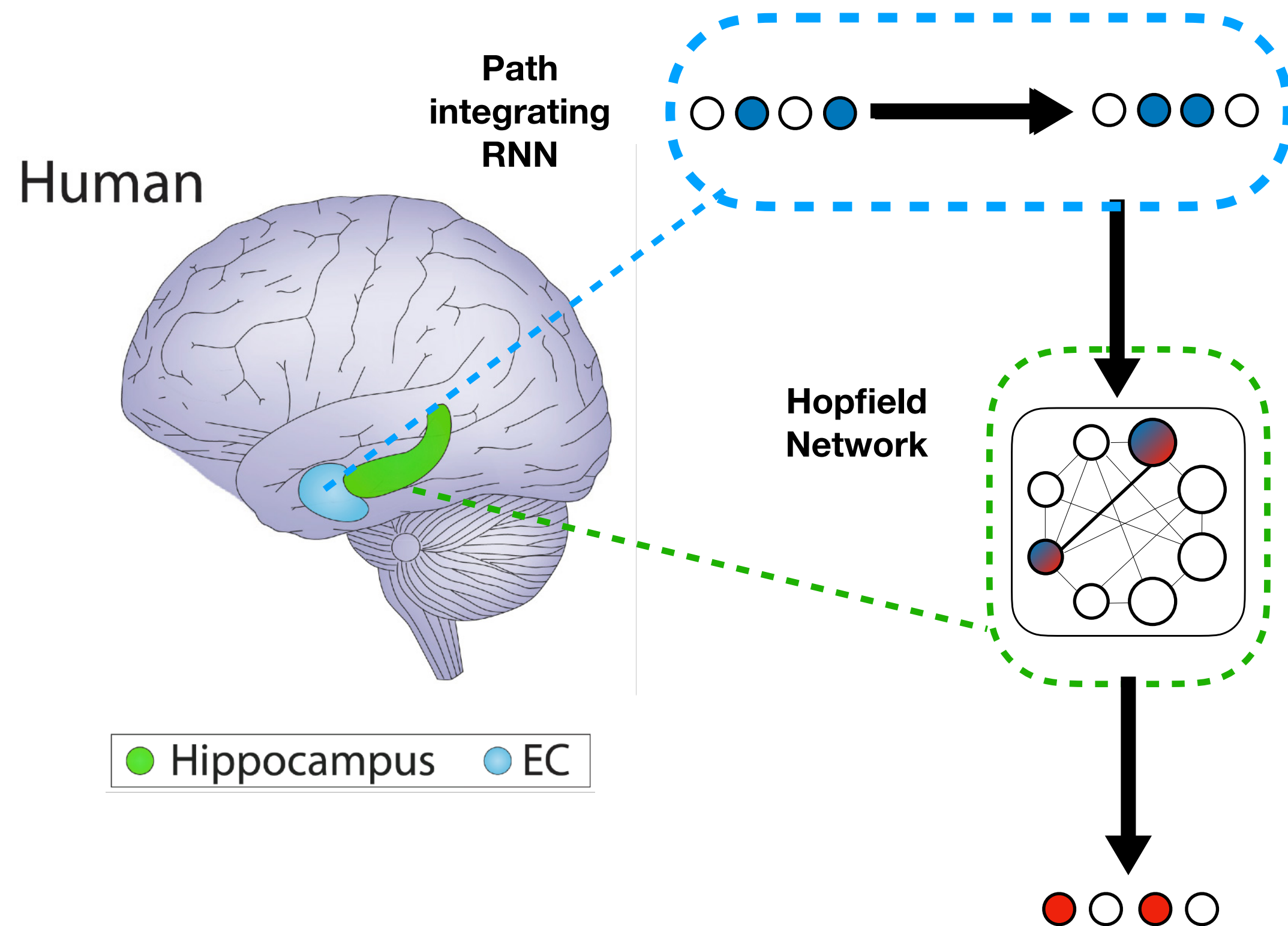
## TEM Model



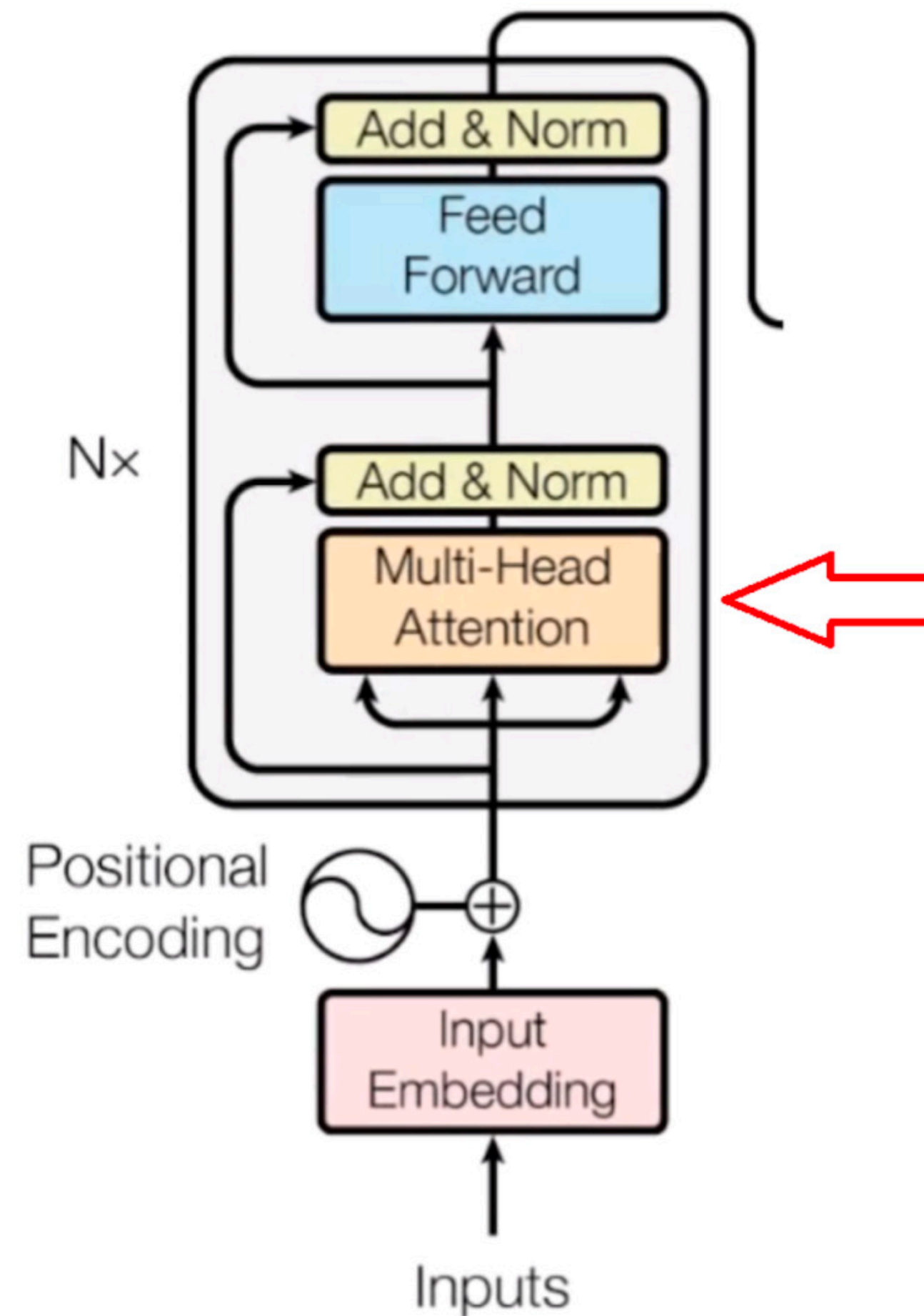


# Side note: TEM is mathematically related to transformers

## TEM Model

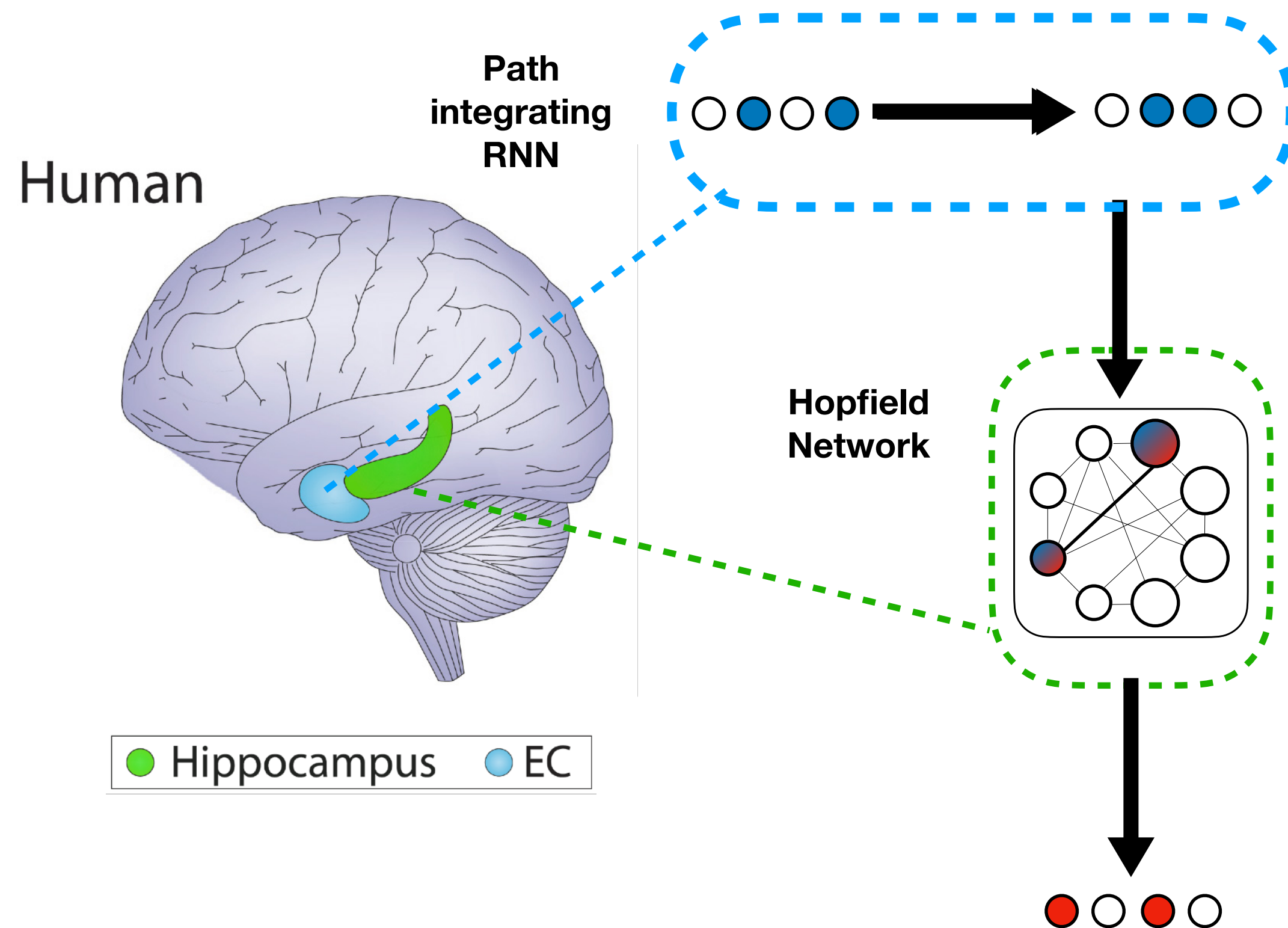


## Transformer neural networks

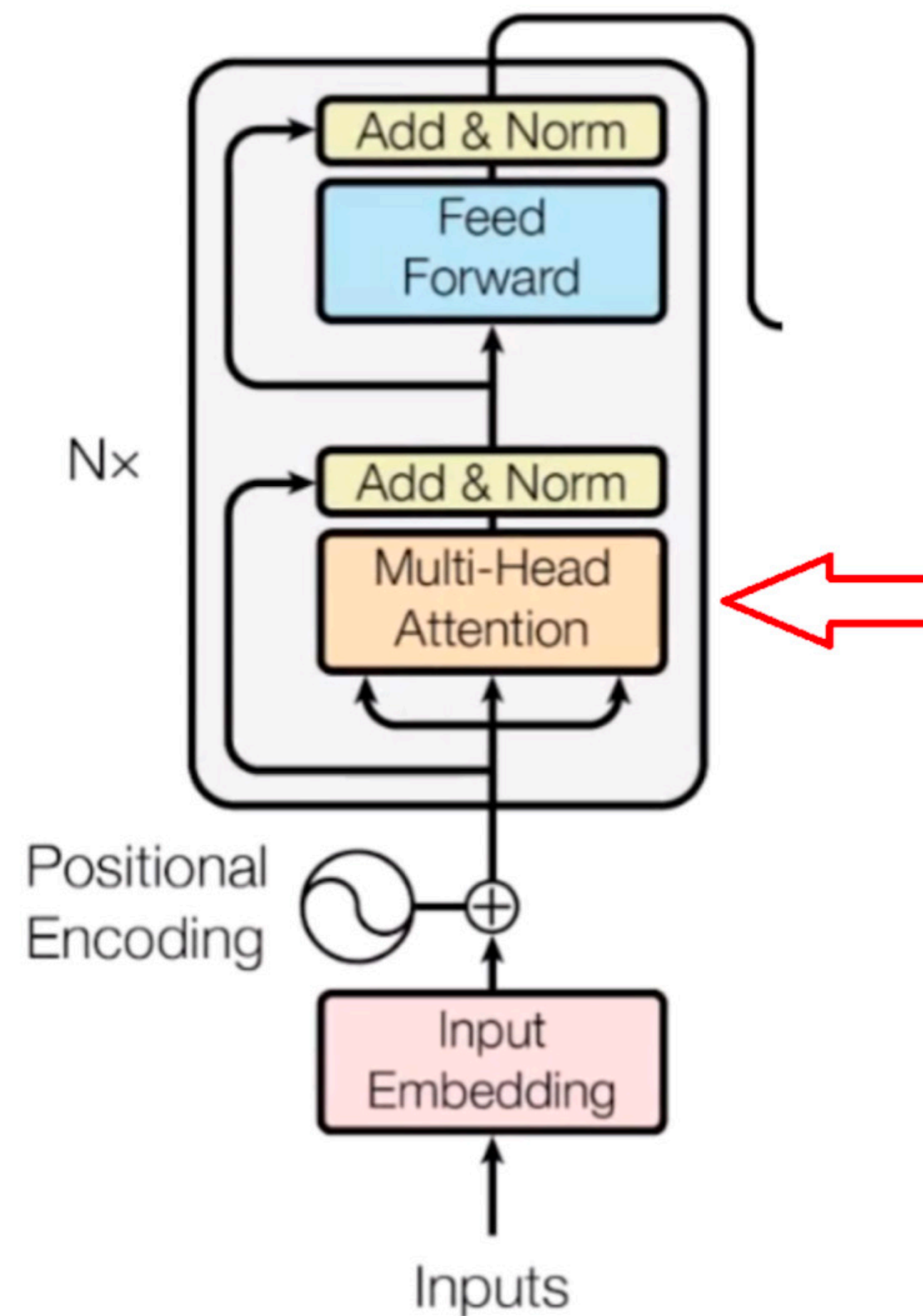


# Side note: TEM is mathematically related to transformers

## TEM Model



## Transformer neural networks

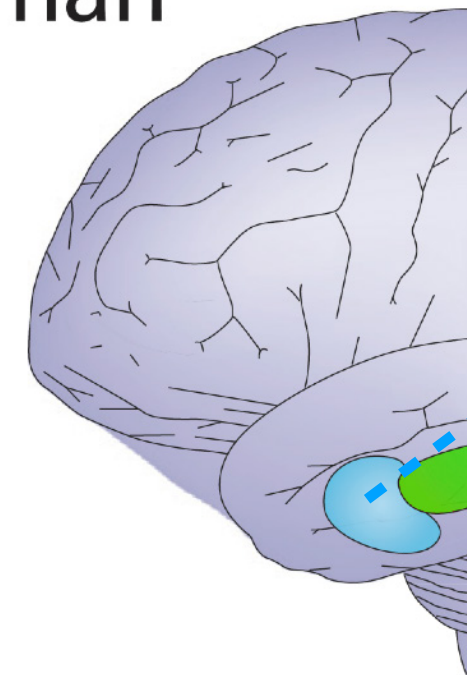


# Side note: TEM is mathematically related to transformers

## TEM Model

## Transformer neural networks

Human



● Hippocamp

Hello, I'm here to find information for you. Start by asking me a question.

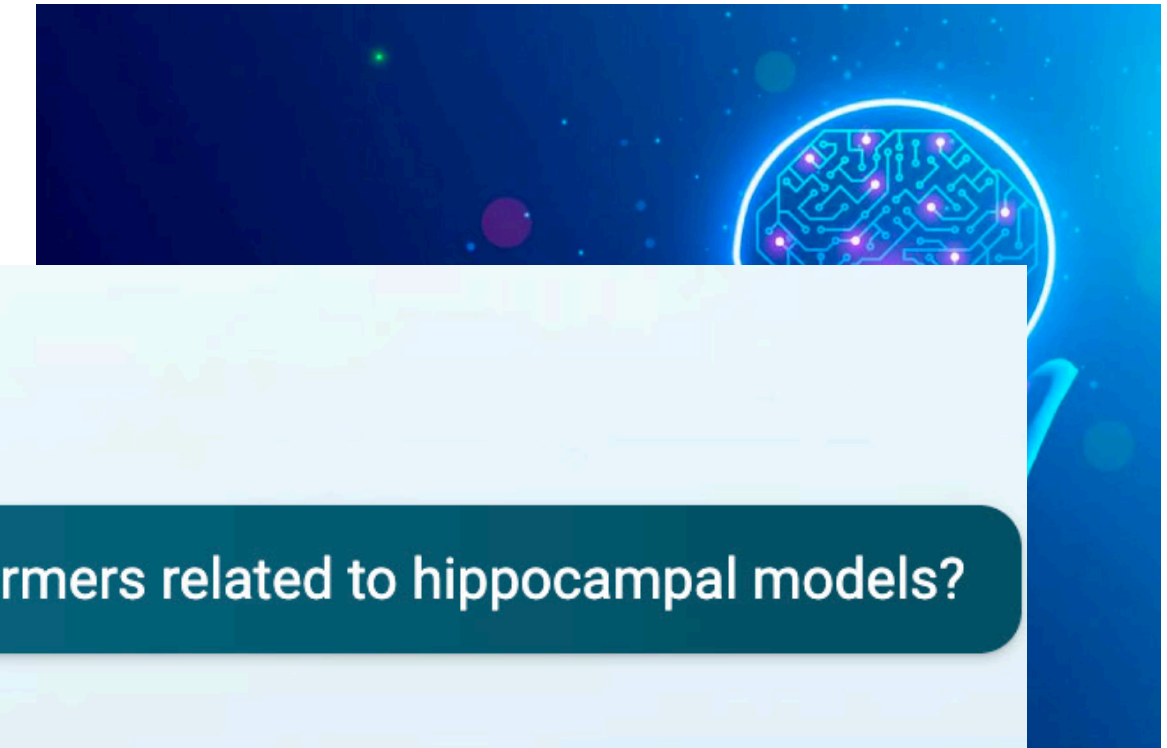
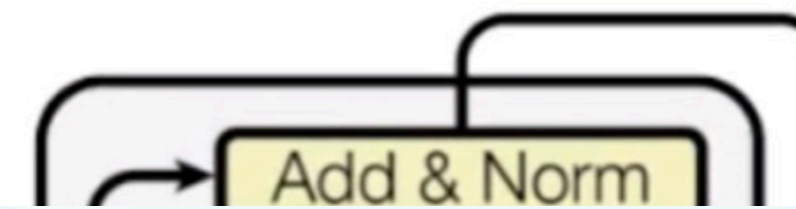
- ✓ Searching for: **transformers related to hippocampal models**
- ✓ Generating answers for you...

Yes, transformers are related to hippocampal models. Researchers have shown that transformers equipped with recurrent position encodings replicate the precisely tuned spatial representations of the hippocampal formation; most notably place and grid cells <sup>1</sup>. This result is closely related to current hippocampal models from neuroscience <sup>1</sup>.

Relating transformers to models and neural representations of the ...  
<https://arxiv.org/abs/2112.04035>

In this work, we show that transformers, when equipped with recurrent position encodings, replic...

1 of 15 ●



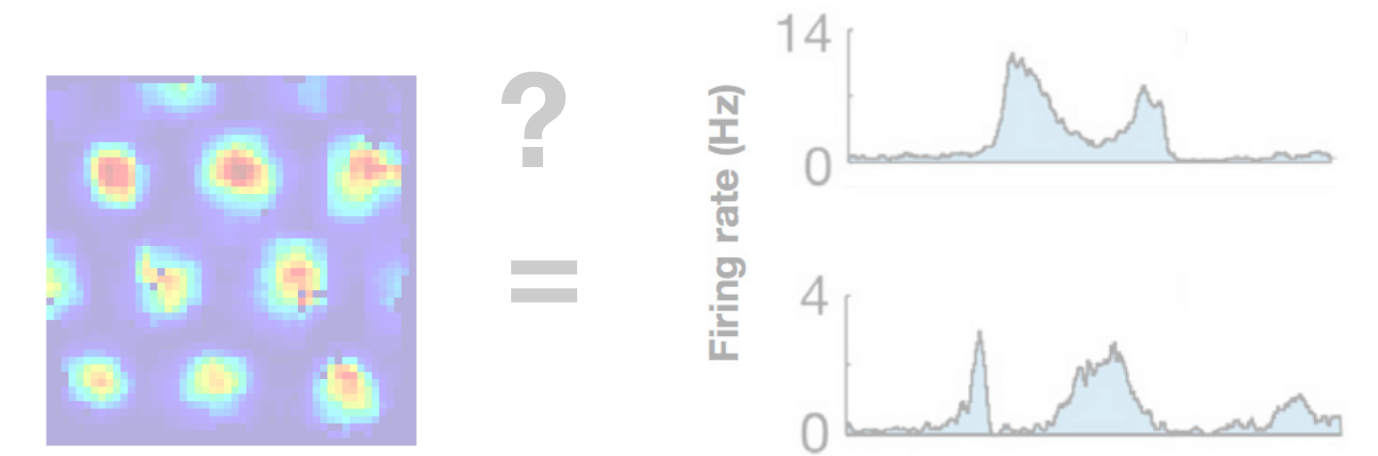
are transformers related to hippocampal models?



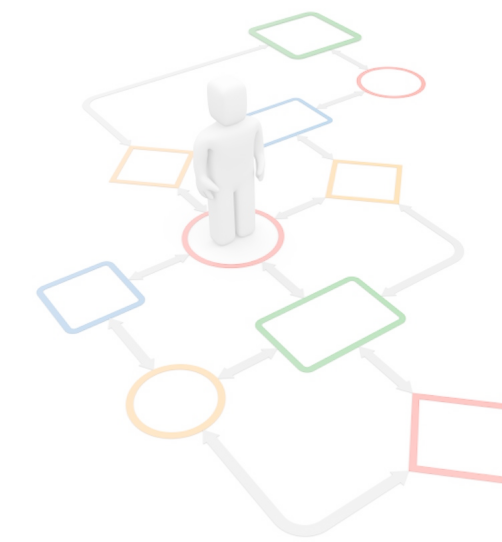
Inputs

# Puzzles of cognitive maps in the brain

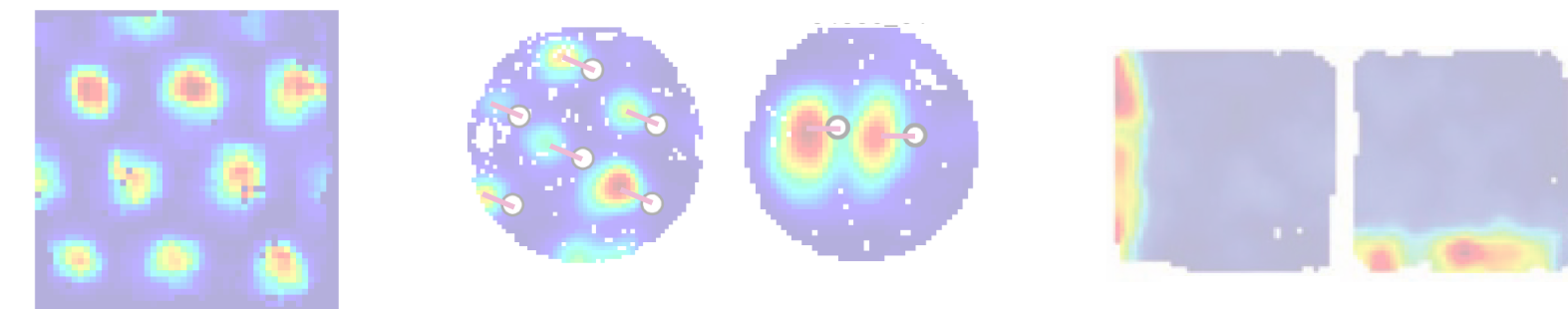
How does the same system do space and non-space?



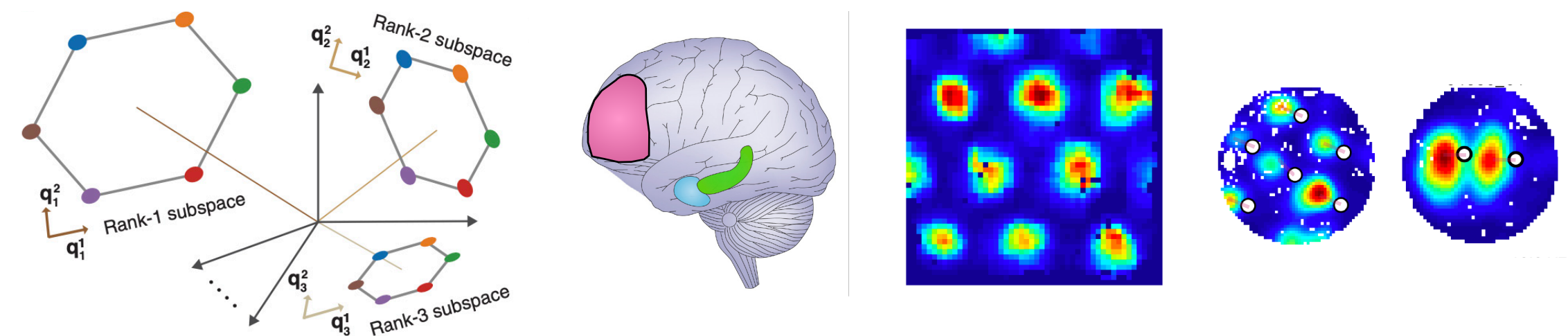
How can brains learn these maps?



Why do the neurons look the ways they do?

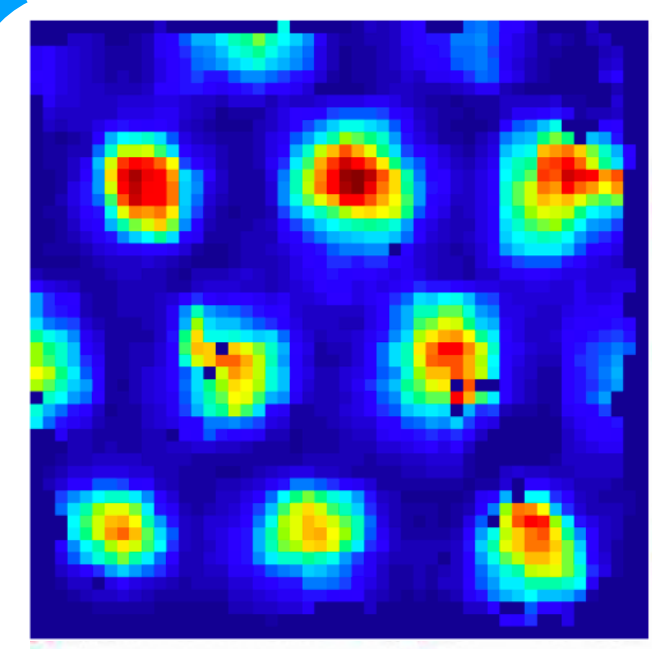
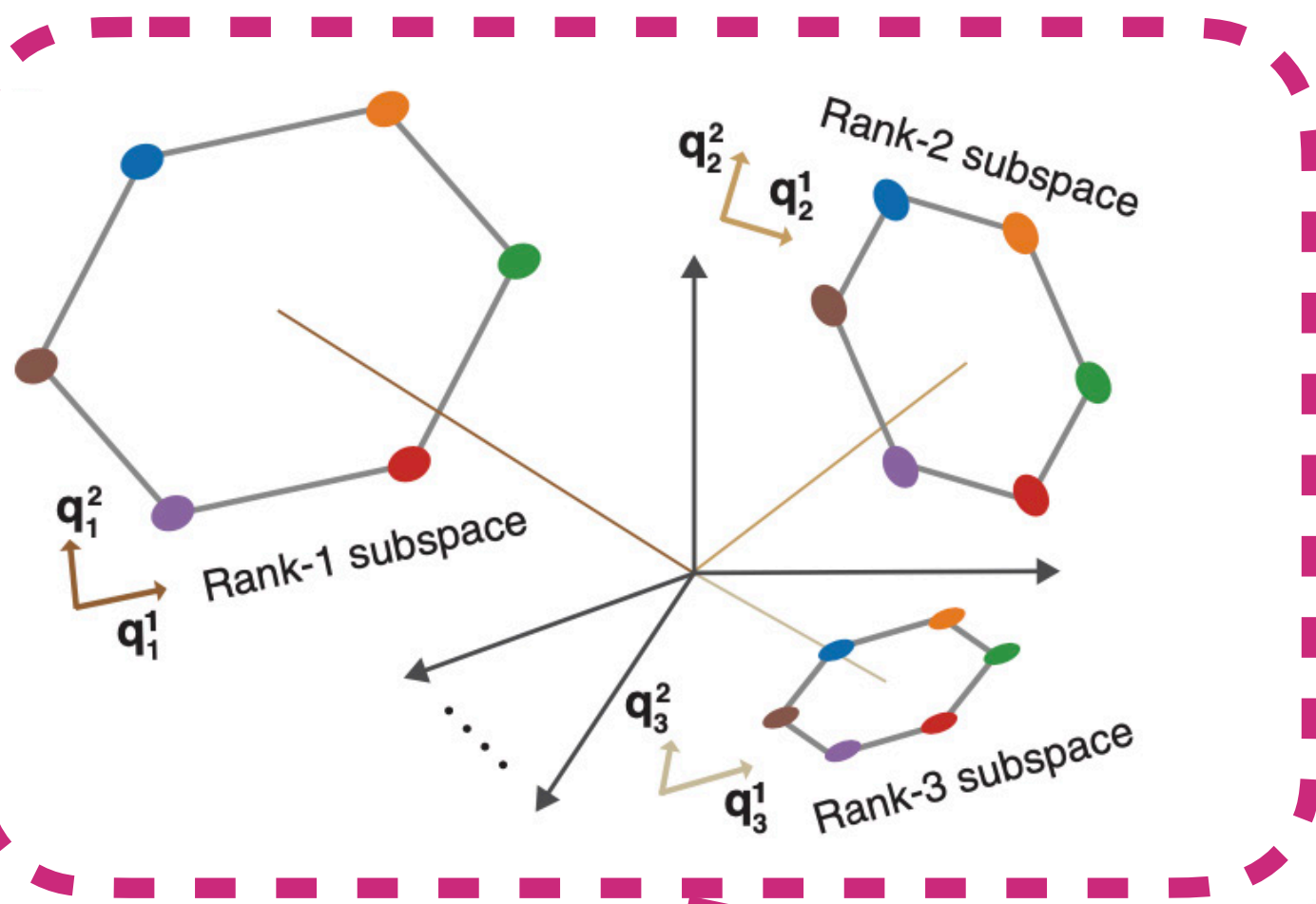


How do different brain regions solve the same problem in different ways?

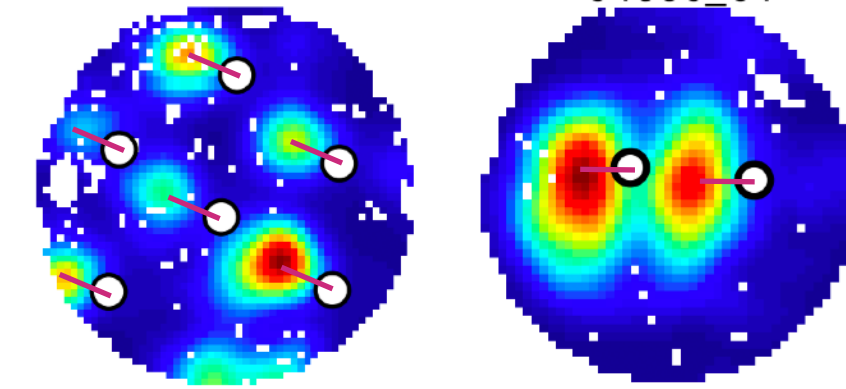


# These models don't capture anything about prefrontal cortex

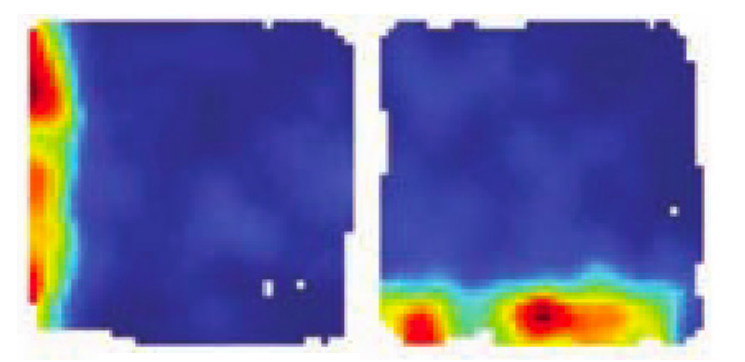
These cells tell you where you are or what you're seeing right now



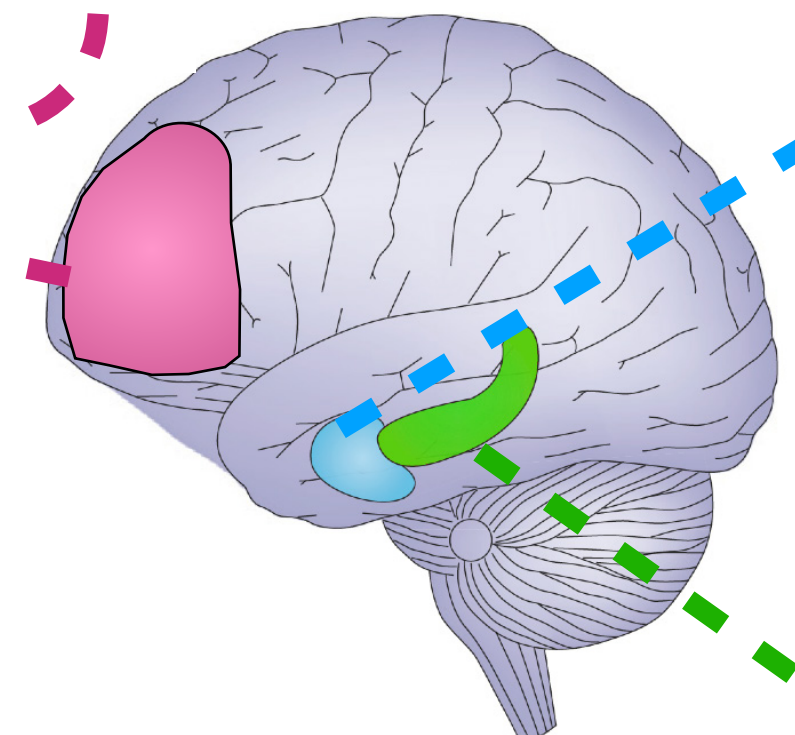
Grid cells  
Hafting et al., 2005



Object vector cells  
Hoydal et al., 2018

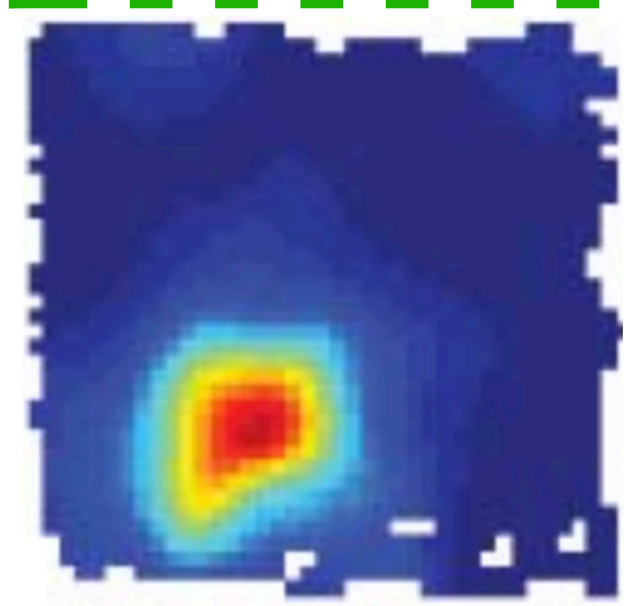


Border cells  
Solstad et al., 2008

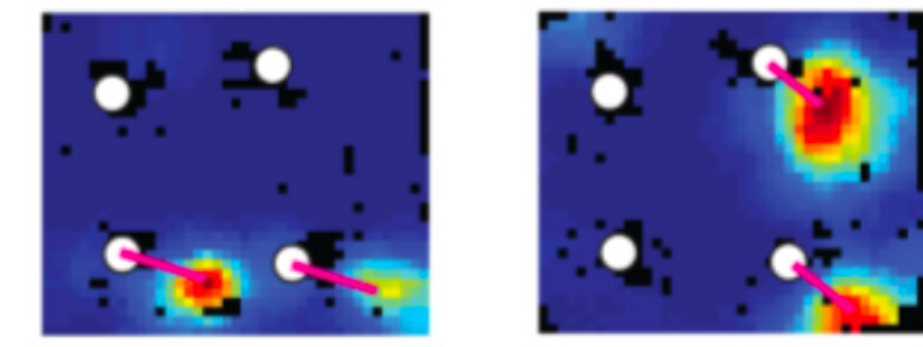


- Hippocampus
- EC
- PFC

These cells include the past and future



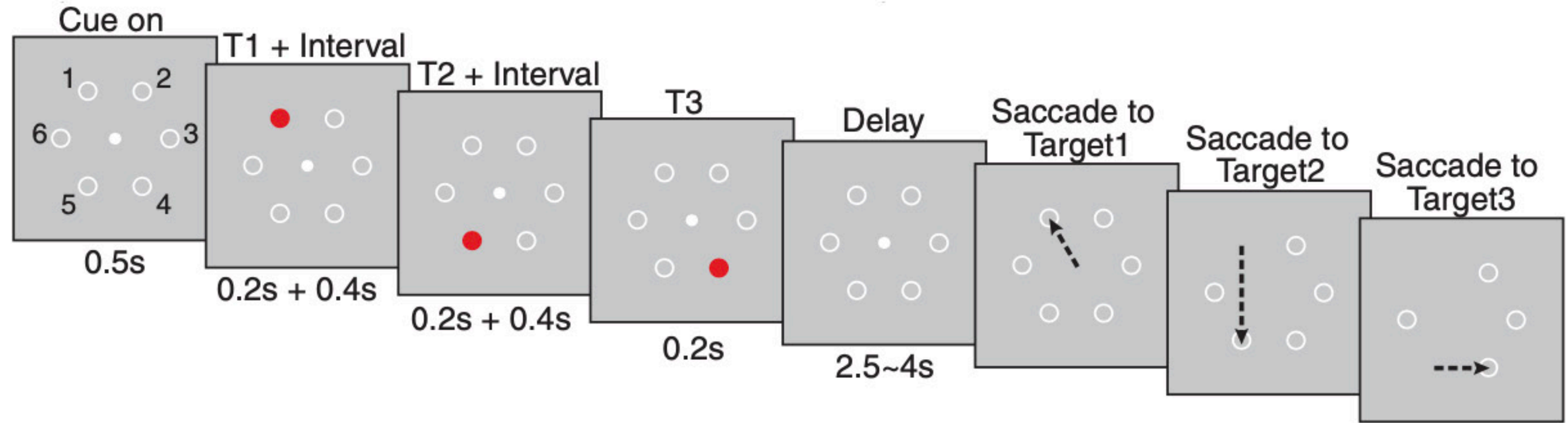
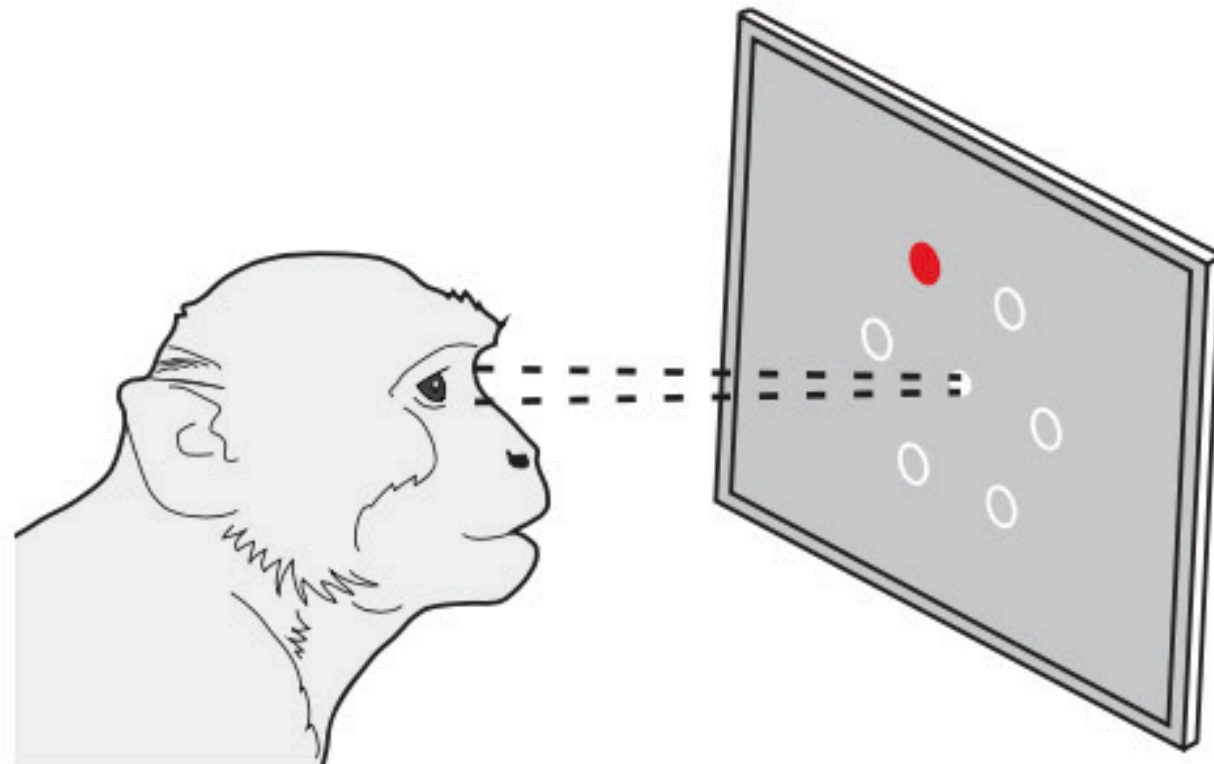
Place cells  
O'Keefe & Dostrovsky, 1971



Landmark cells  
Deshmukh & Knierim, 2013

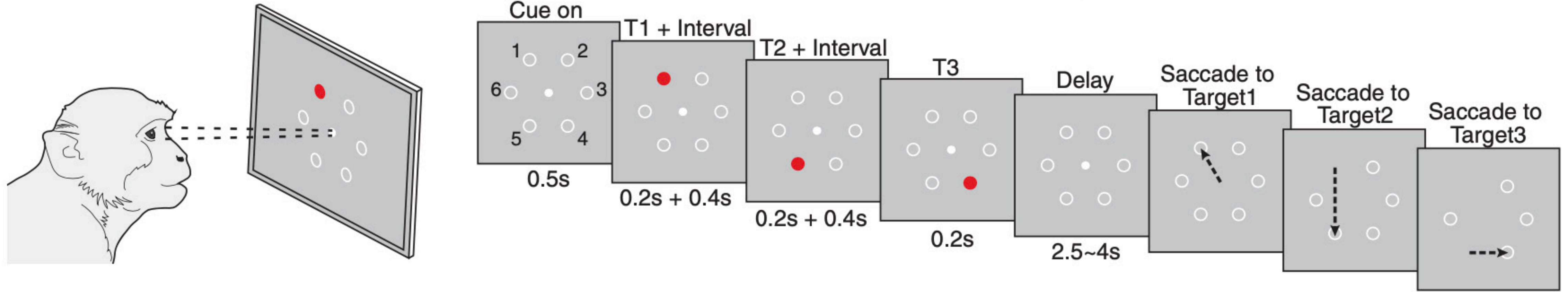
**Prefrontal cells represent the whole history at any given time**

# Prefrontal cells represent the whole history at any given time

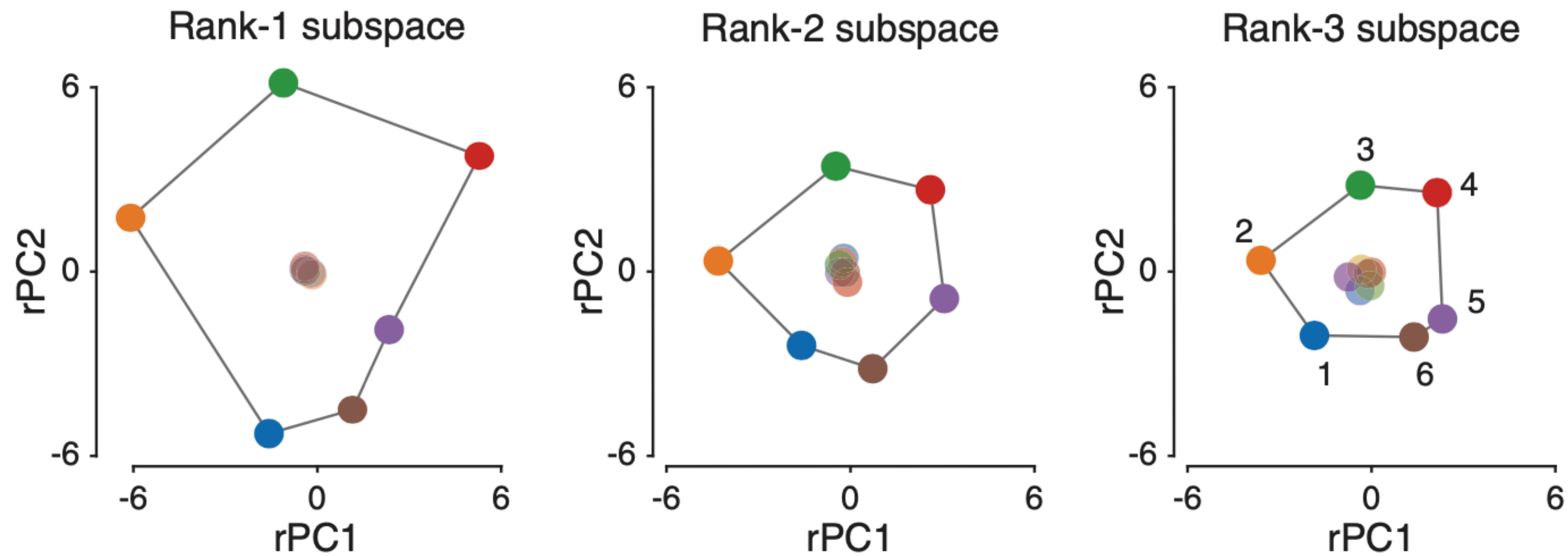


Xie et al., 2022

# Prefrontal cells represent the whole history at any given time

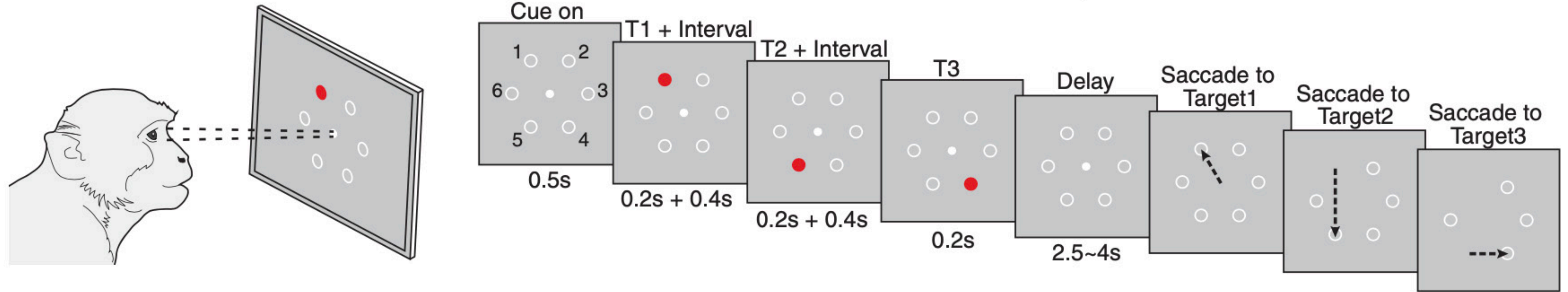


Xie et al., 2022

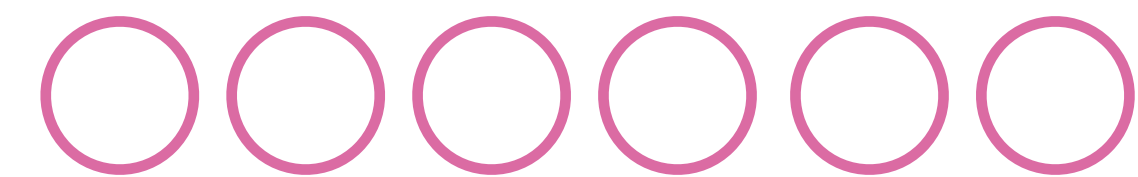
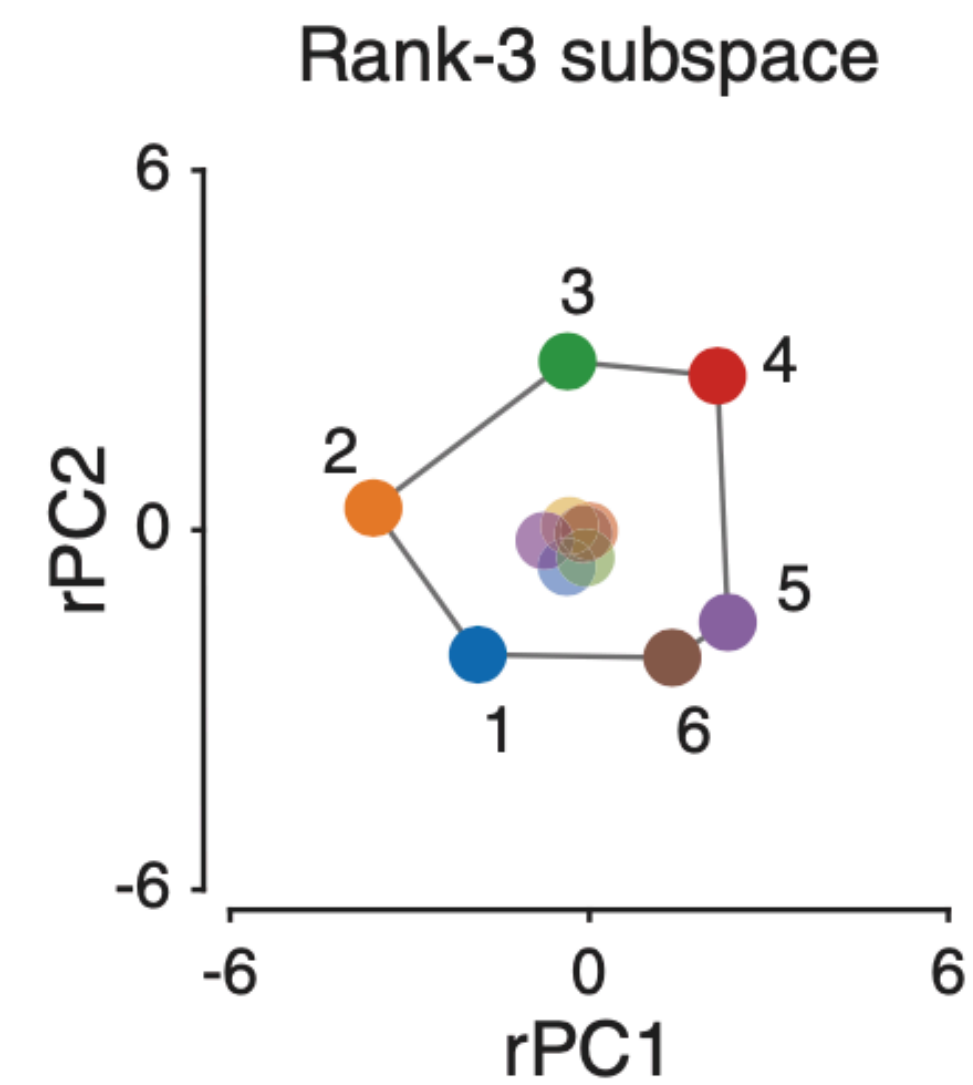
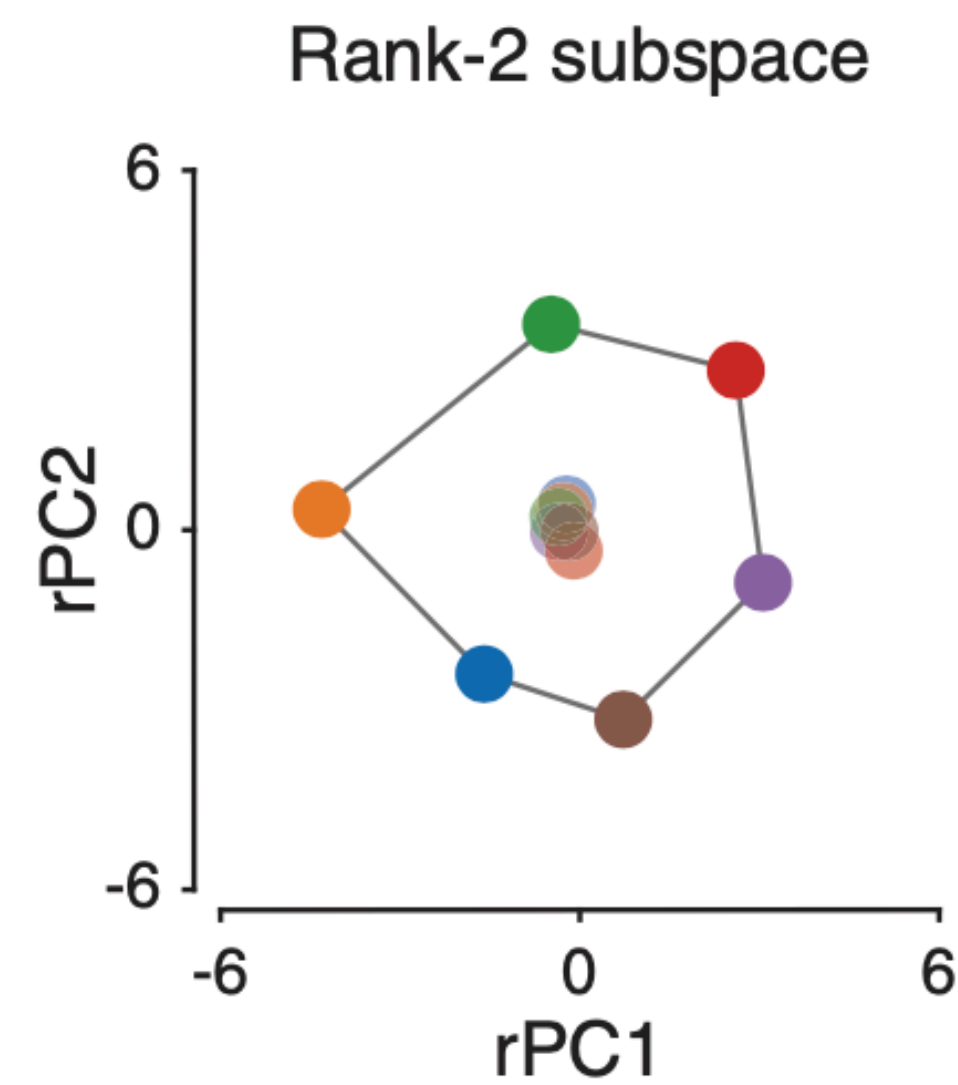
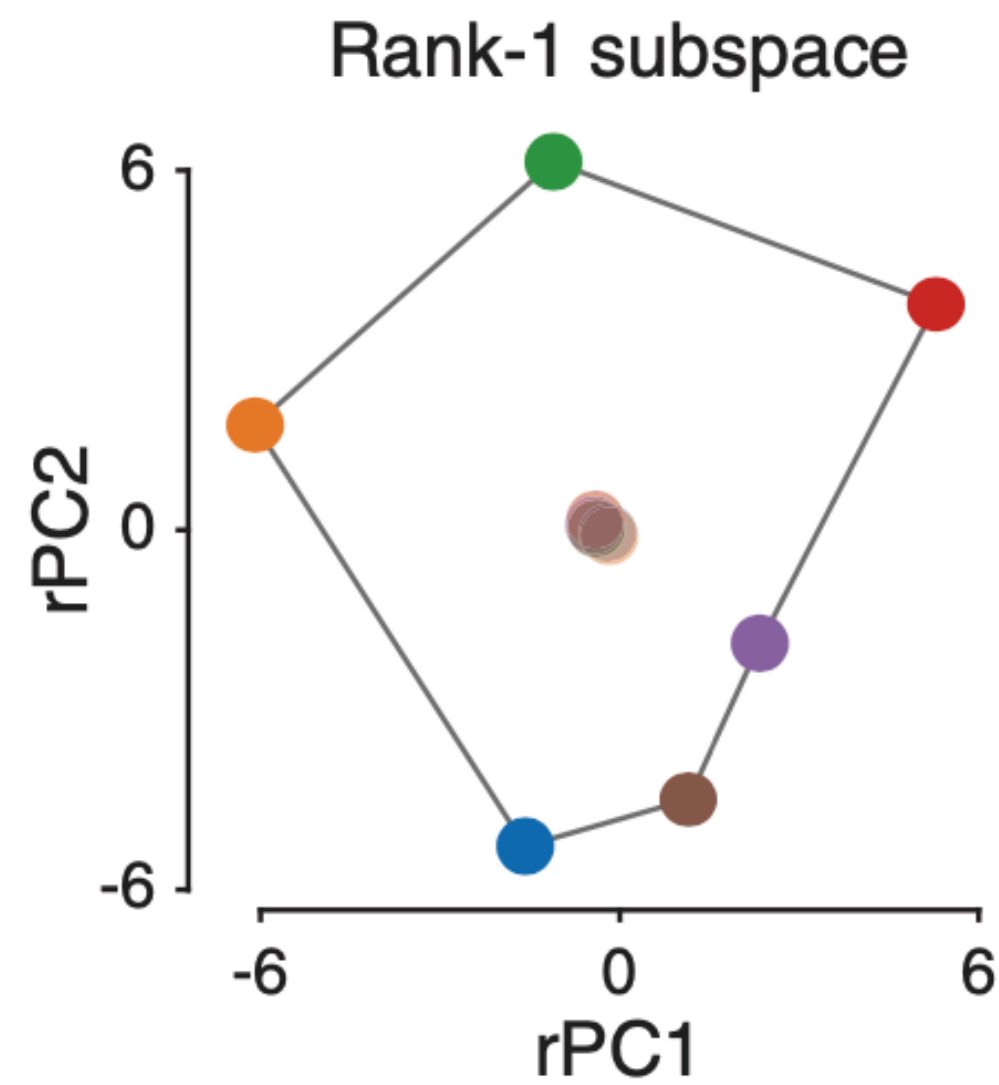




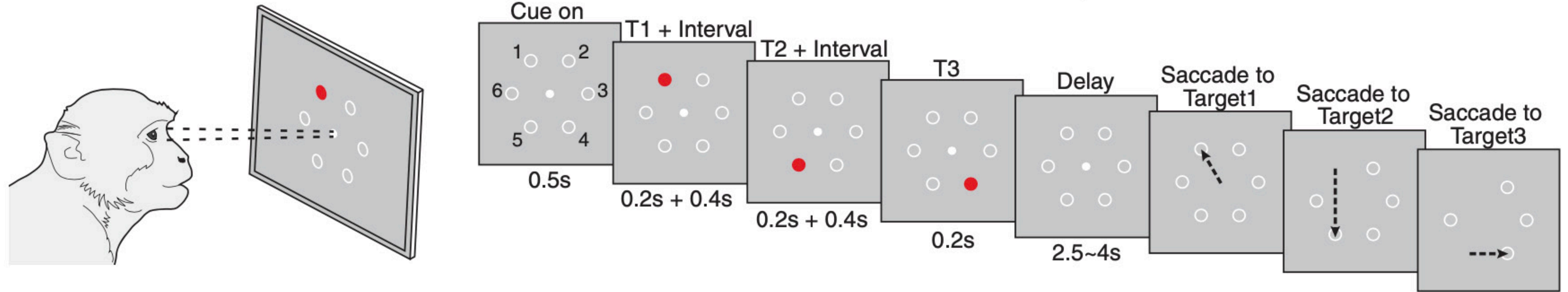
# Prefrontal cells represent the whole history at any given time



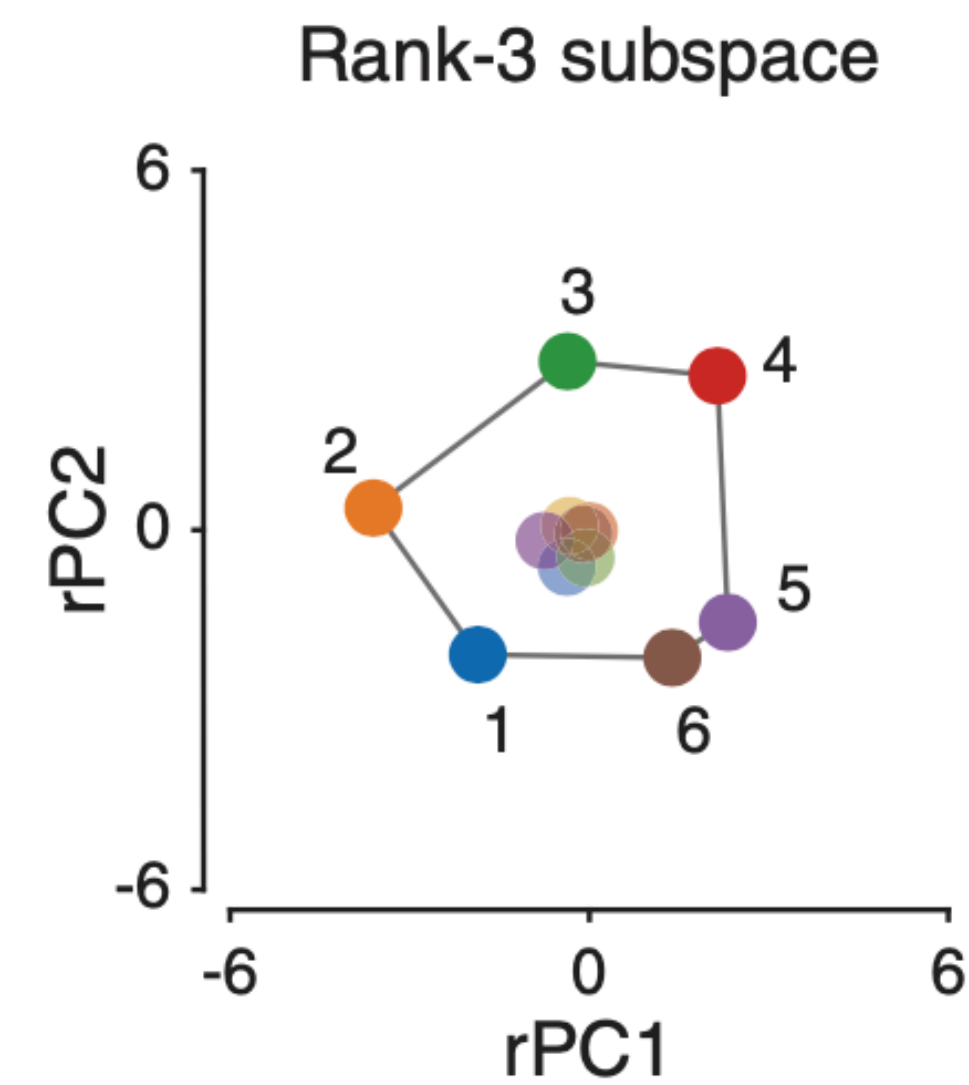
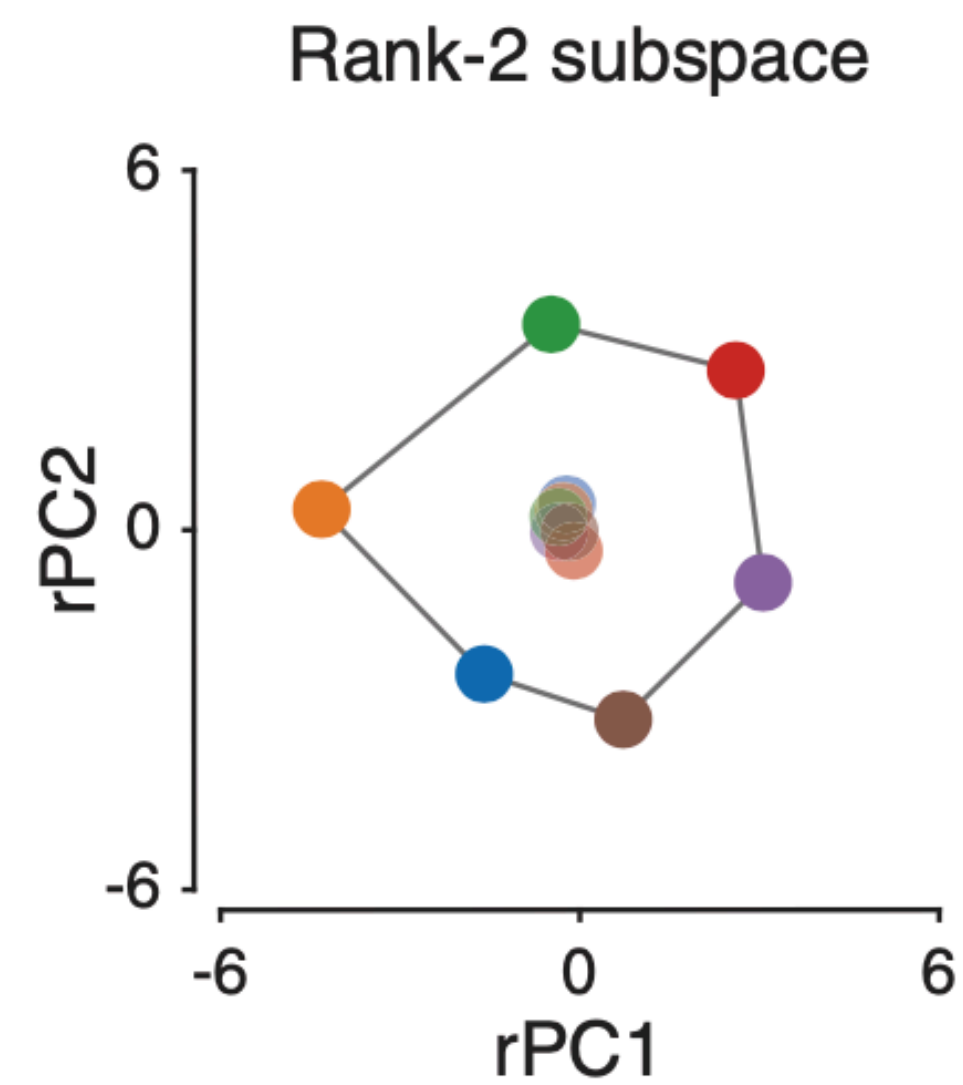
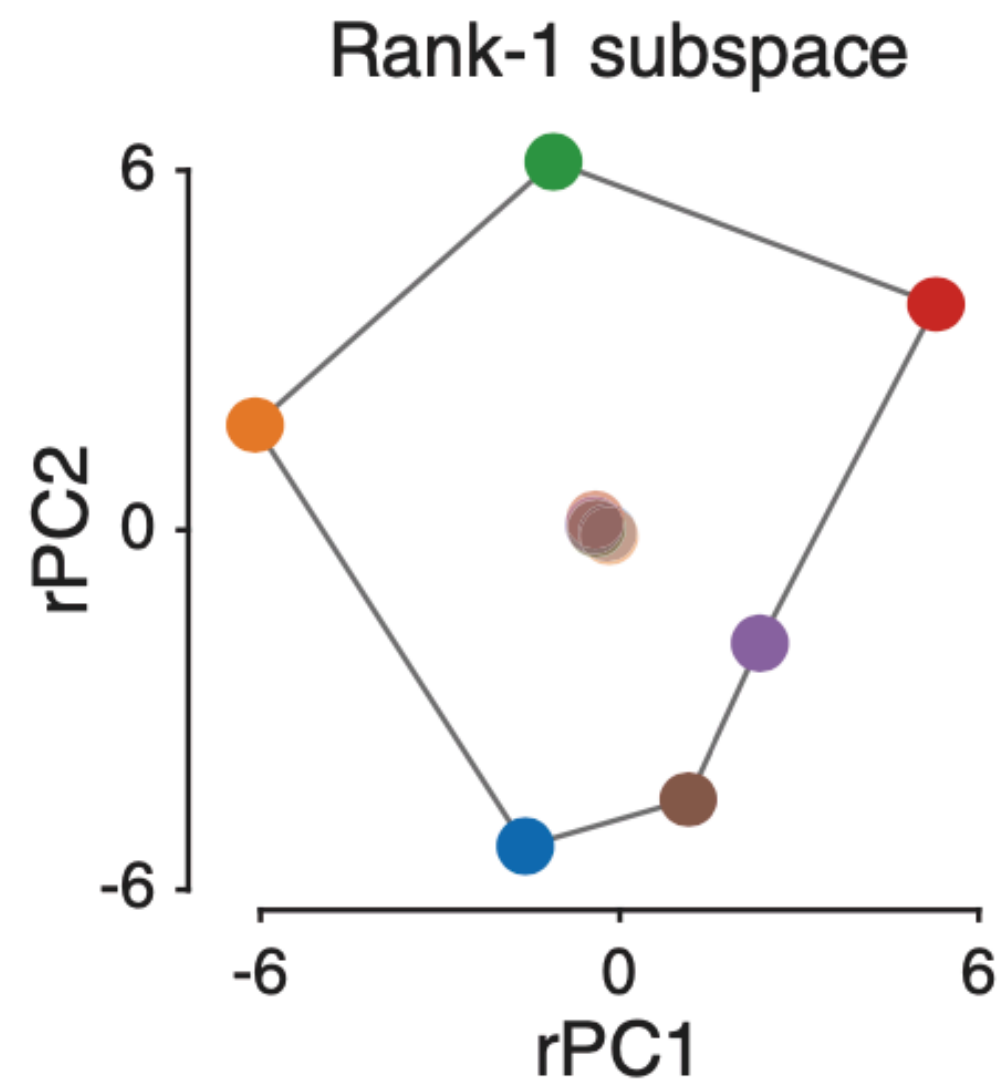
Xie et al., 2022



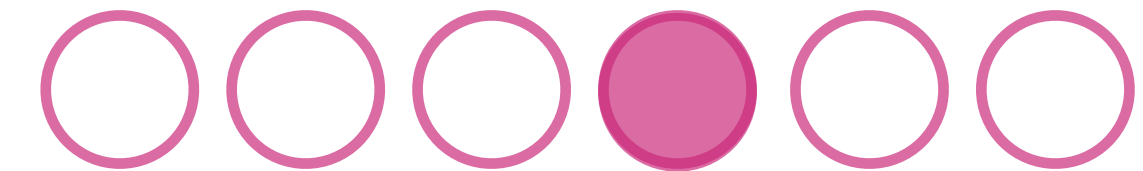
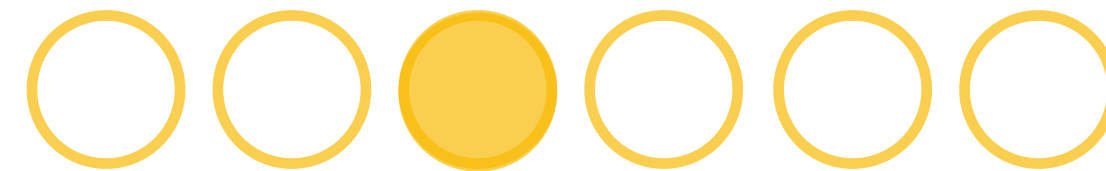
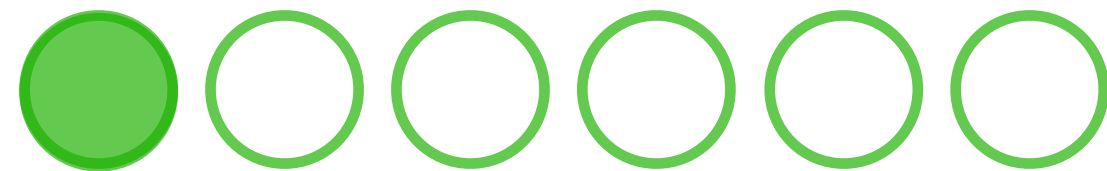
# Prefrontal cells represent the whole history at any given time



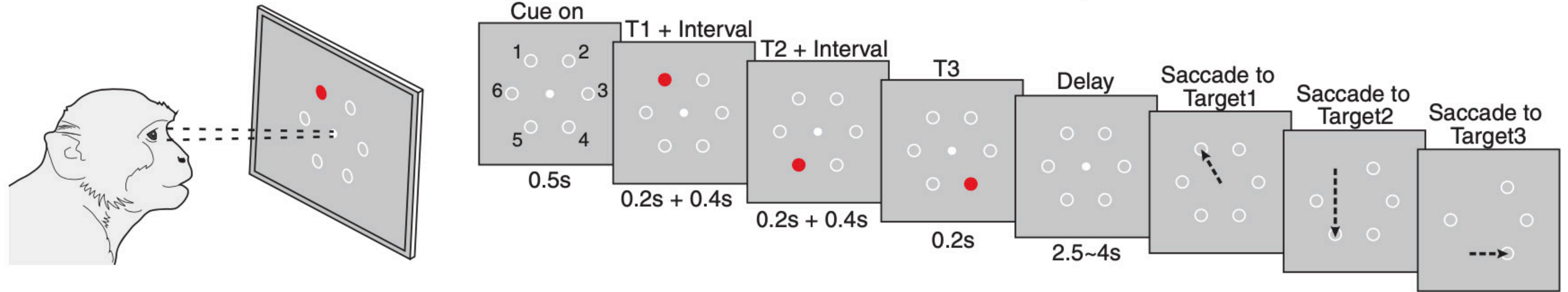
Xie et al., 2022



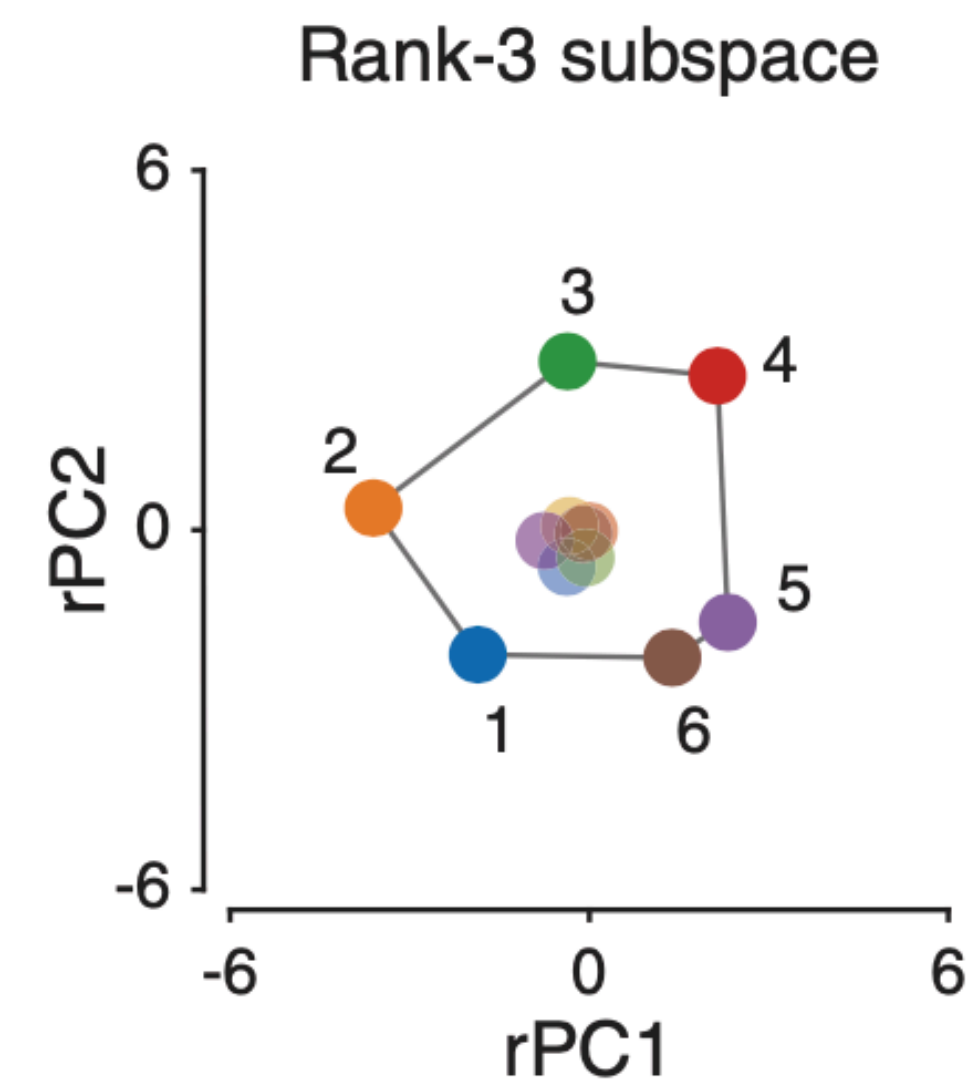
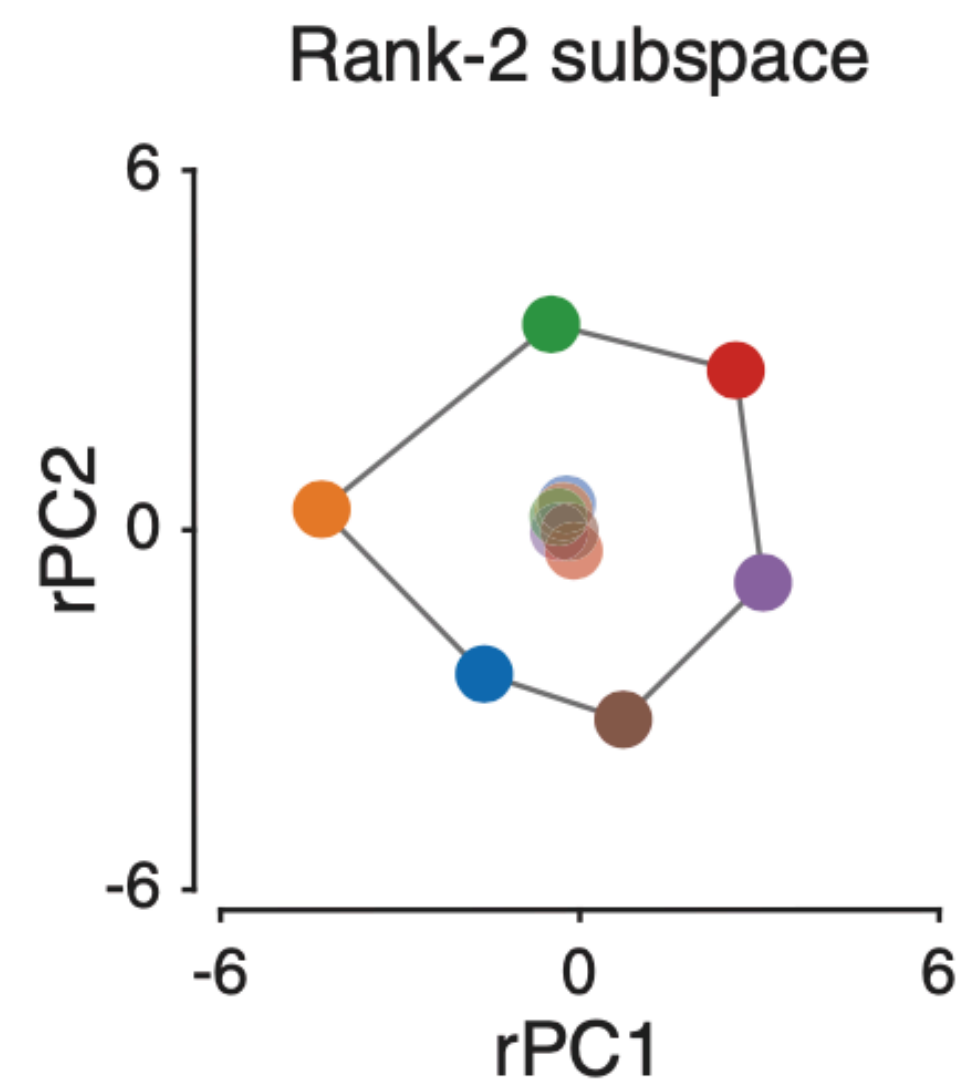
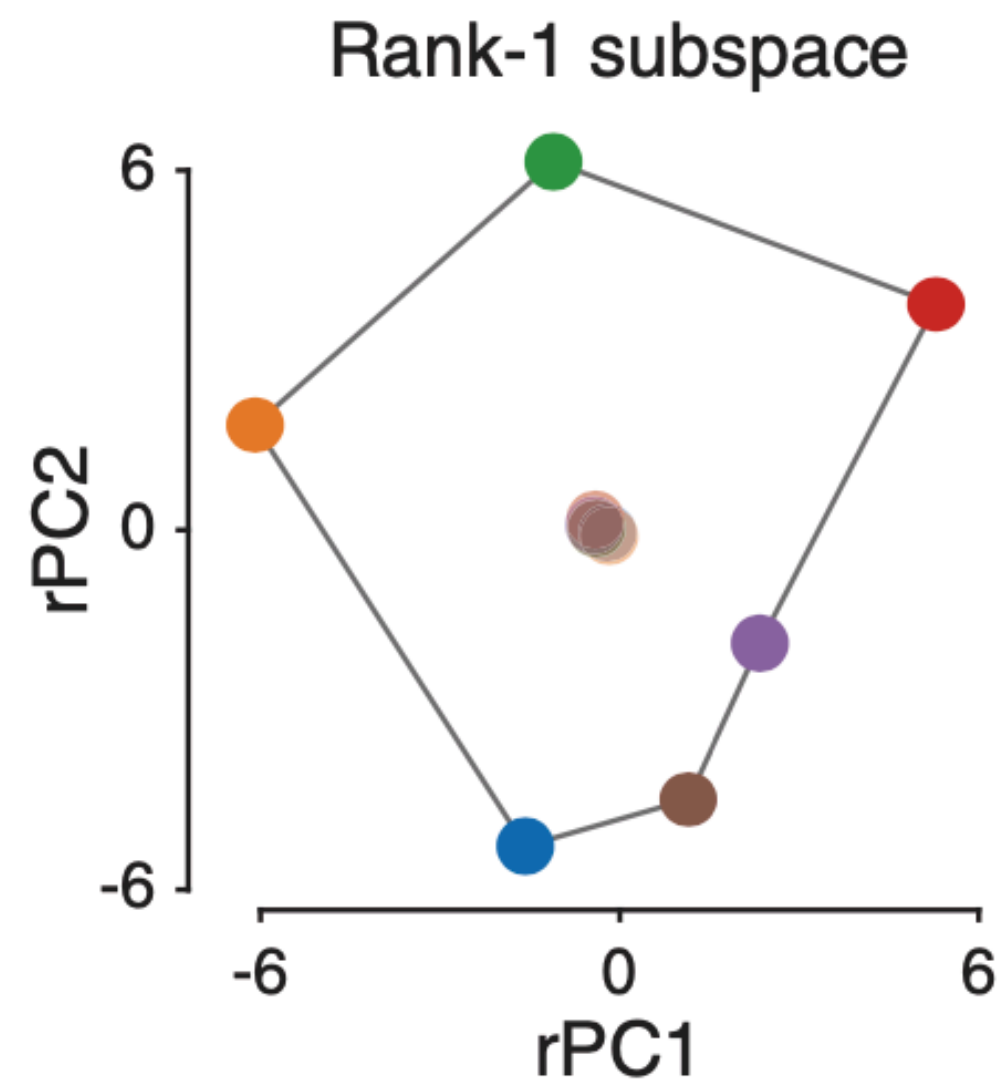
For sequence 1,3,4



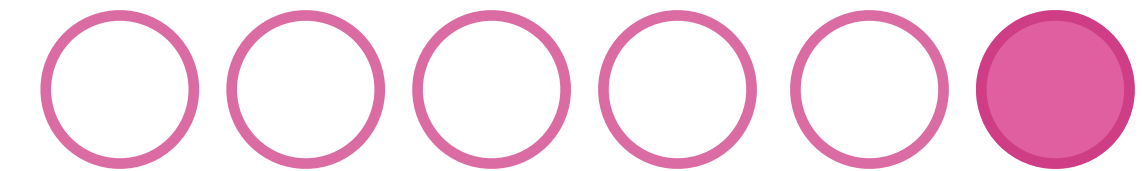
# Prefrontal cells represent the whole history at any given time



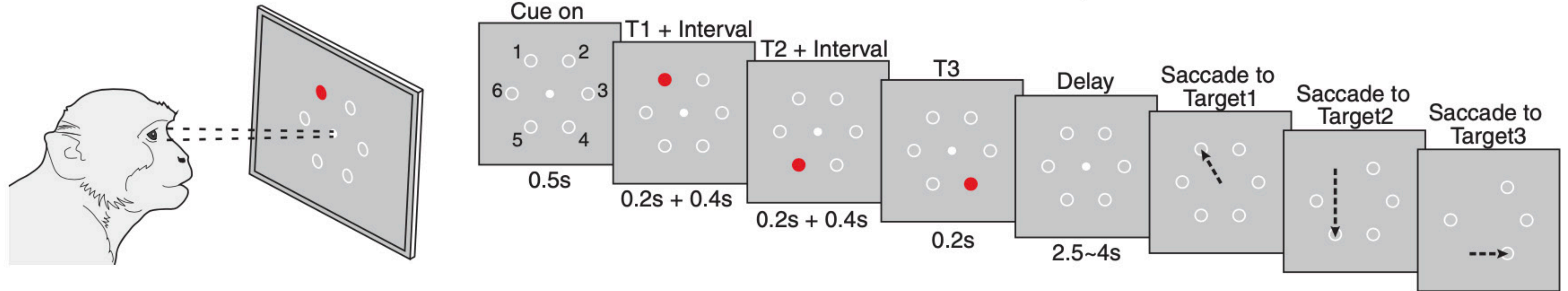
Xie et al., 2022



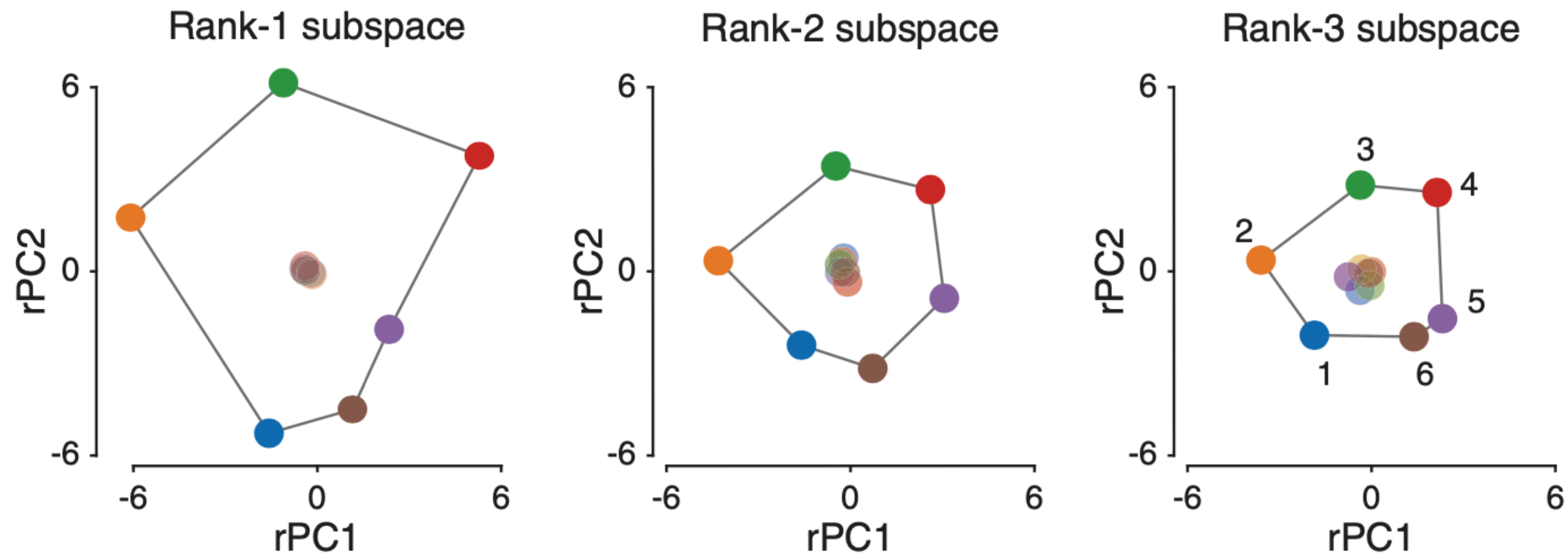
For sequence 5,2,6



# Prefrontal cells represent the whole history at any given time



Xie et al., 2022



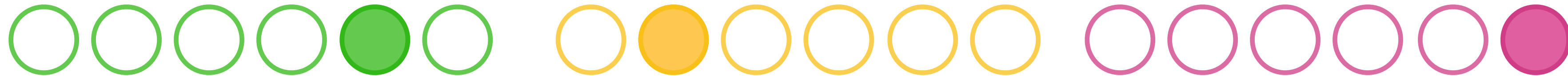
For sequence 5,2,6



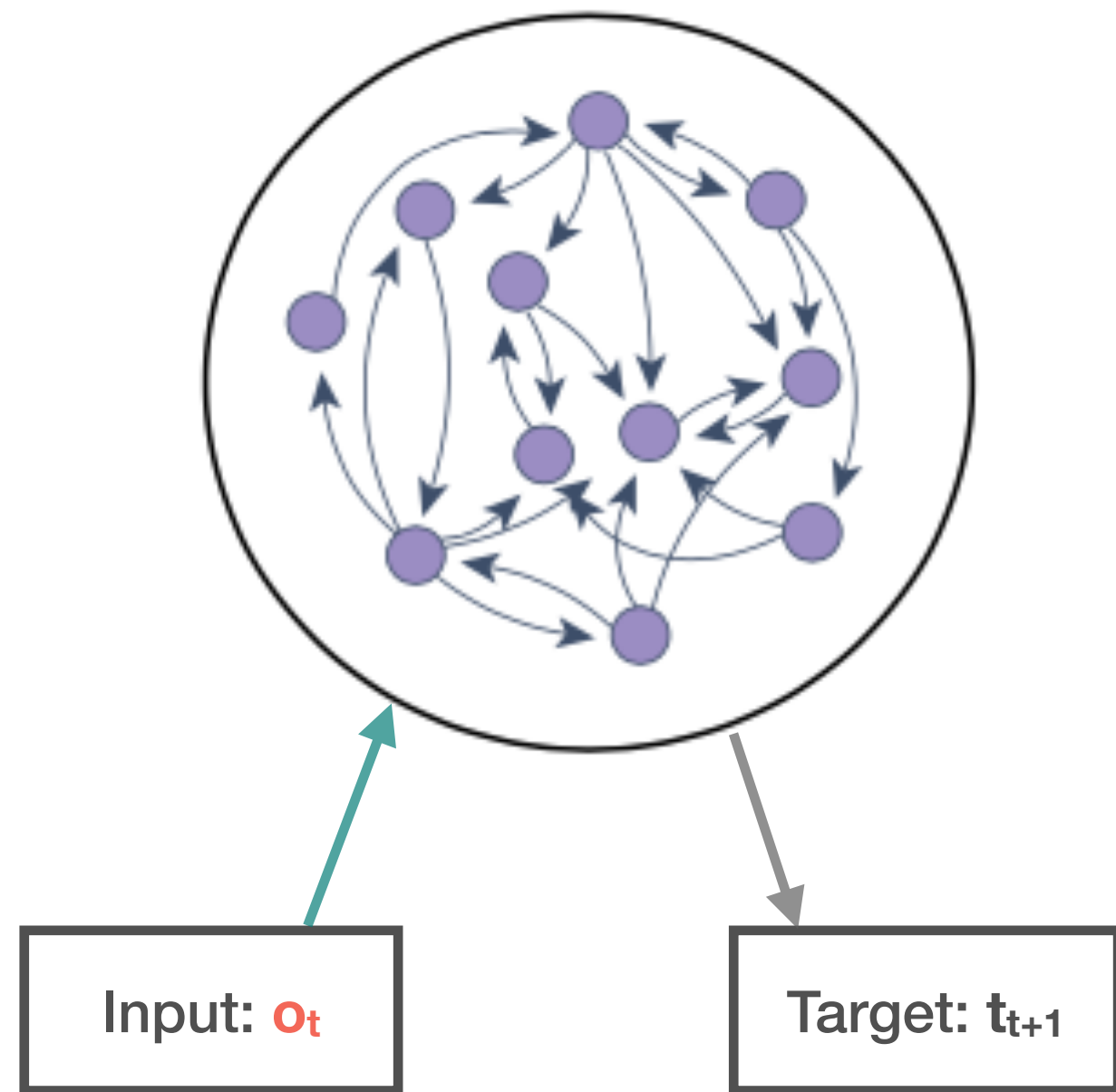
**These neurons keep a memory of the past!**

# Recurrent neural networks have internal memory

For sequence 5,2,6



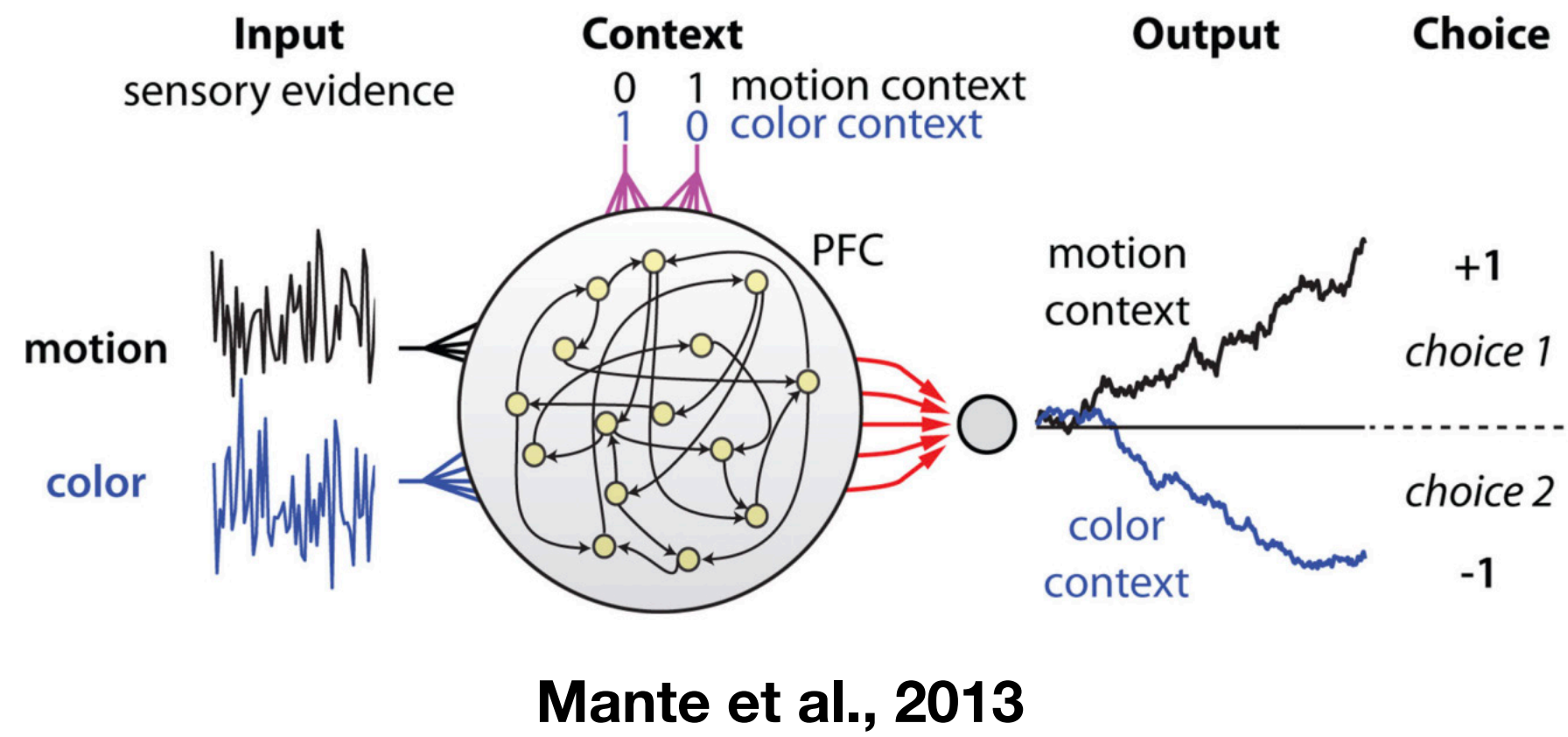
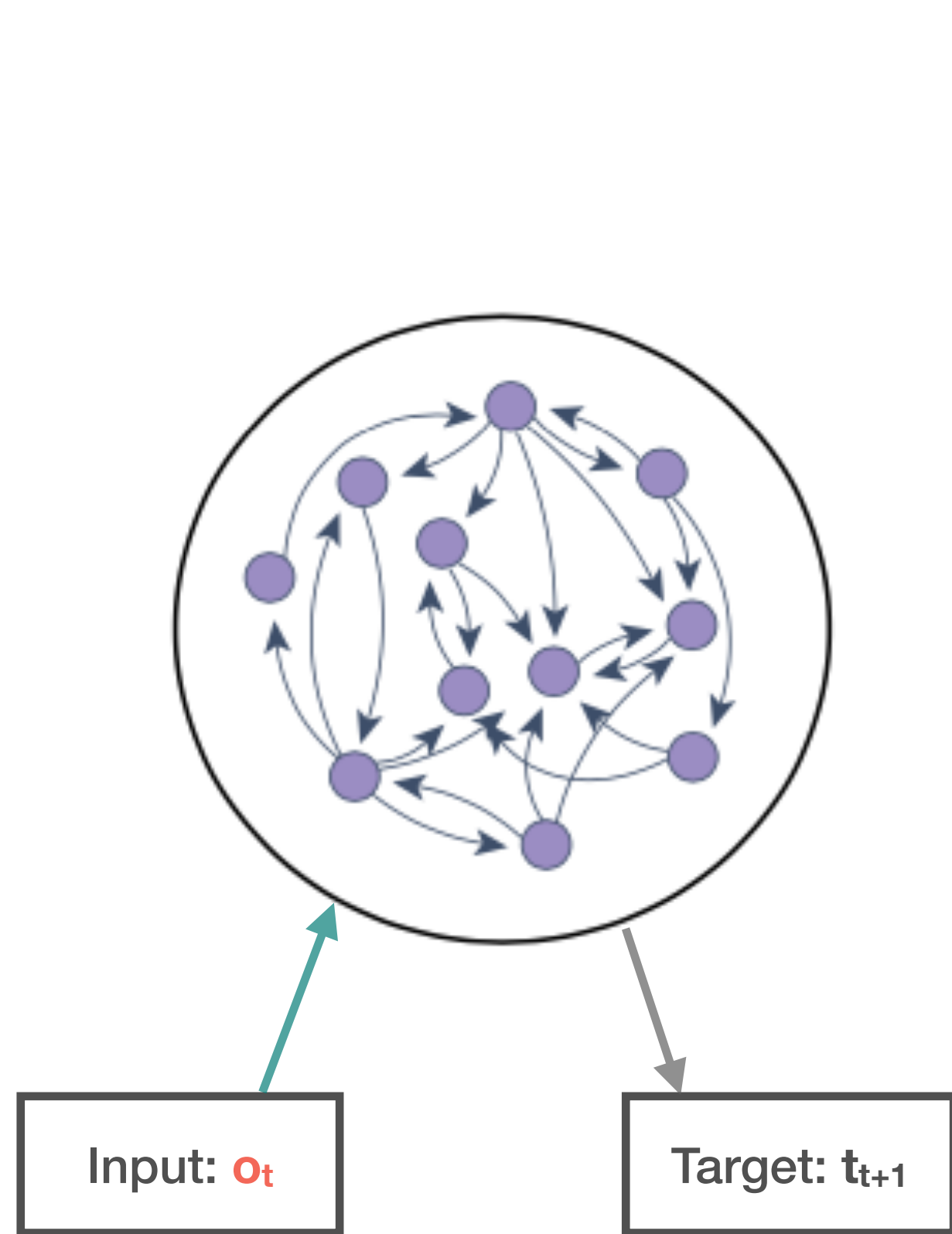
# Recurrent neural networks have internal memory



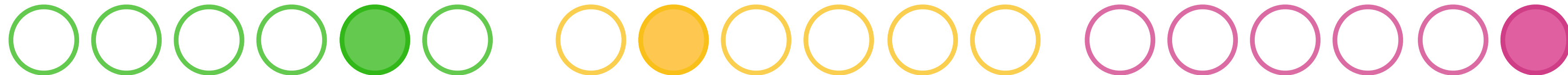
For sequence 5,2,6



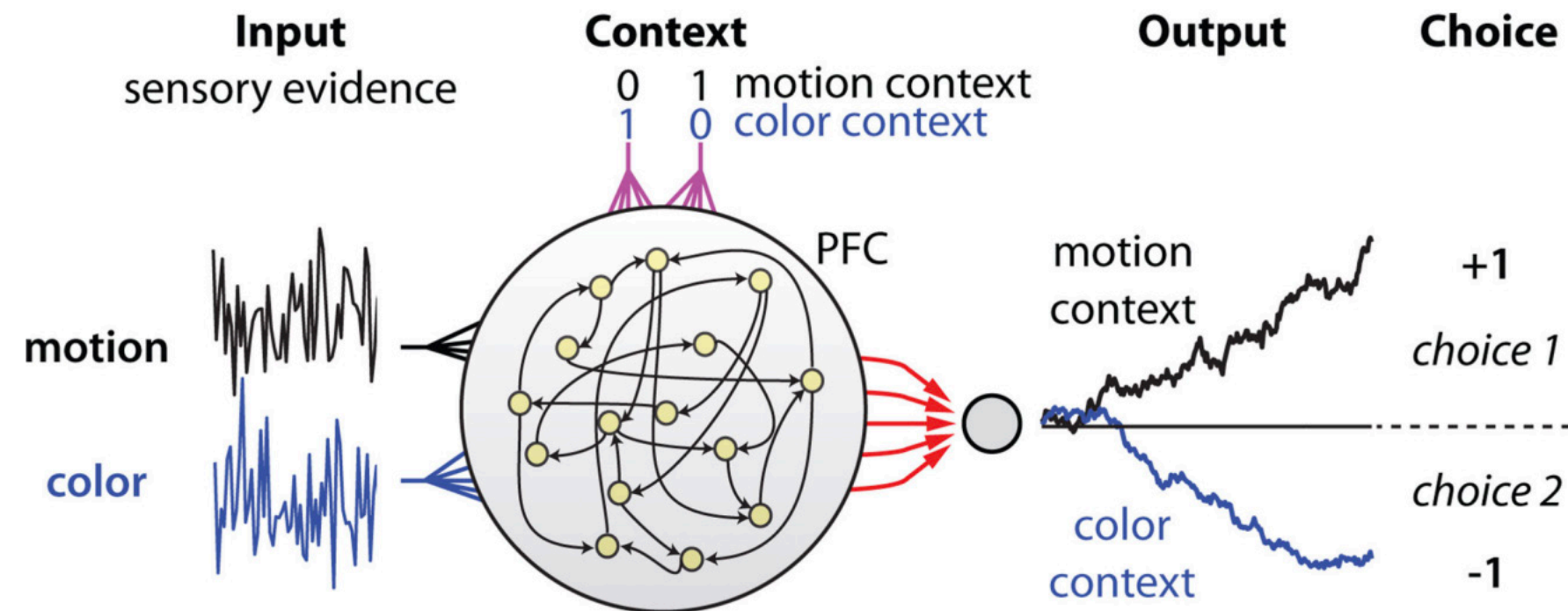
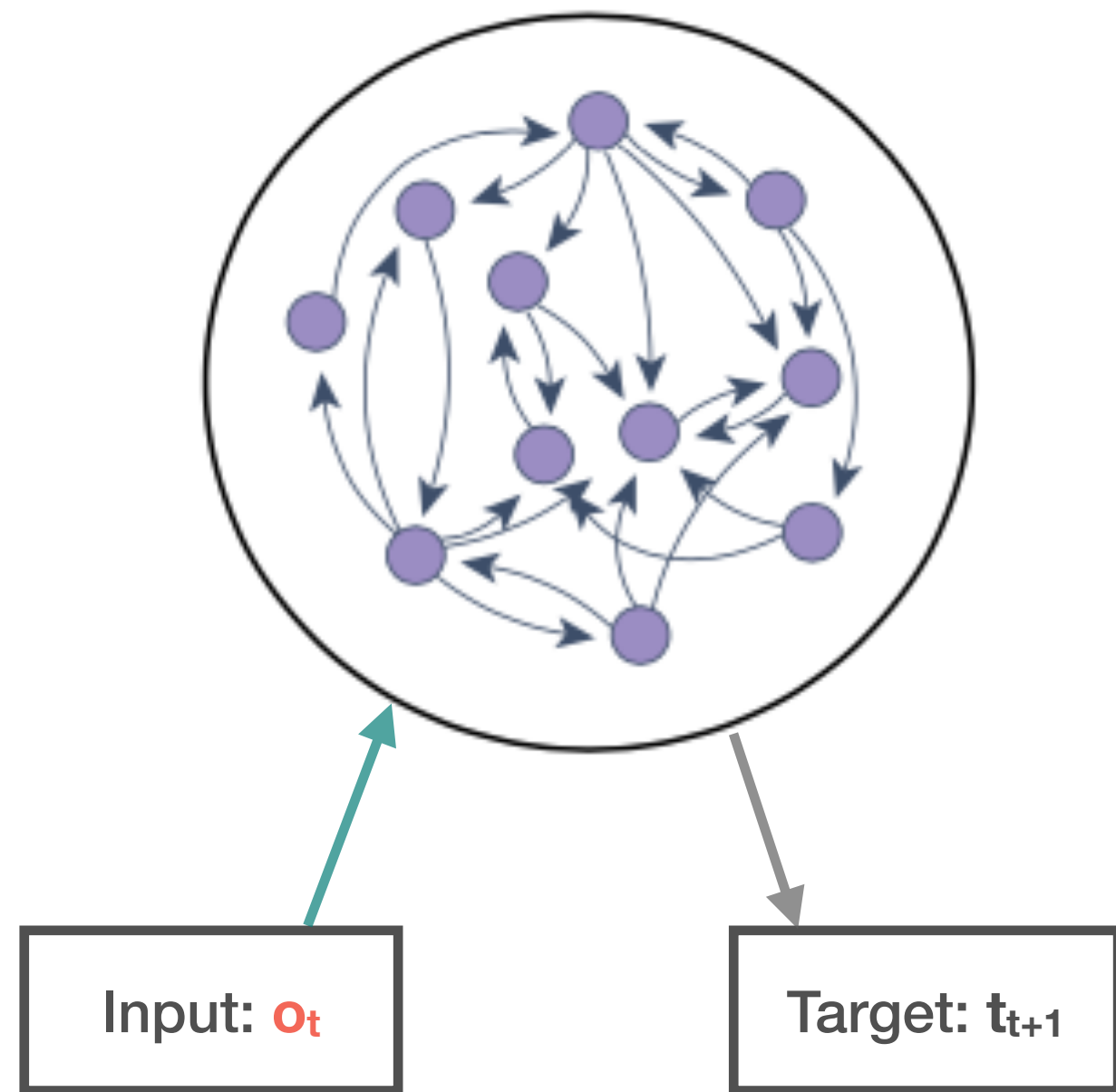
# Recurrent neural networks have internal memory



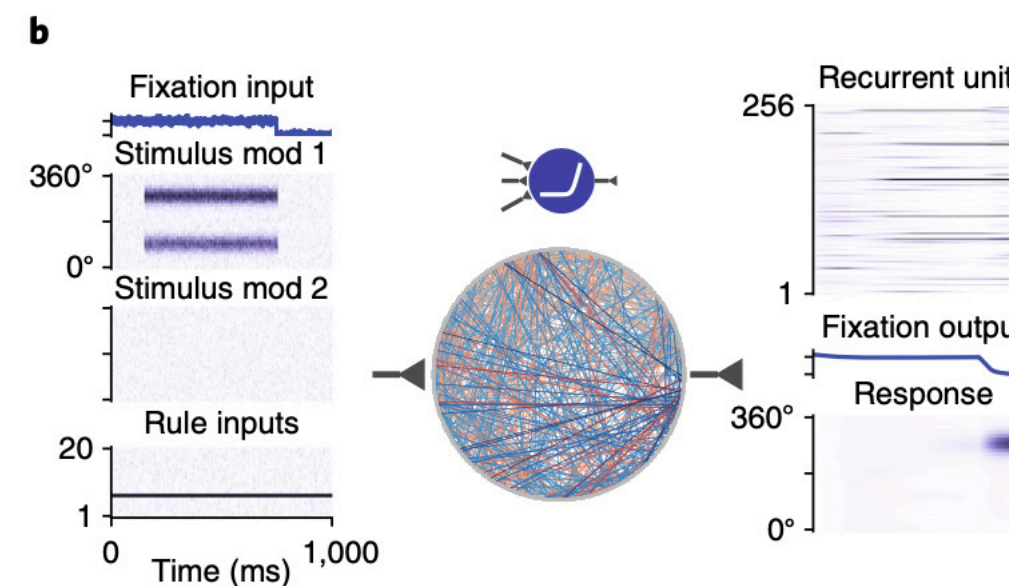
For sequence 5,2,6



# Recurrent neural networks have internal memory



Mante et al., 2013



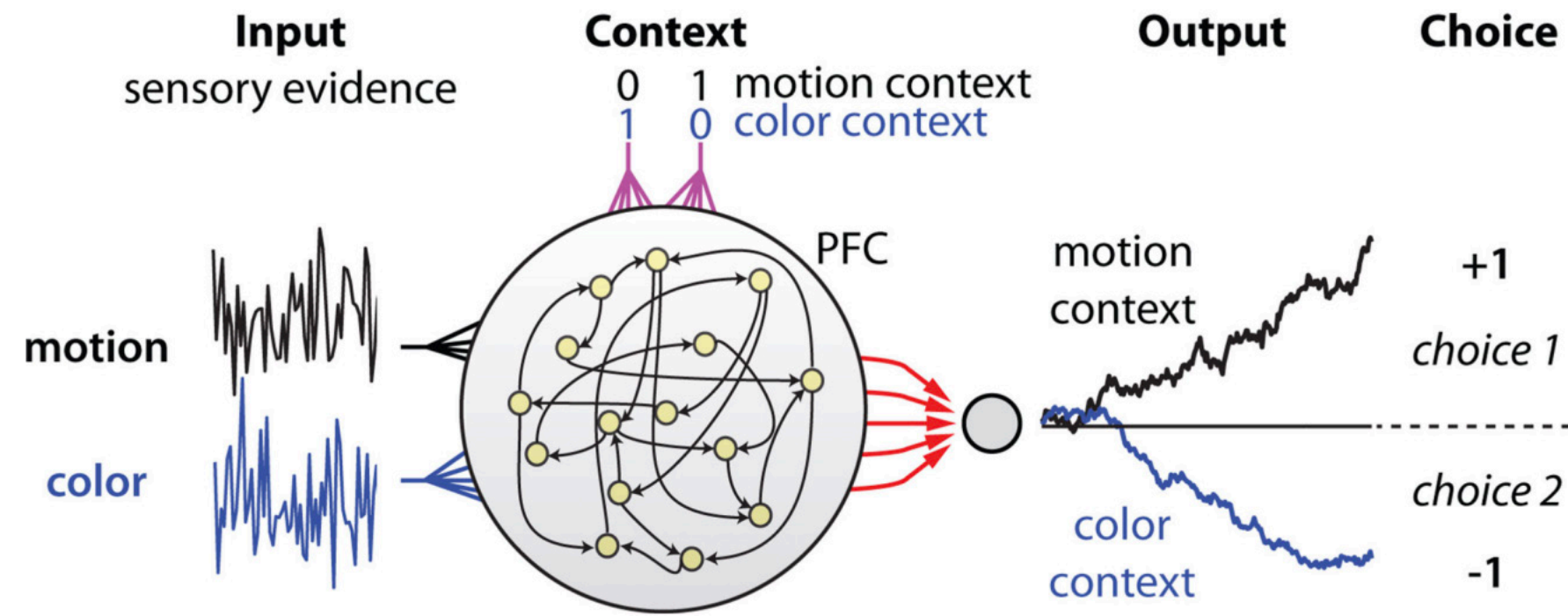
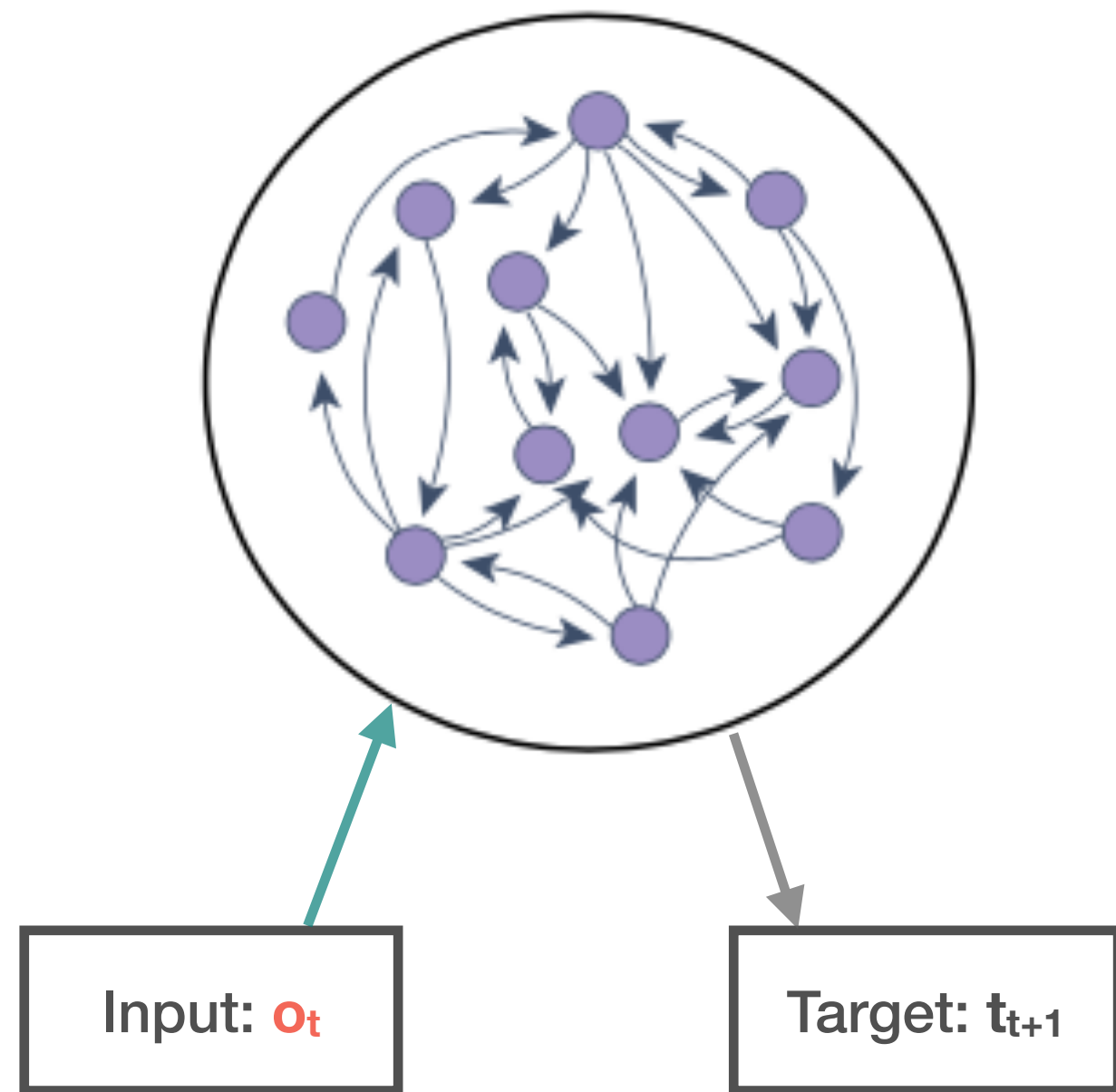
Yang et al., 2019

For sequence 5,2,6

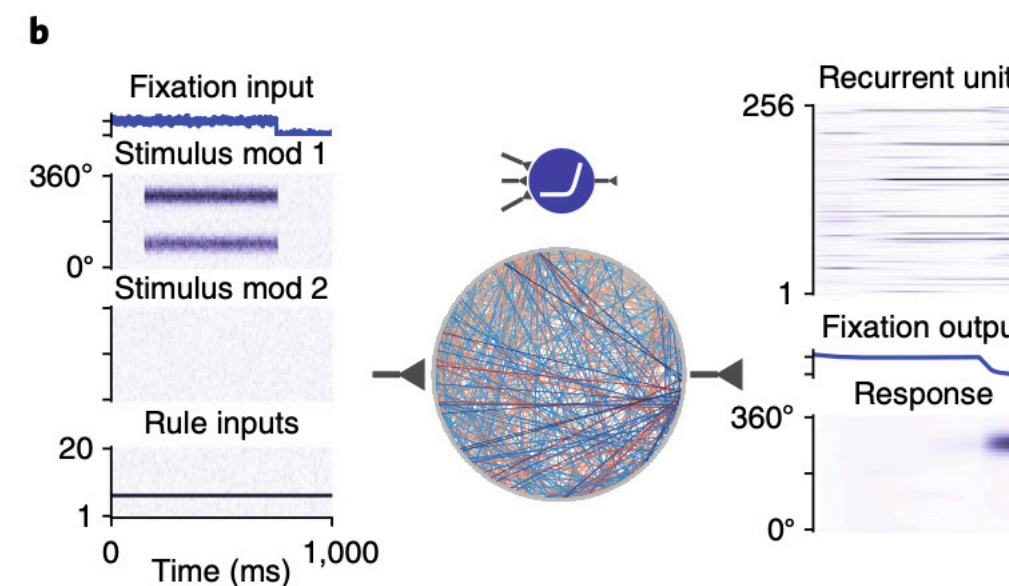




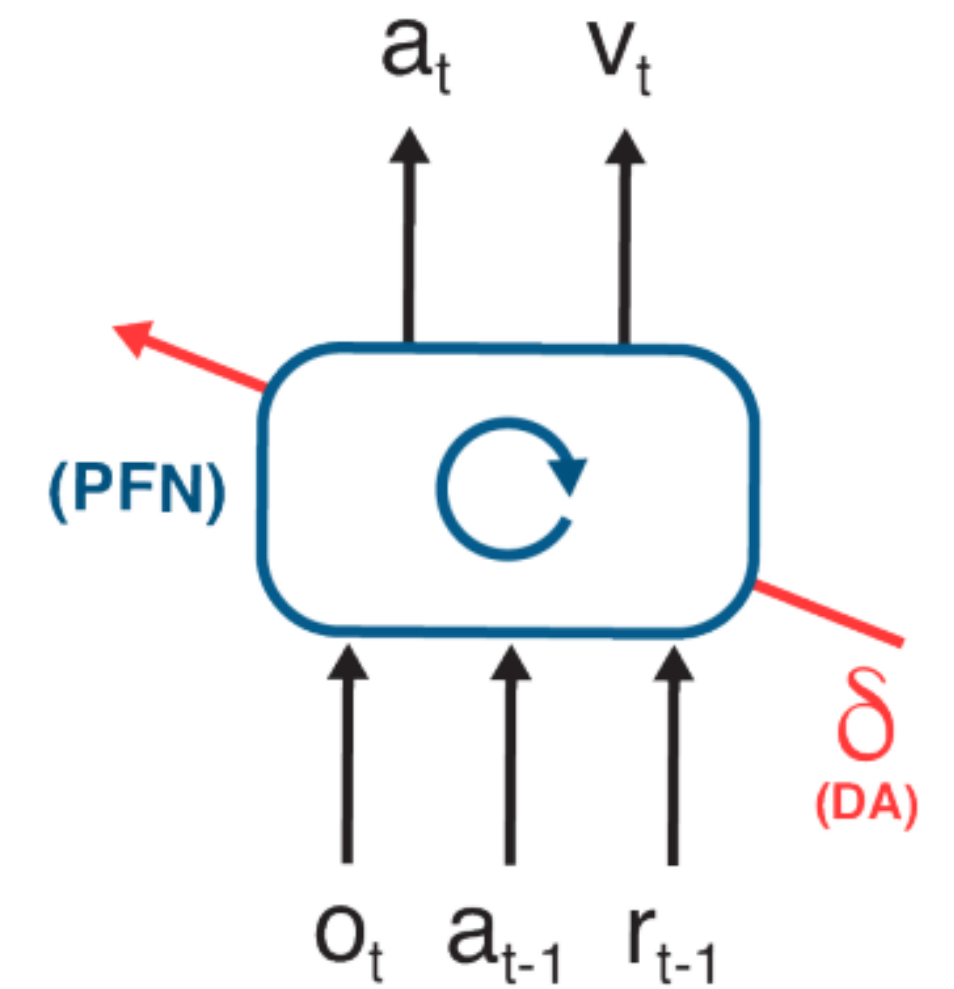
# Recurrent neural networks have internal memory



Mante et al., 2013



Yang et al., 2019

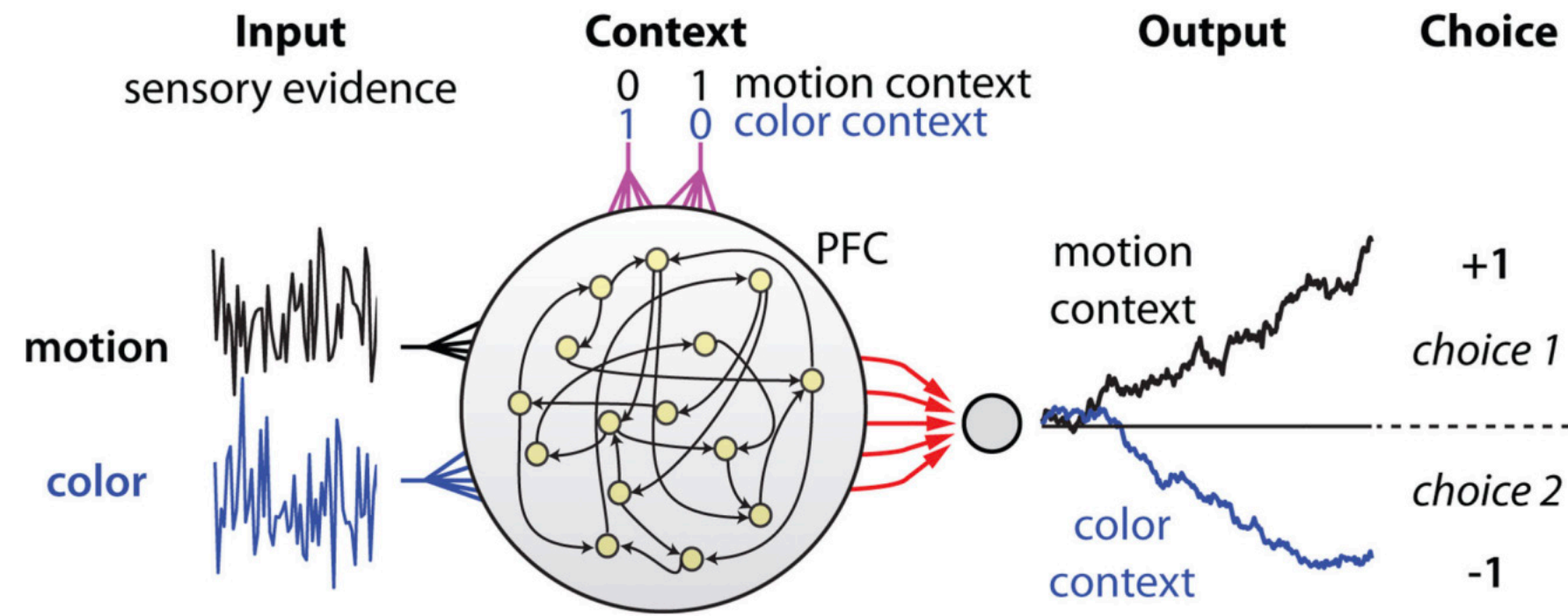
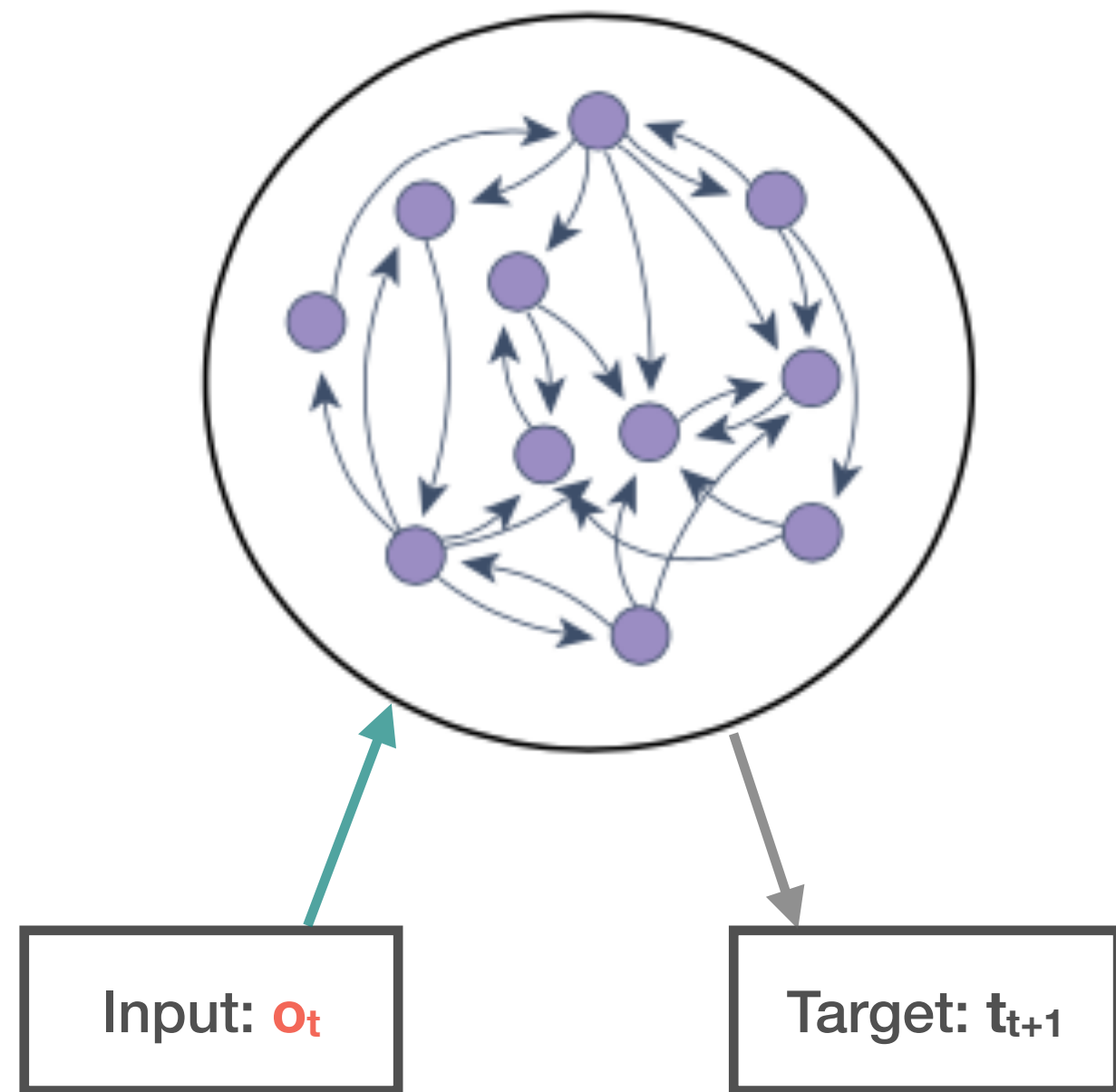


Wang et al., 2018

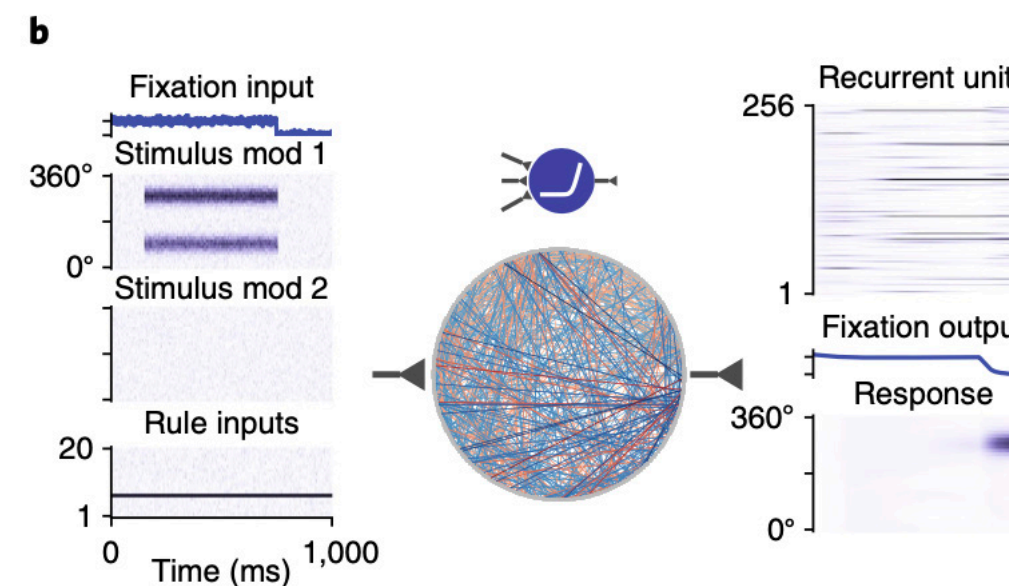
For sequence 5,2,6



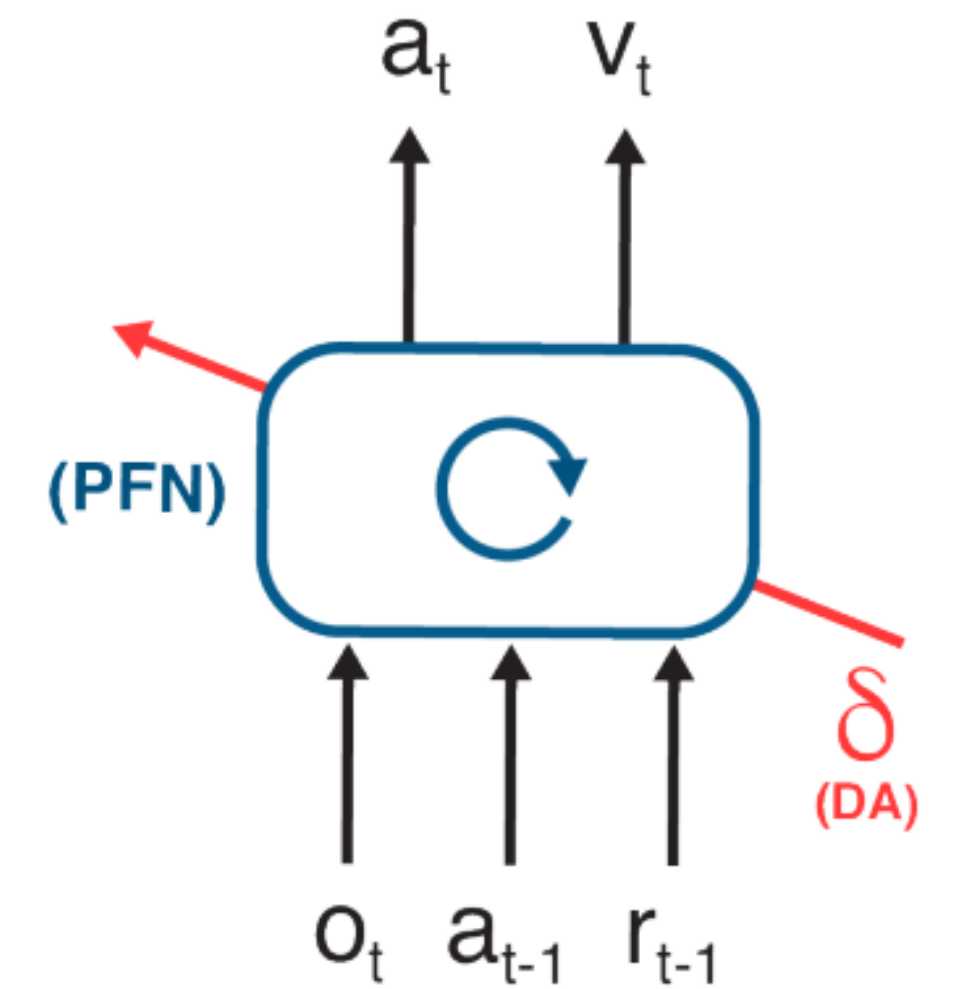
# Recurrent neural networks have internal memory



Mante et al., 2013



Yang et al., 2019



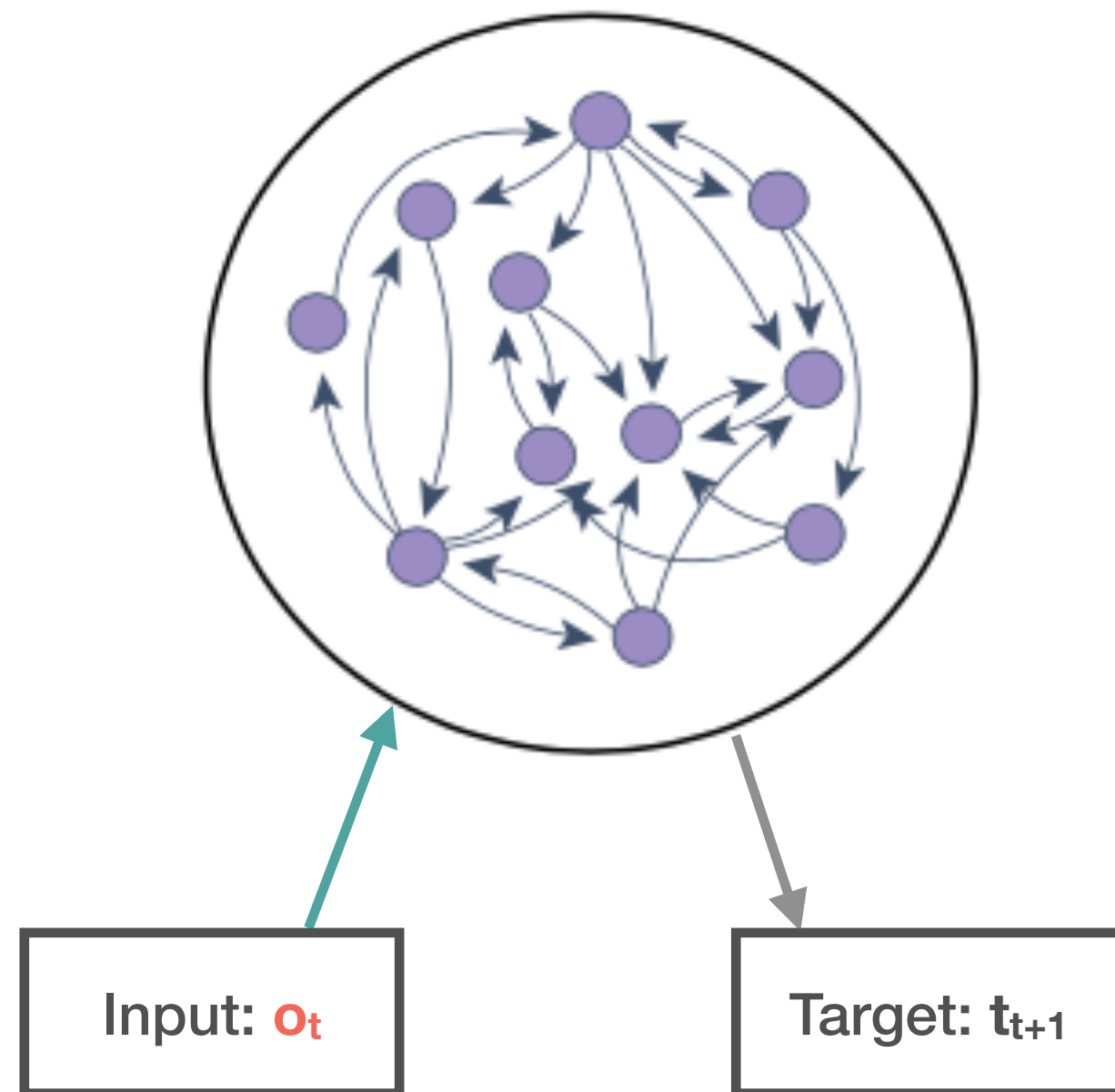
Wang et al., 2018

**We want to understand the mechanism of these RNNs!**

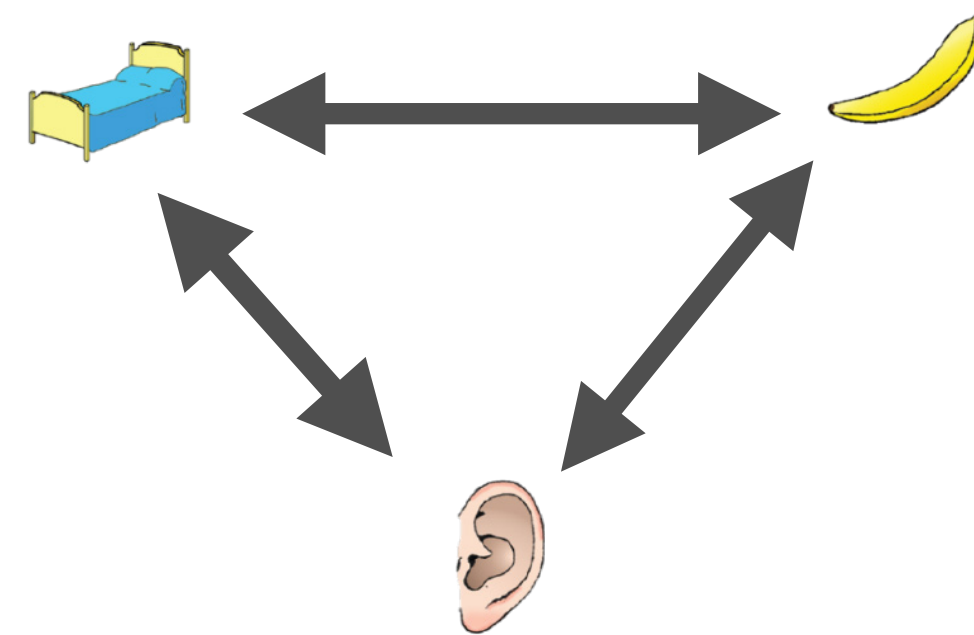
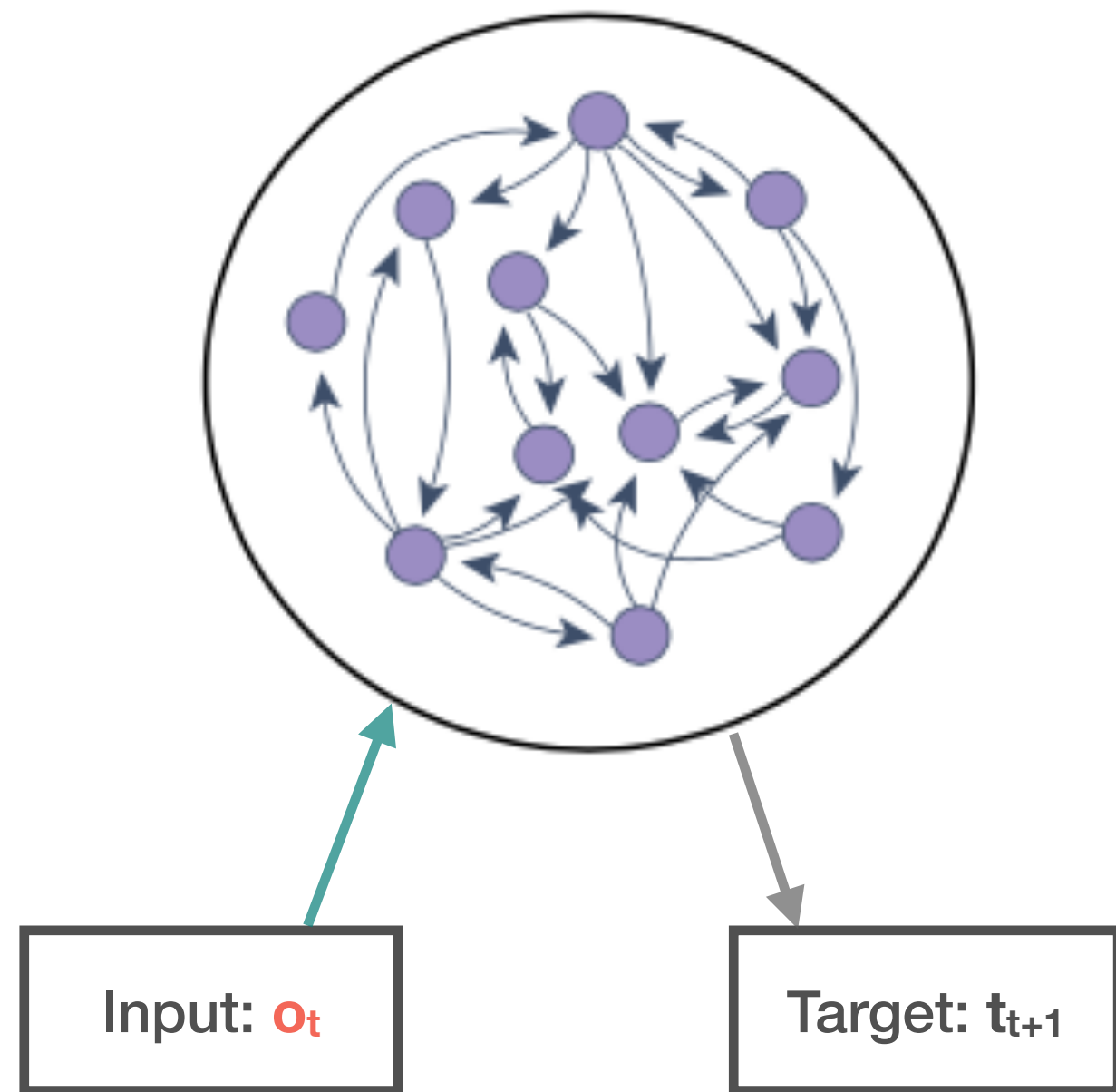
For sequence 5,2,6



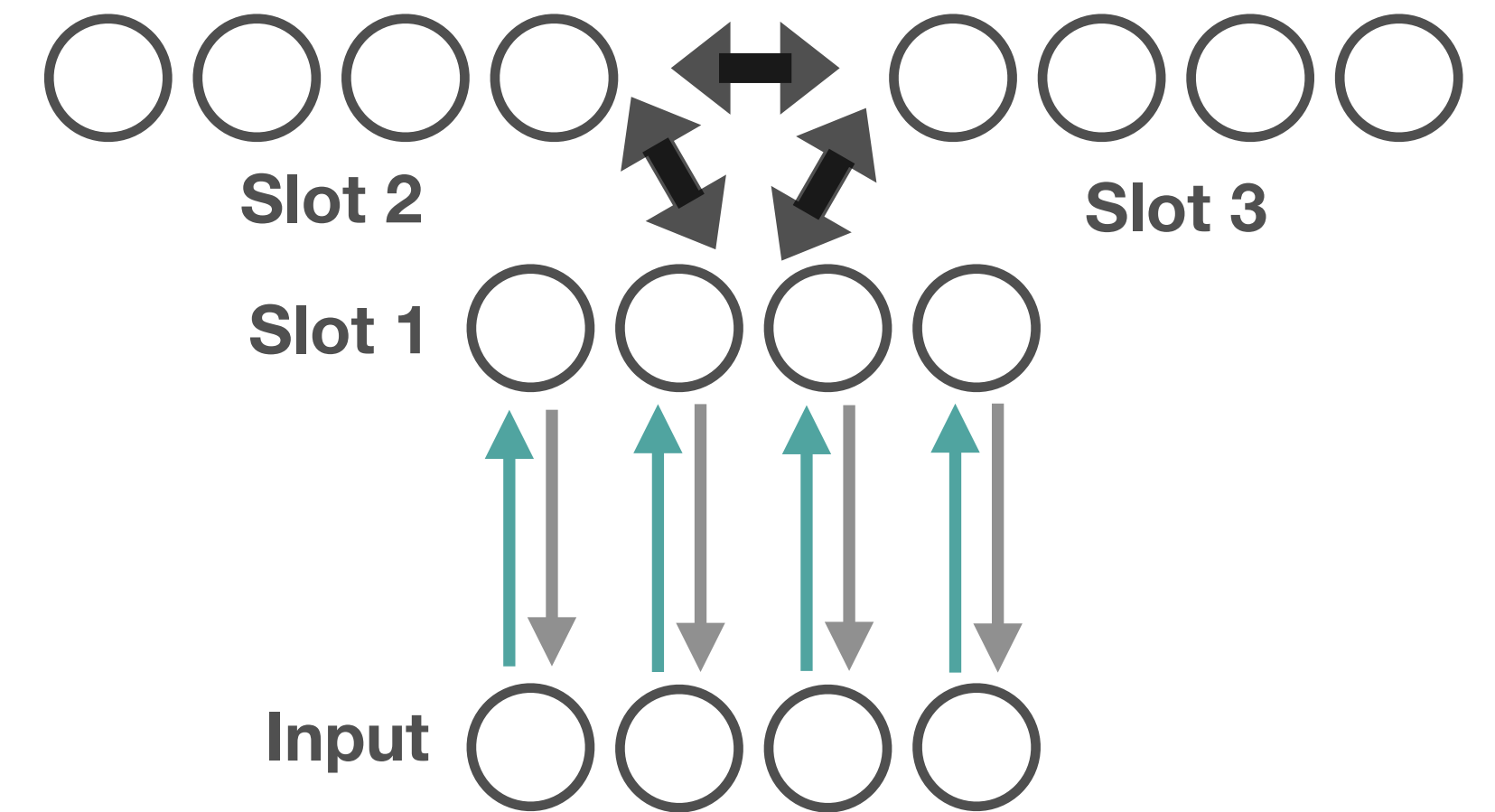
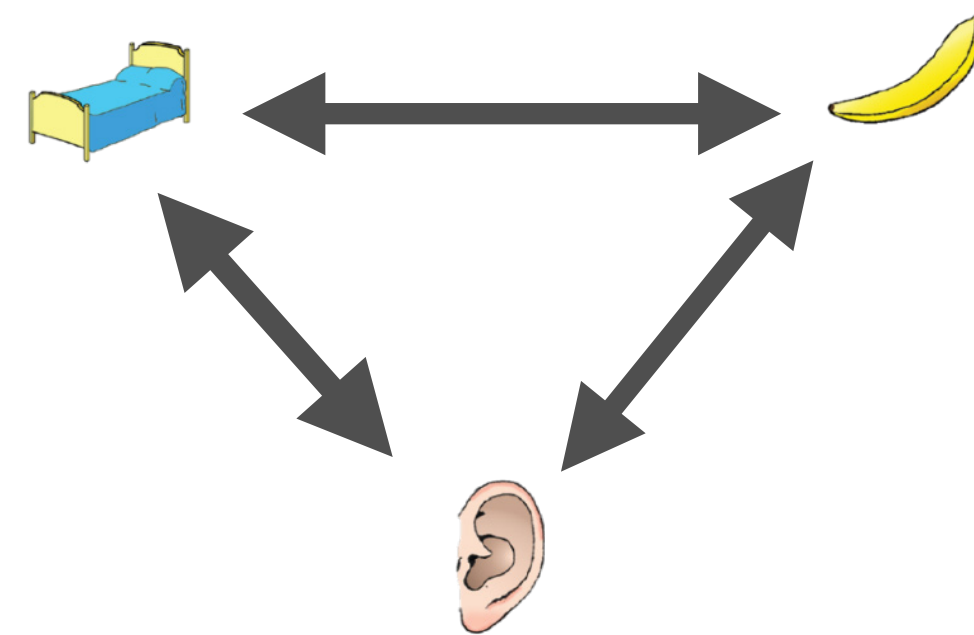
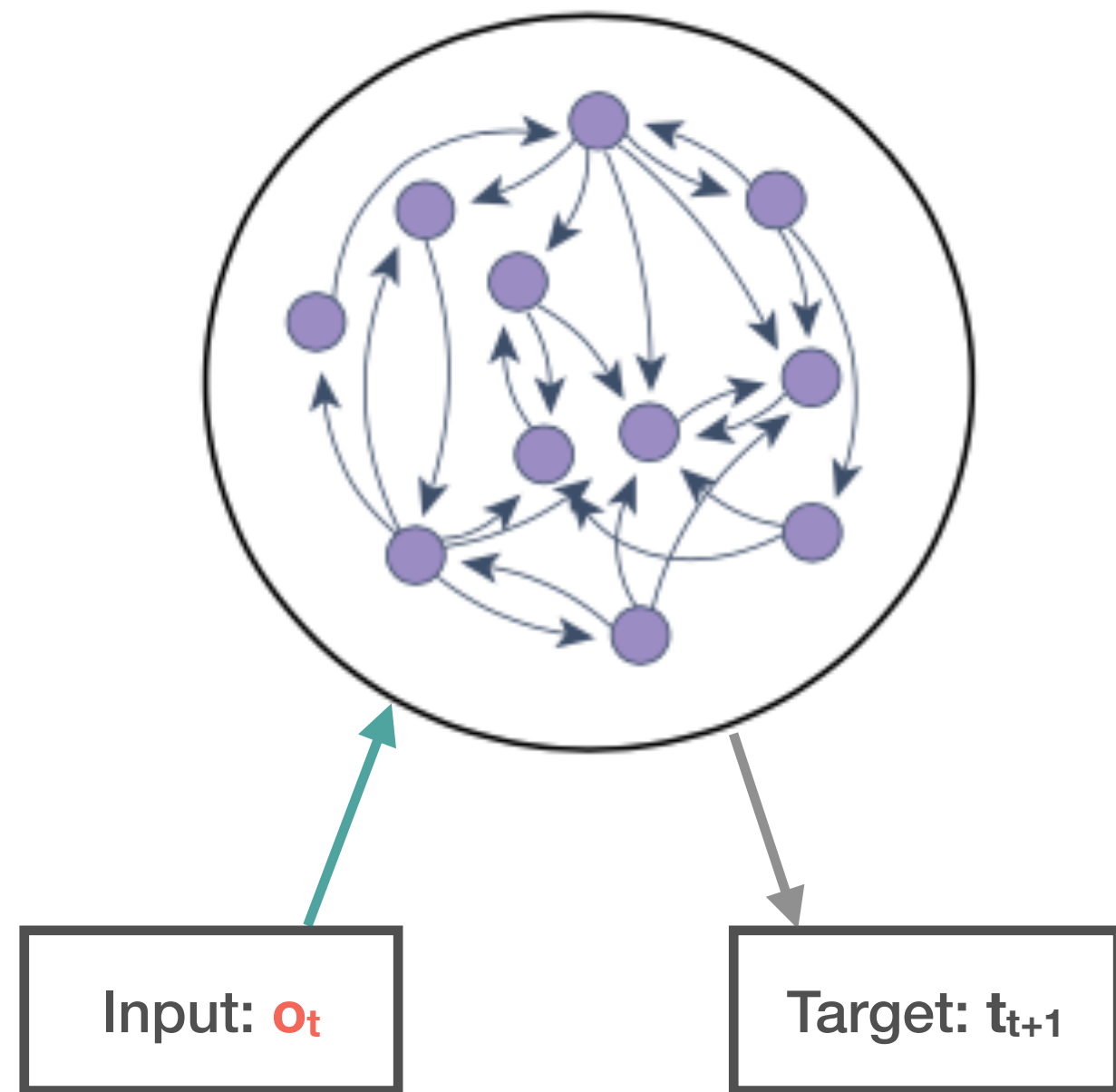
# Hypothesis: Sequence working memory using neural circuits of structured slots



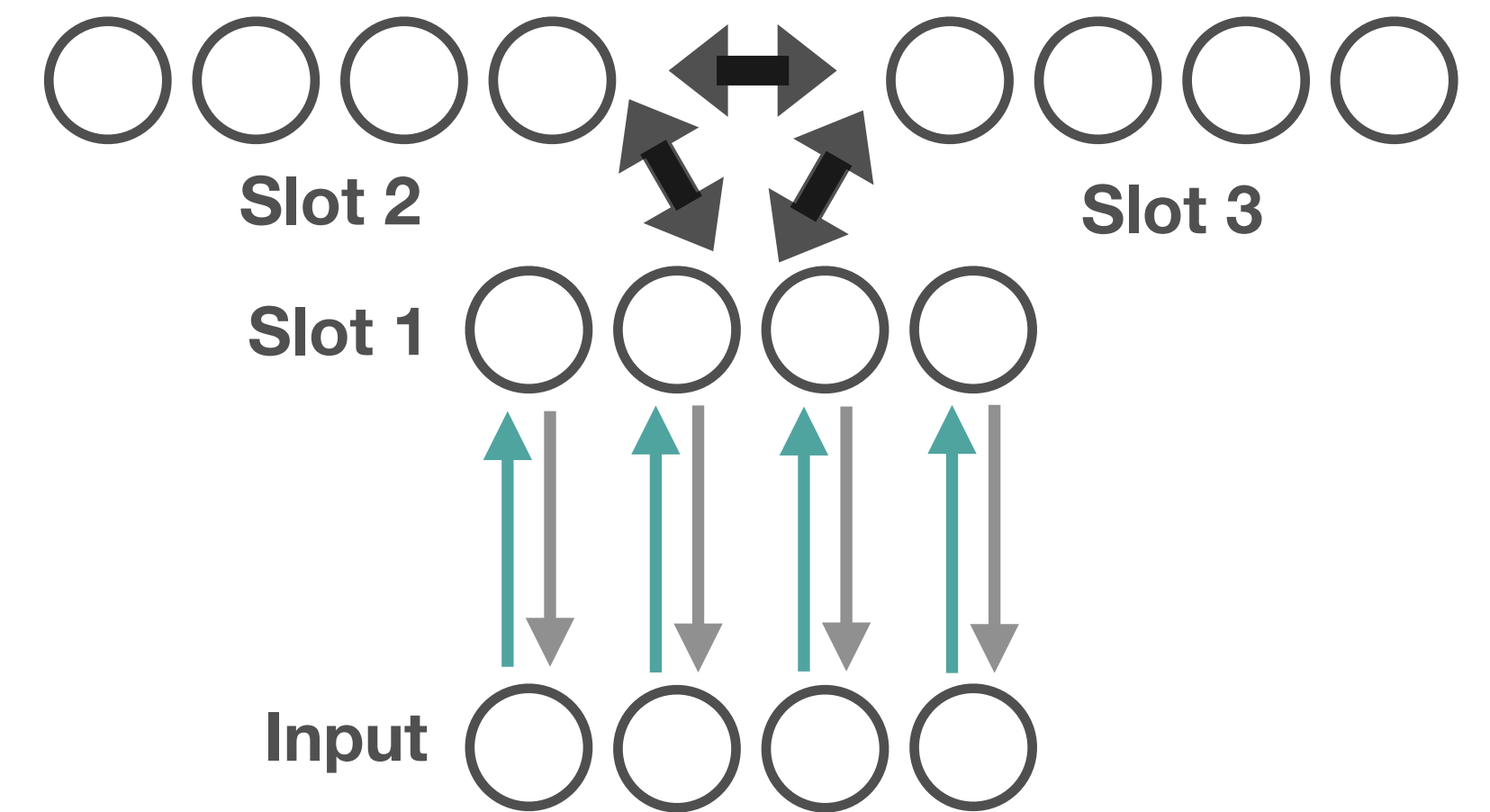
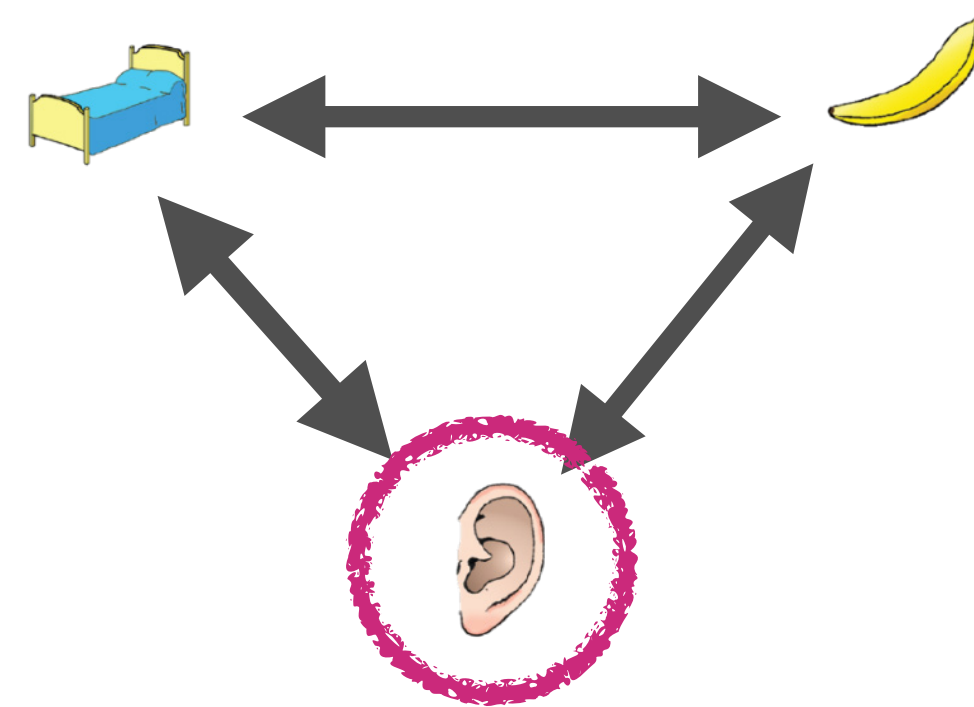
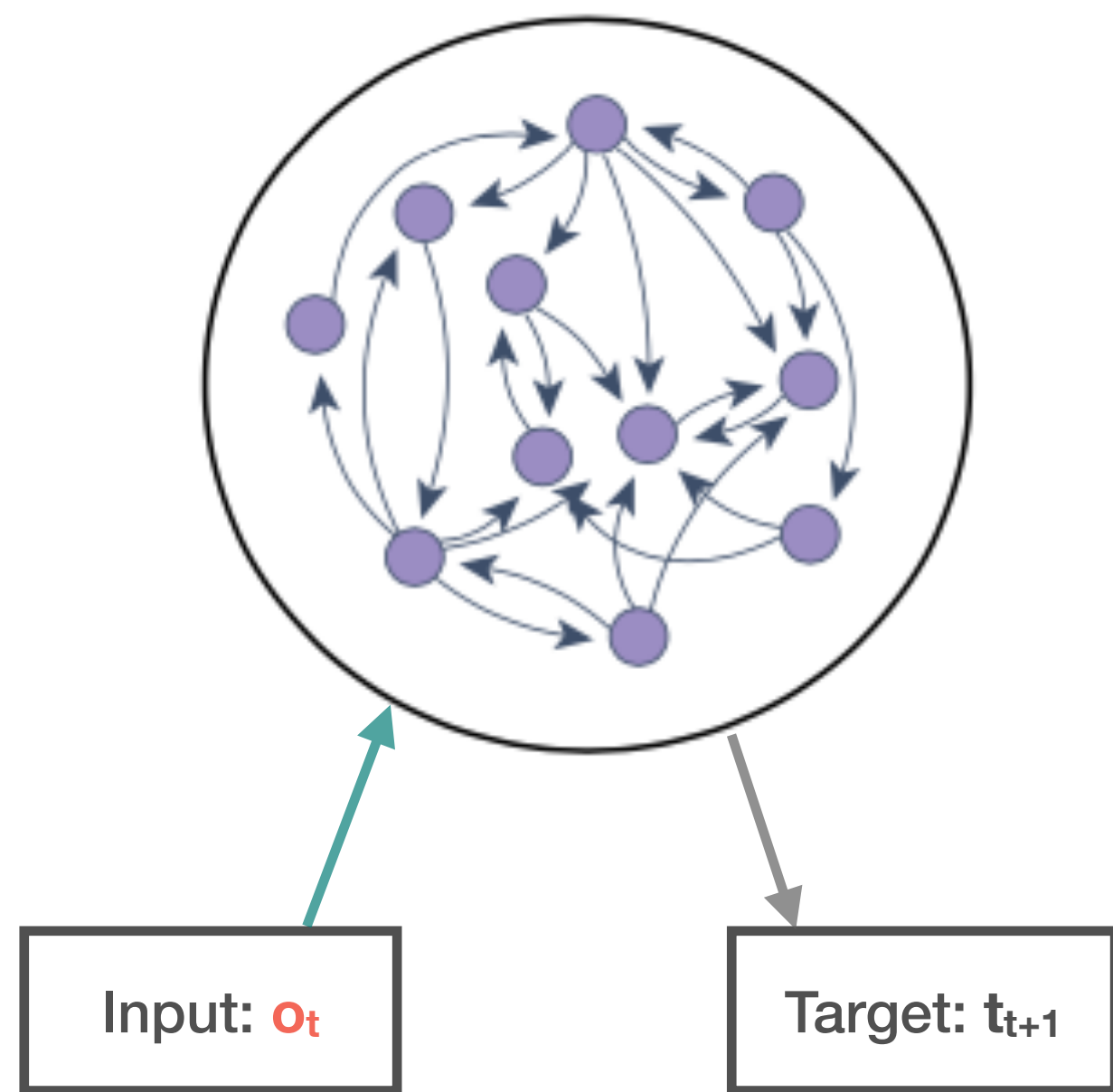
# Hypothesis: Sequence working memory using neural circuits of structured slots



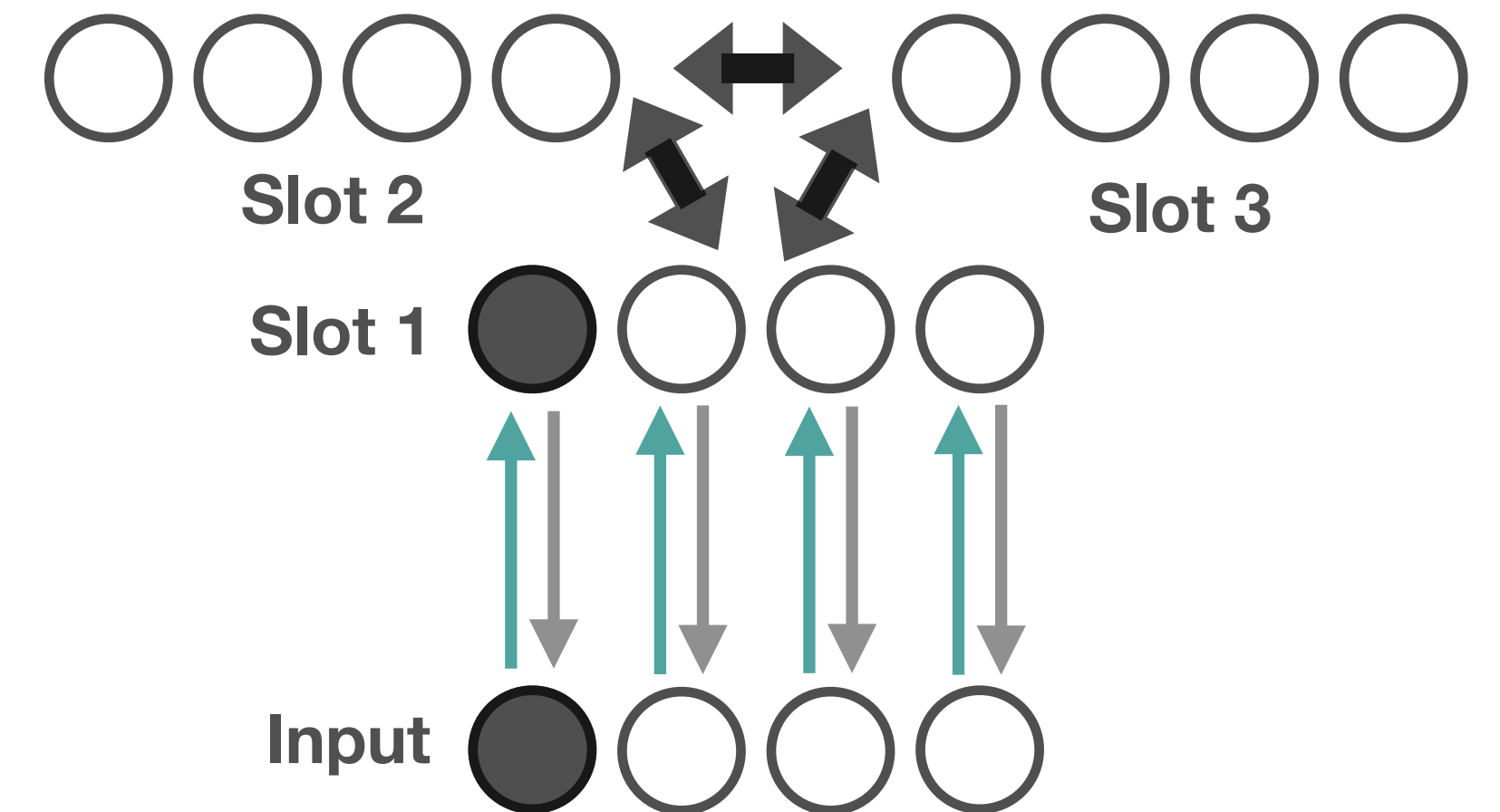
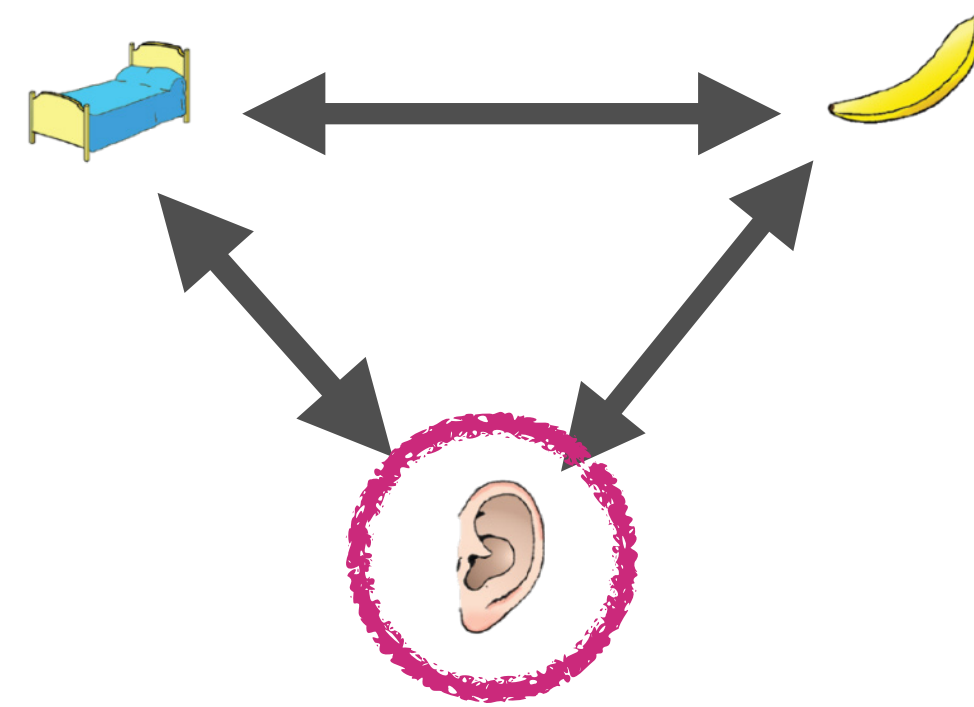
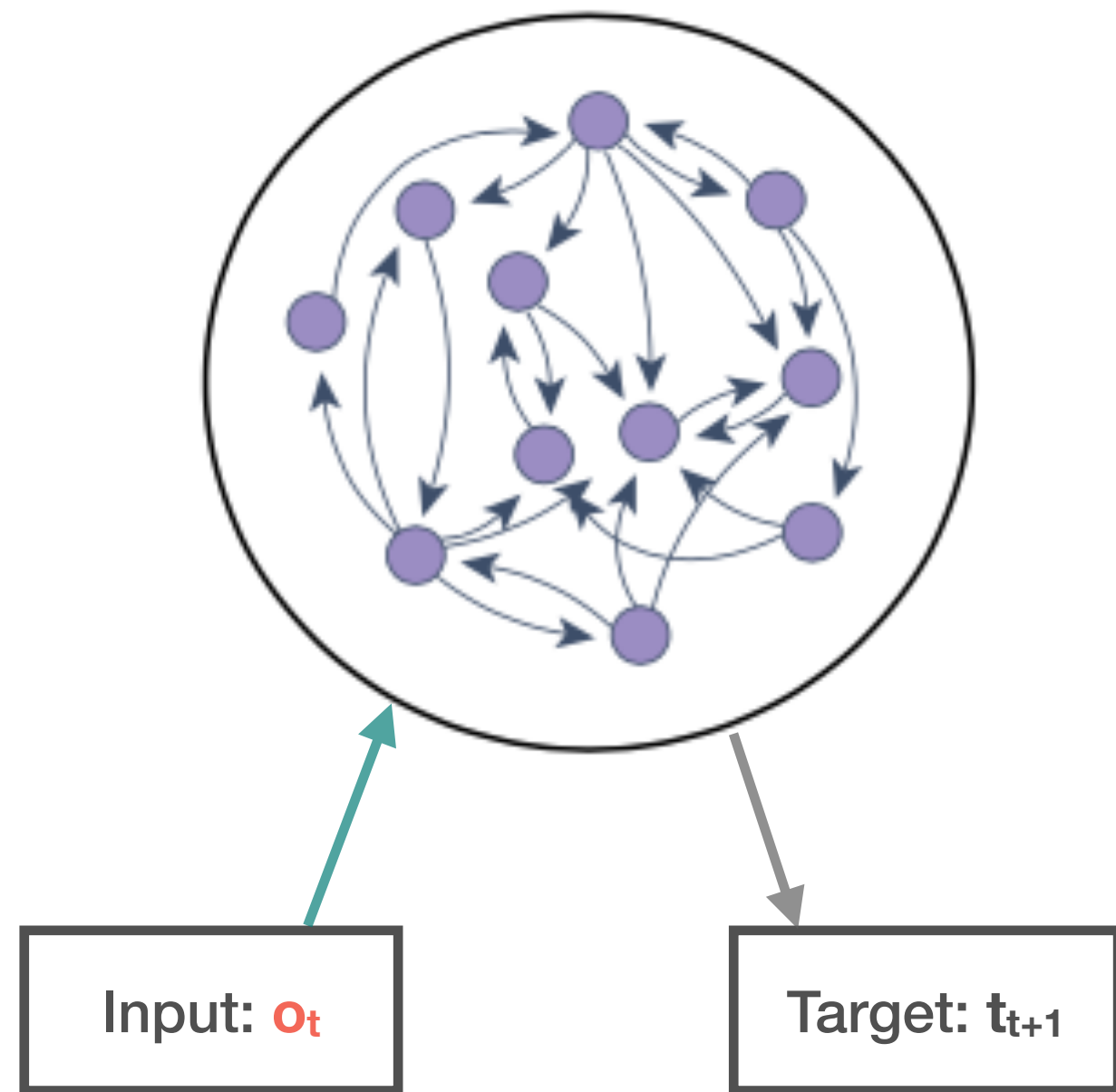
# Hypothesis: Sequence working memory using neural circuits of structured slots



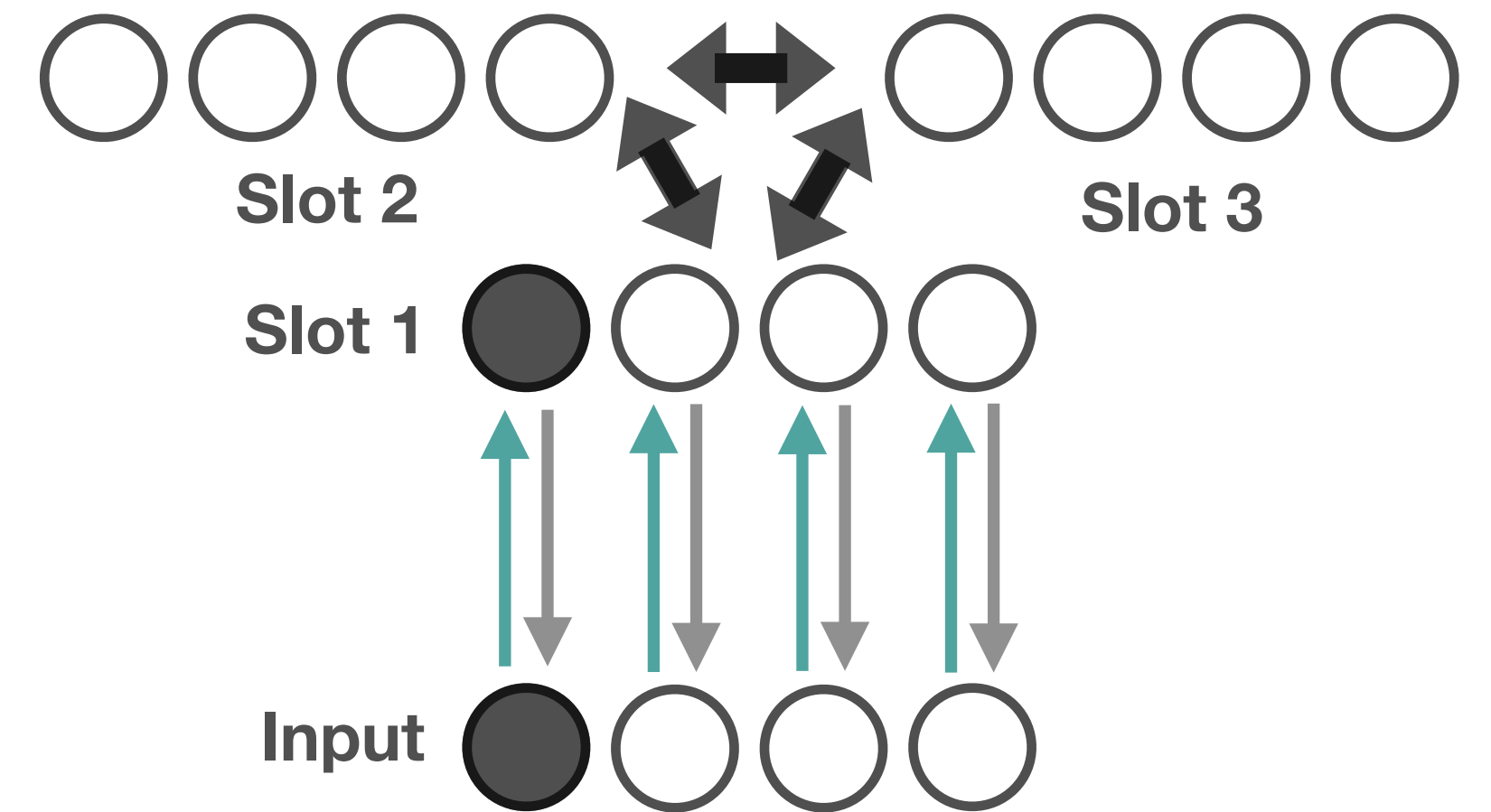
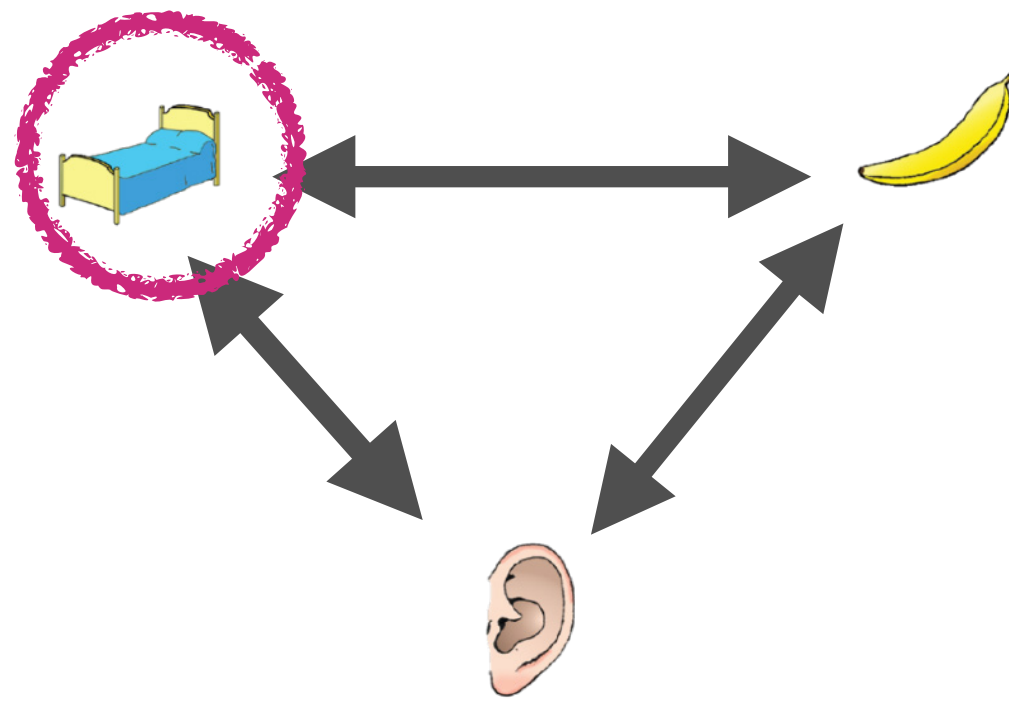
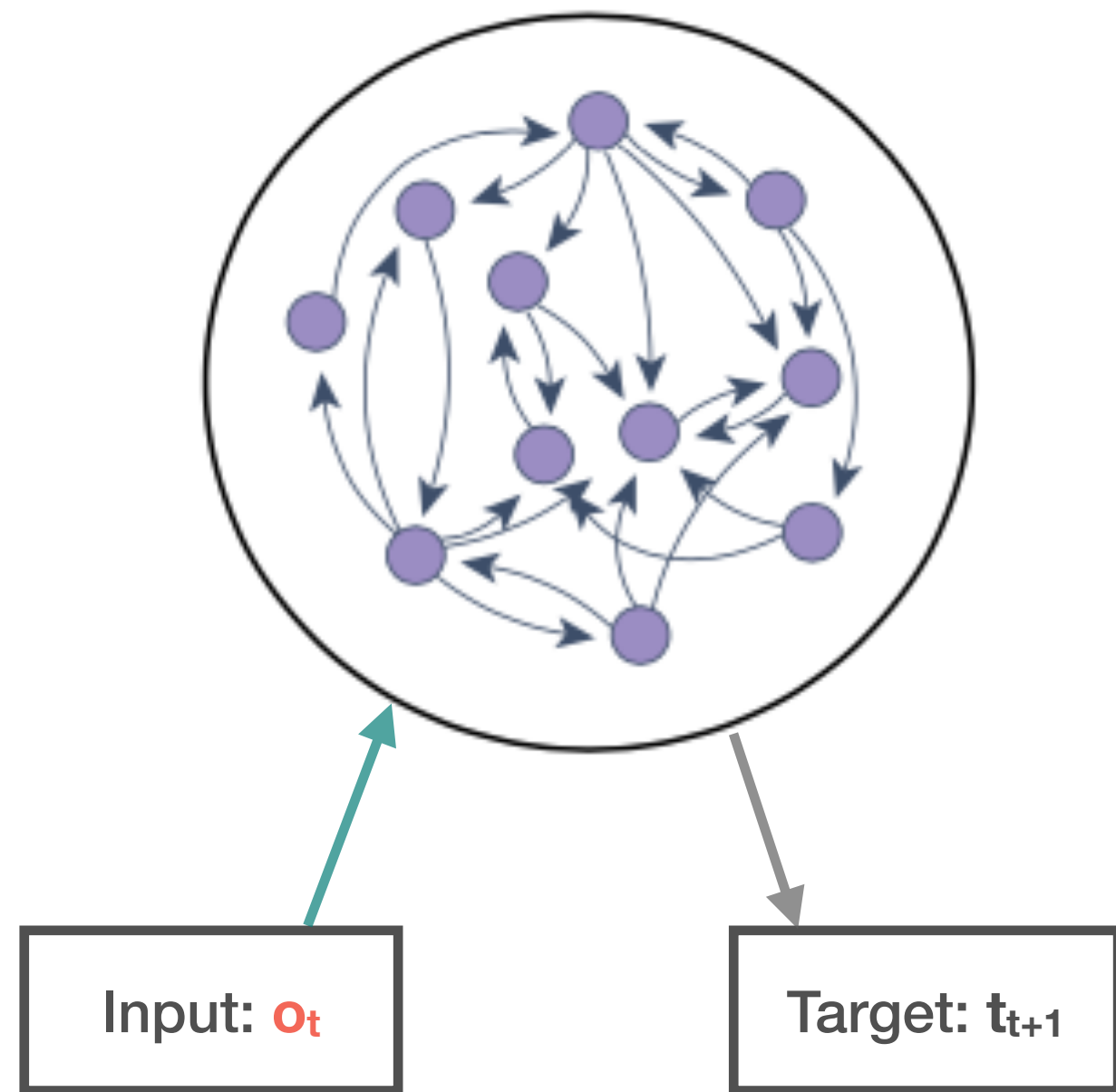
# Hypothesis: Sequence working memory using neural circuits of structured slots



# Hypothesis: Sequence working memory using neural circuits of structured slots

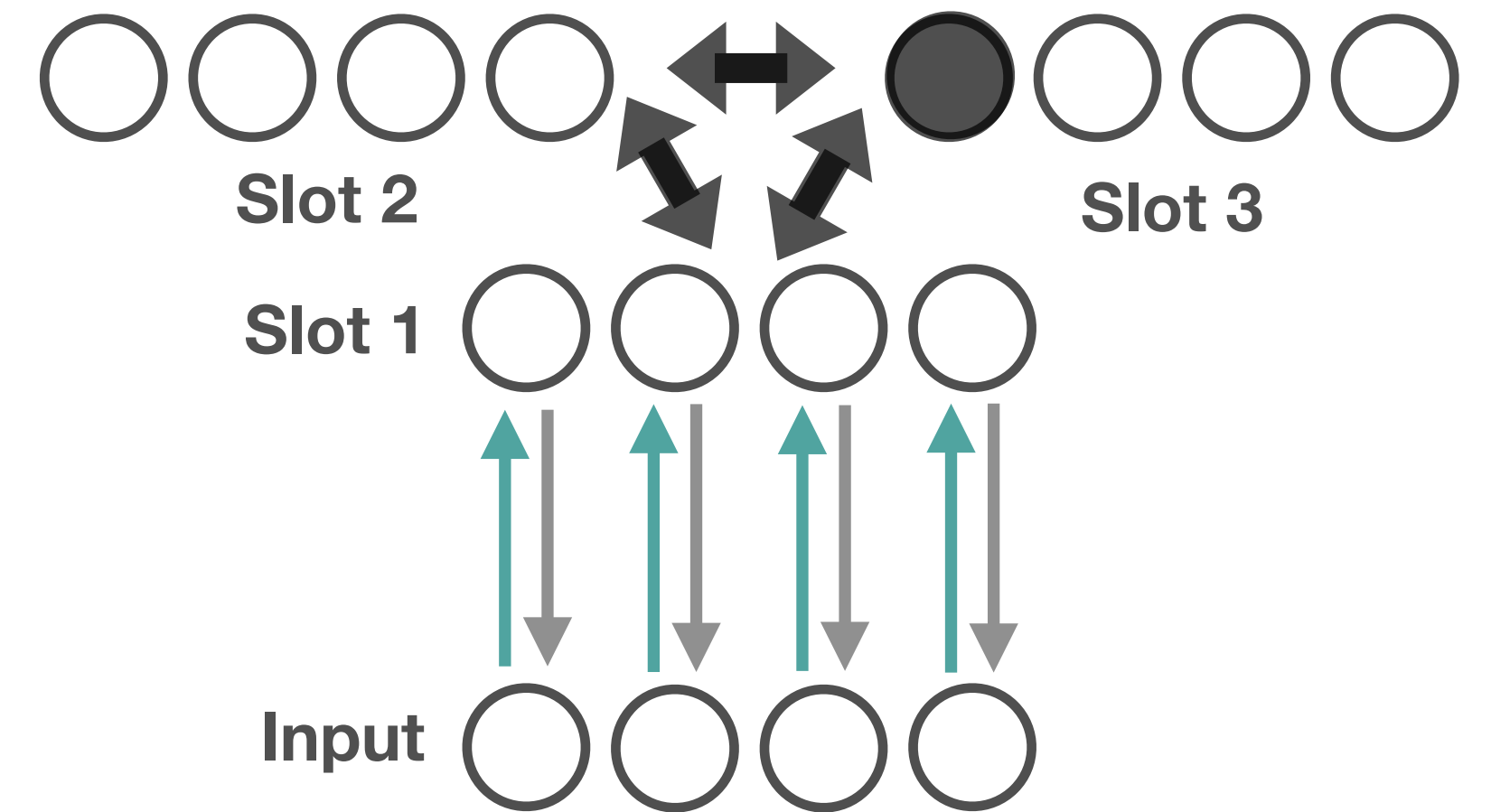
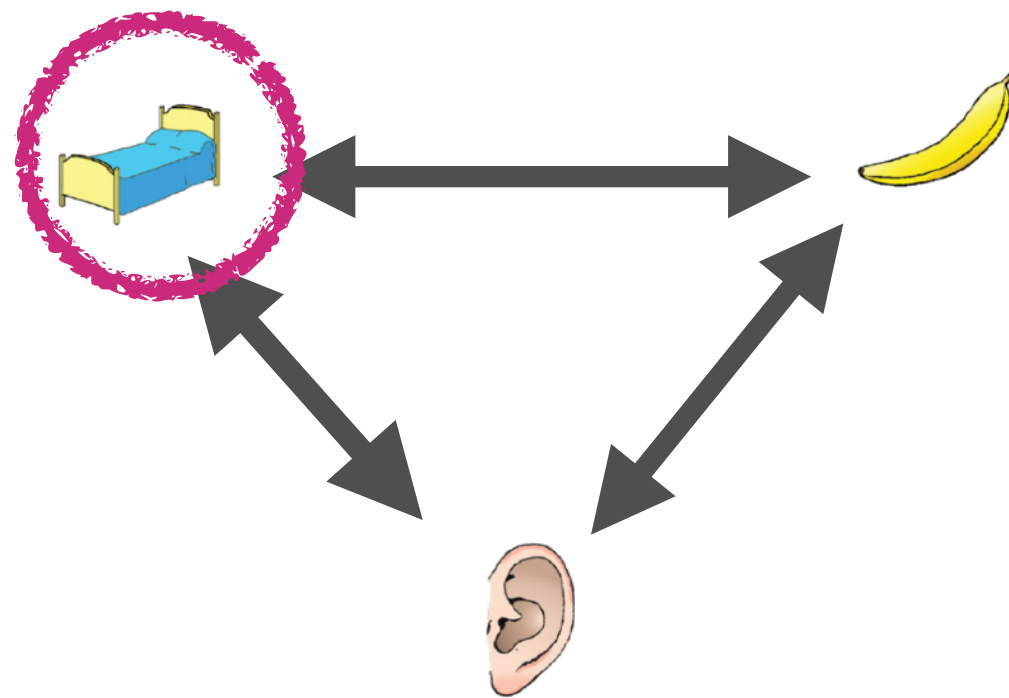
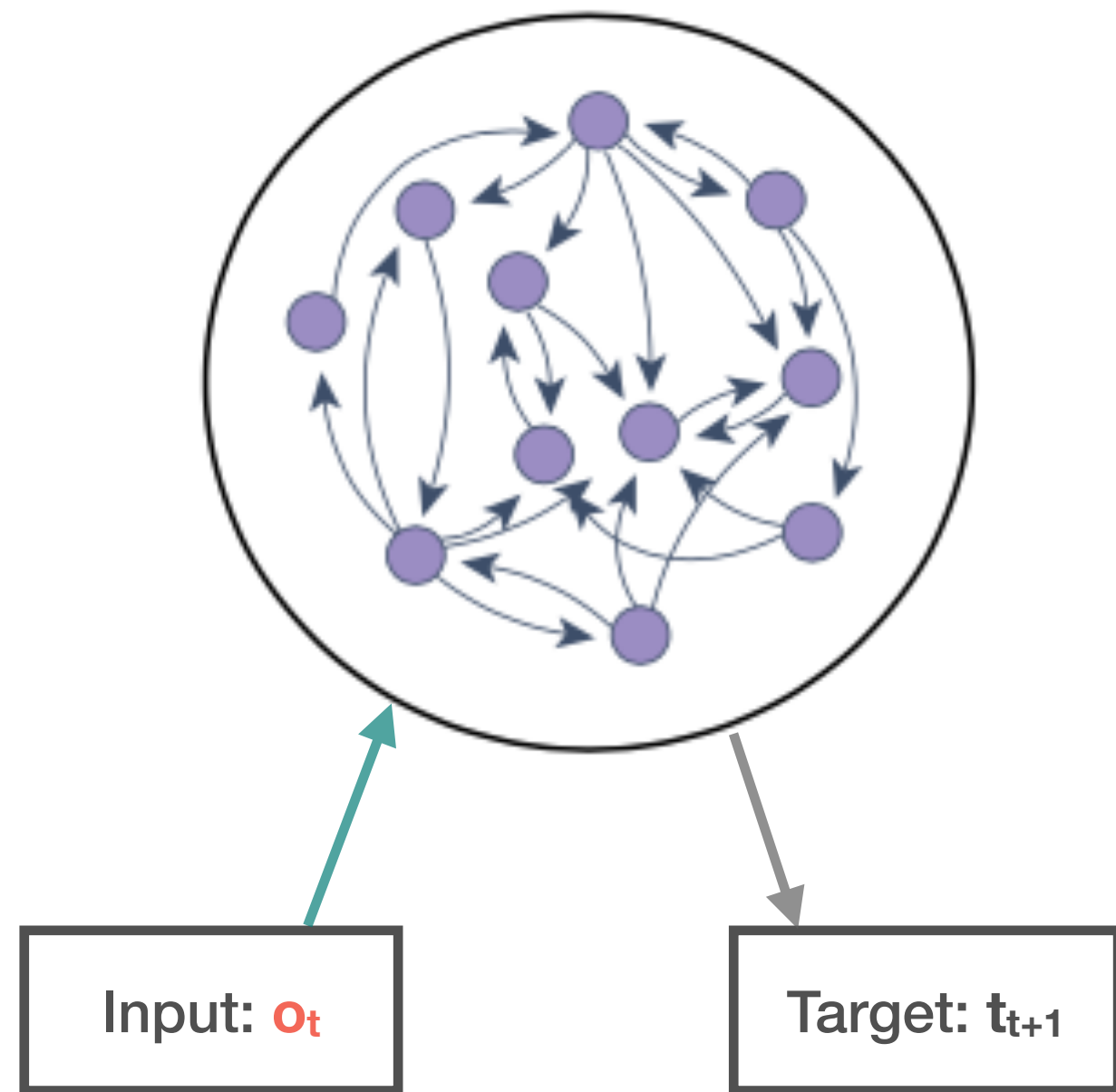


# Hypothesis: Sequence working memory using neural circuits of structured slots

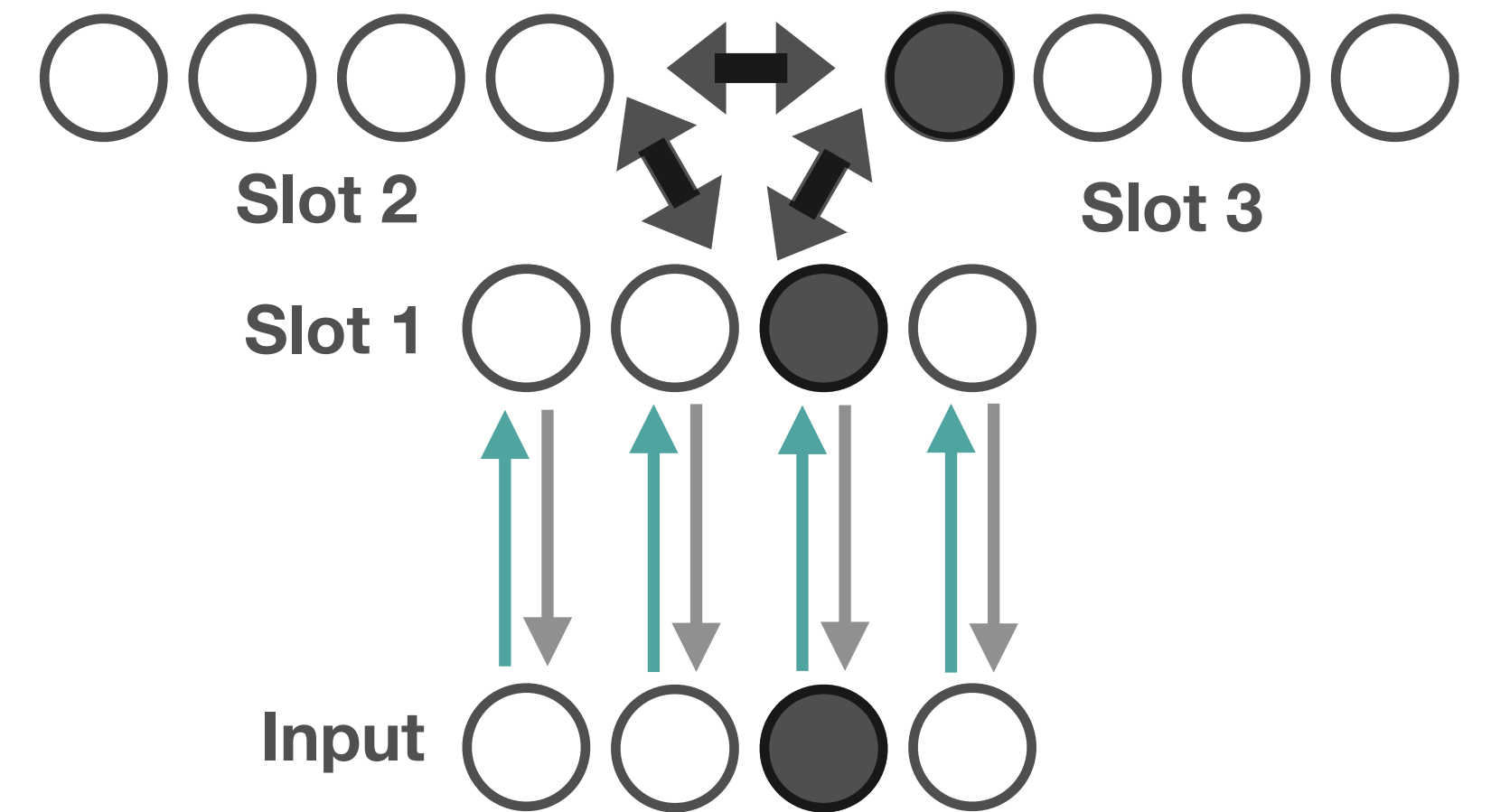
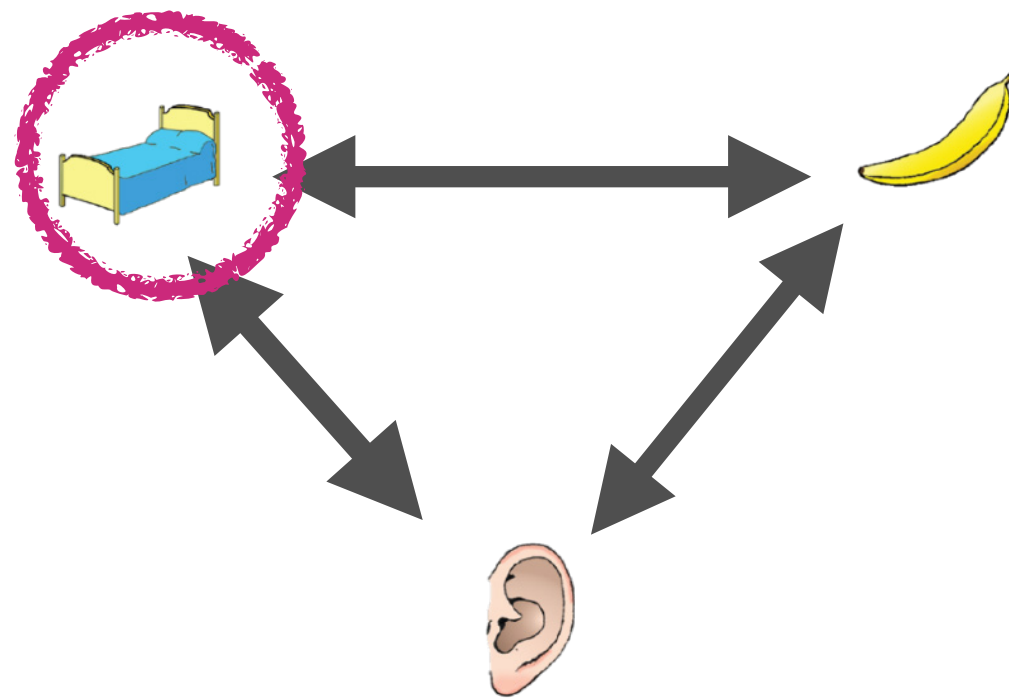
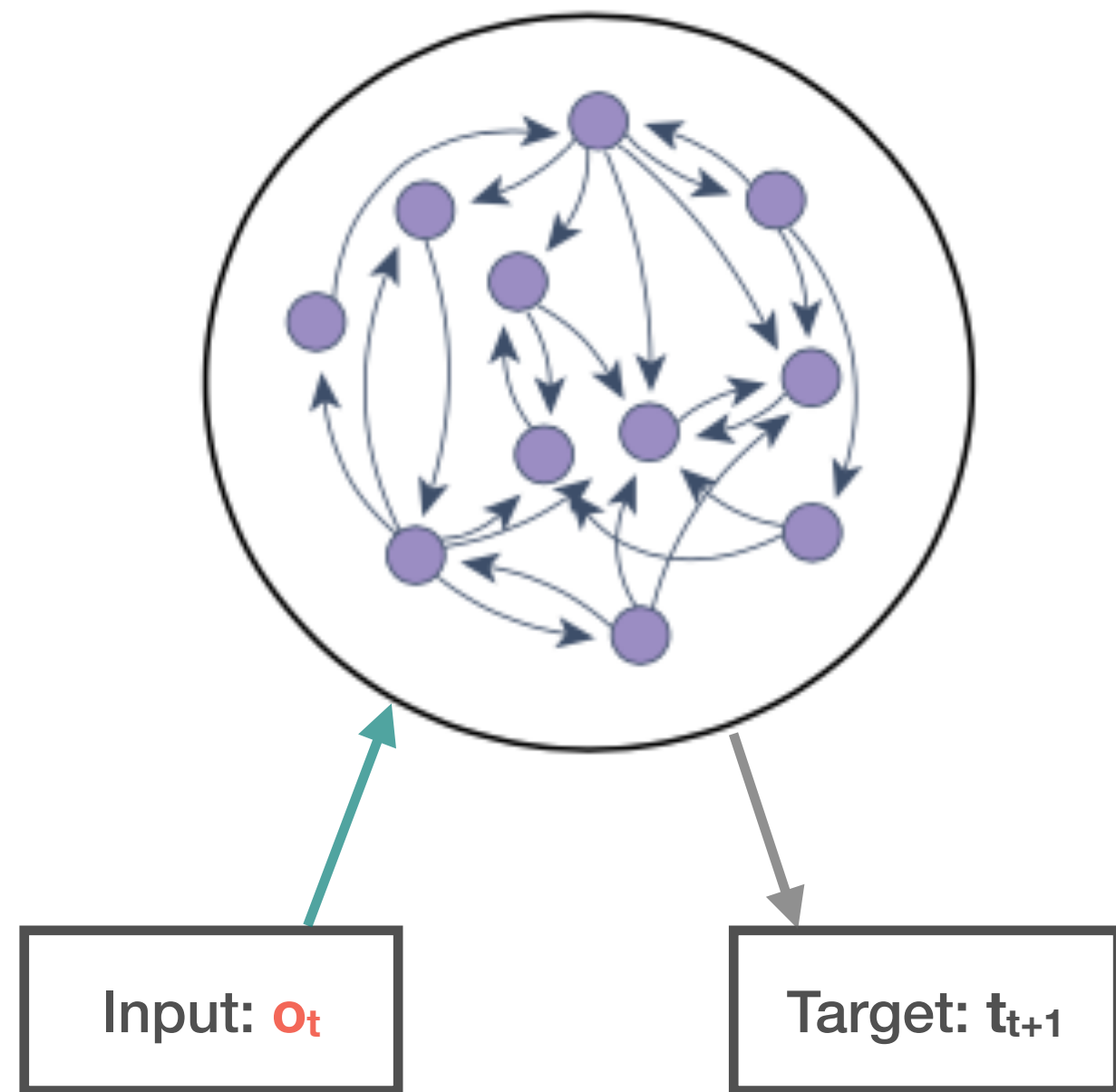




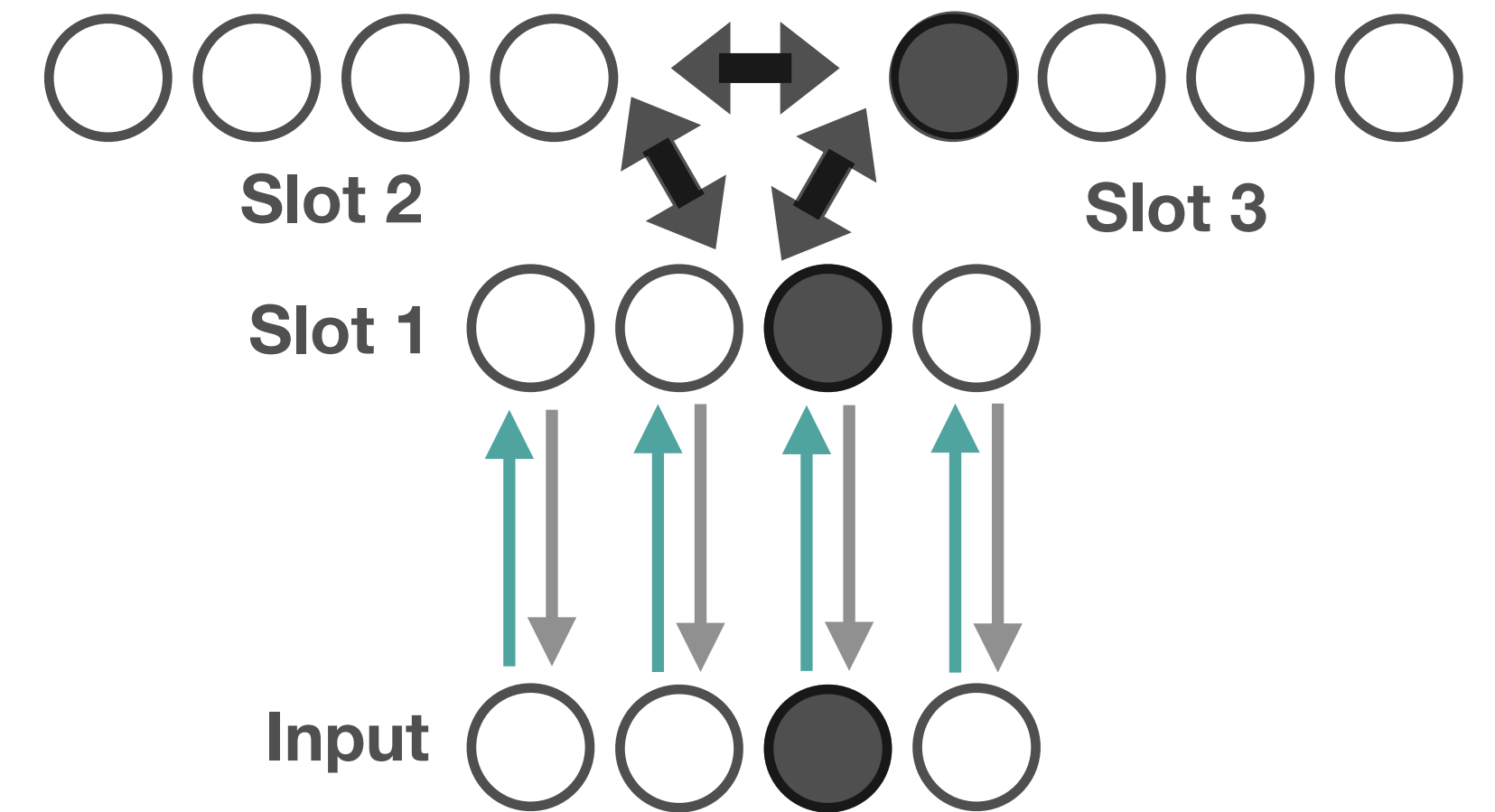
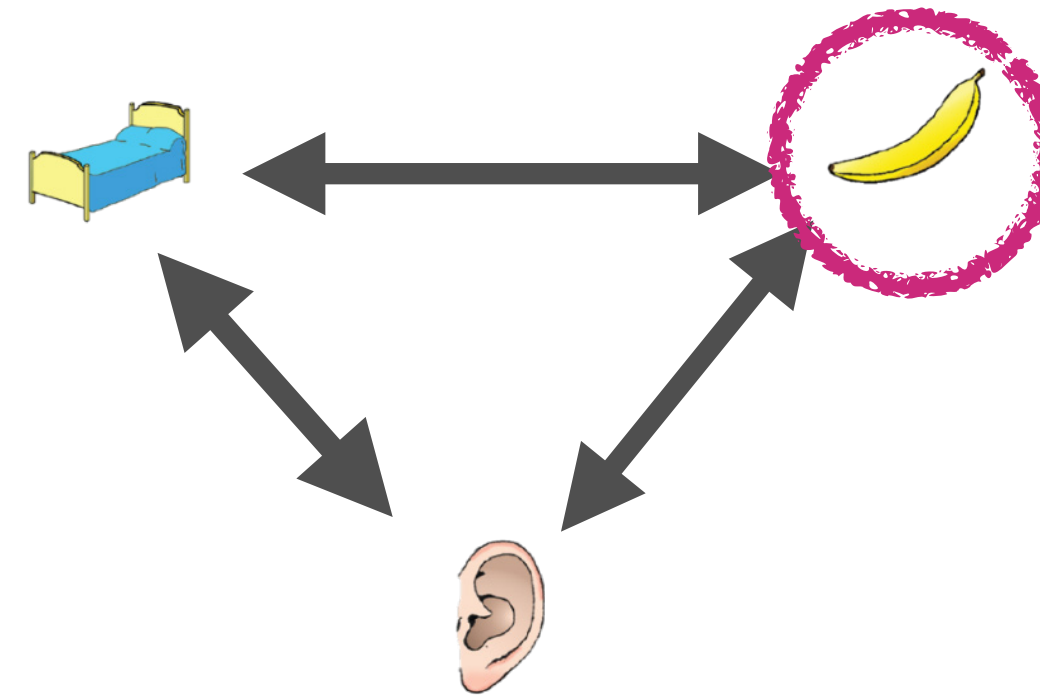
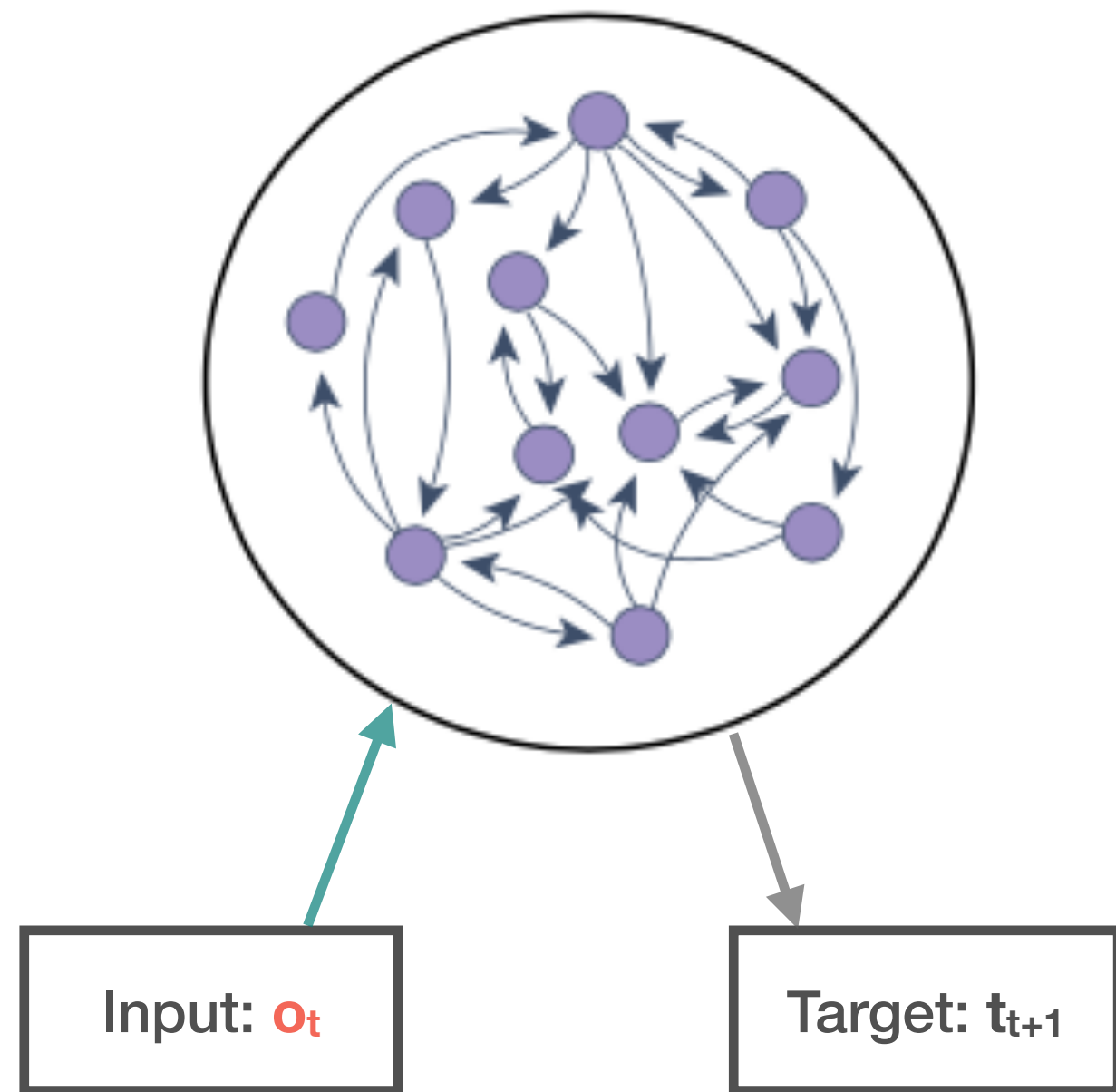
# Hypothesis: Sequence working memory using neural circuits of structured slots



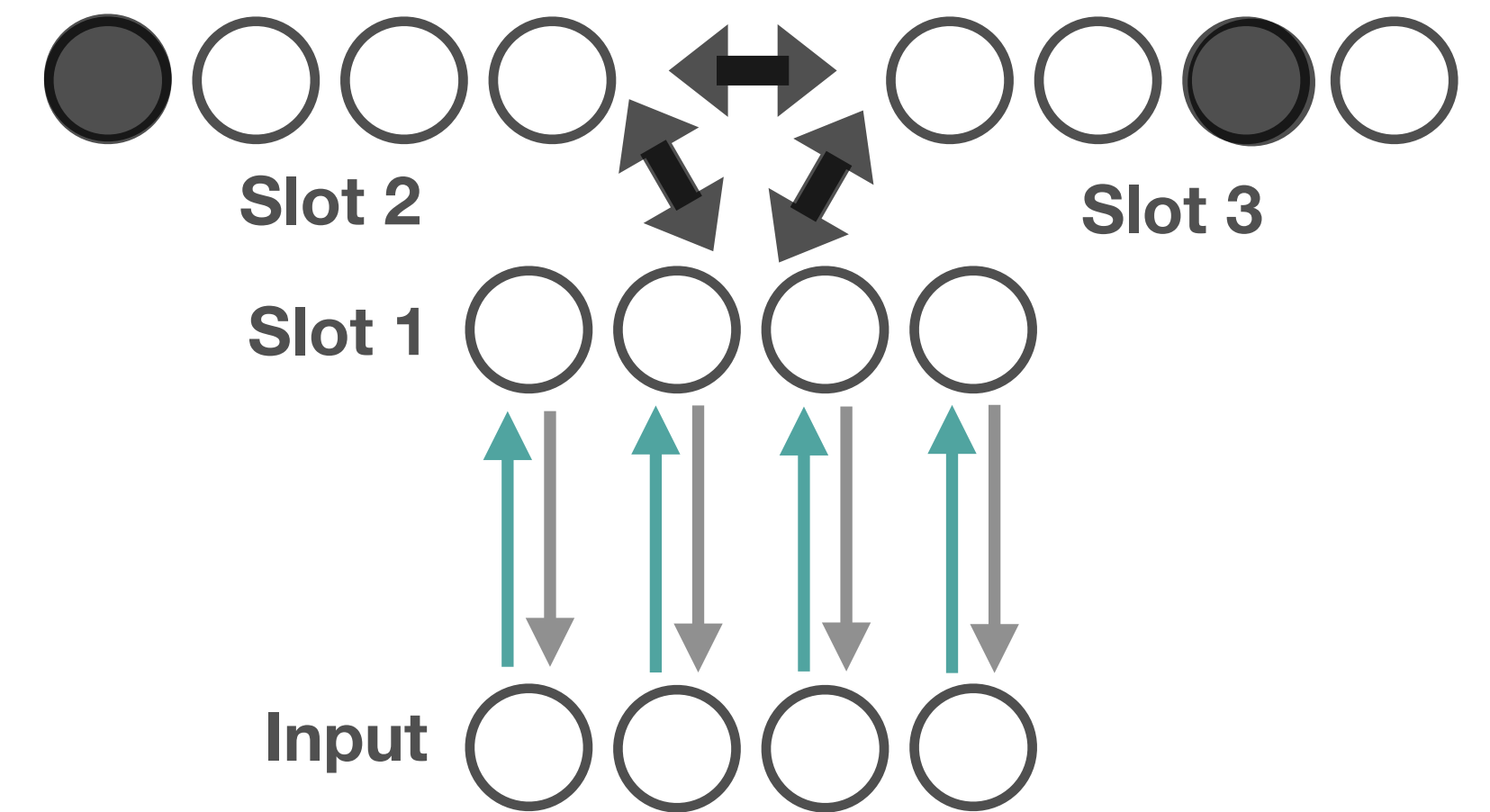
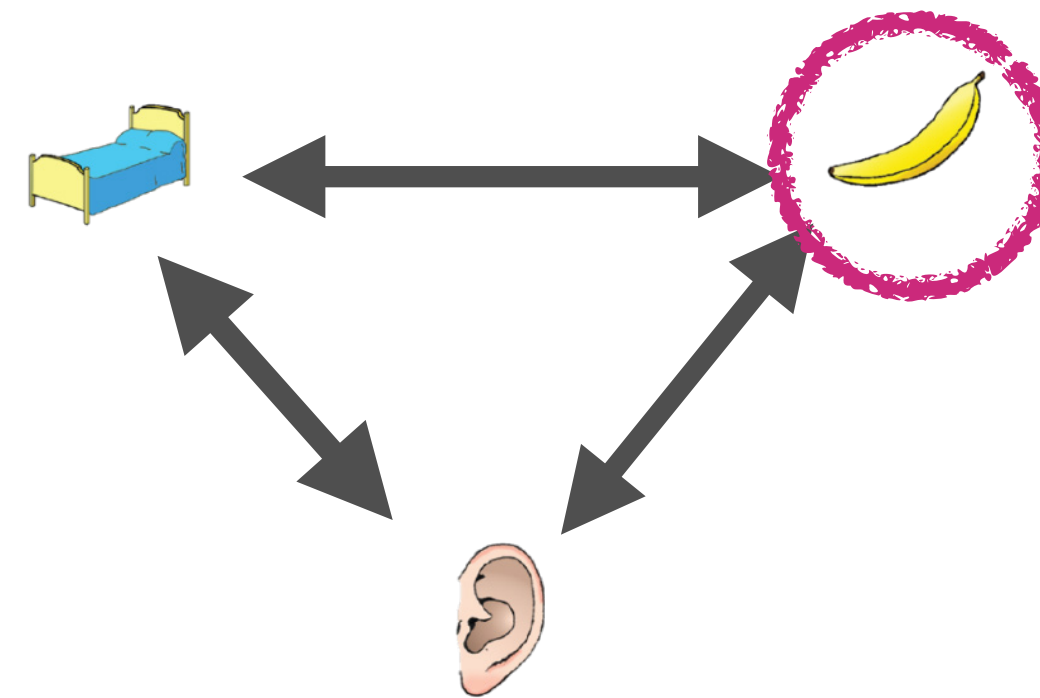
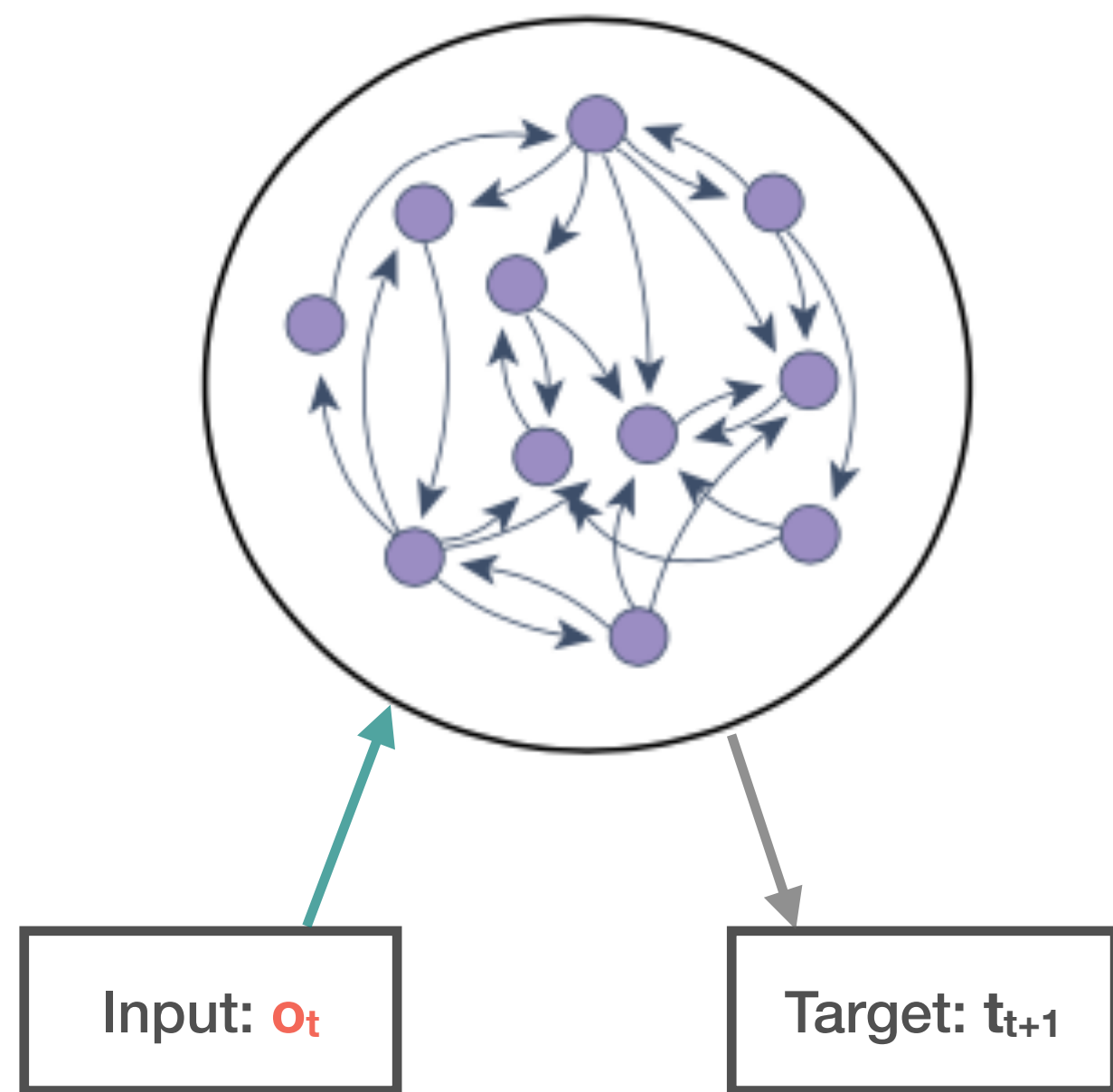
# Hypothesis: Sequence working memory using neural circuits of structured slots



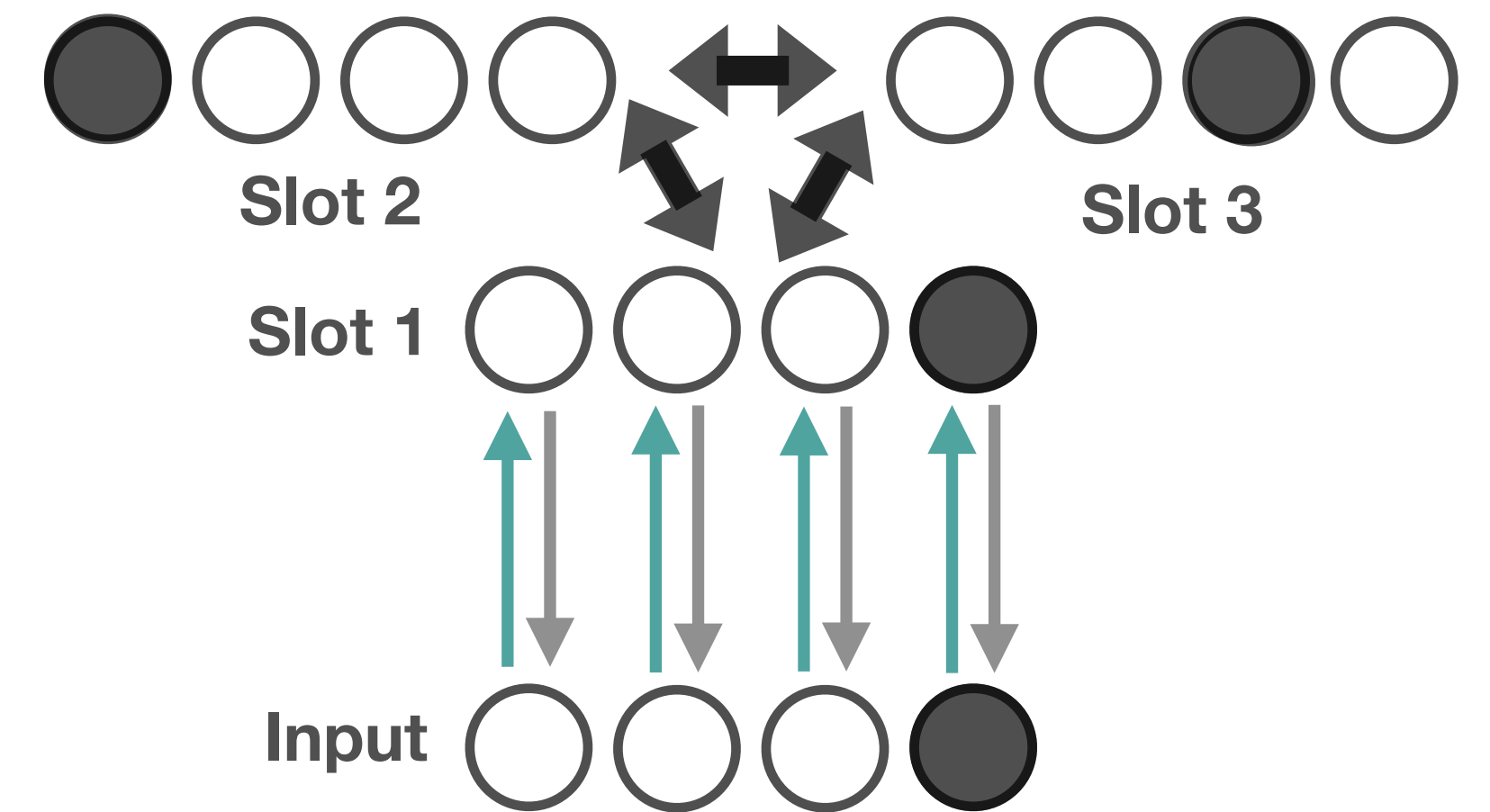
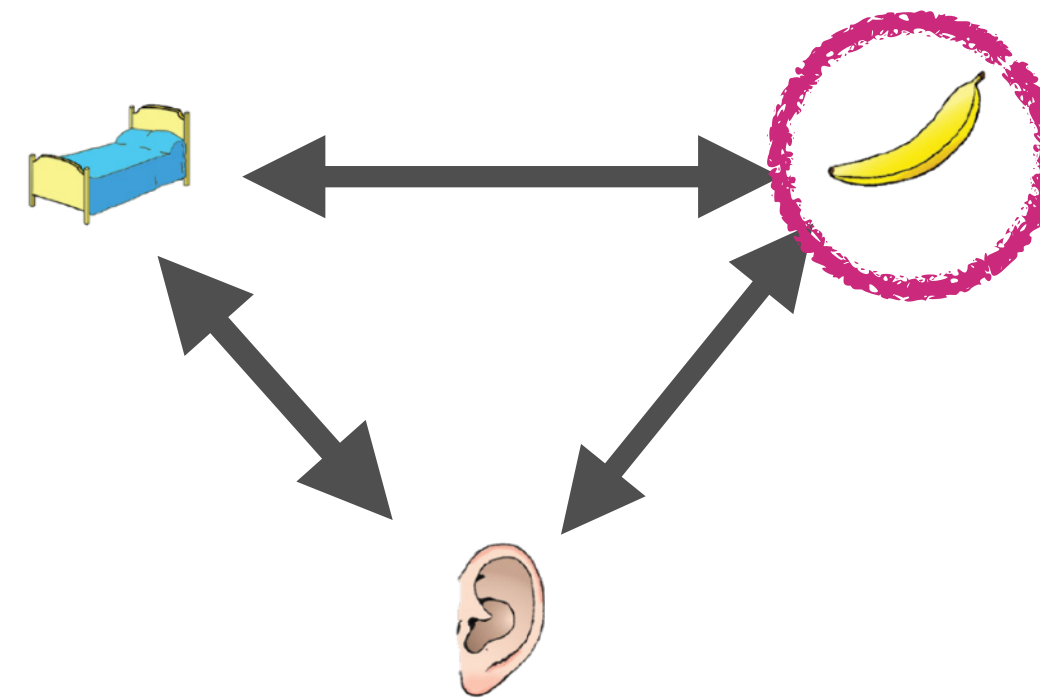
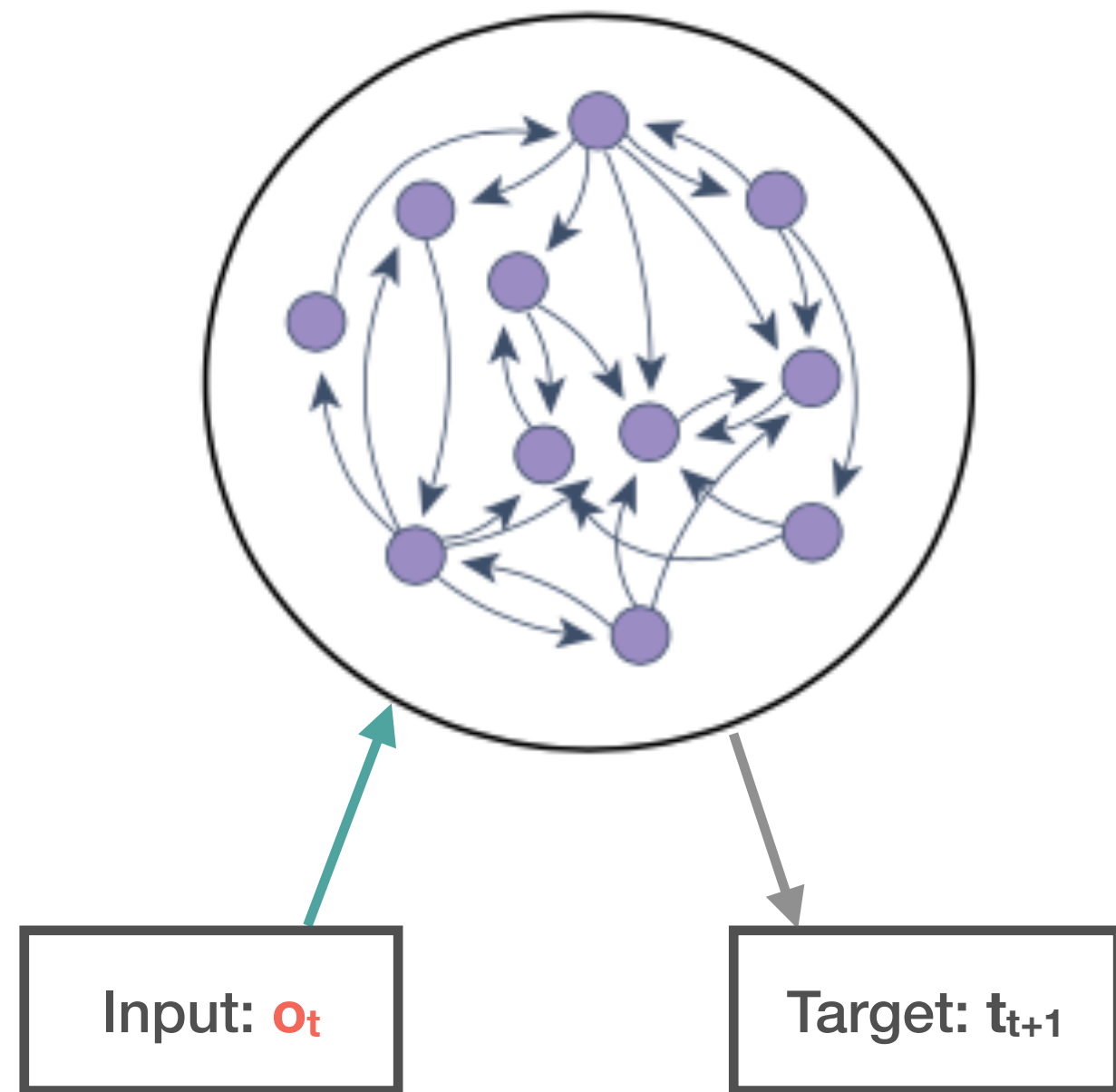
# Hypothesis: Sequence working memory using neural circuits of structured slots



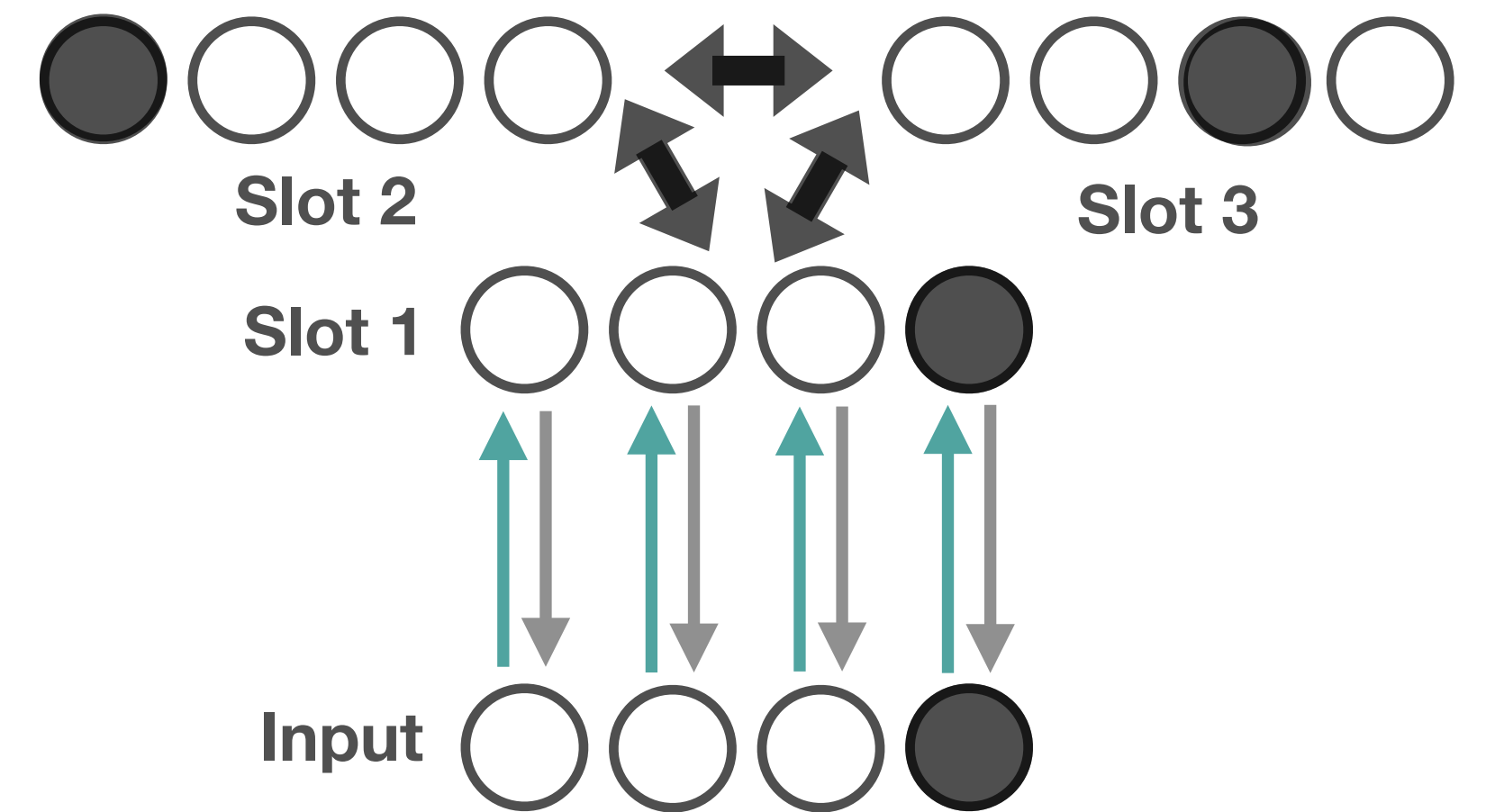
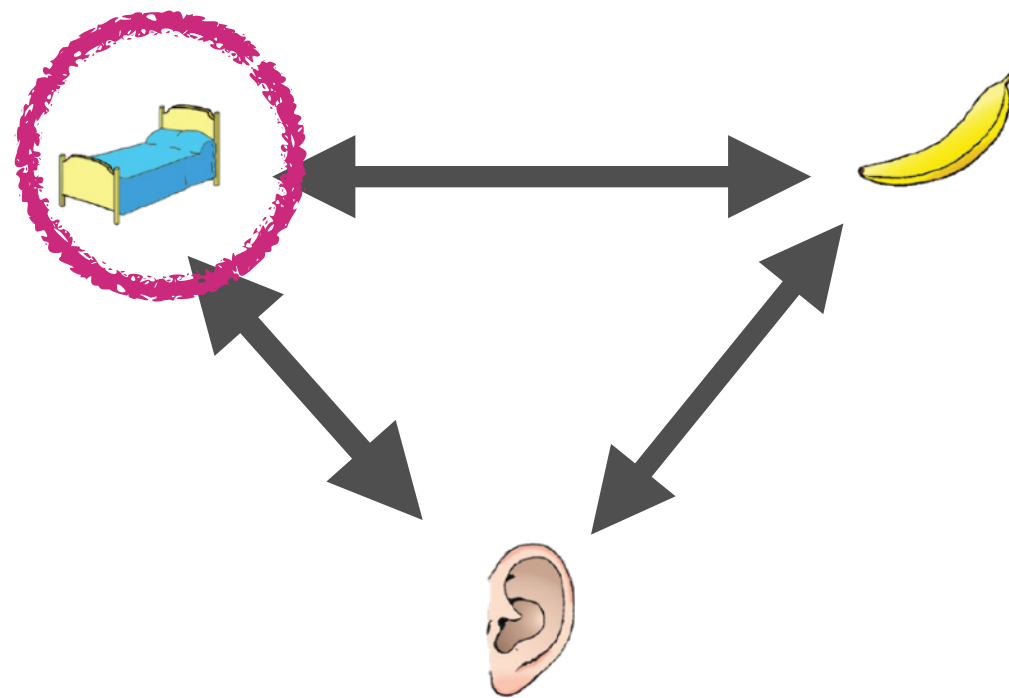
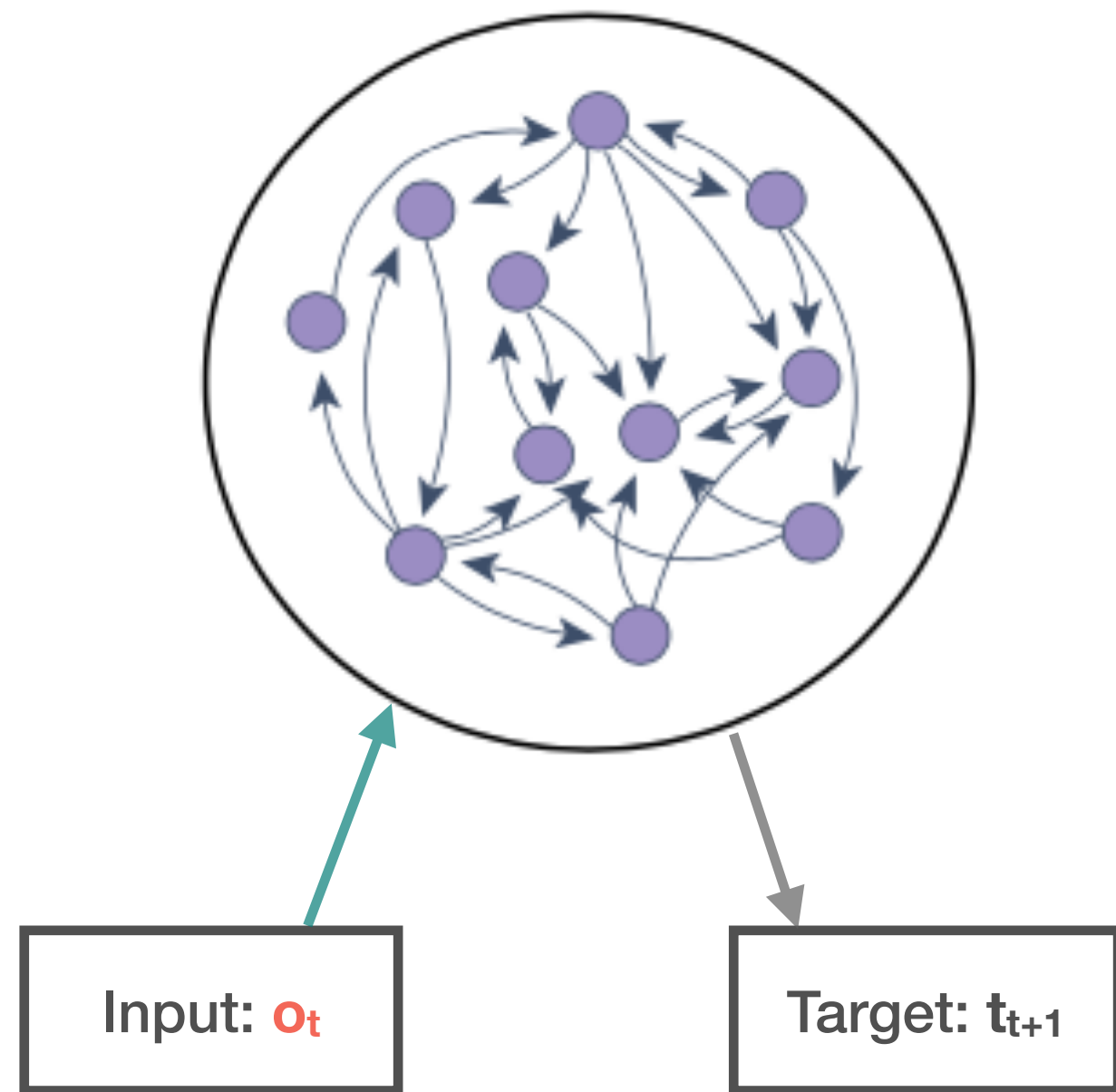
# Hypothesis: Sequence working memory using neural circuits of structured slots



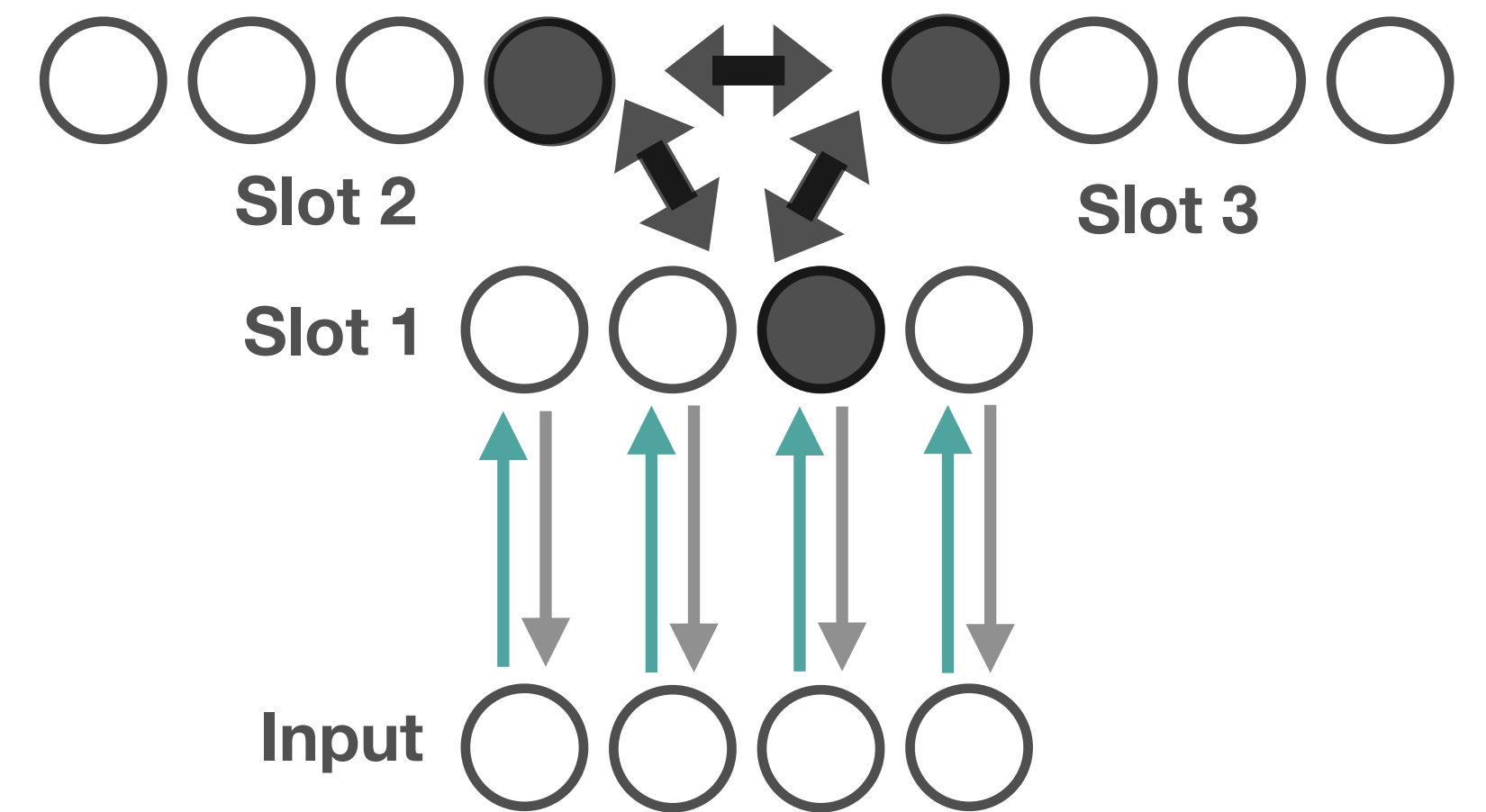
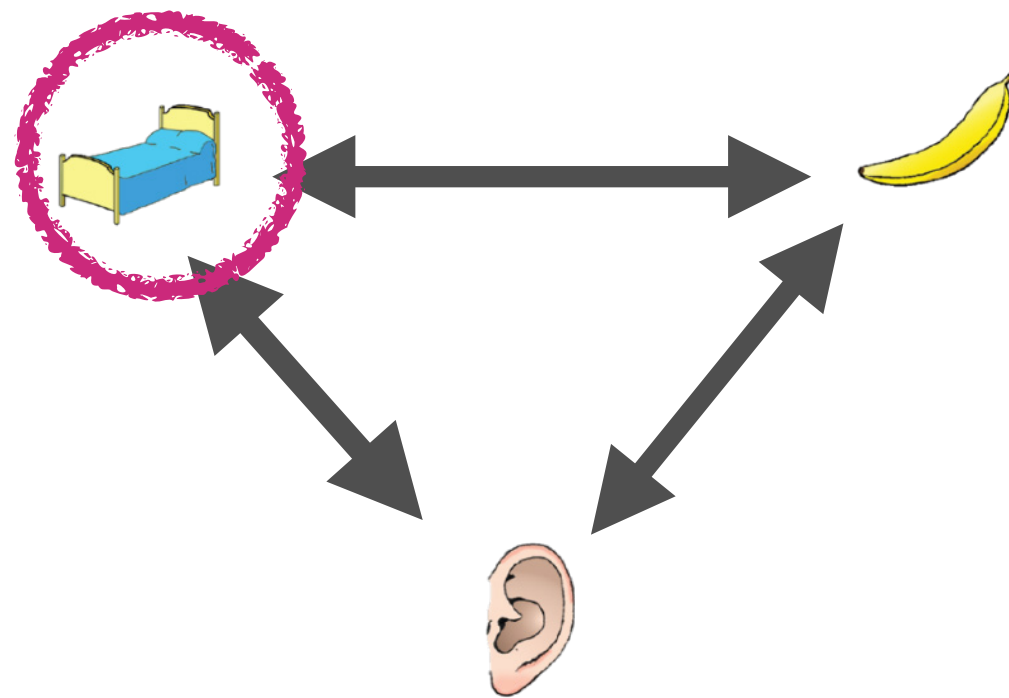
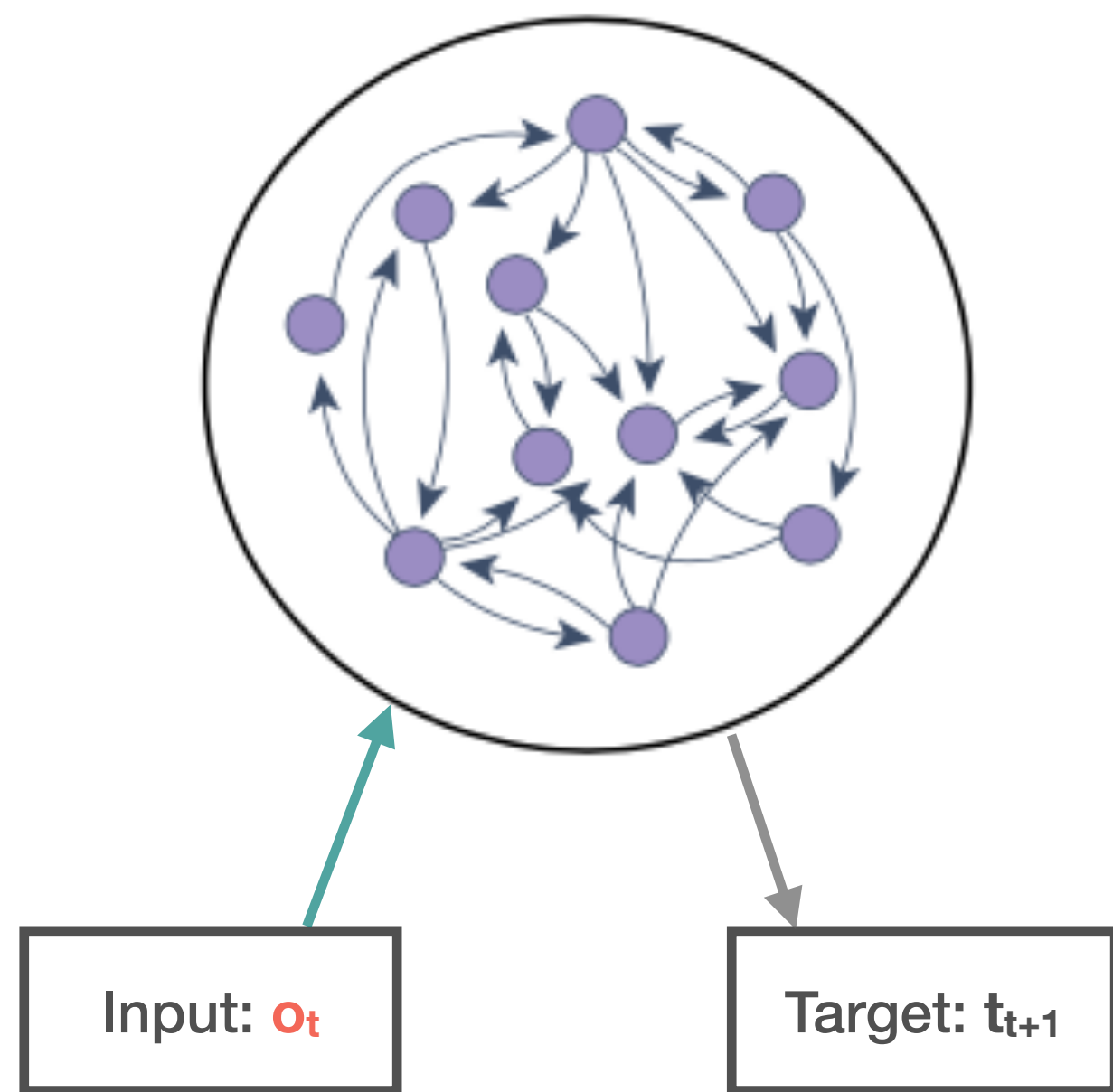
# Hypothesis: Sequence working memory using neural circuits of structured slots



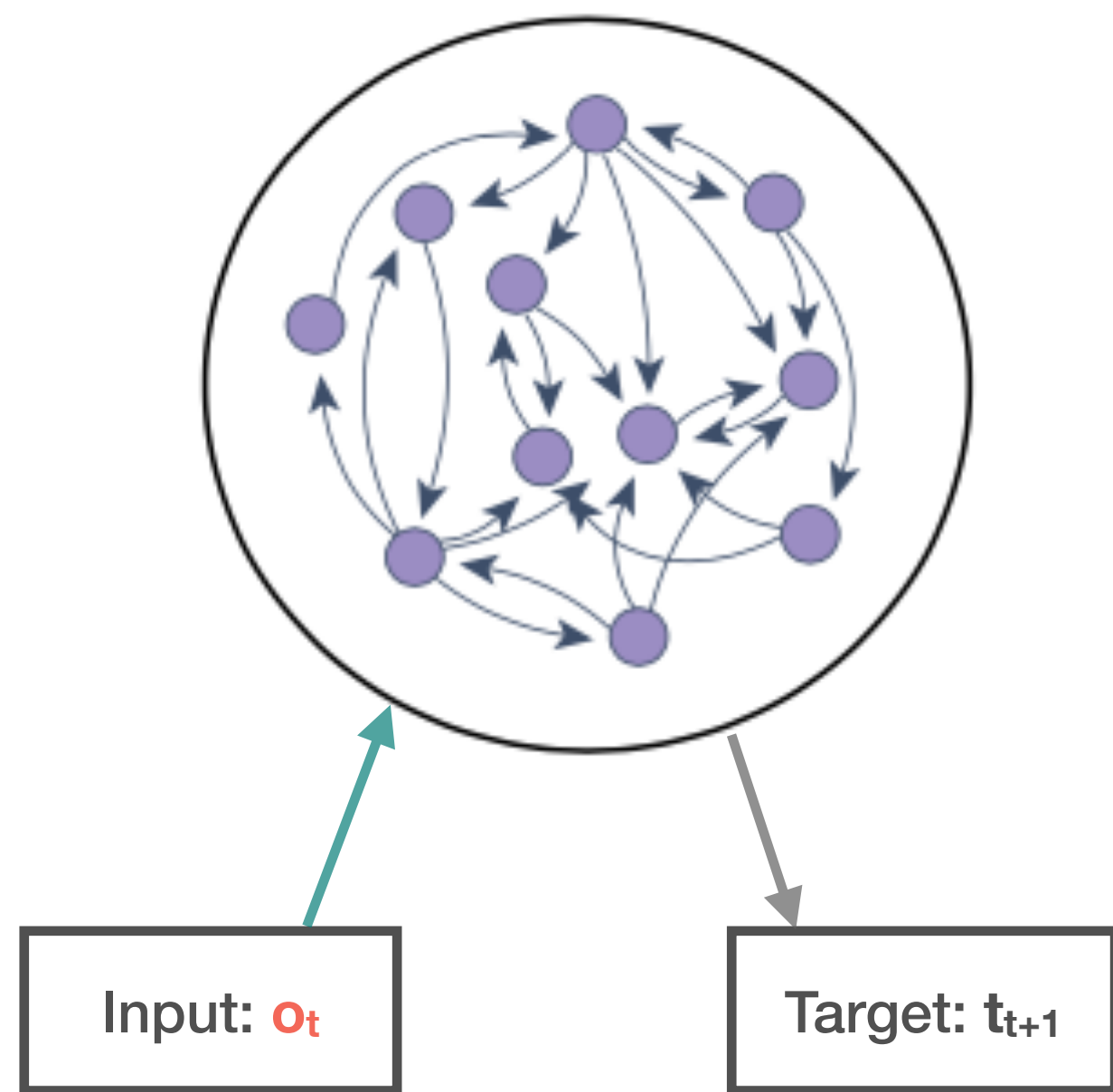
# Hypothesis: Sequence working memory using neural circuits of structured slots



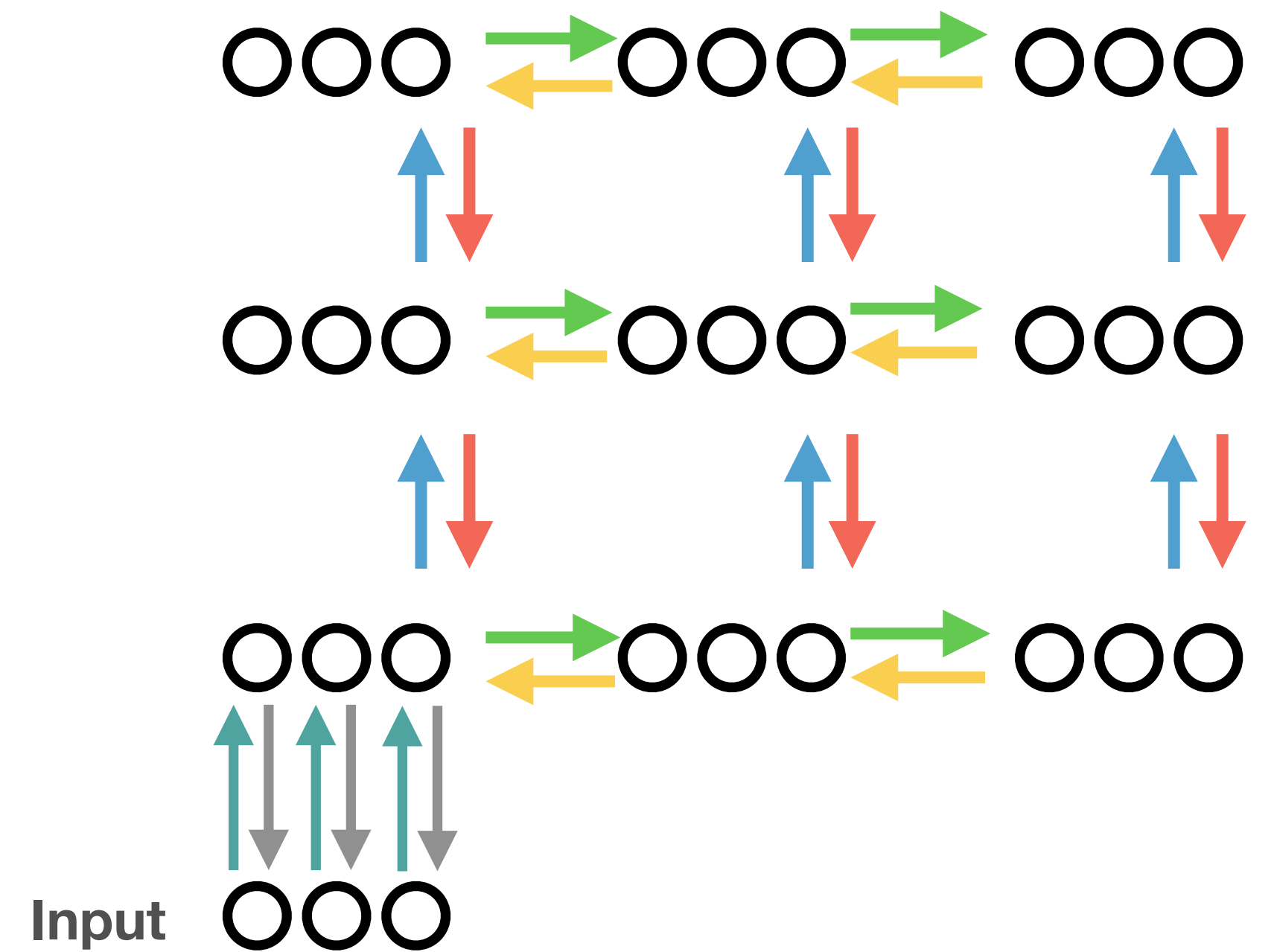
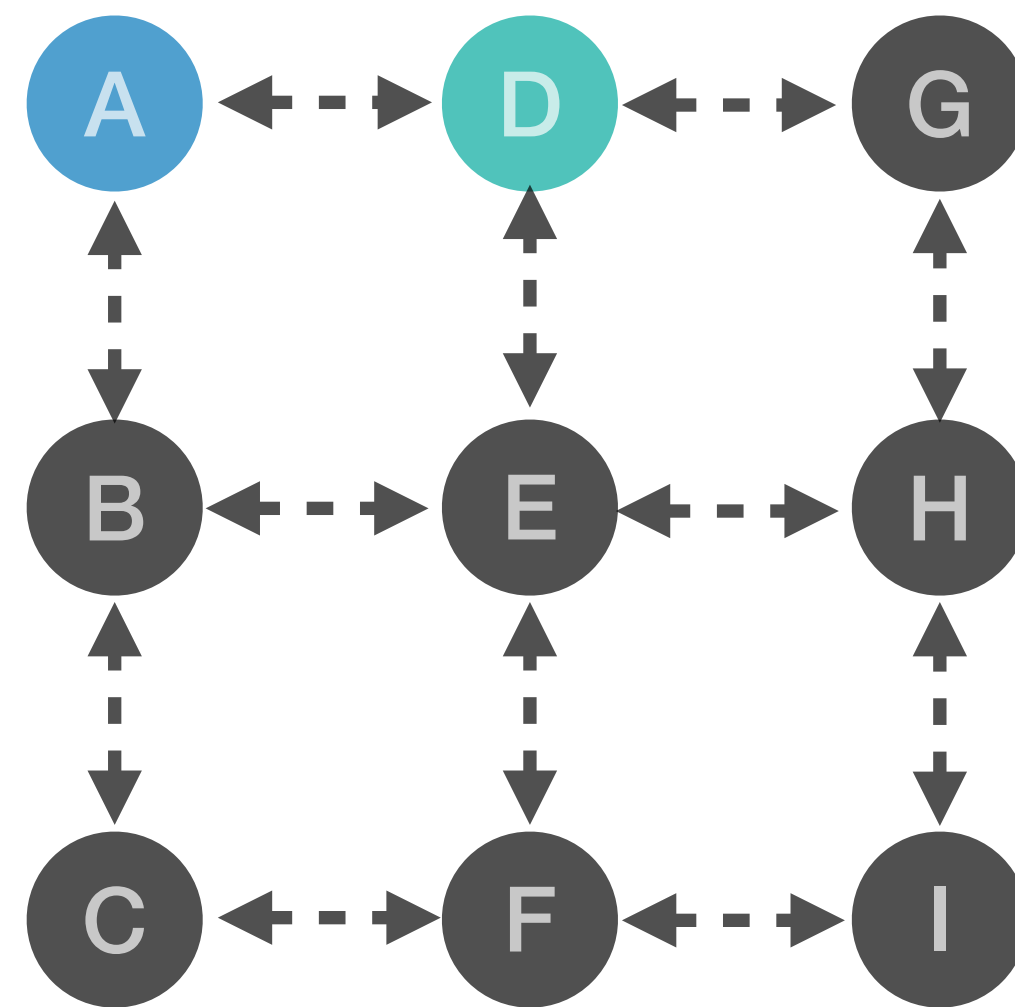
# Hypothesis: Sequence working memory using neural circuits of structured slots



# Different problems have different slot structure

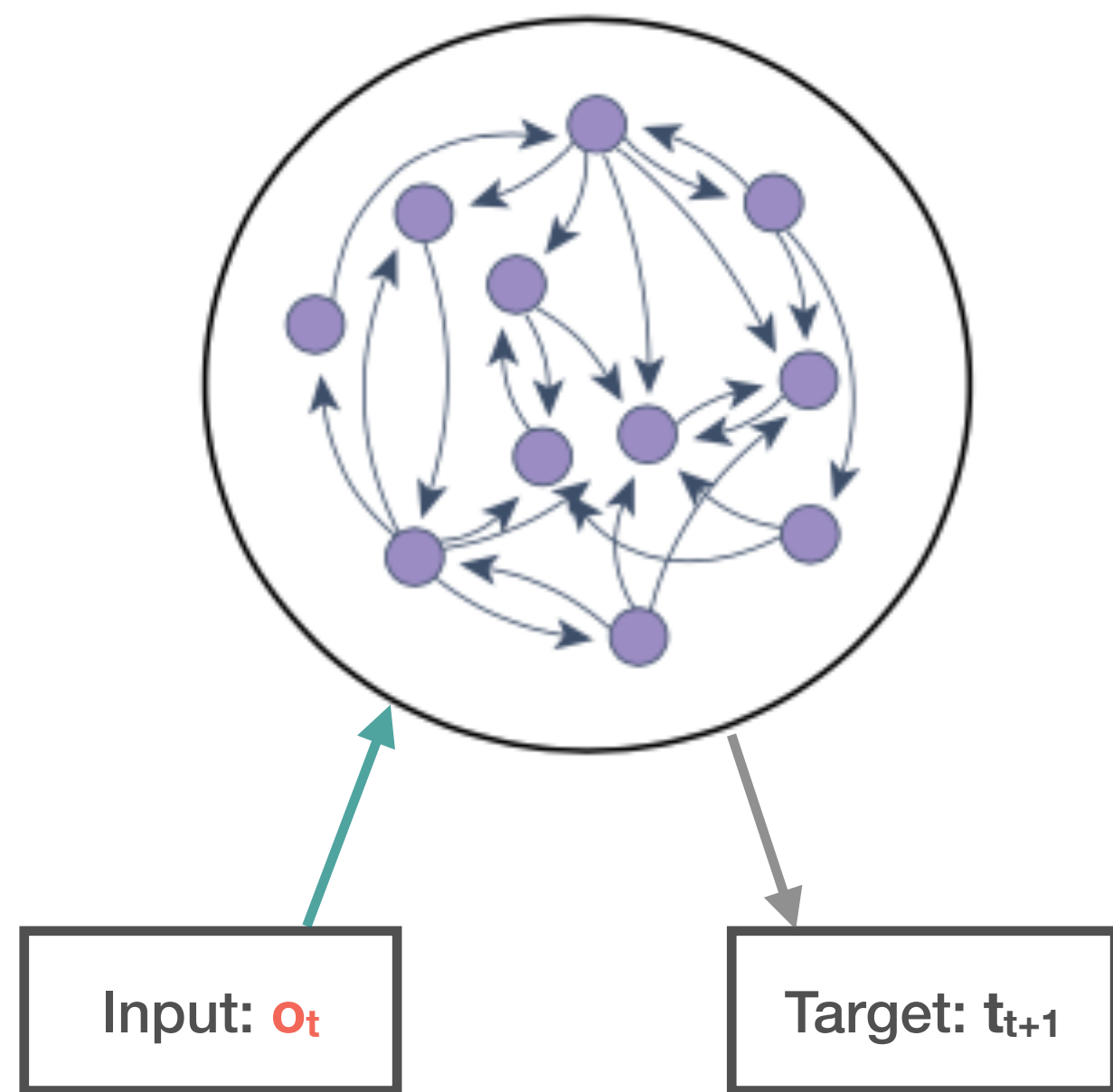


2D Navigation

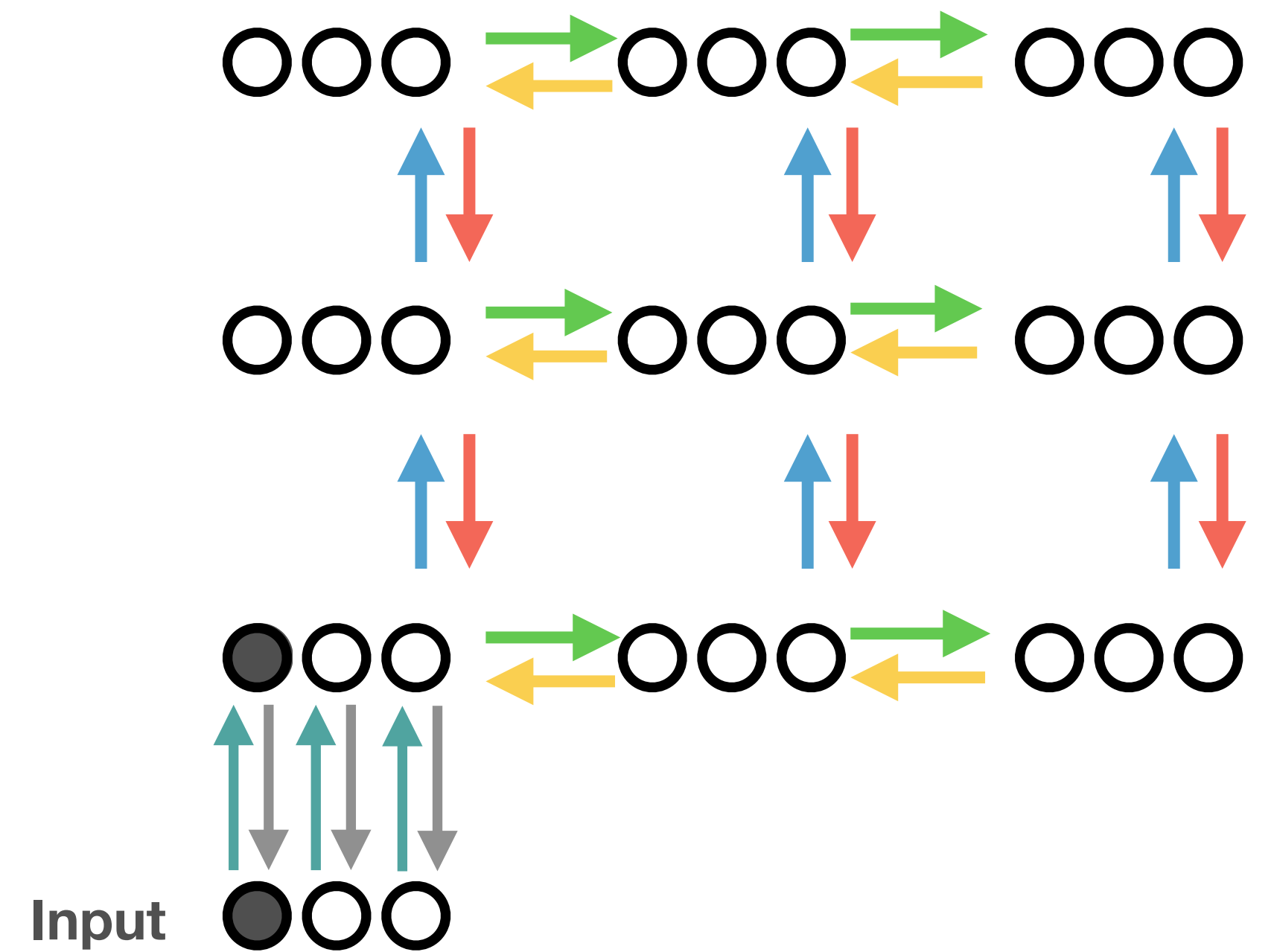
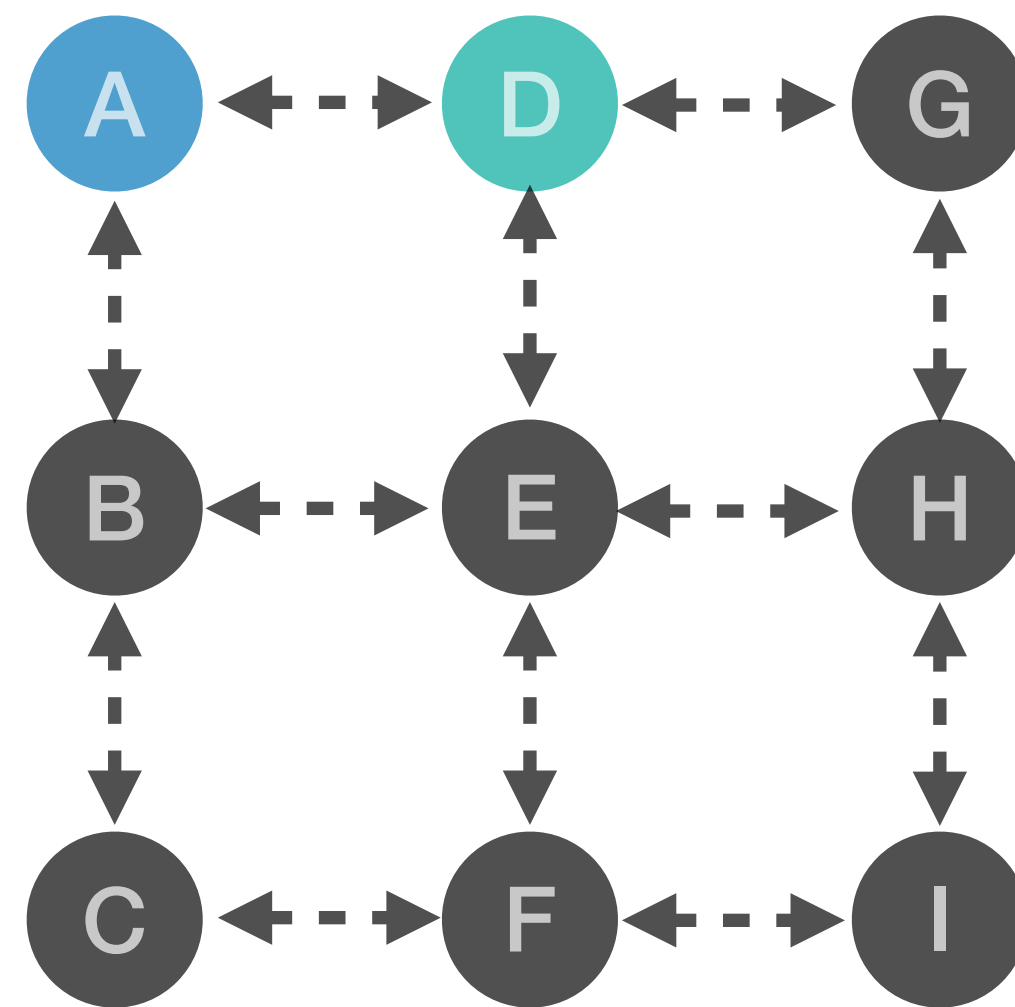




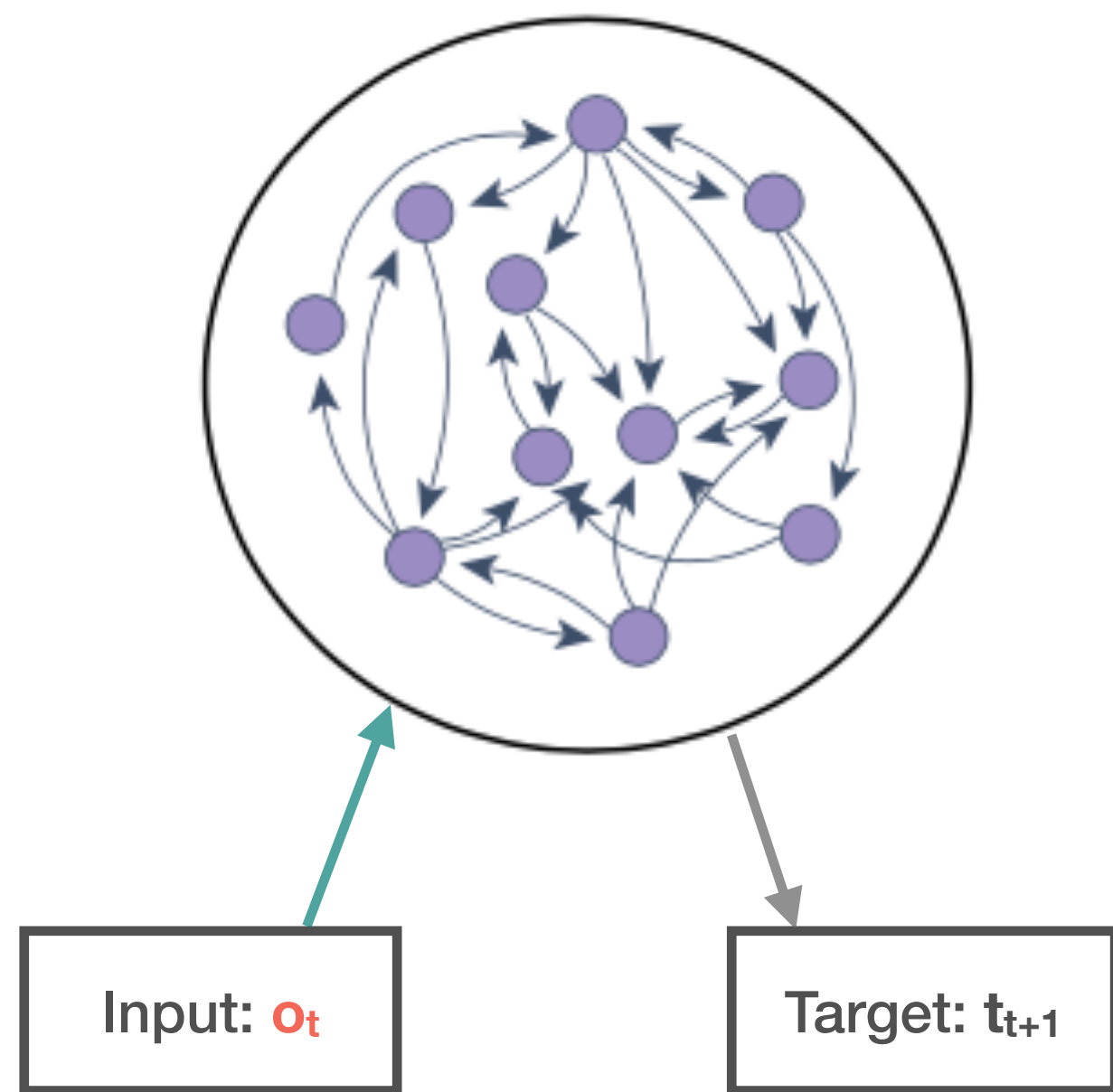
# Different problems have different slot structure



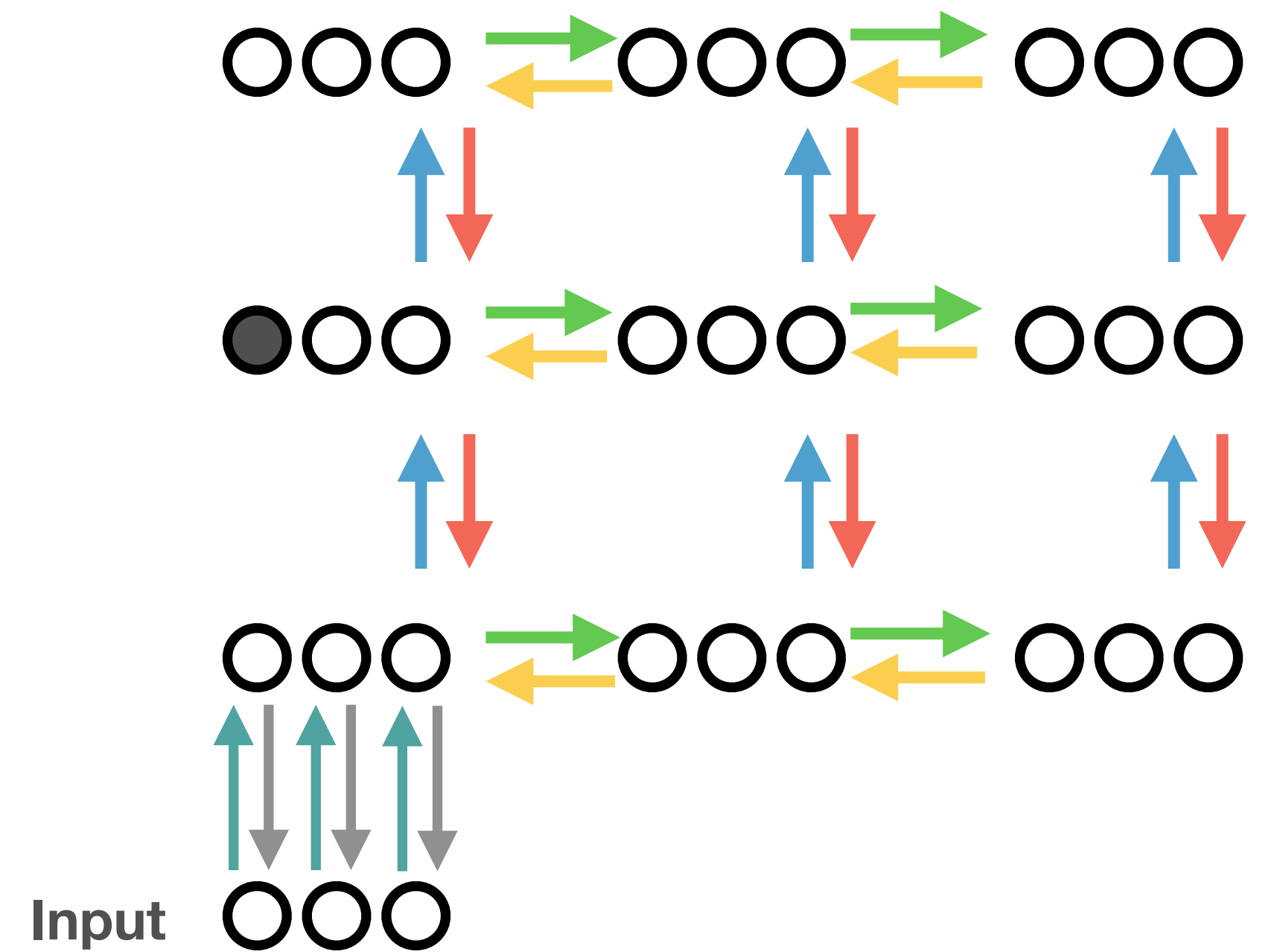
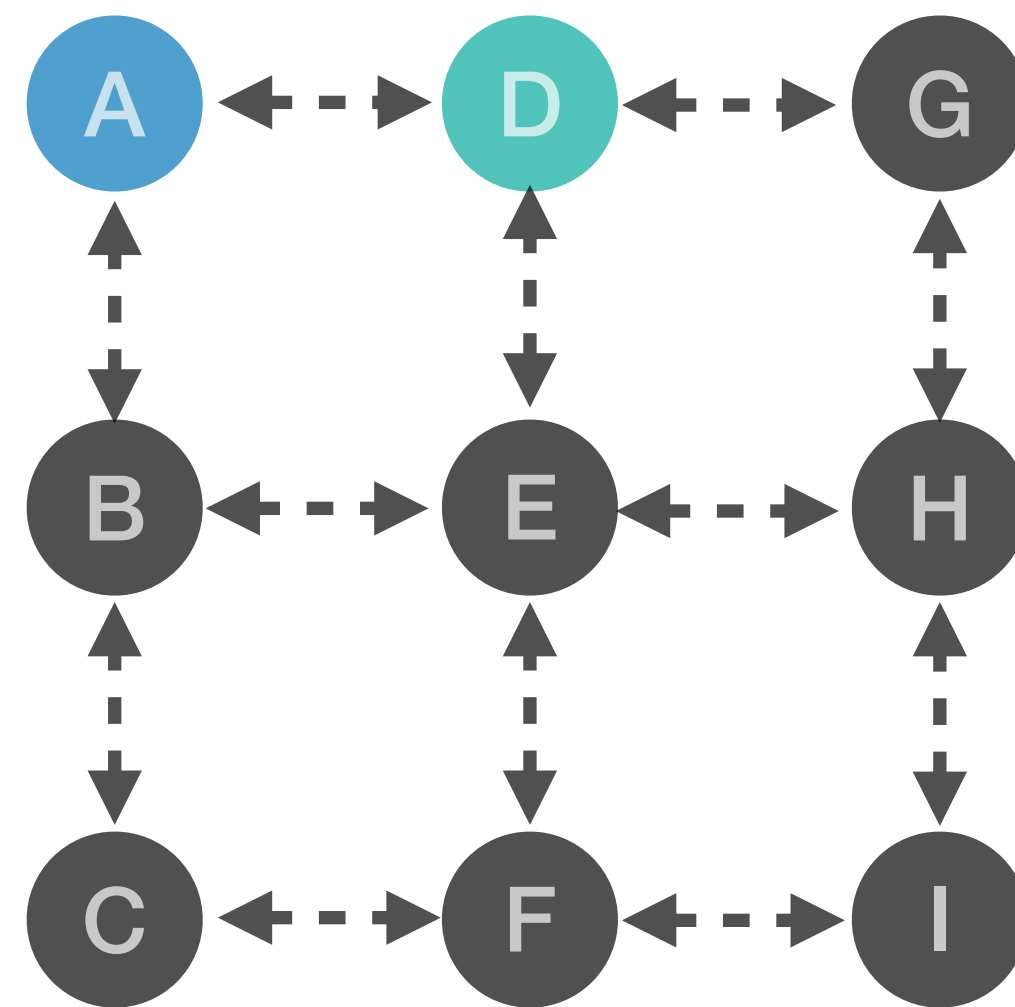
2D Navigation



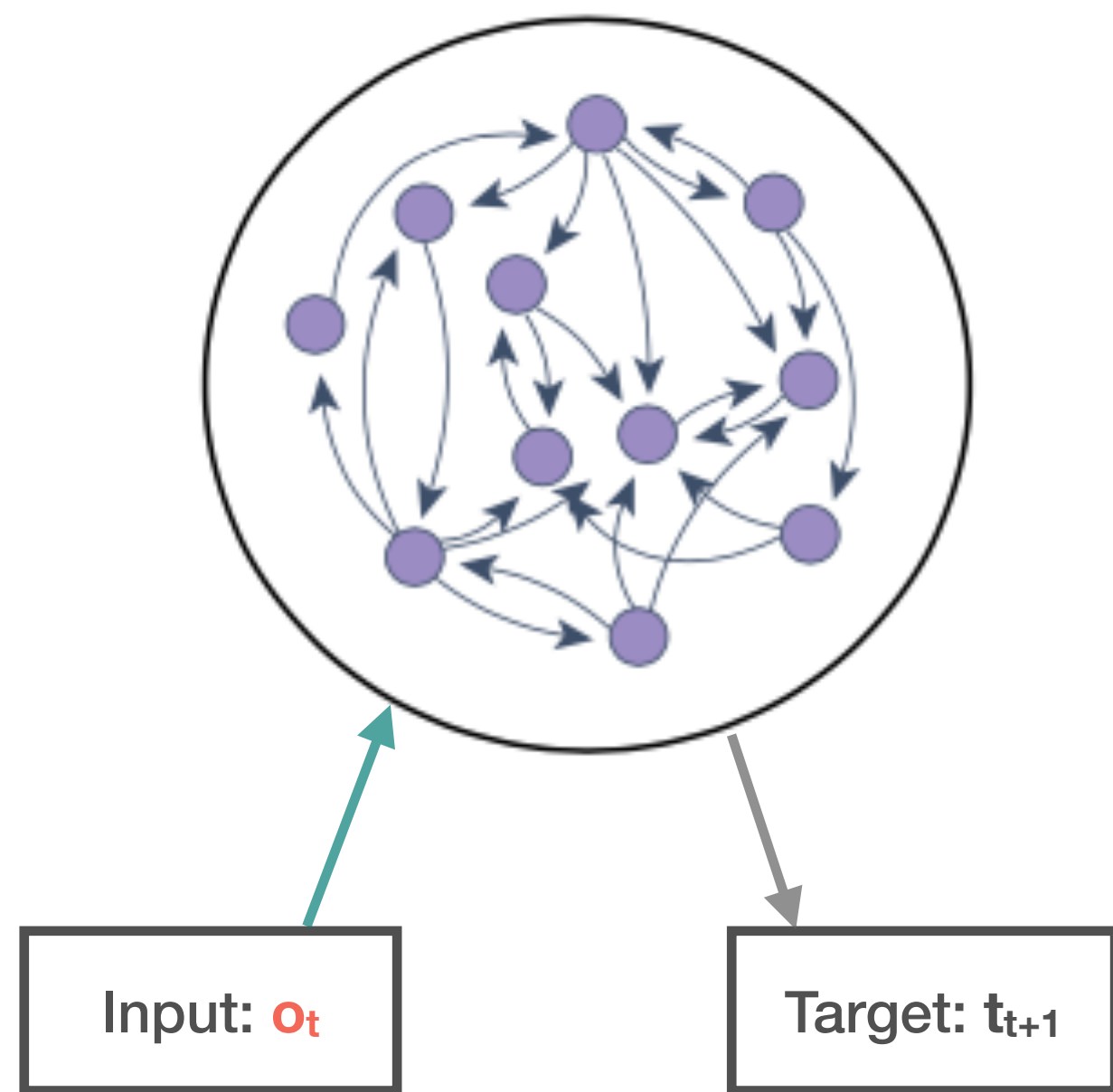
# Different problems have different slot structure



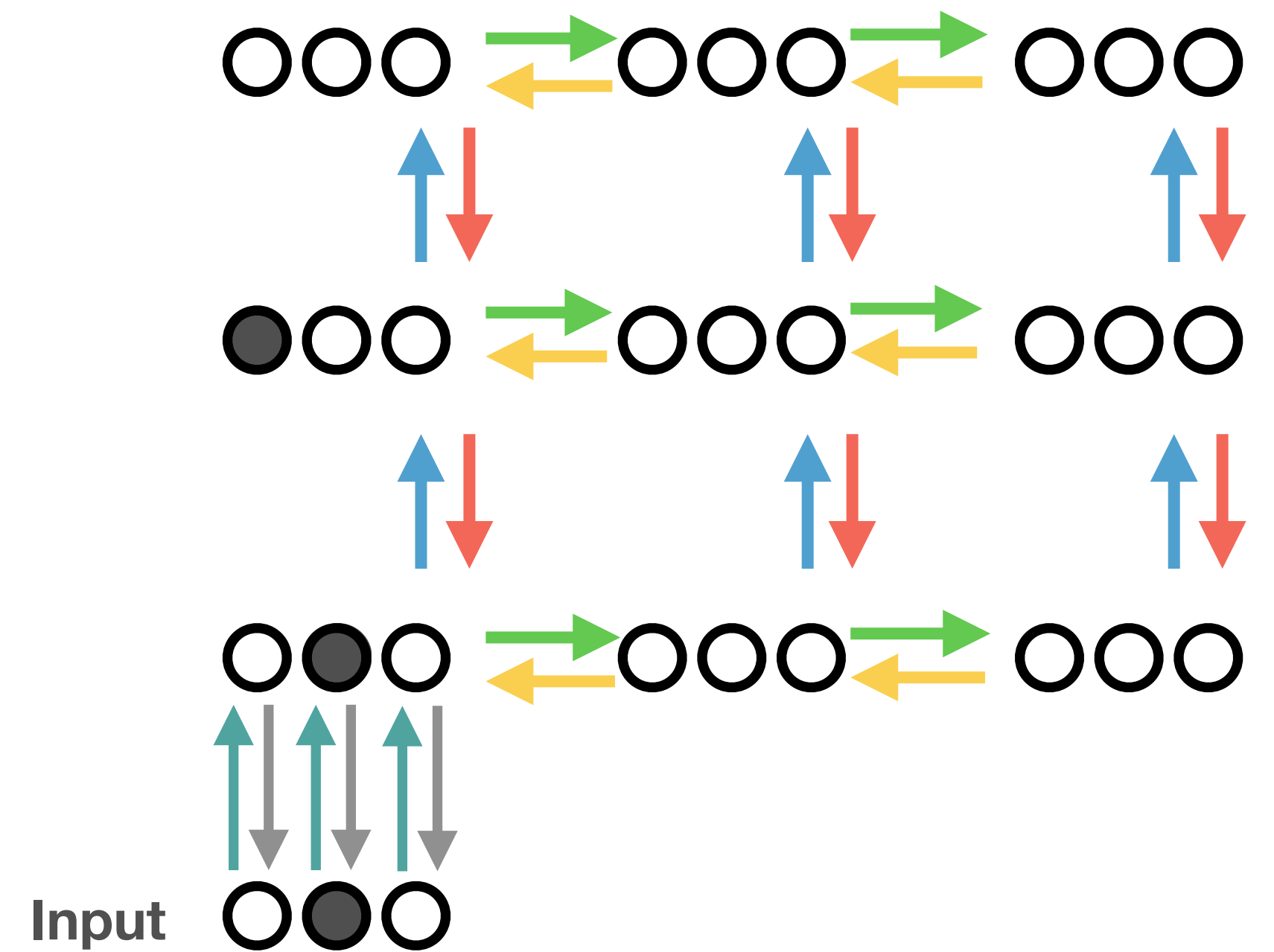
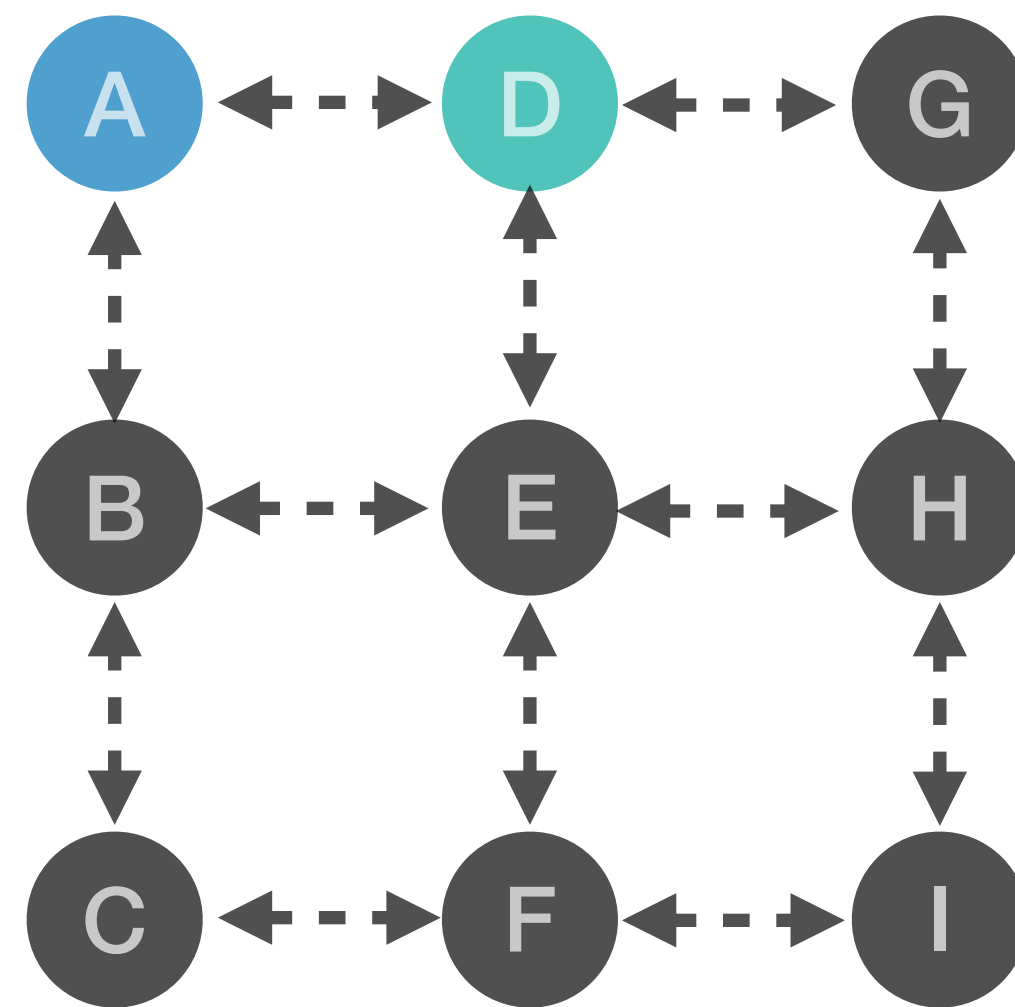
2D Navigation



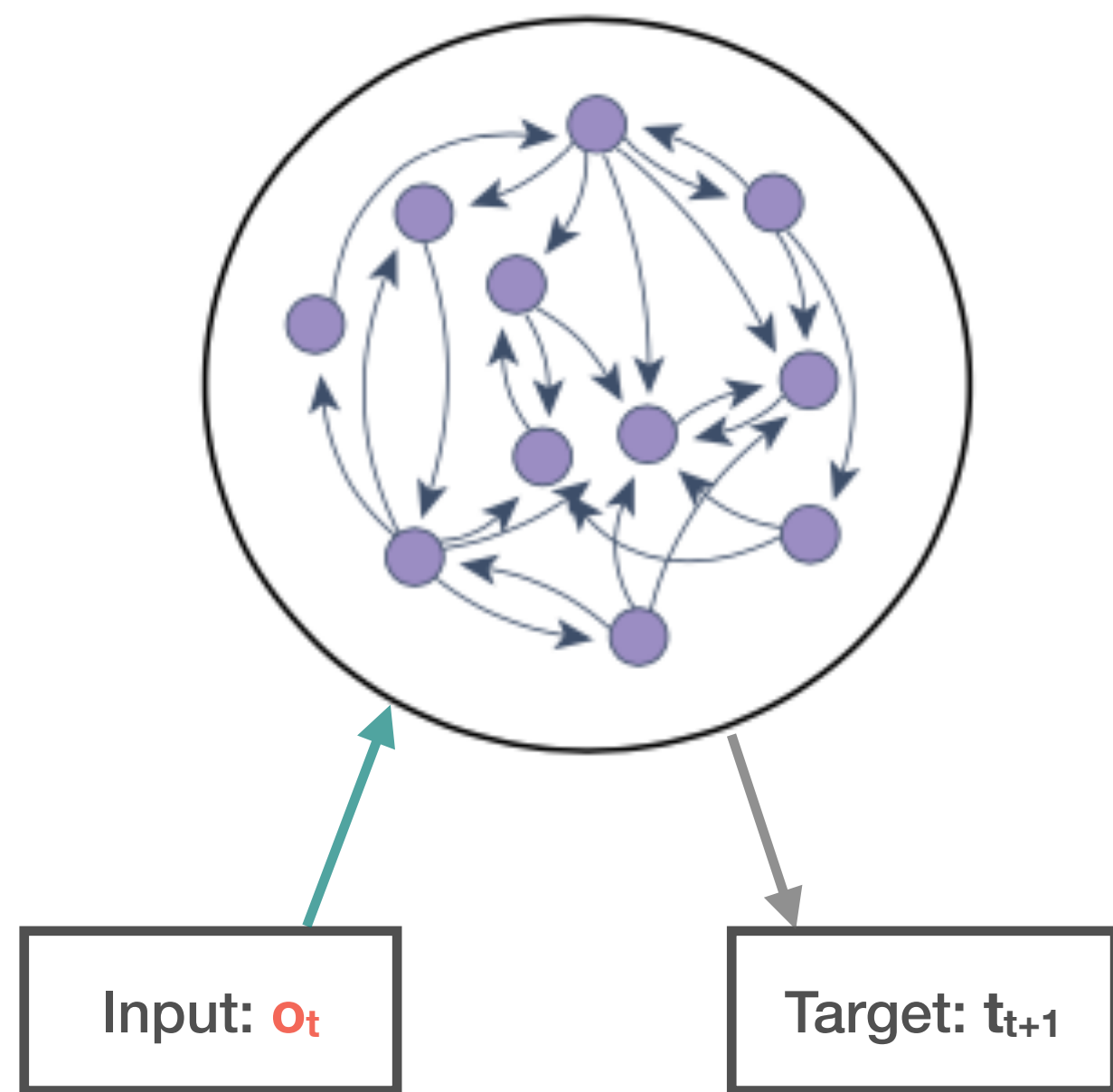
# Different problems have different slot structure



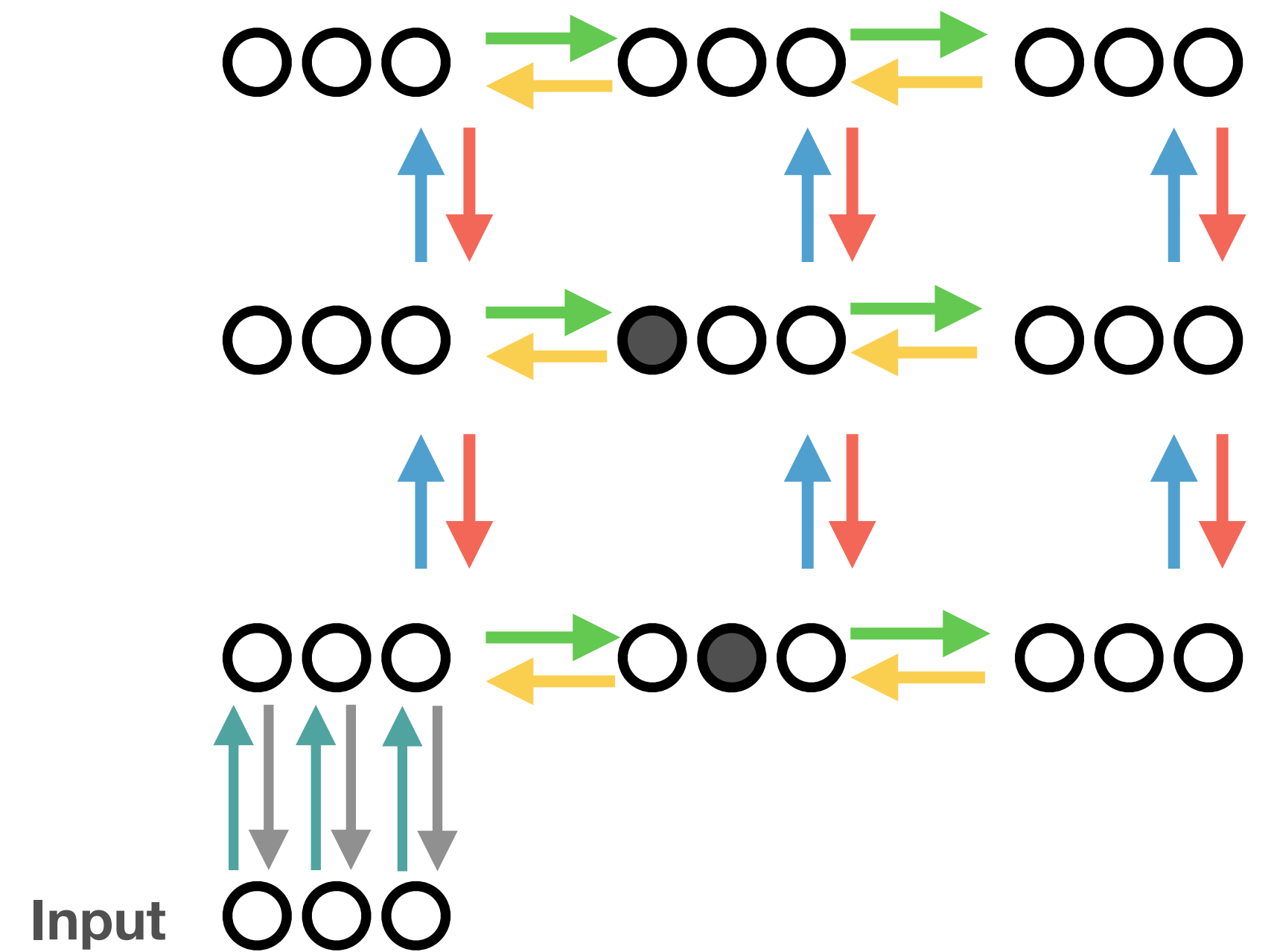
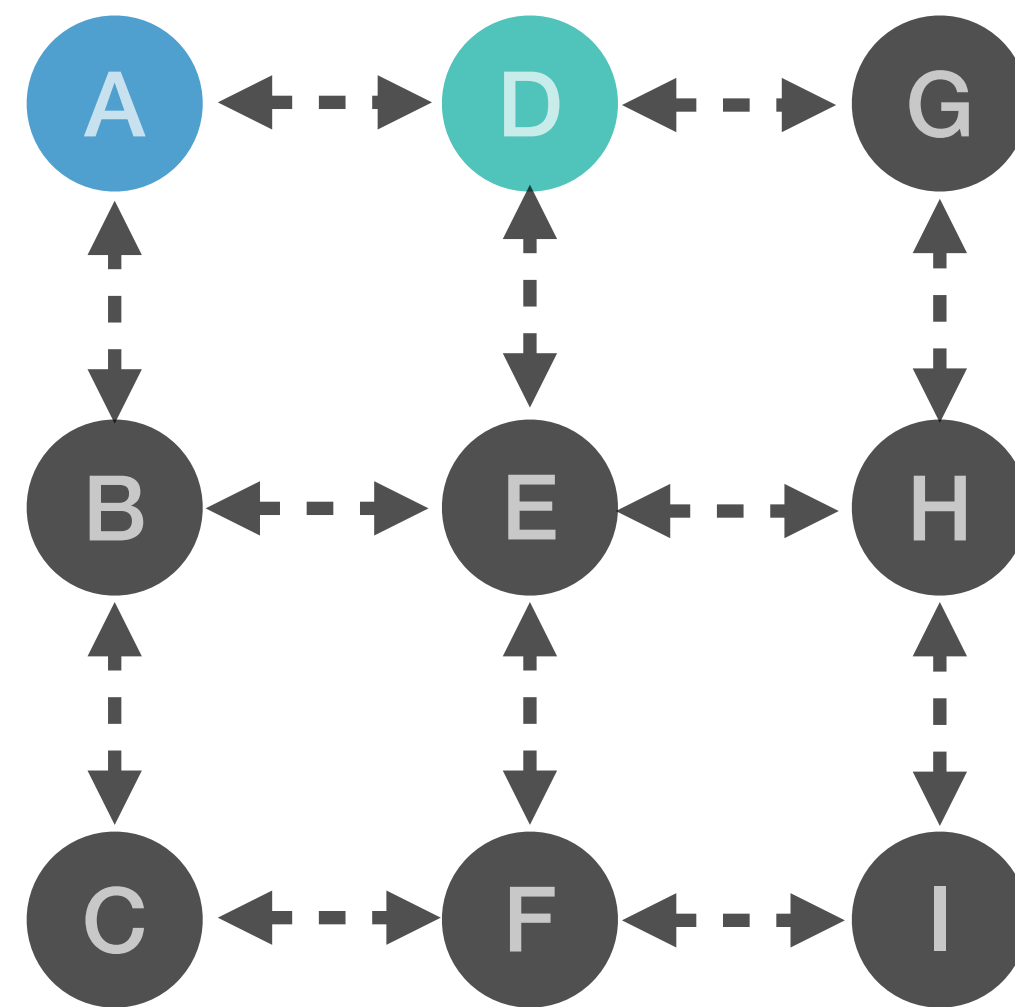
2D Navigation



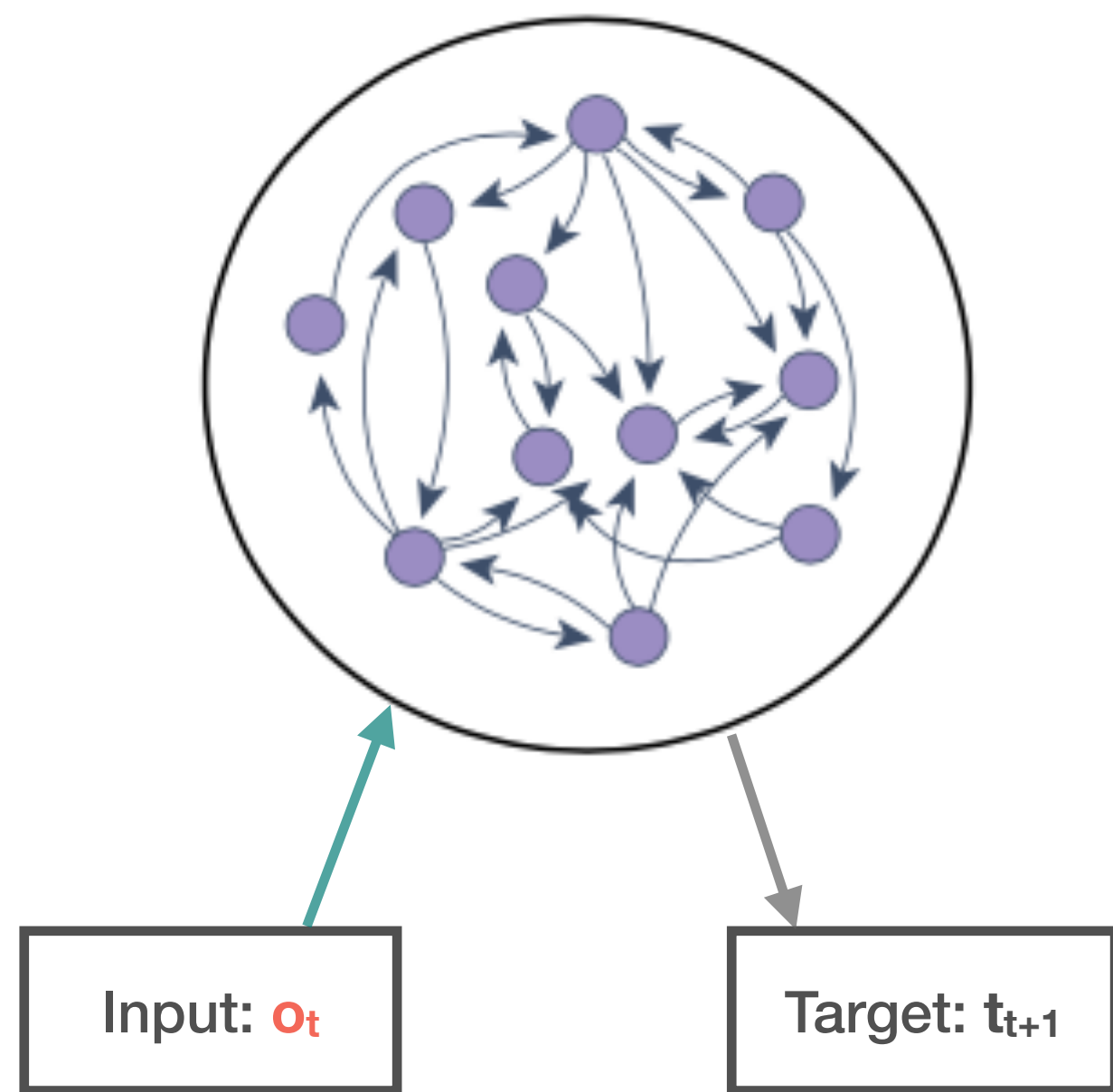
# Different problems have different slot structure



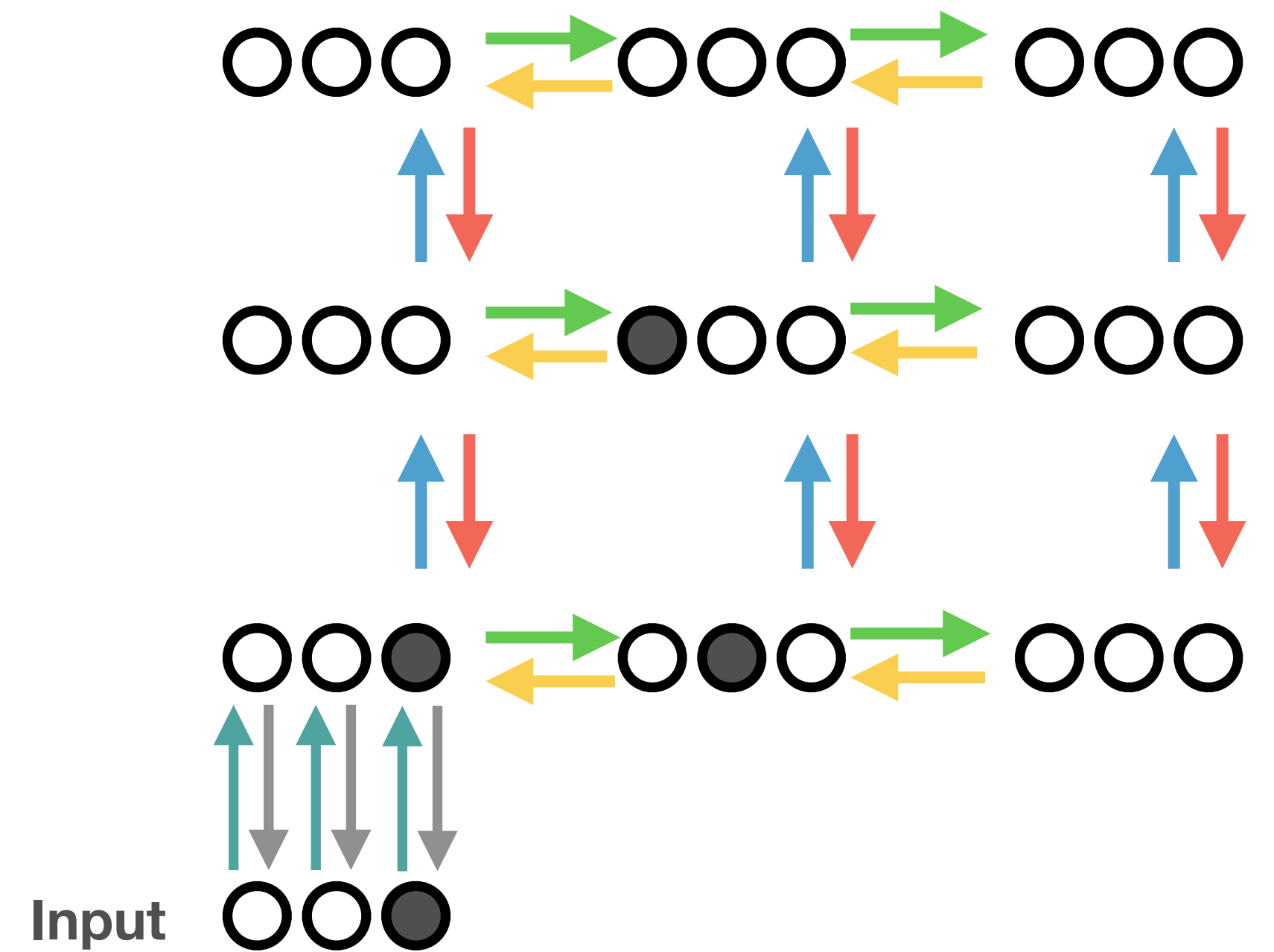
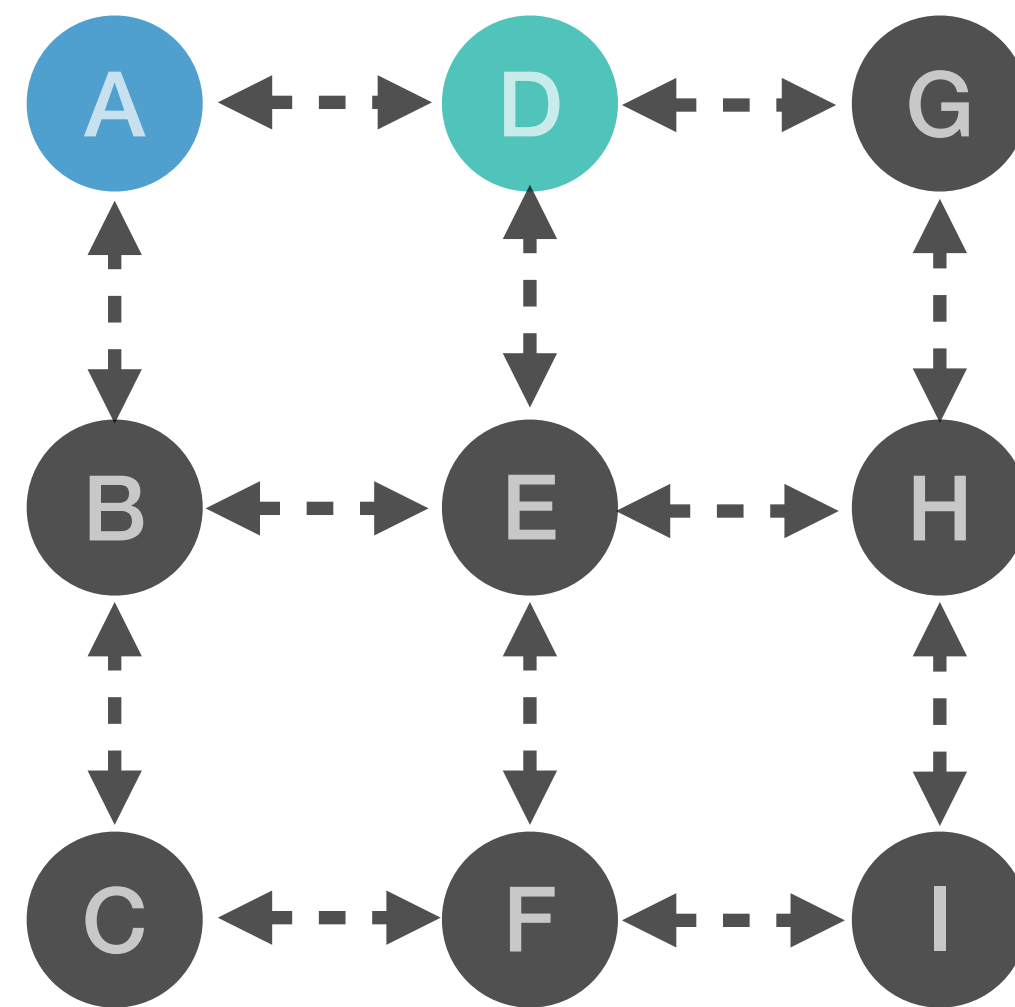
2D Navigation



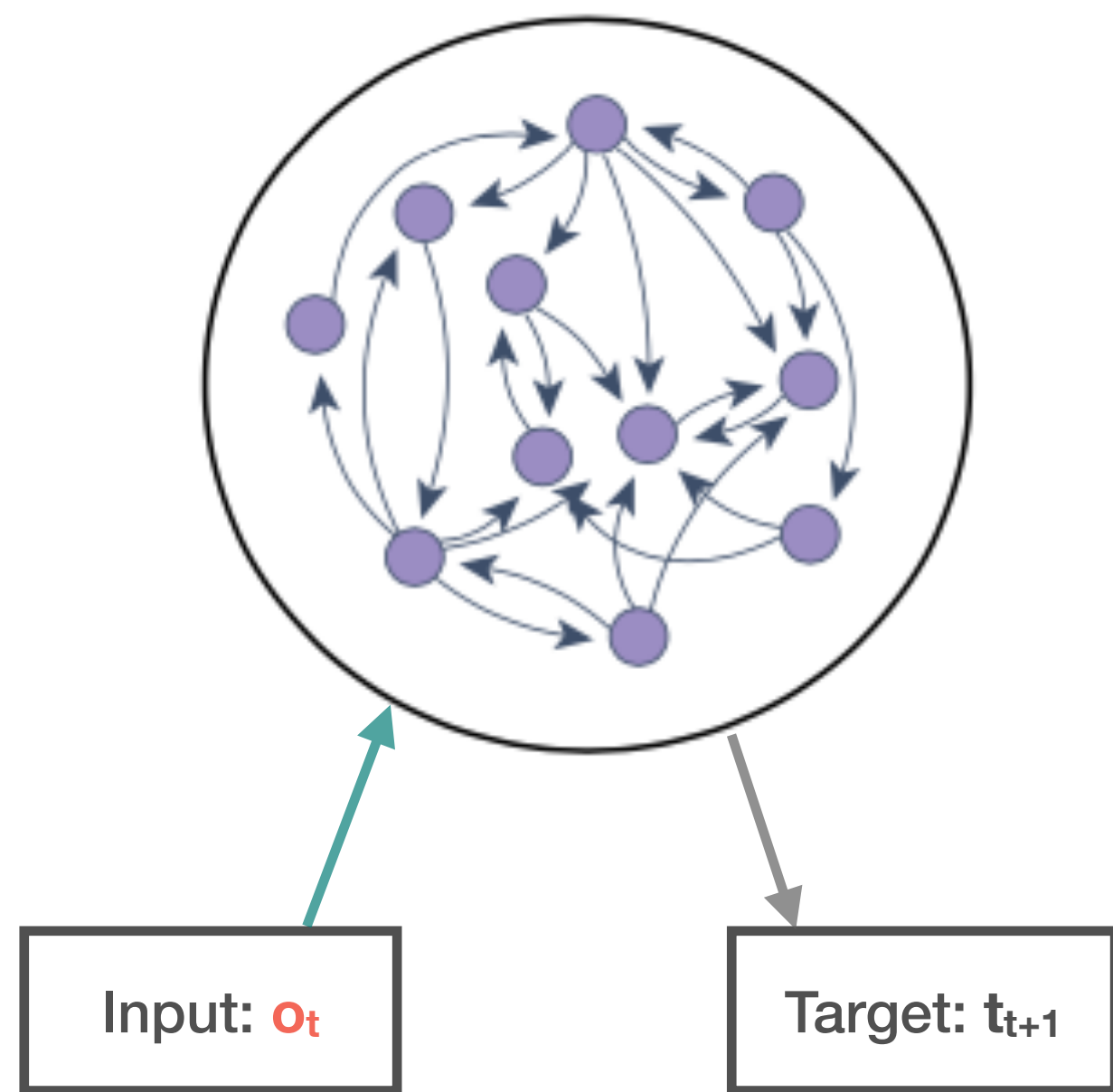
# Different problems have different slot structure



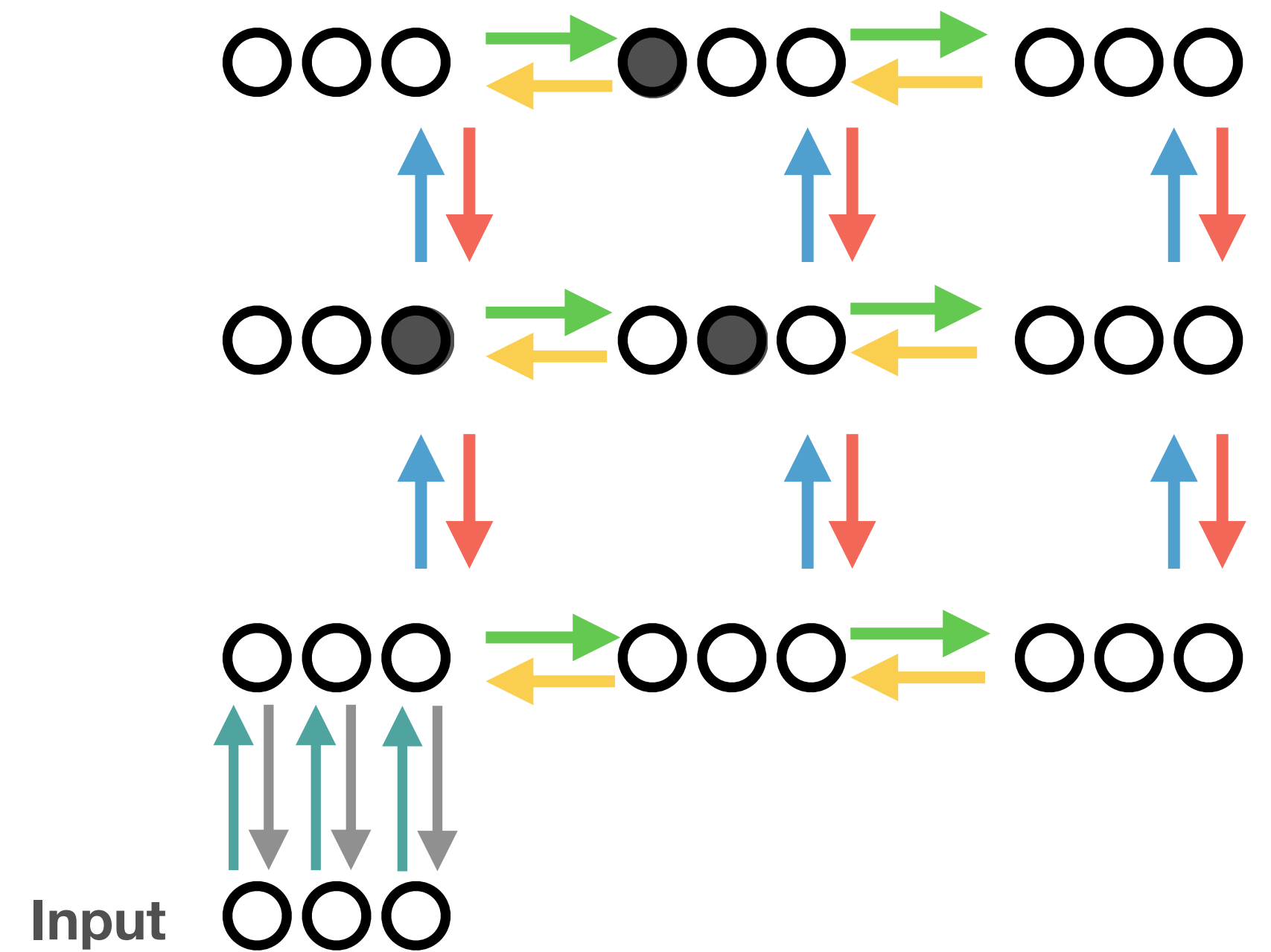
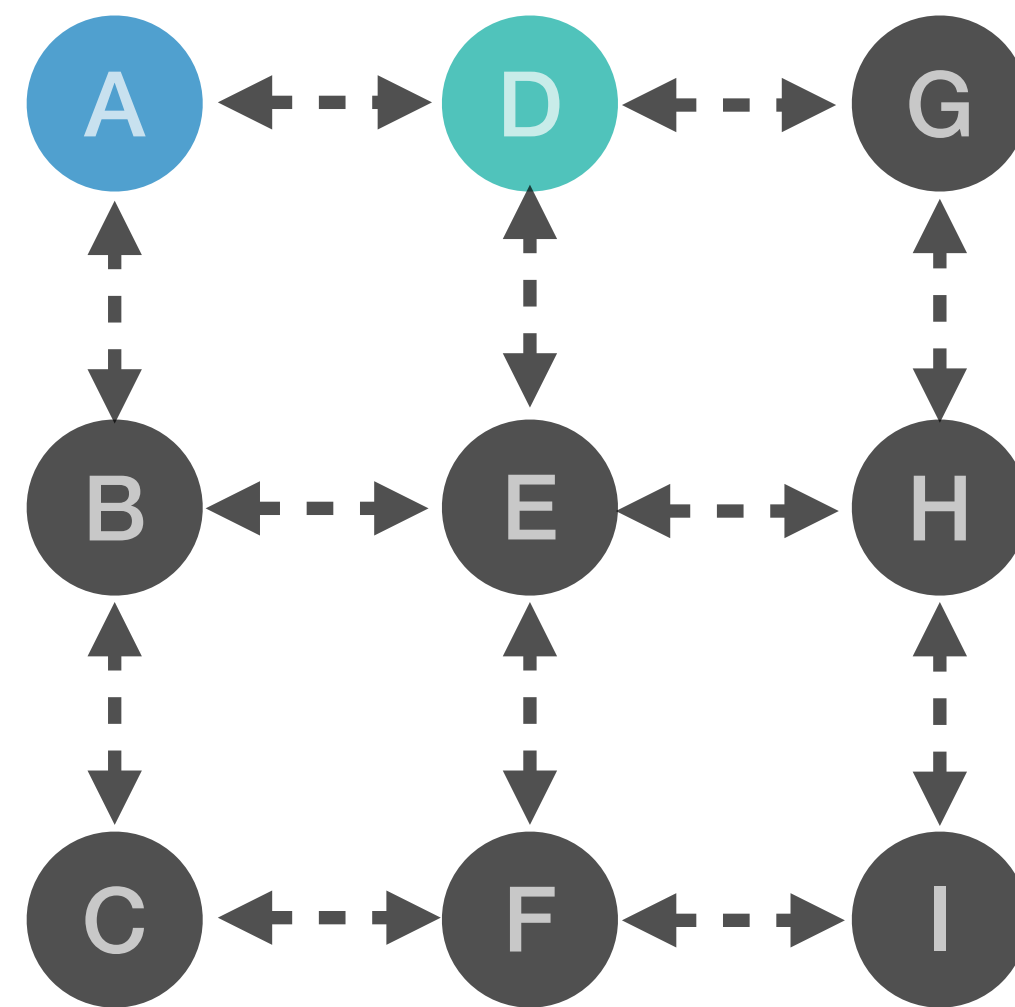
2D Navigation



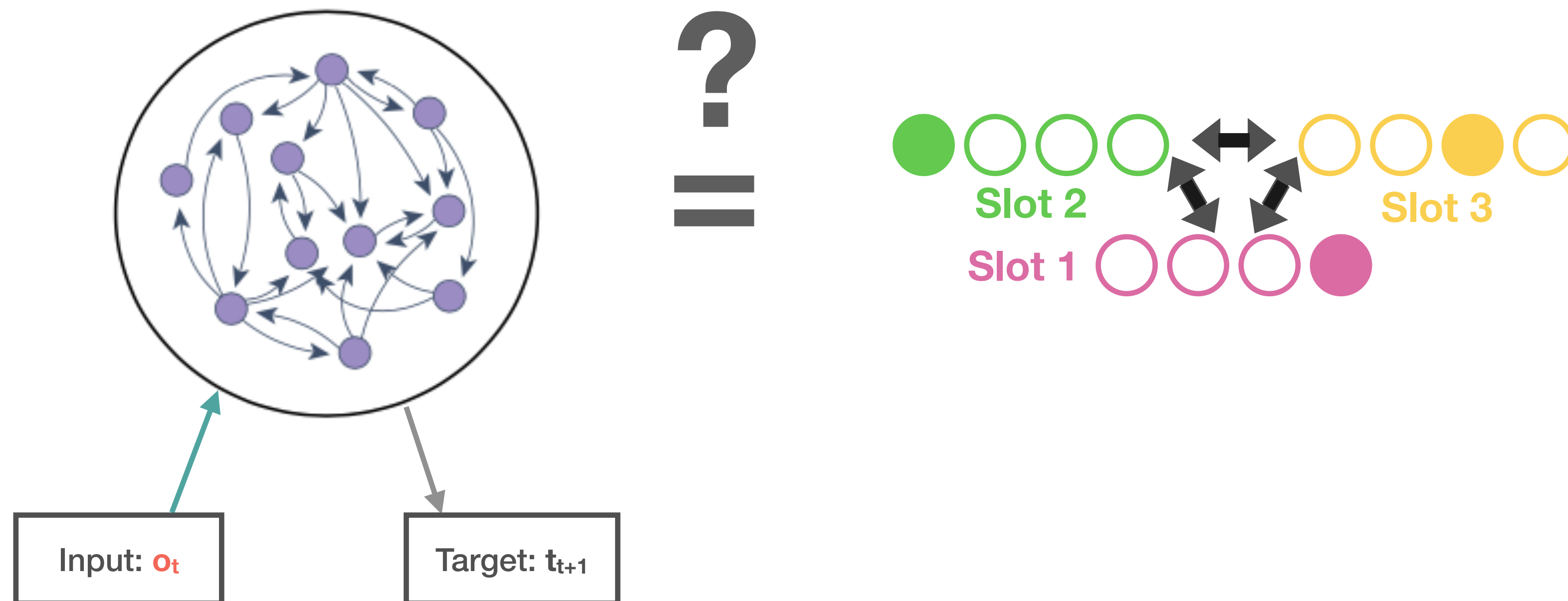
# Different problems have different slot structure



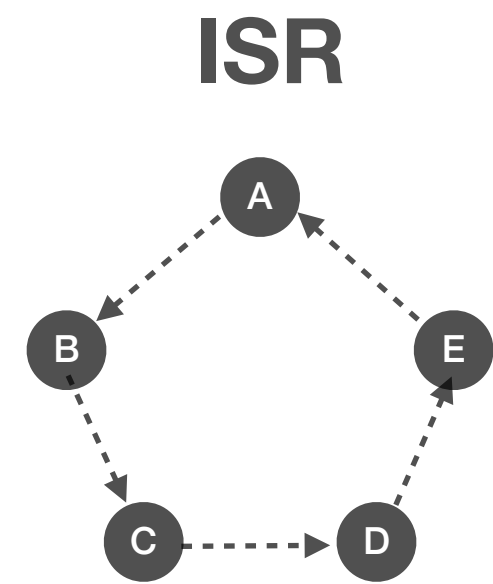
2D Navigation



# Do RNNs actually learn this mechanistic model?

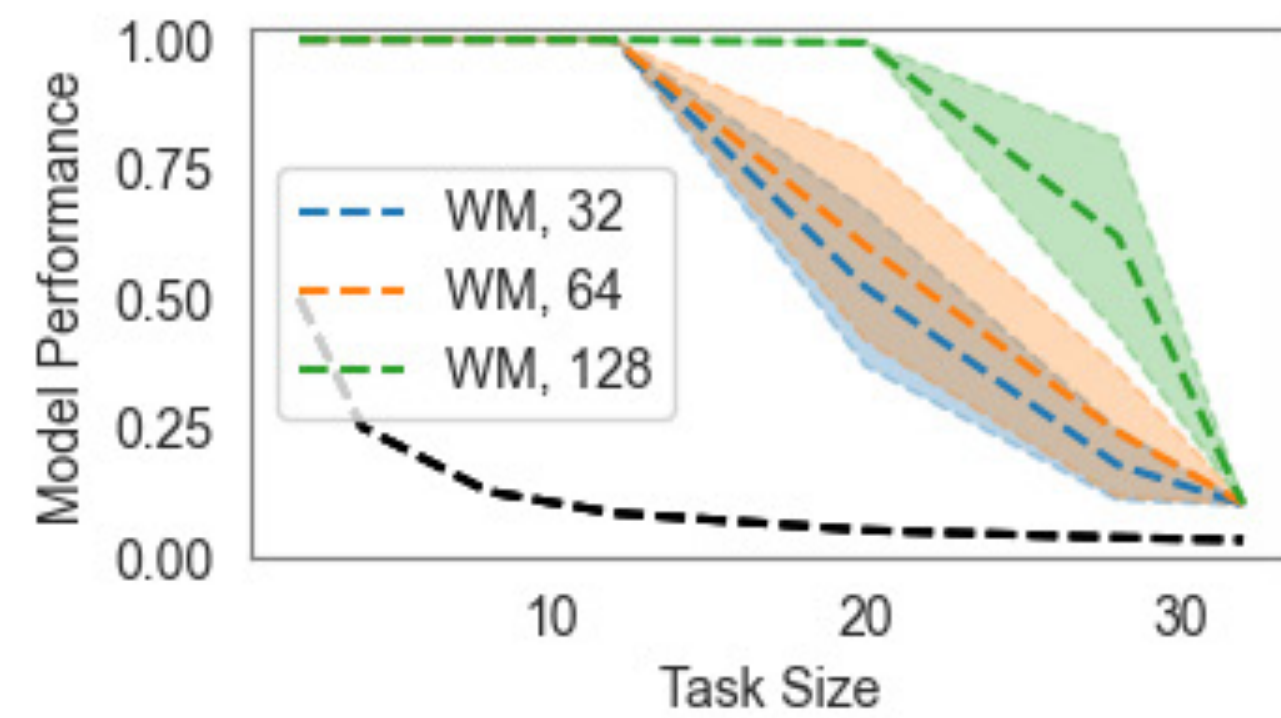
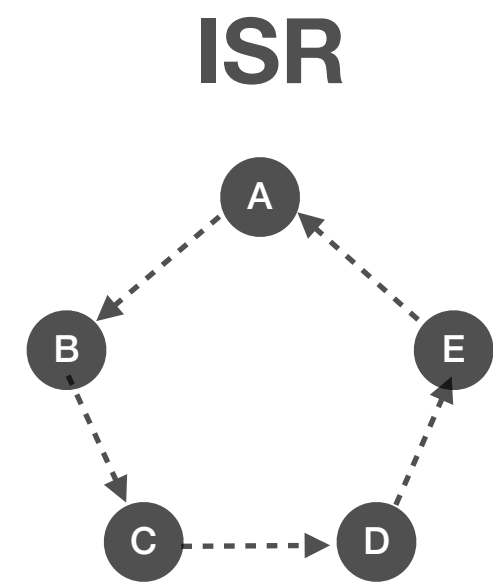


# Frist, RNNs can actually solve the task

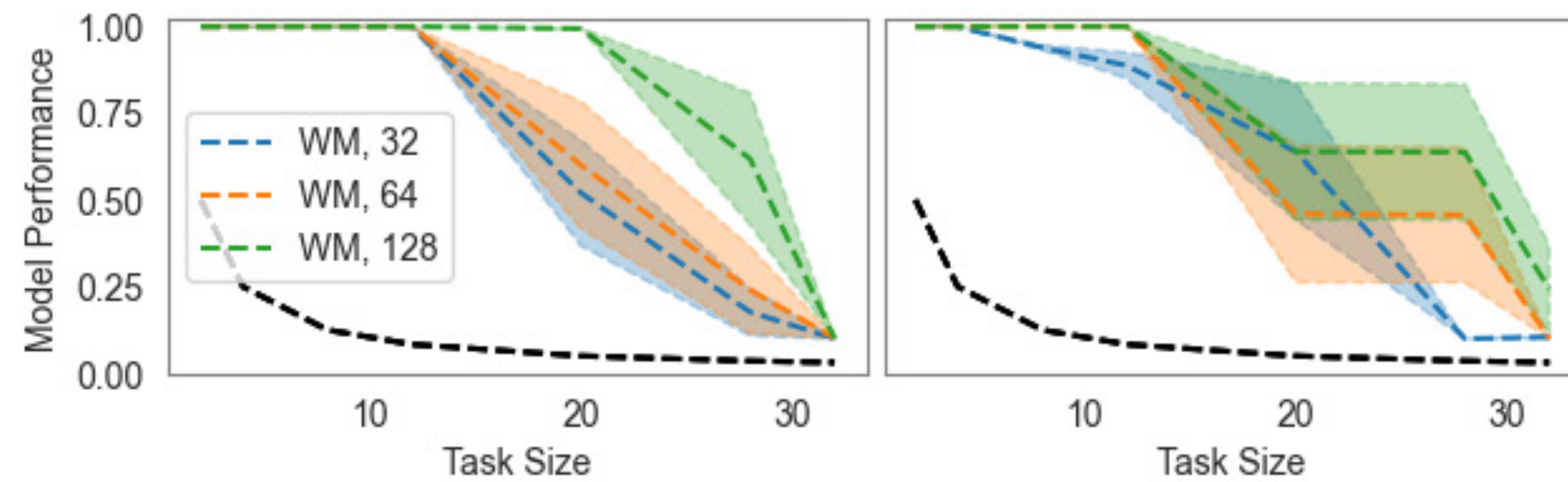
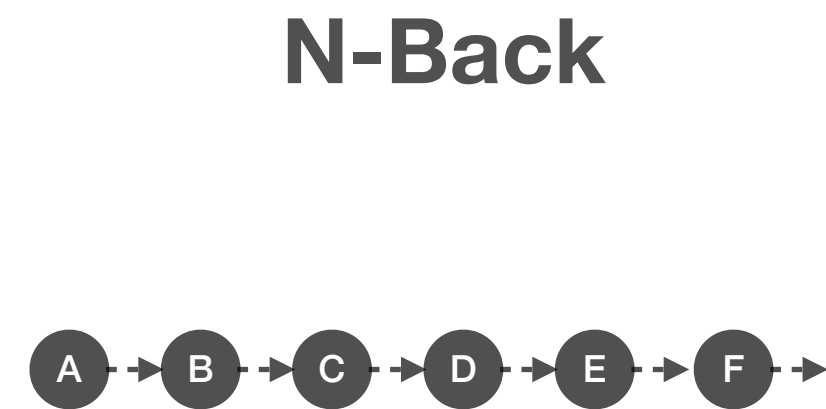
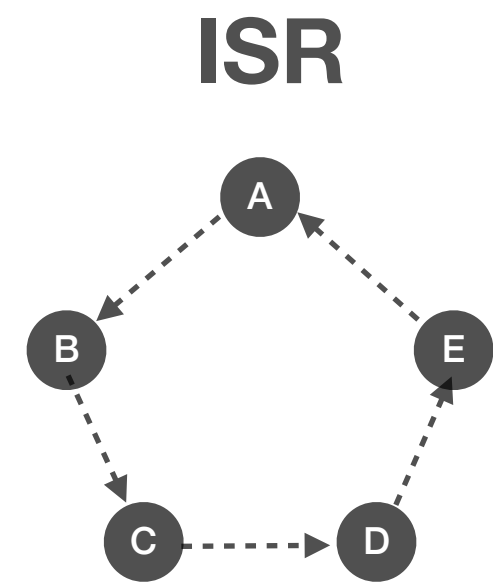




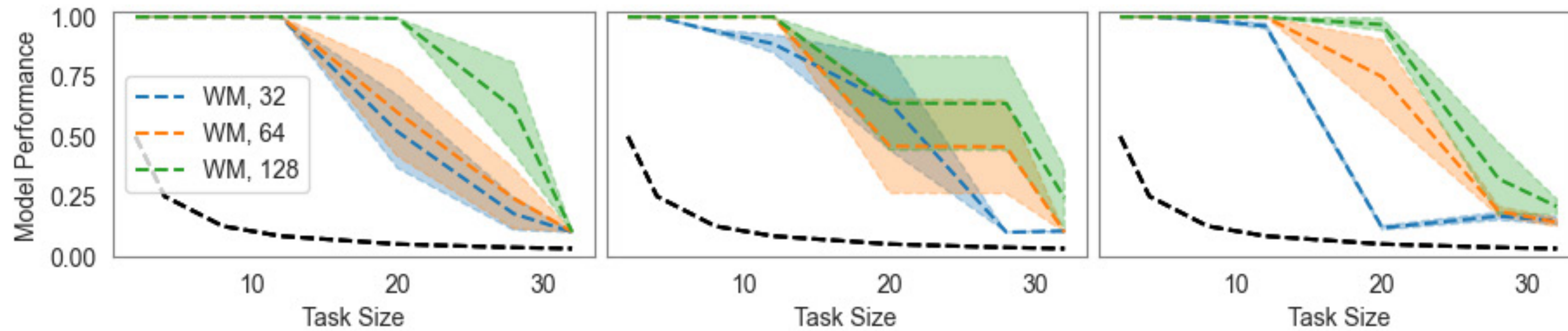
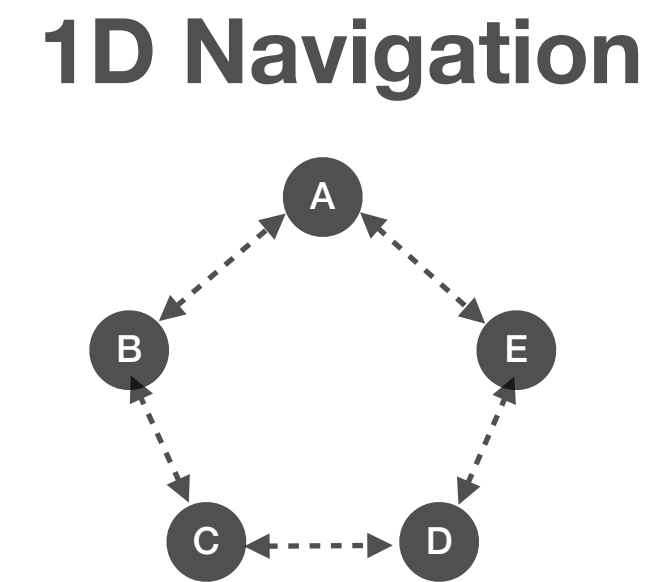
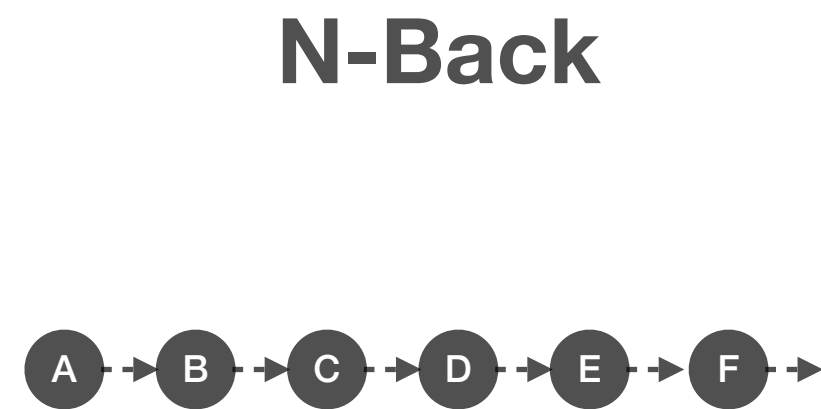
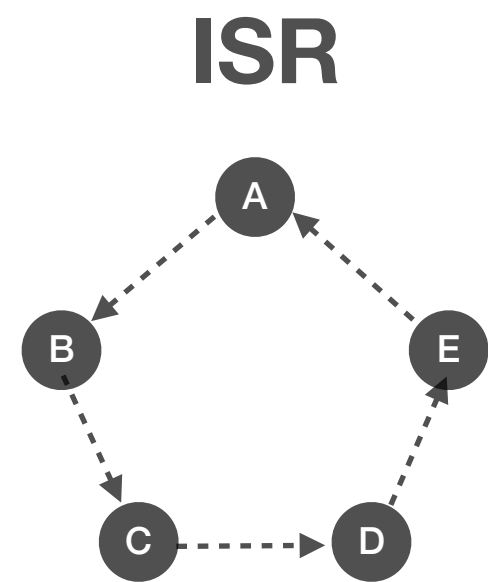
# Frist, RNNs can actually solve the task



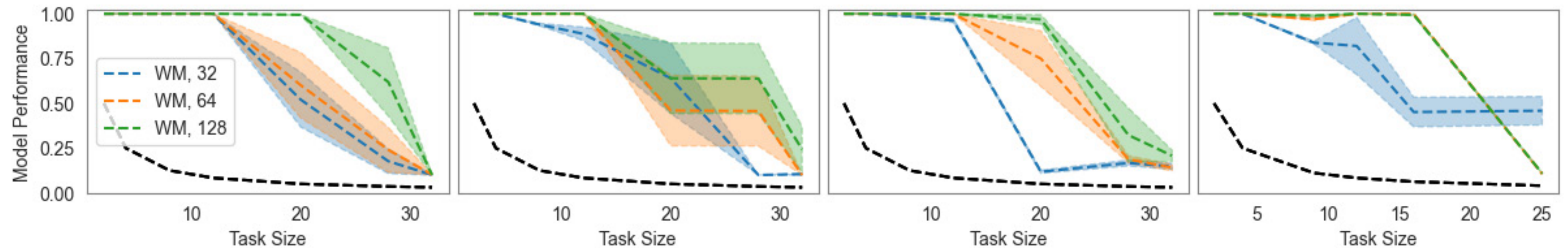
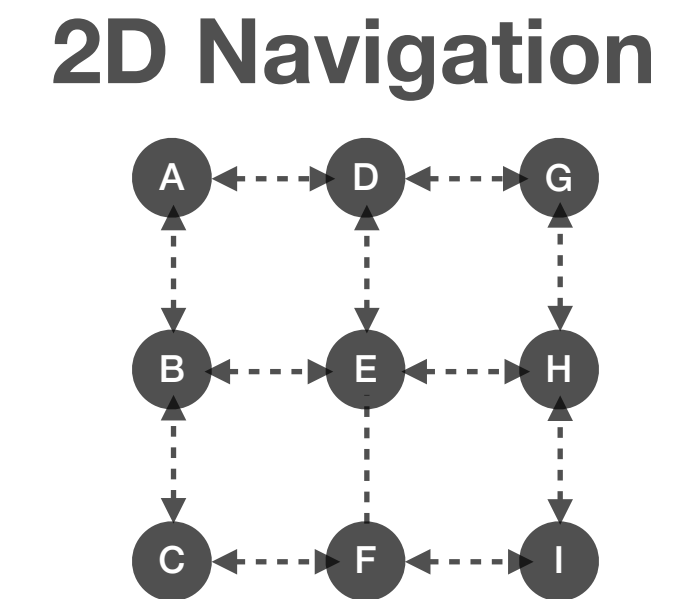
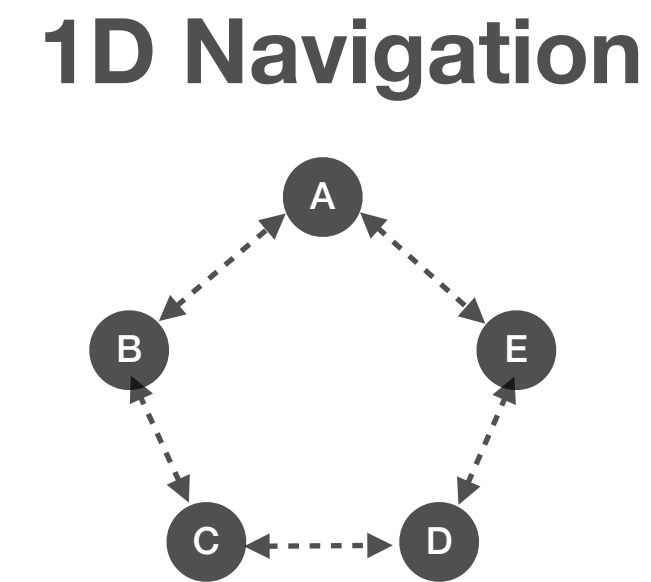
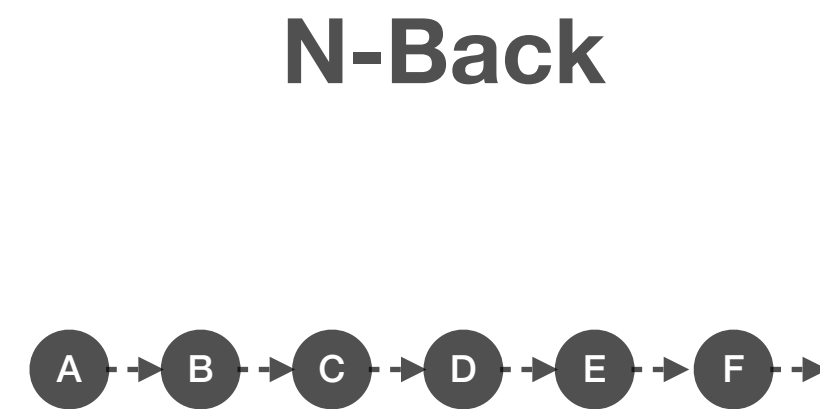
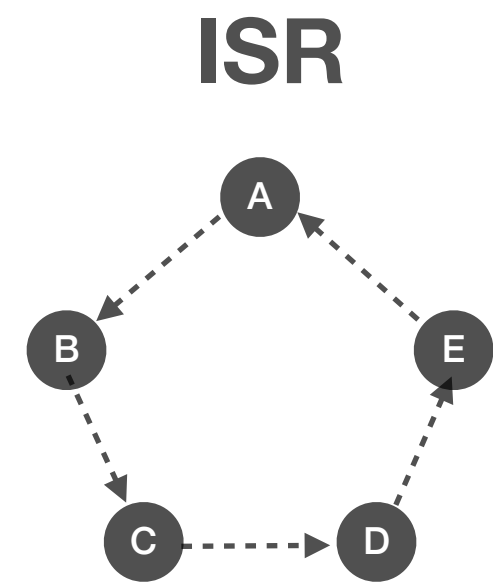
# Frist, RNNs can actually solve the task



# Frist, RNNs can actually solve the task



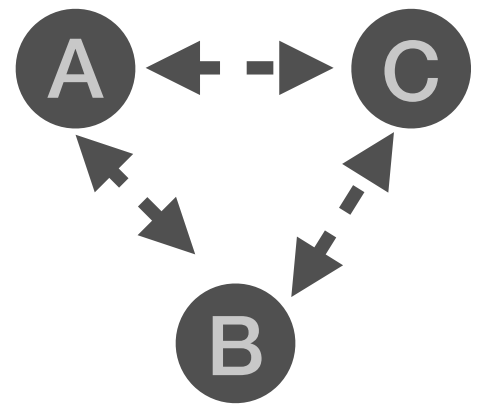
# Frist, RNNs can actually solve the task



# RNNs do use structured slots

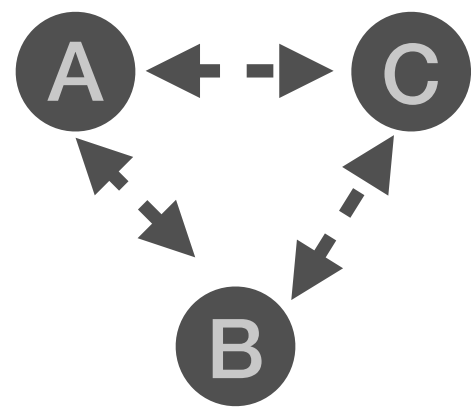
# RNNs do use structured slots

1D Navigation

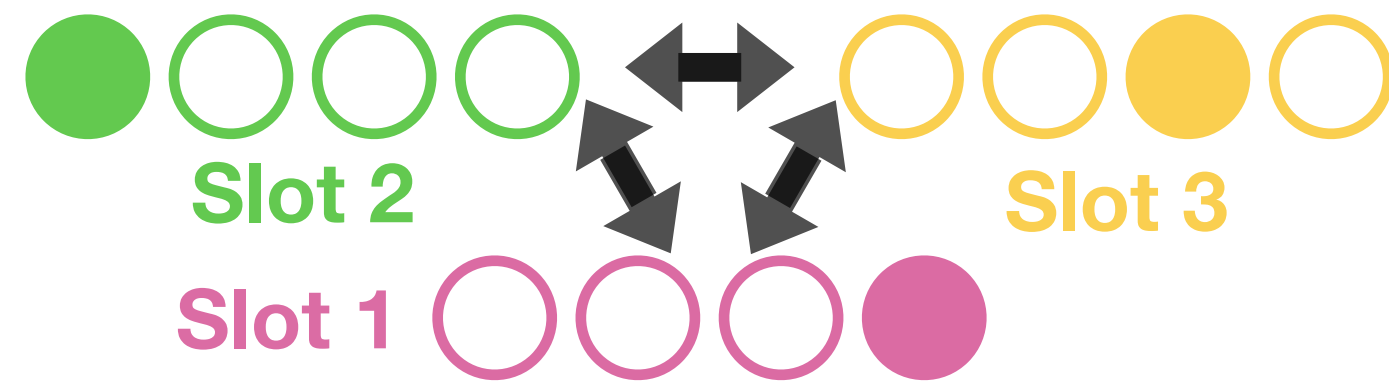


# RNNs do use structured slots

1D Navigation

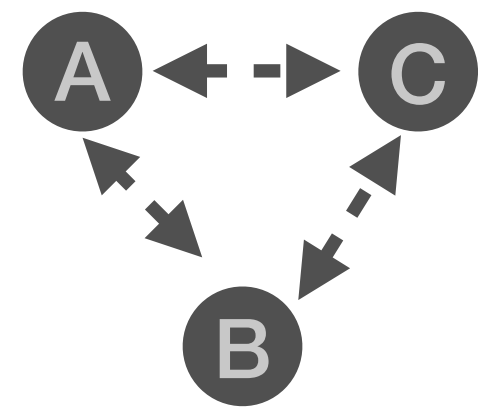


1D slot prediction

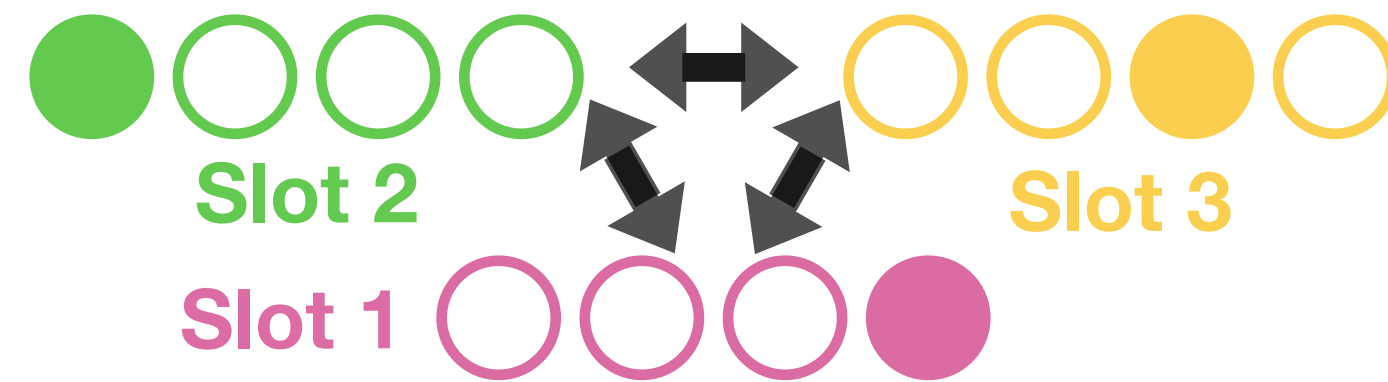


# RNNs do use structured slots

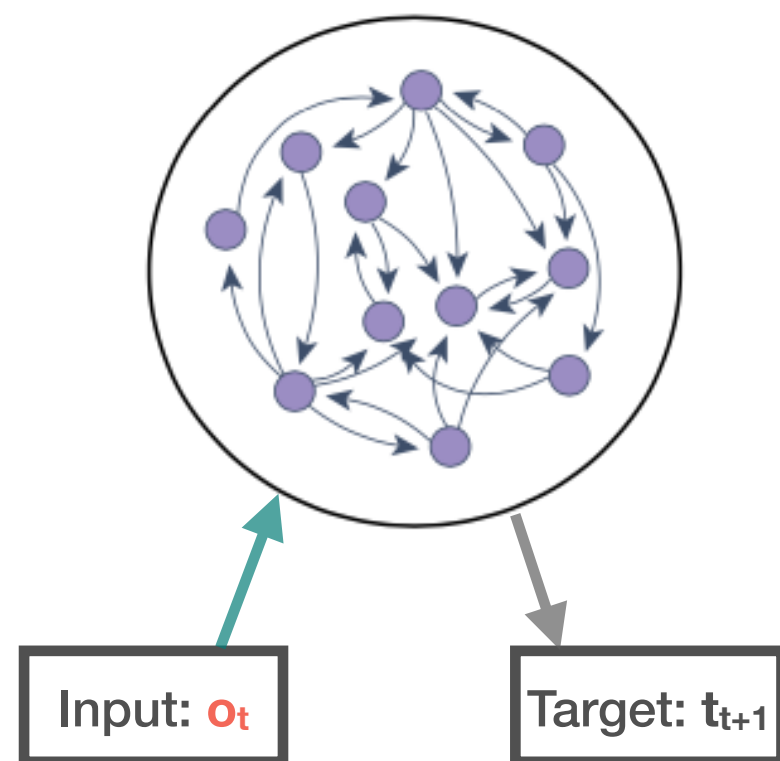
1D Navigation



1D slot prediction



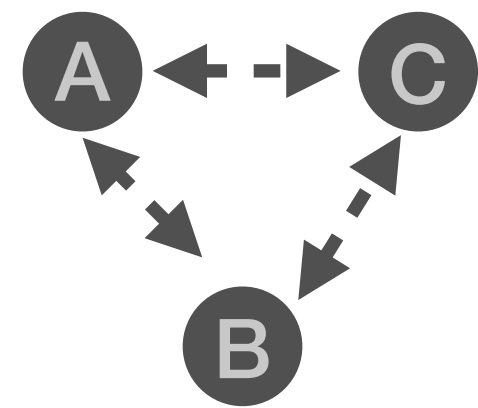
Should be able to decode the predicted contents of each slot at every tilmestep



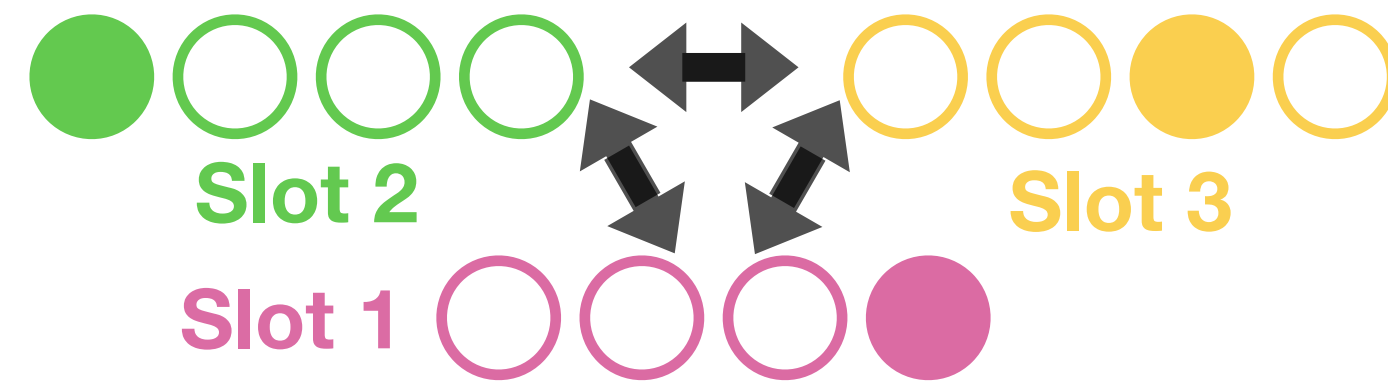


# RNNs do use structured slots

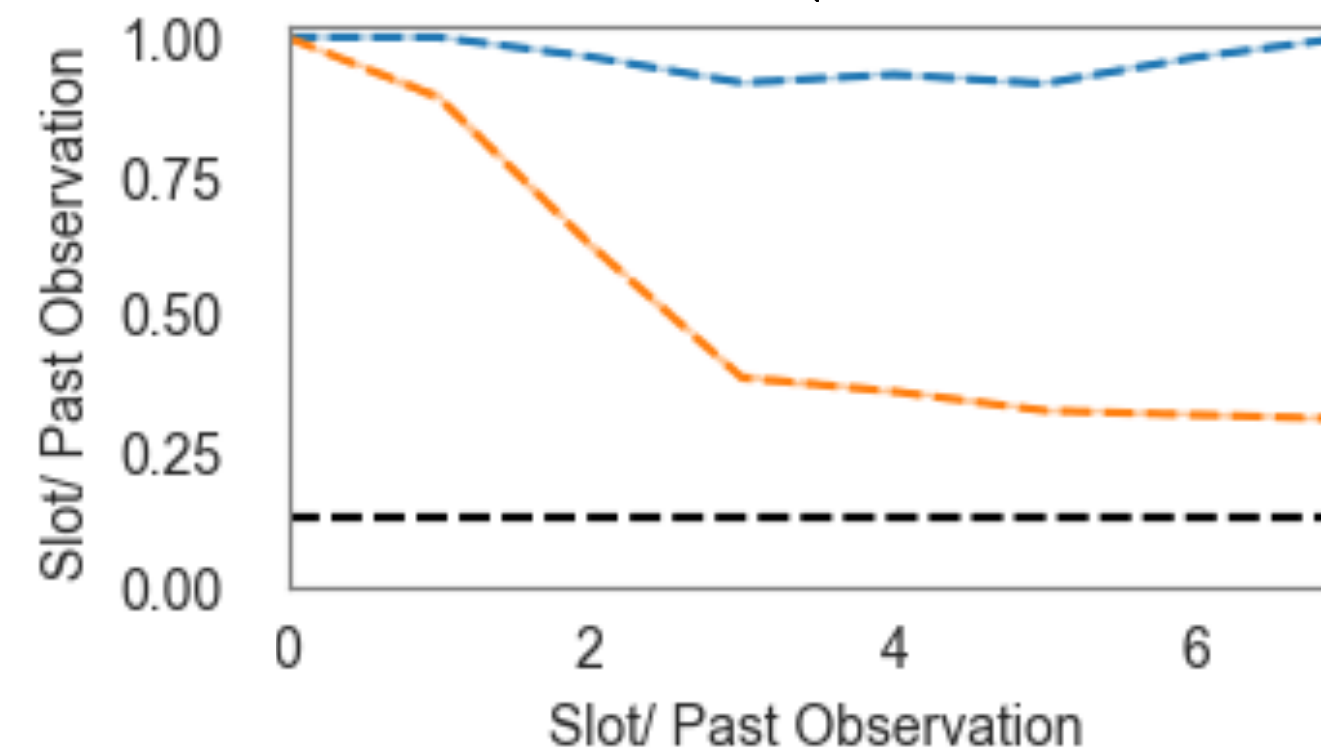
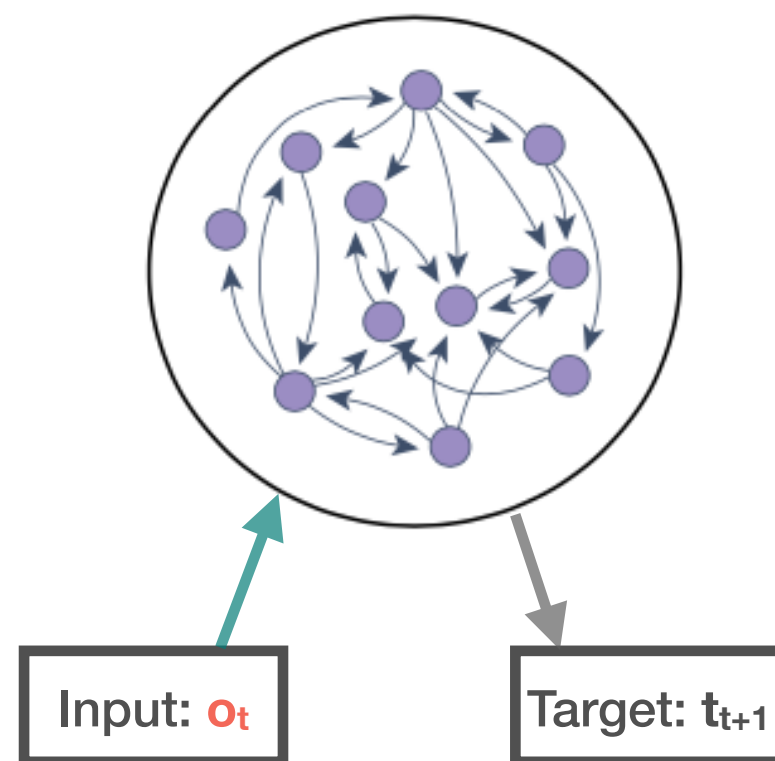
1D Navigation



1D slot prediction

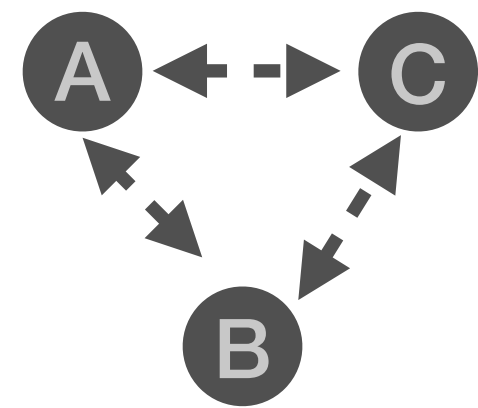


Should be able to decode the predicted contents of each slot at every timesteps

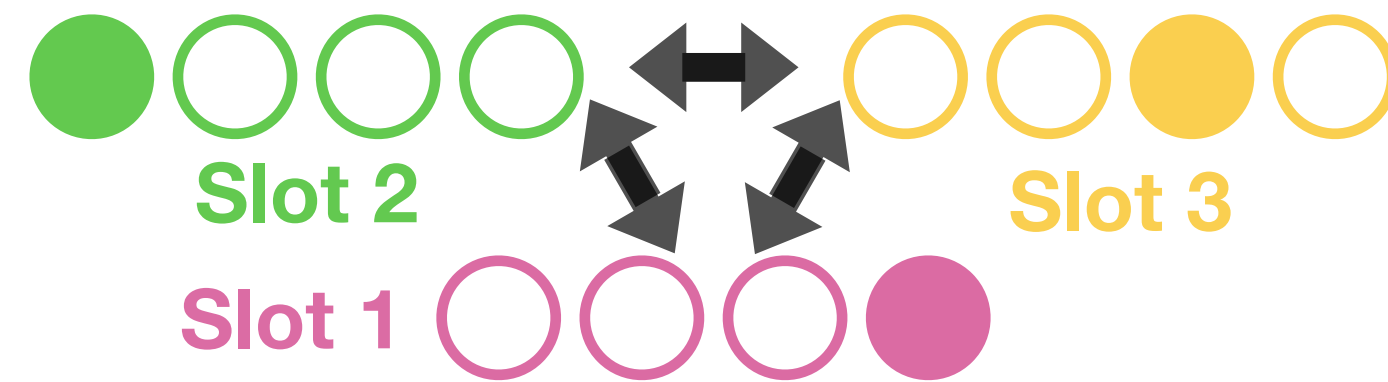


# RNNs do use structured slots

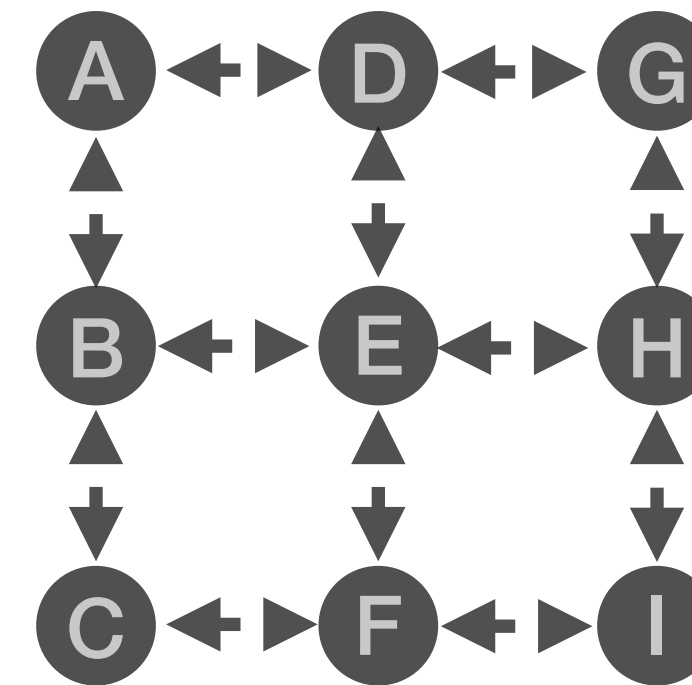
## 1D Navigation



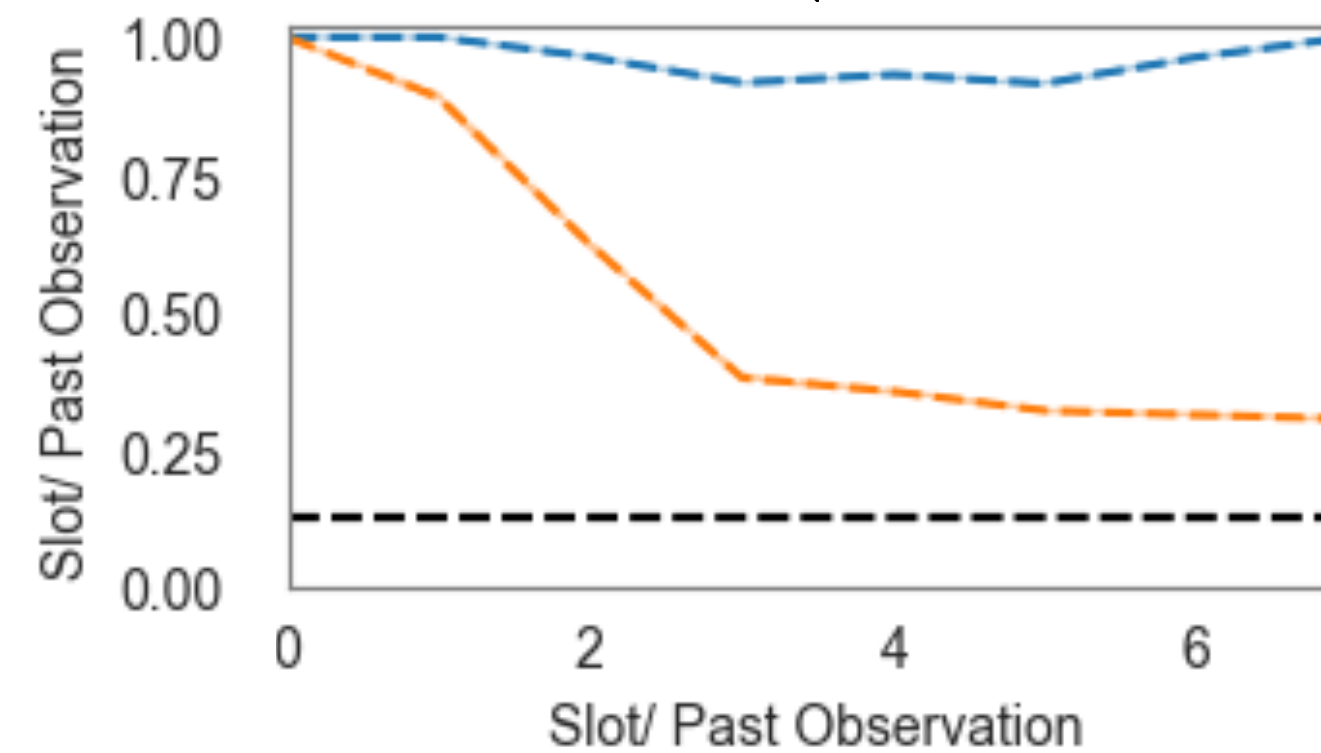
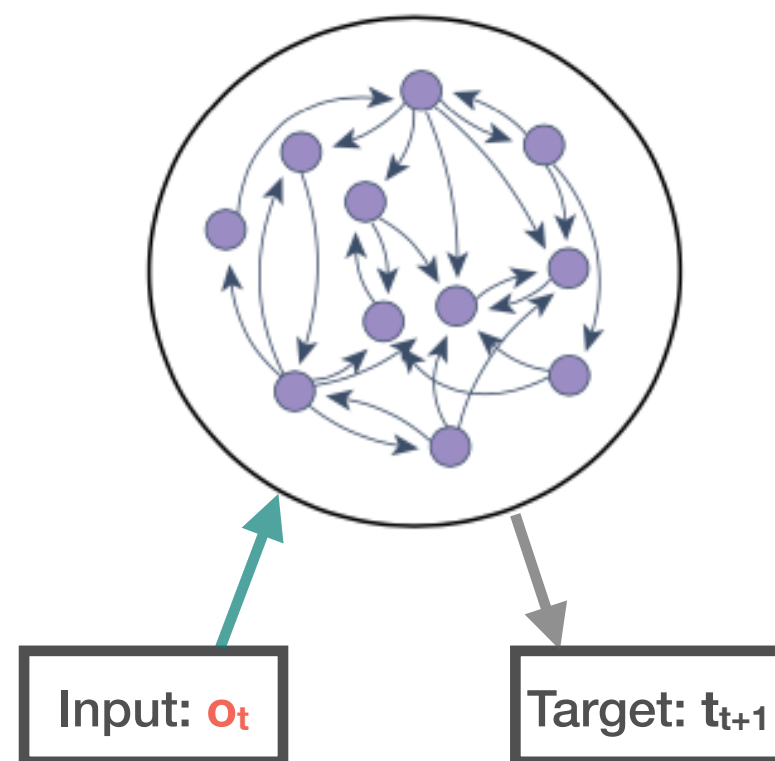
## 1D slot prediction



## 2D Navigation

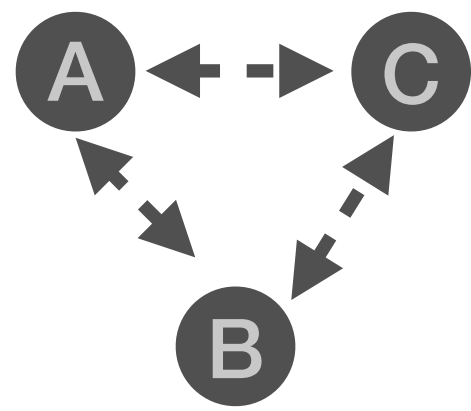


Should be able to decode the predicted contents of each slot at every timesteps

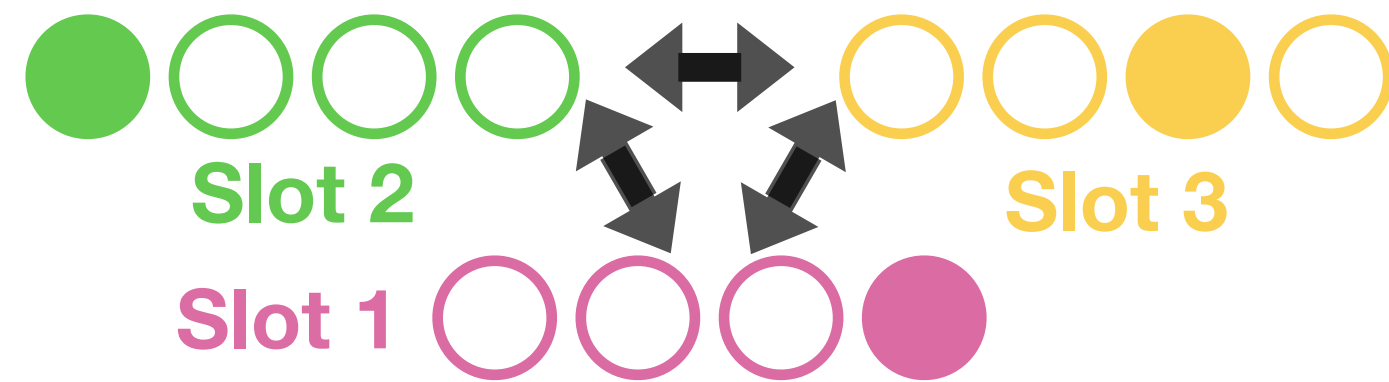


# RNNs do use structured slots

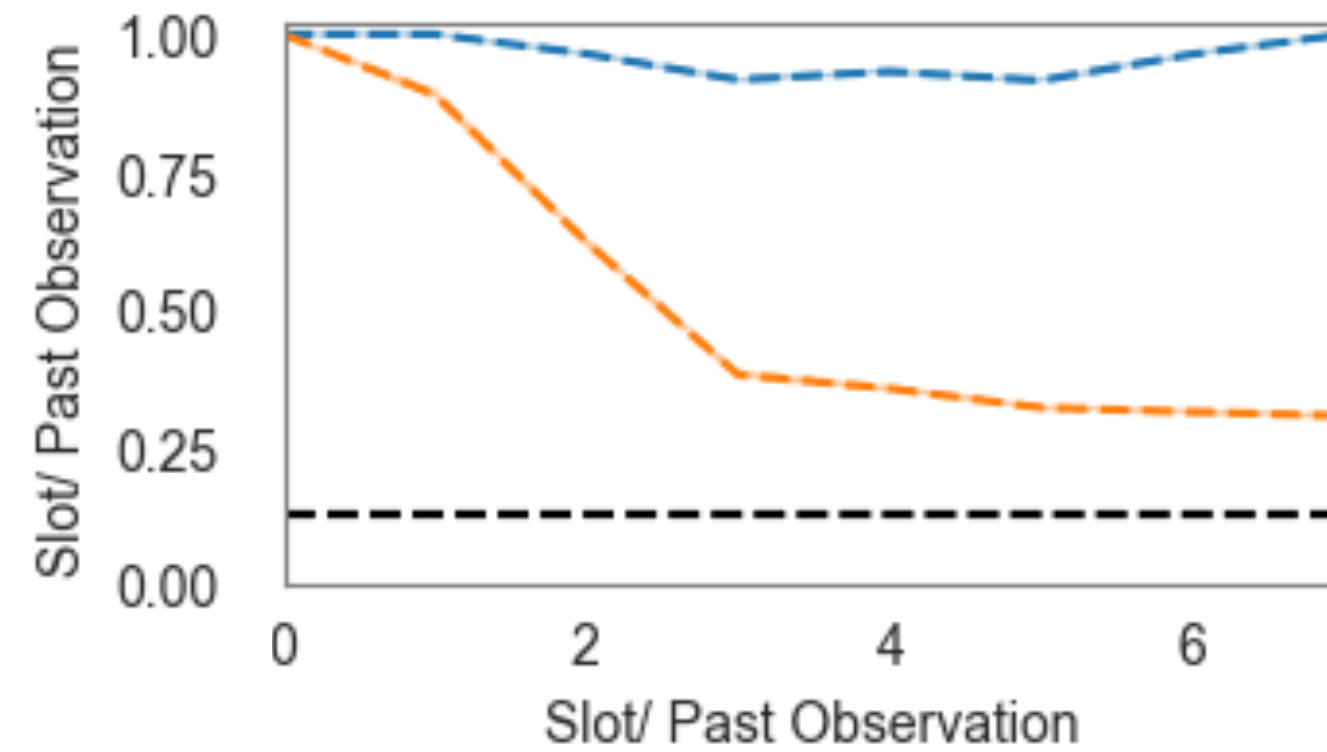
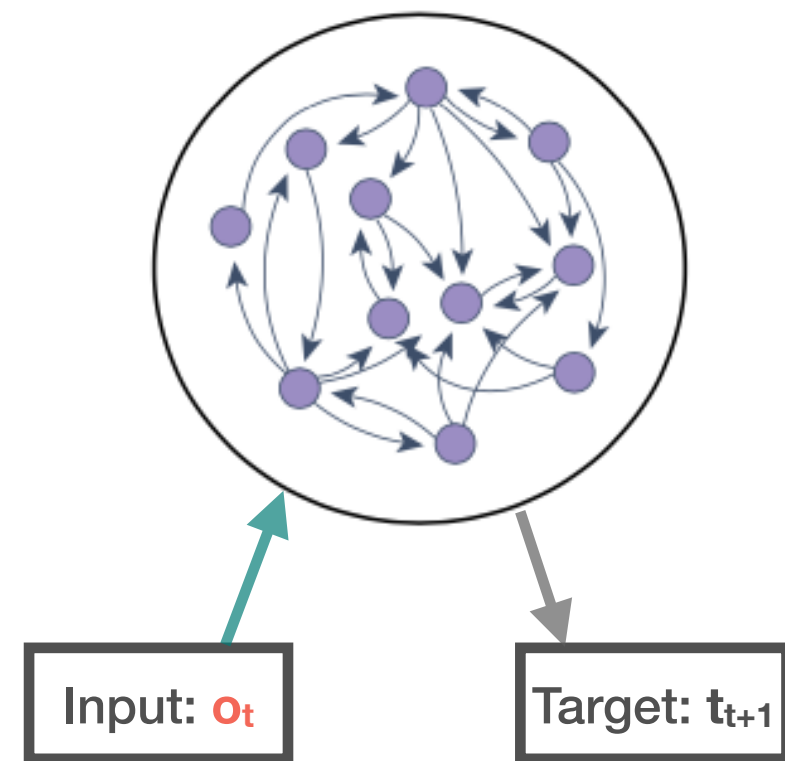
## 1D Navigation



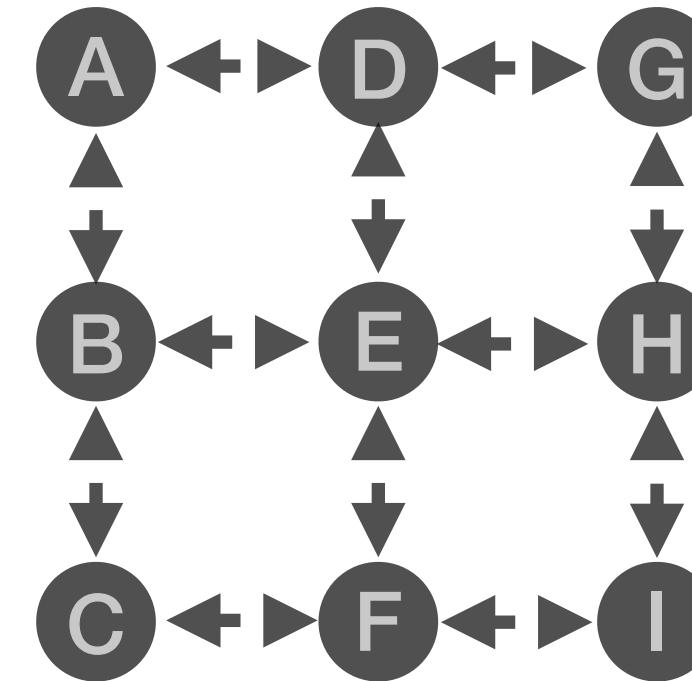
## 1D slot prediction



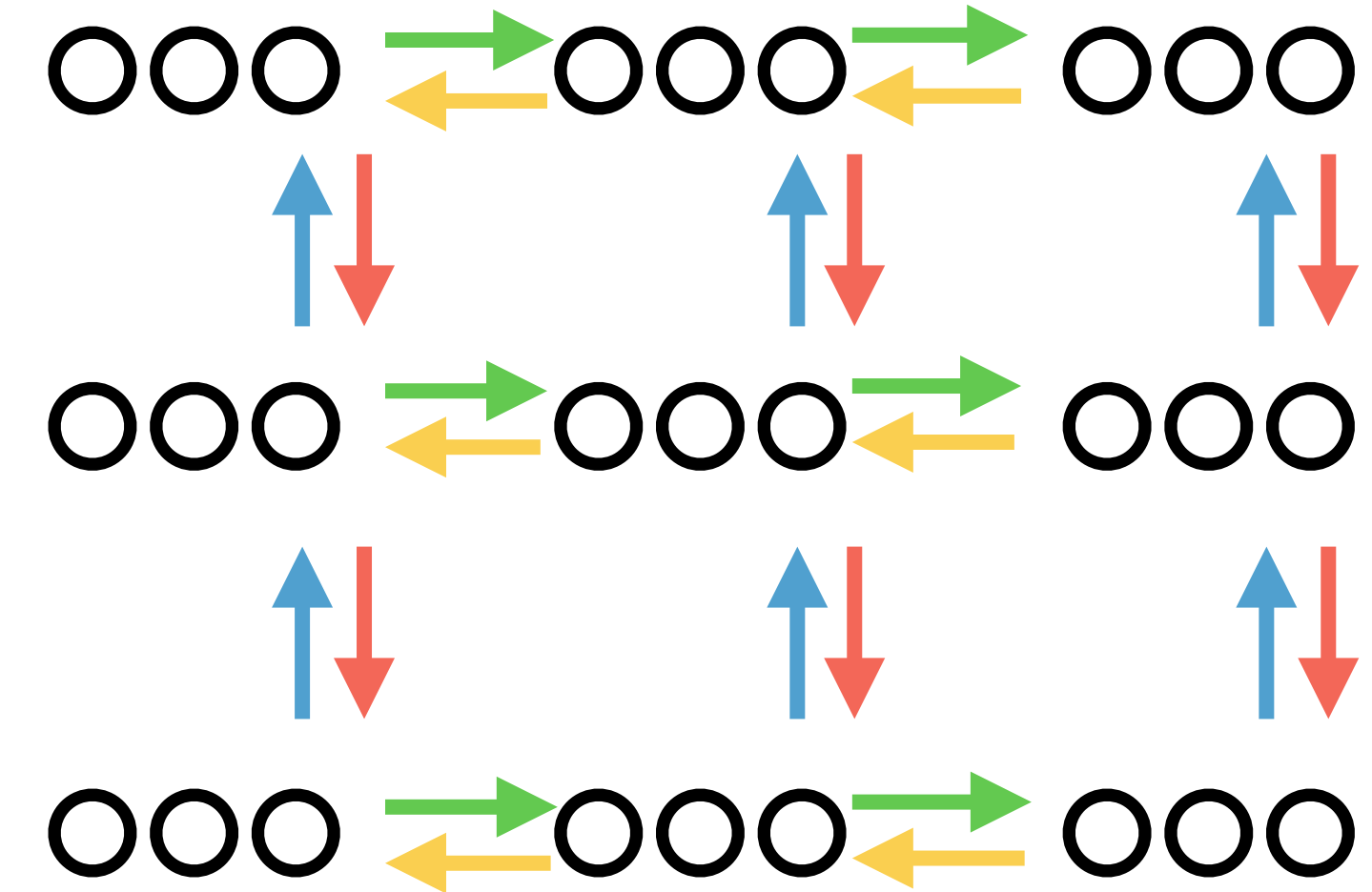
Should be able to decode the predicted contents of each slot at every timesteps



## 2D Navigation

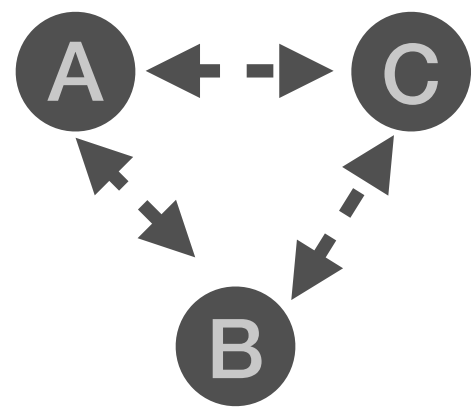


## 2D slot prediction

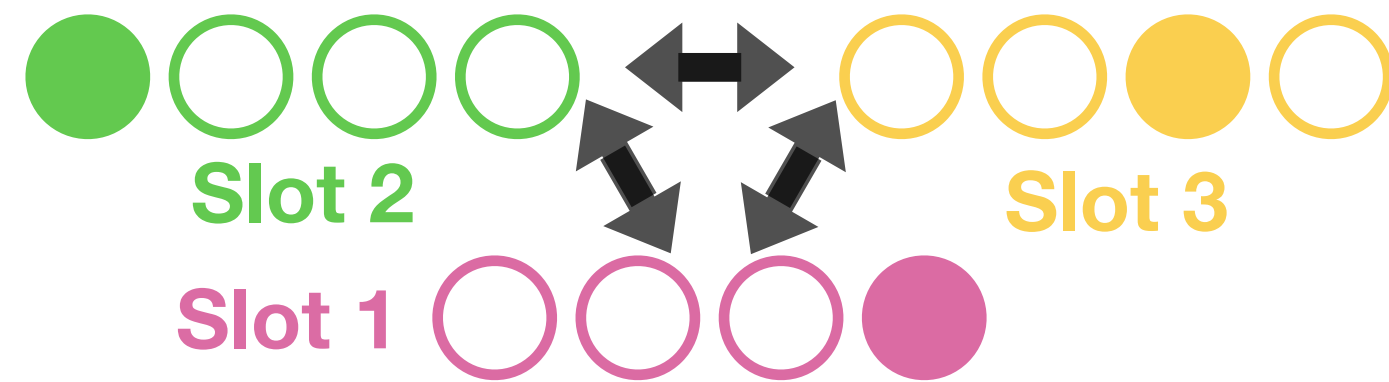


# RNNs do use structured slots

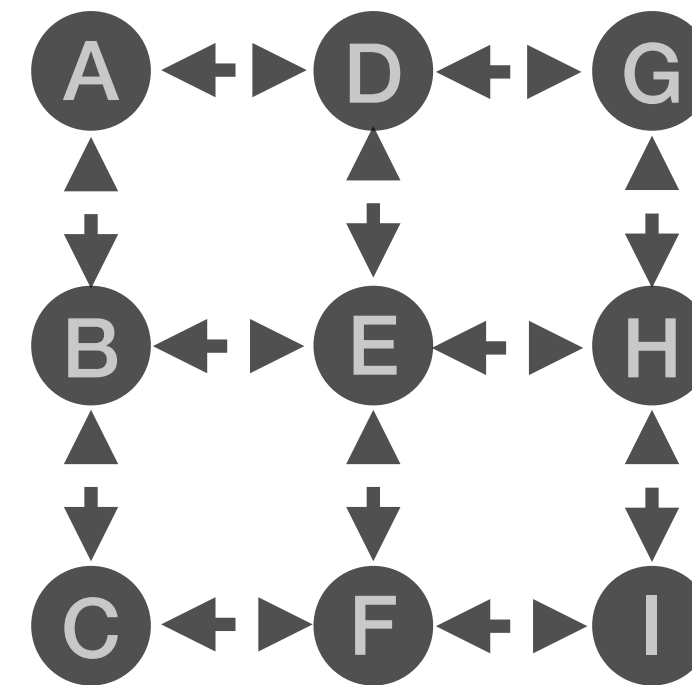
### 1D Navigation



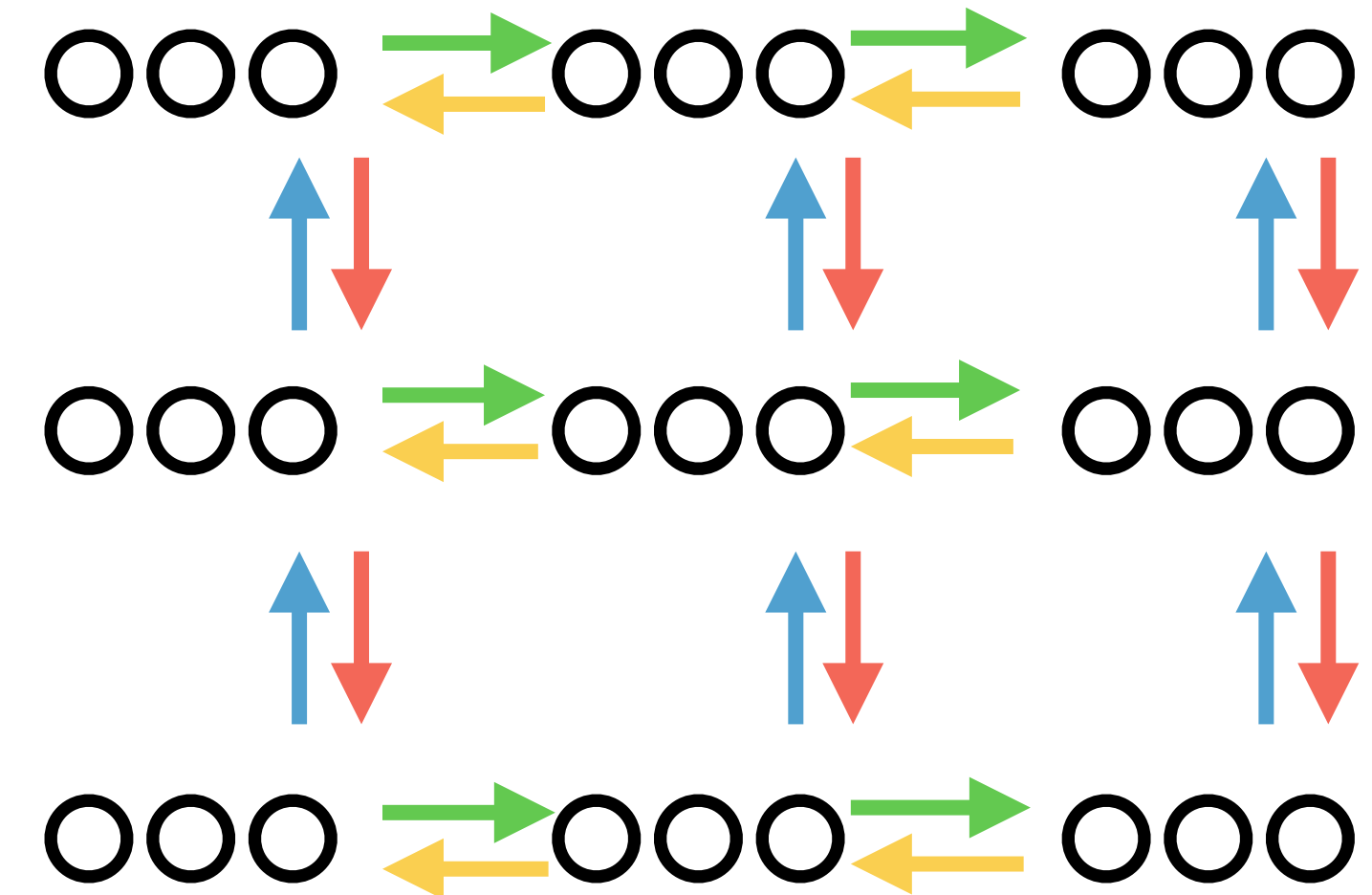
### 1D slot prediction



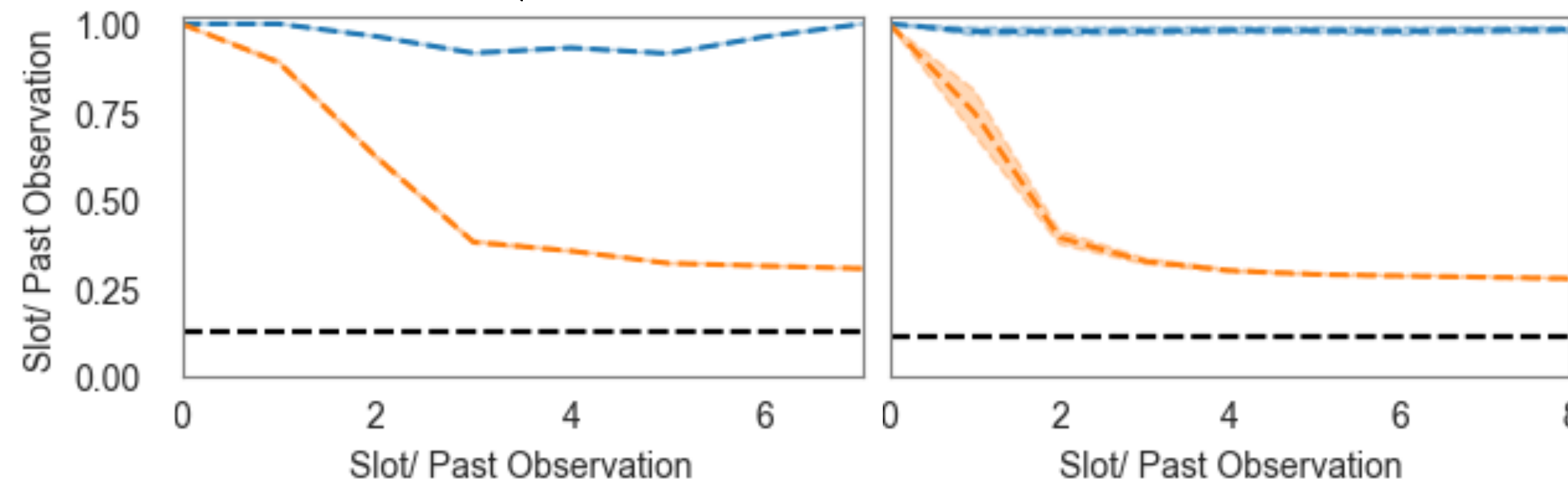
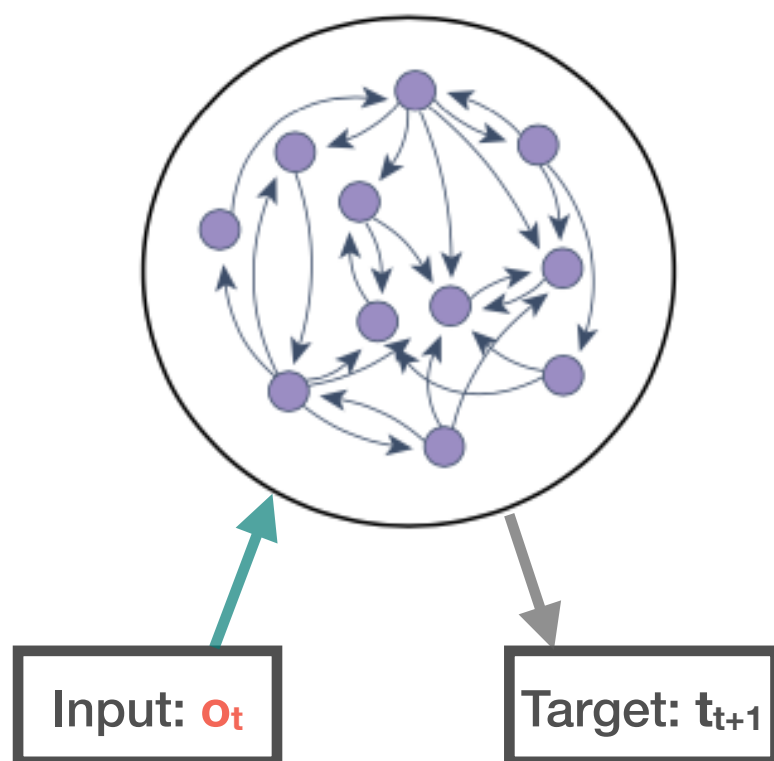
### 2D Navigation



### 2D slot prediction



Should be able to decode the predicted contents of each slot at every timestep



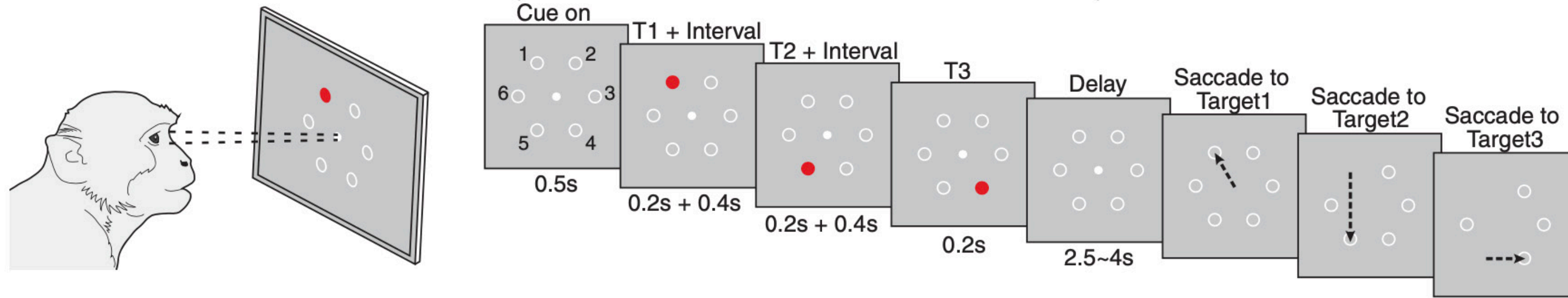
**Does this match mechanistic interpretation match PFC data?**

# A slot based understanding unifies many representations in prefrontal cortex

Xie et al., 2022

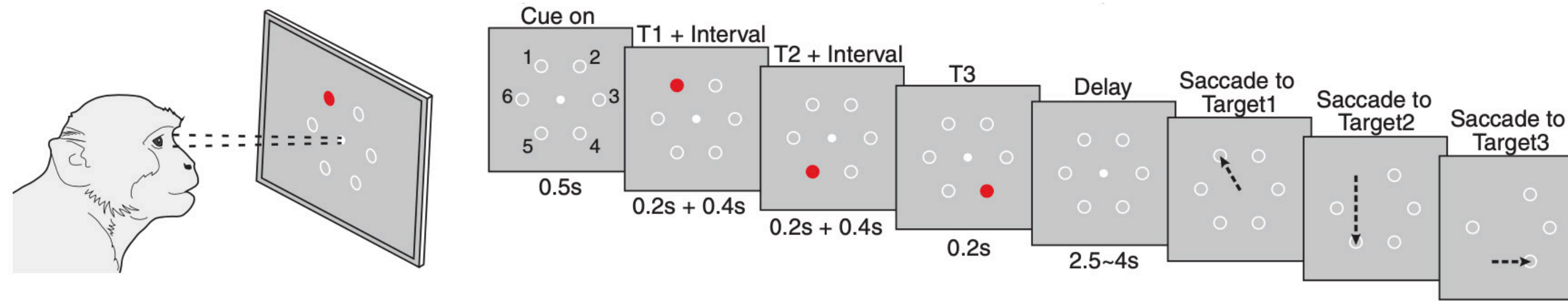
Whittington et al., 2023, *bioRxiv*

# A slot based understanding unifies many representations in prefrontal cortex

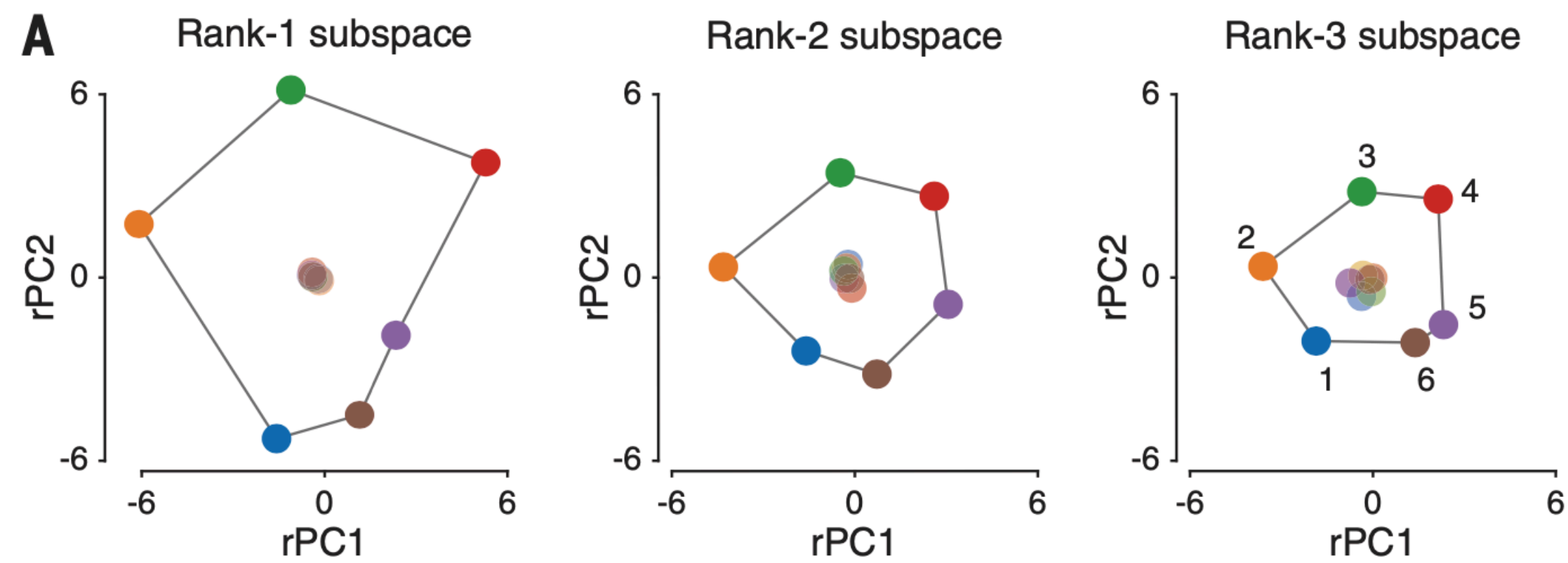


Xie et al., 2022

# A slot based understanding unifies many representations in prefrontal cortex

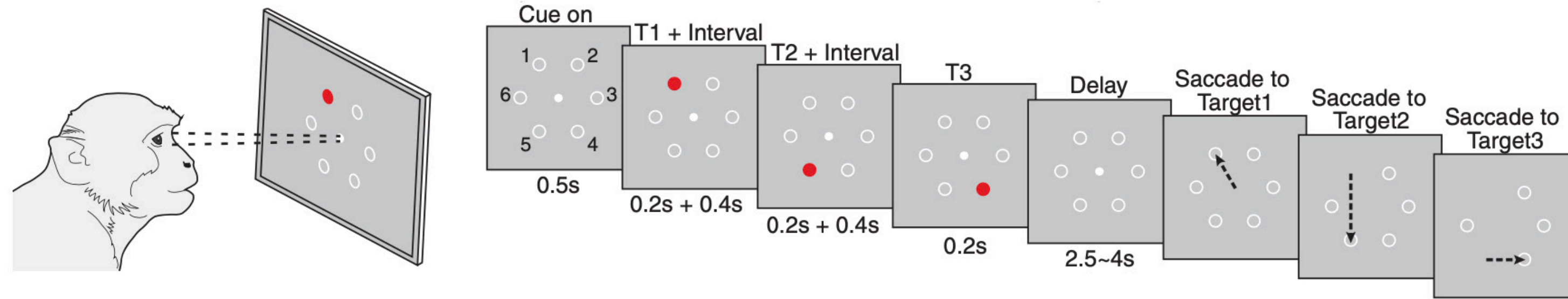


Xie et al., 2022

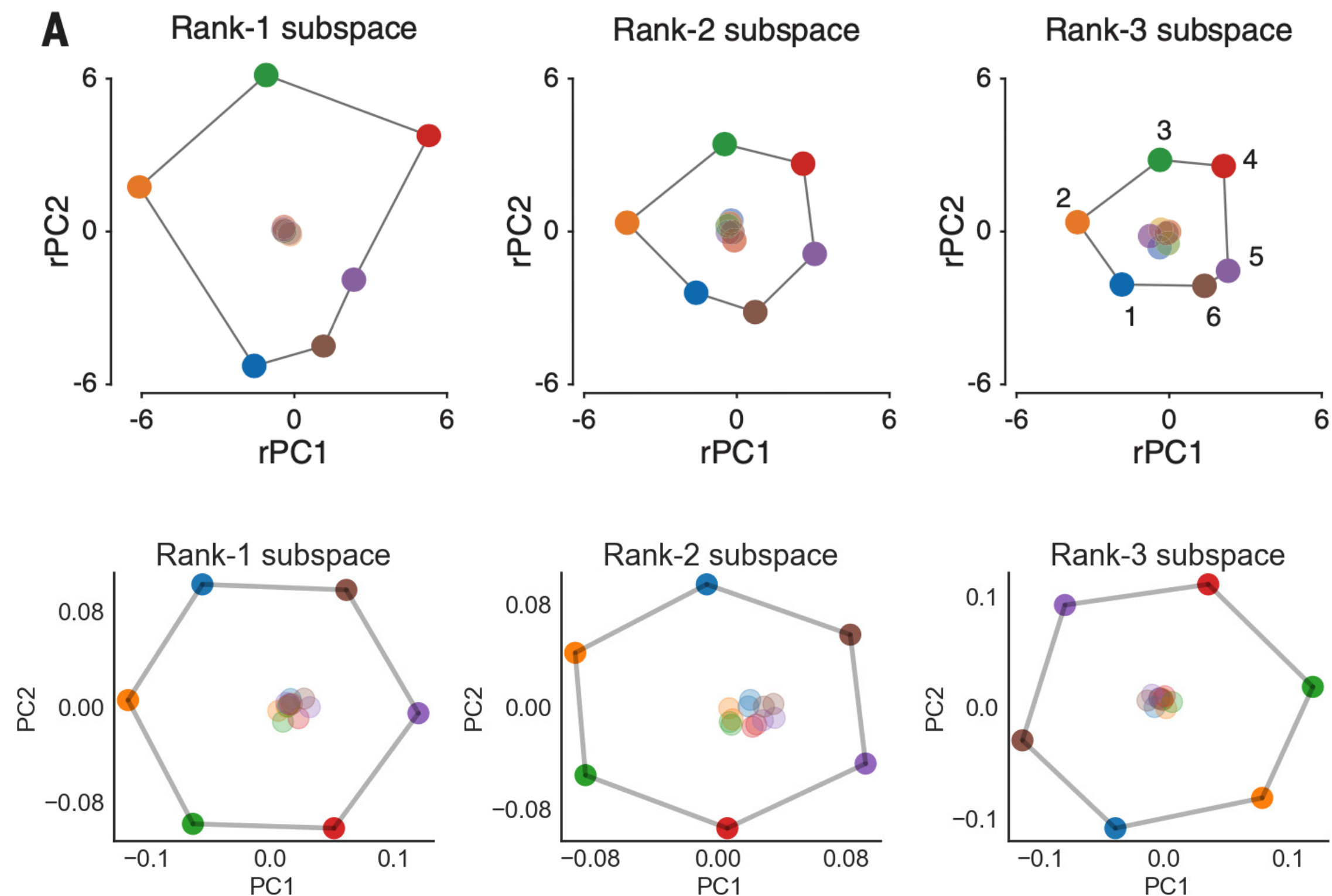




# A slot based understanding unifies many representations in prefrontal cortex

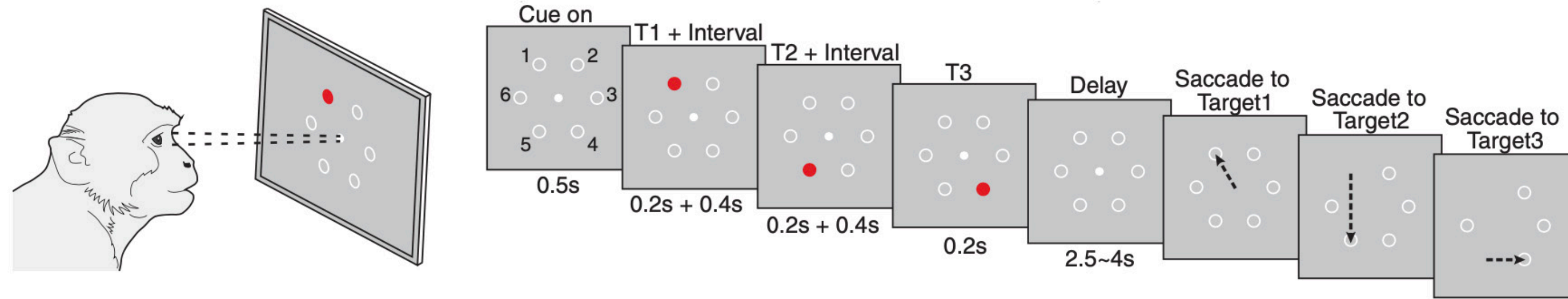


Xie et al., 2022

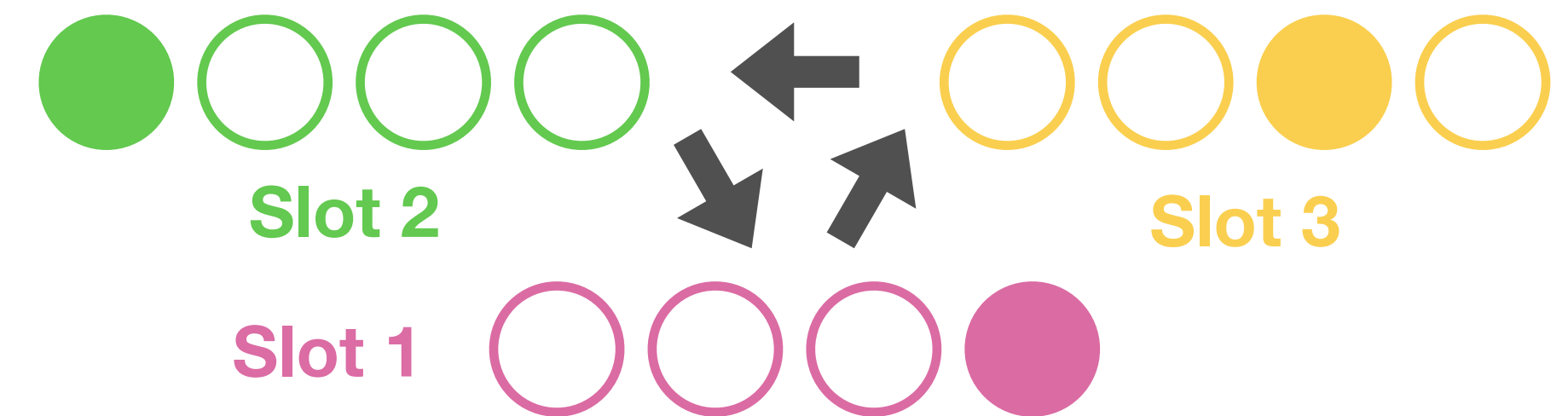
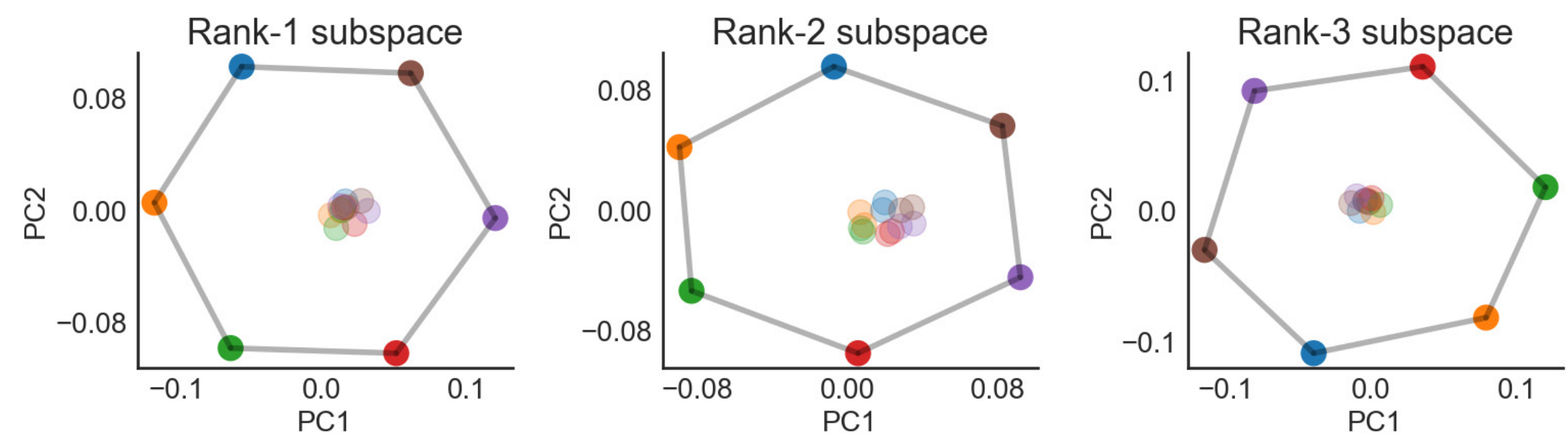
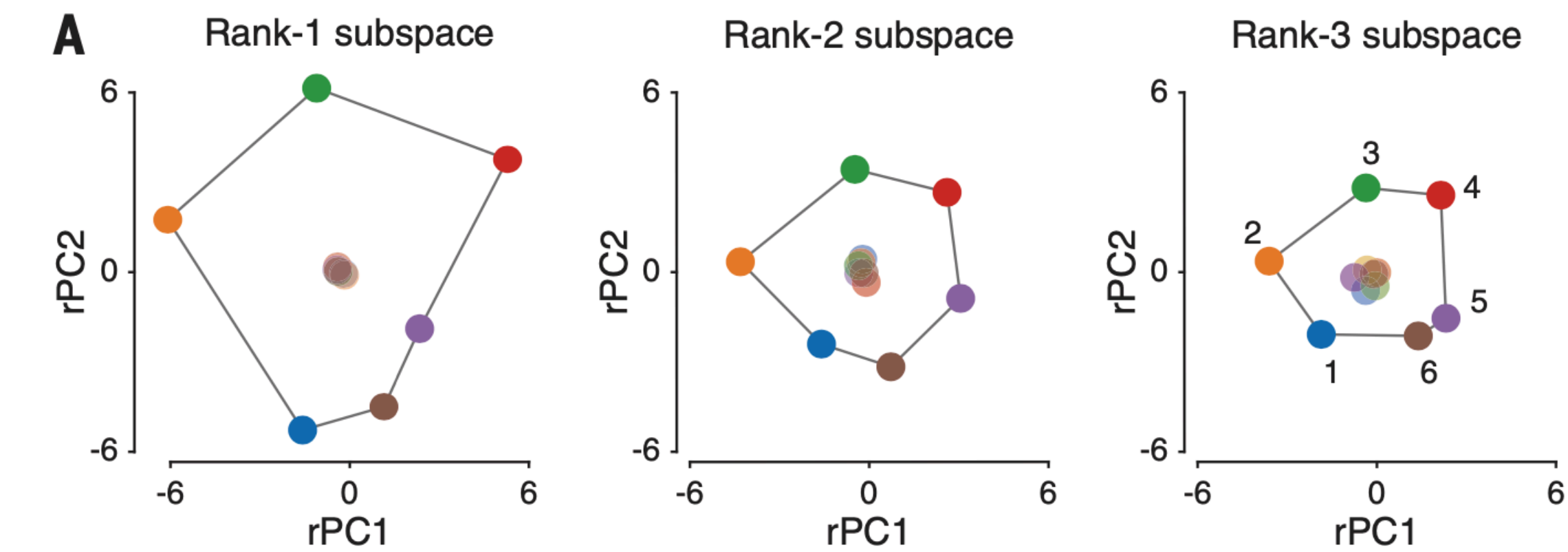


Whittington et al., 2023, *bioRxiv*

# A slot based understanding unifies many representations in prefrontal cortex



Xie et al., 2022



Whittington et al., 2023, *bioRxiv*

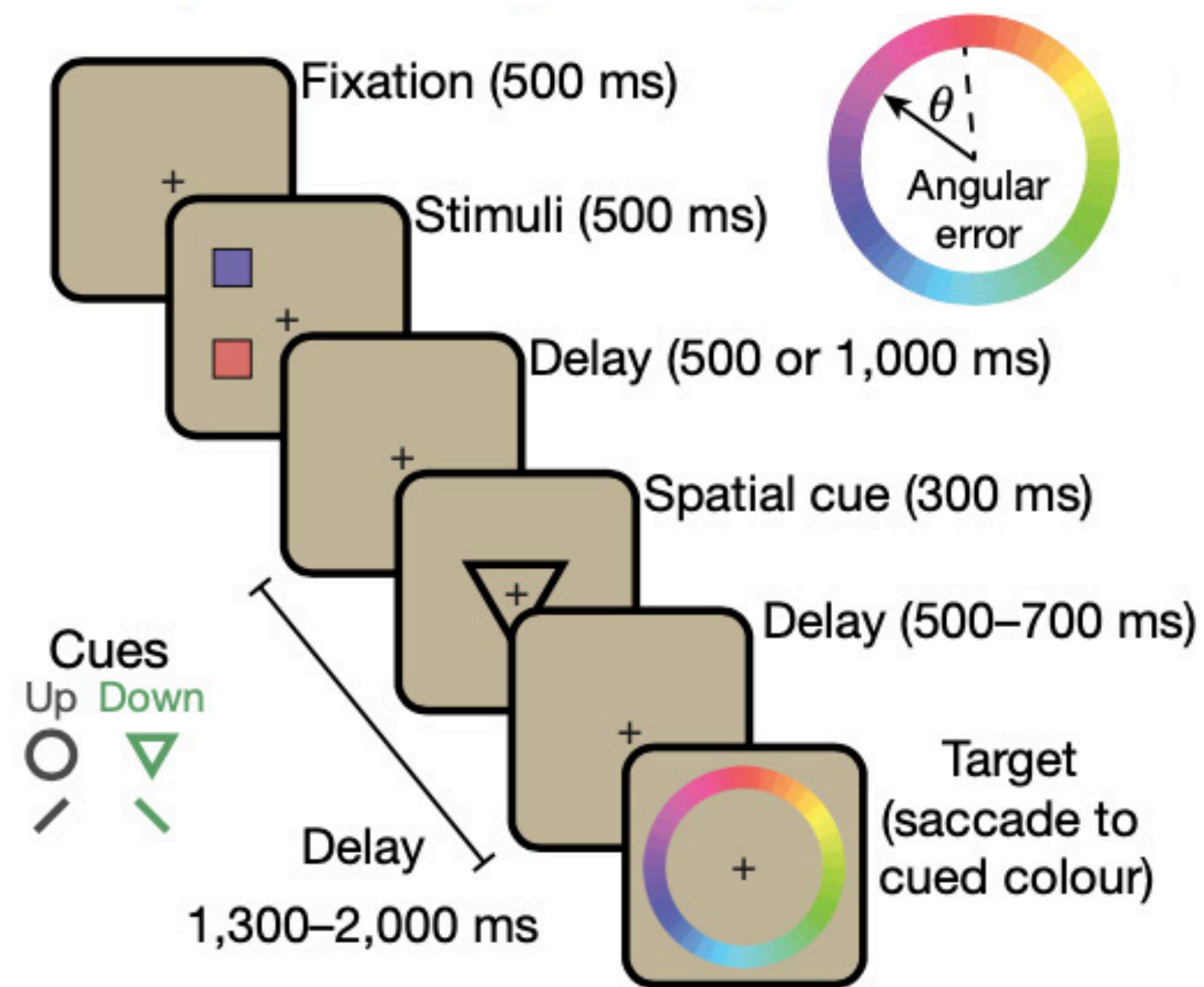
# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021

Whittington et al., 2023, *bioRxiv*

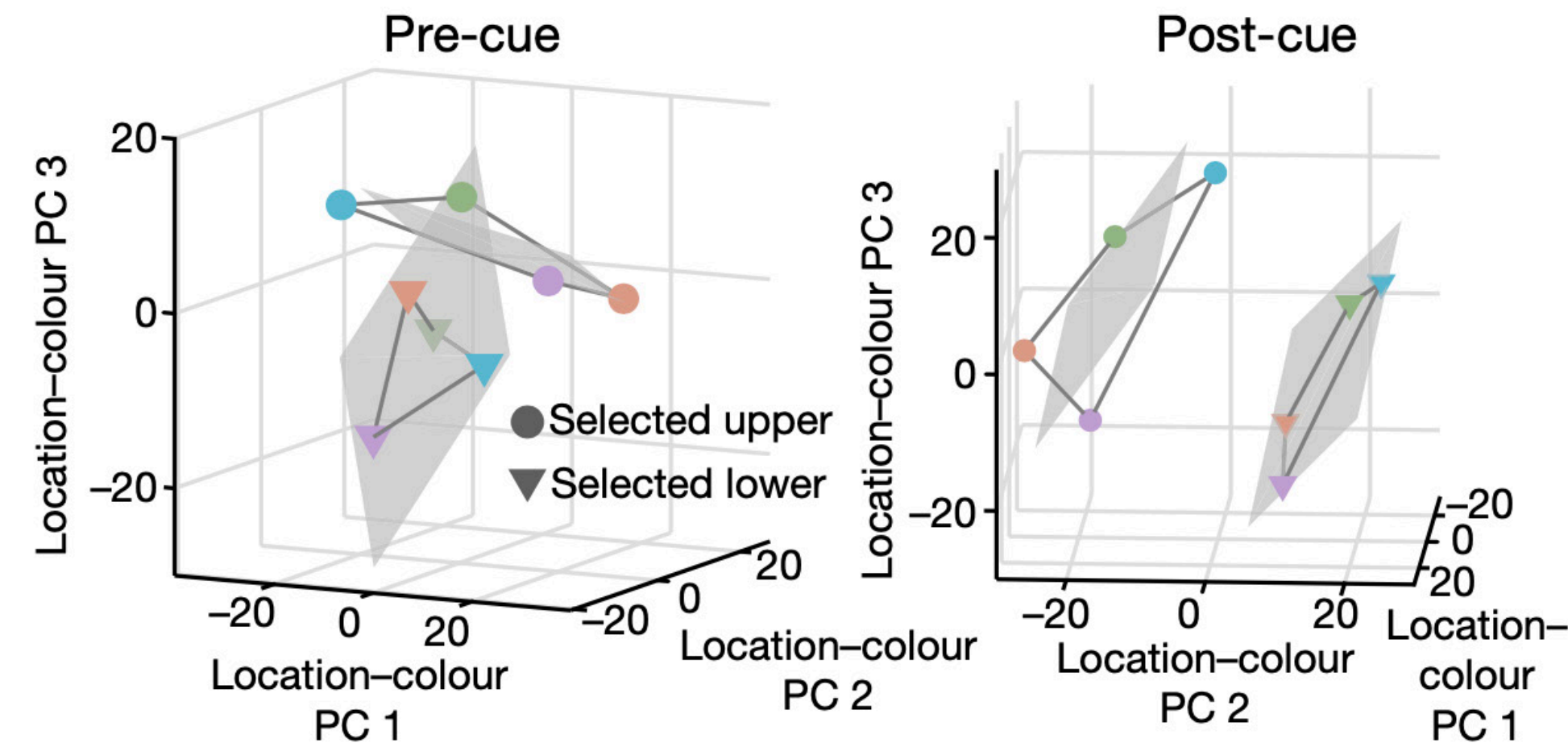
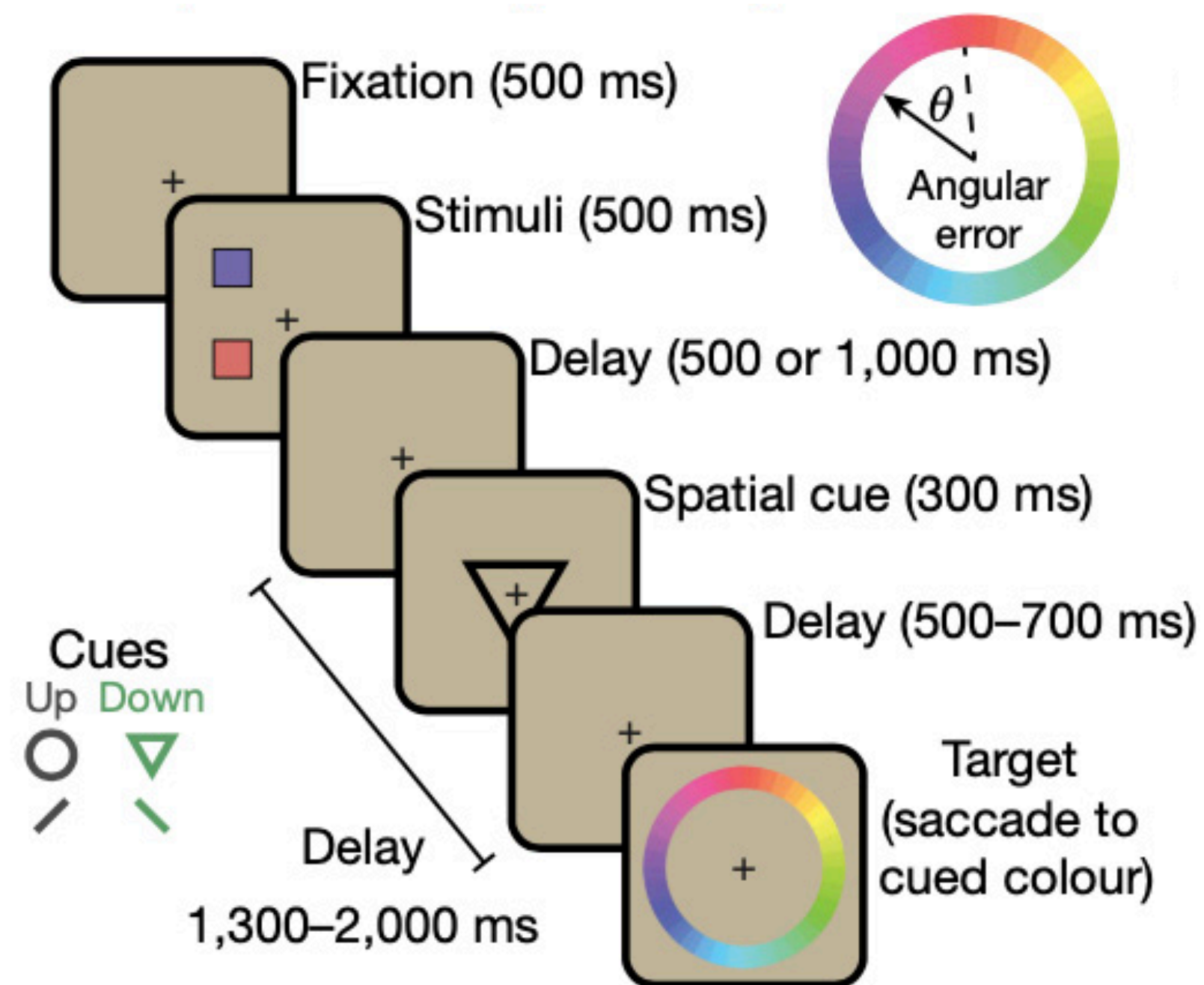
# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021



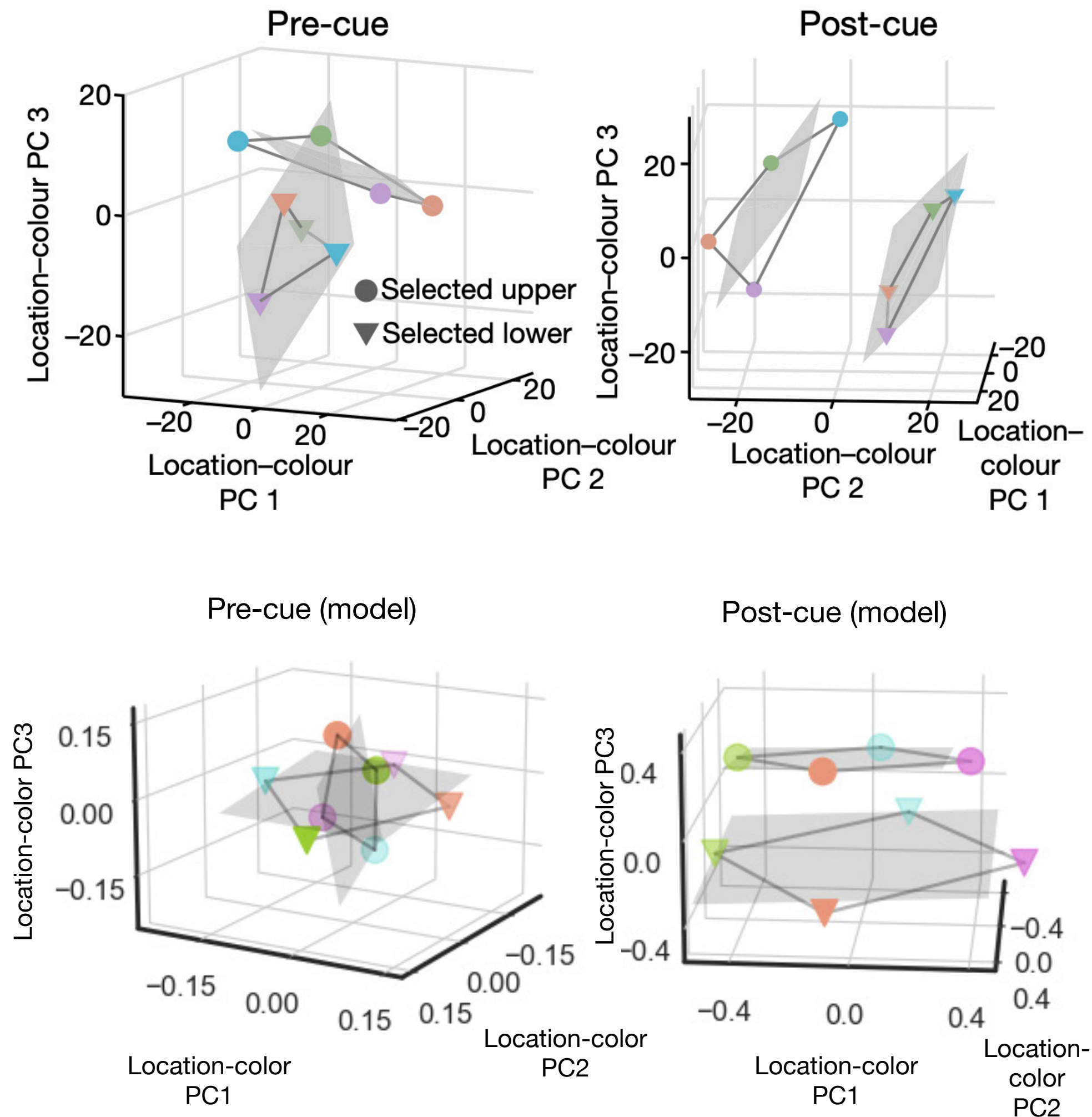
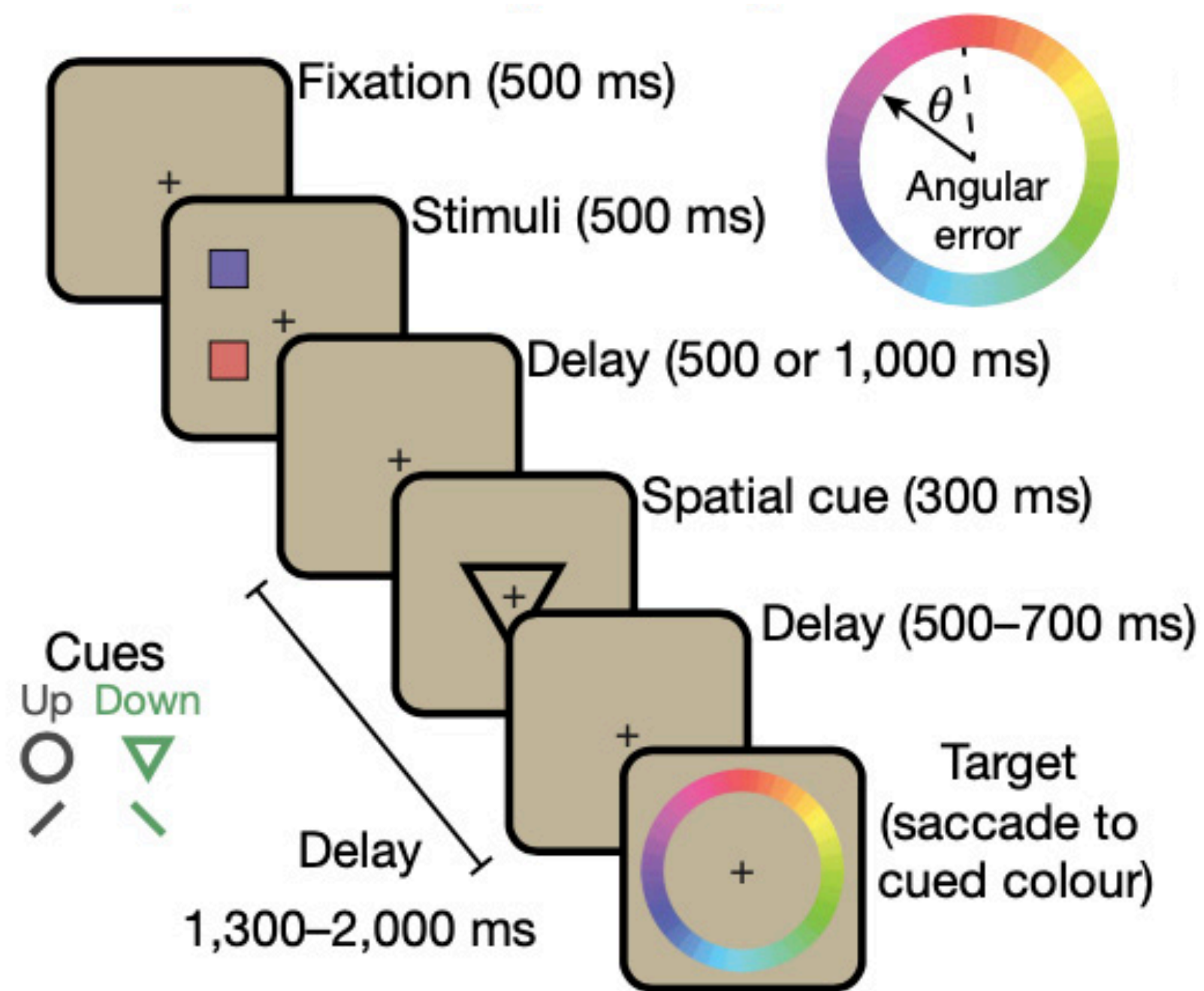
# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021



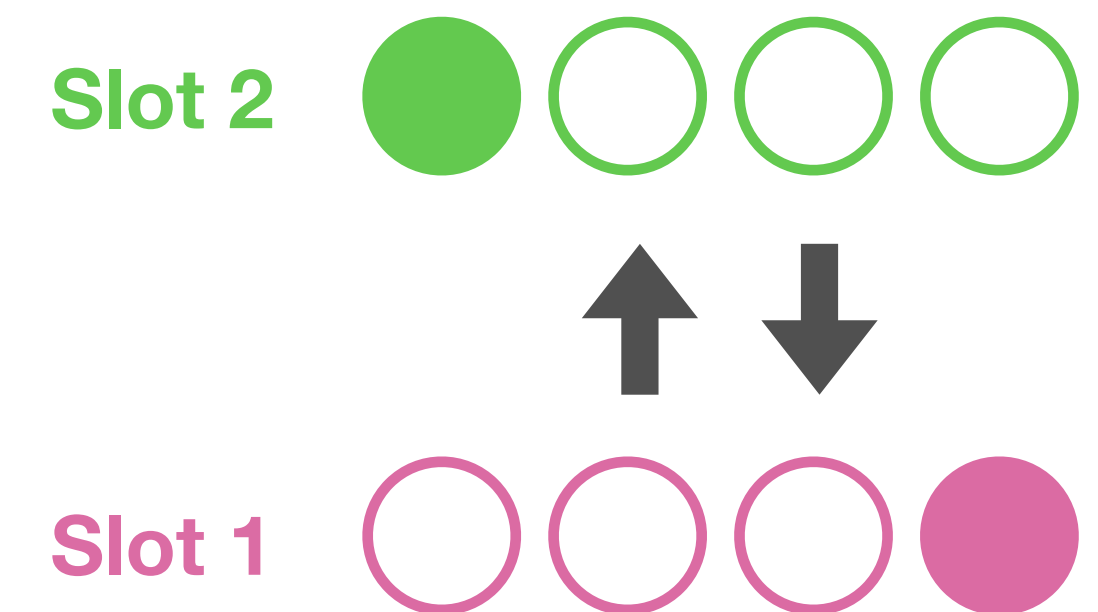
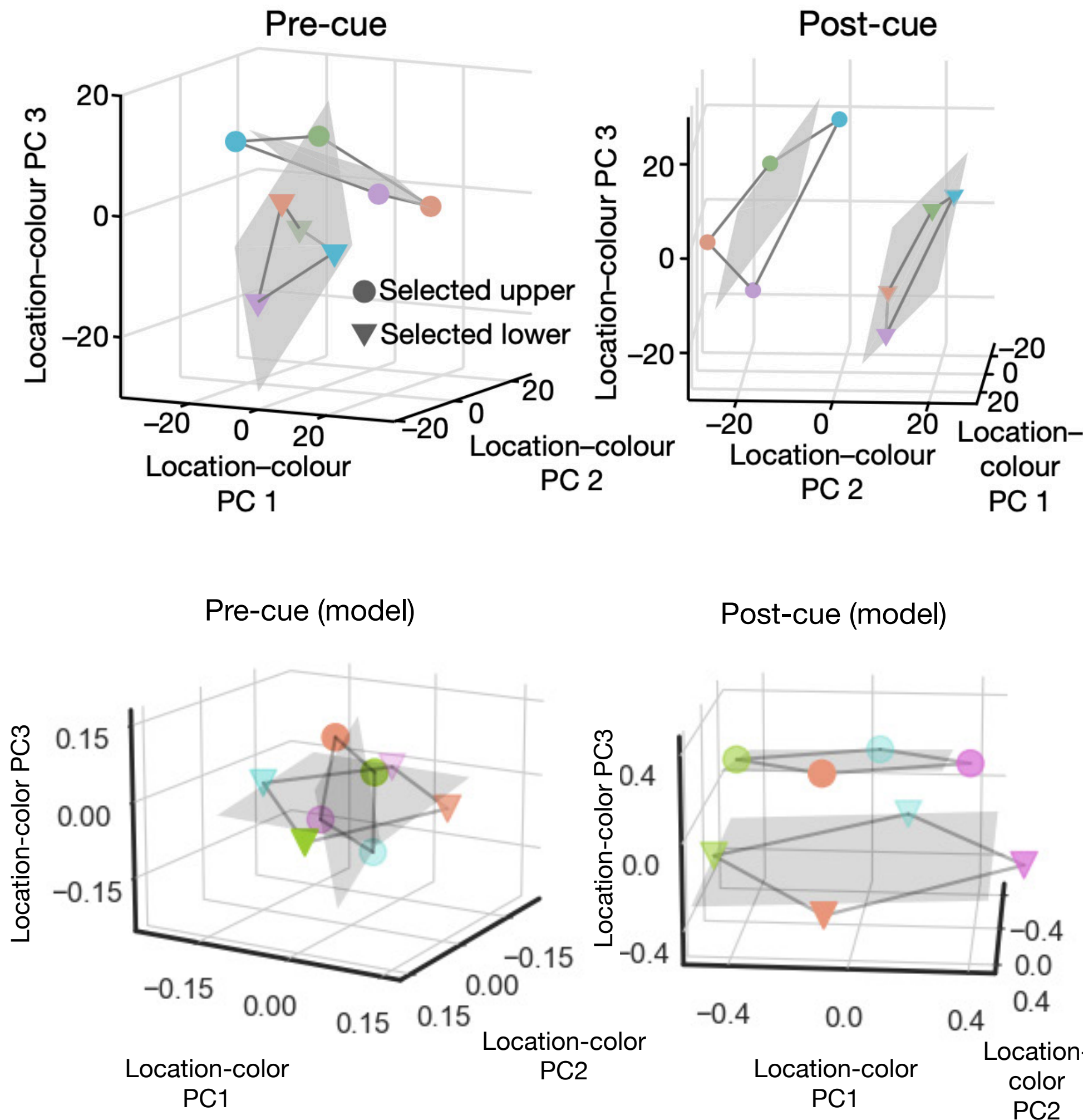
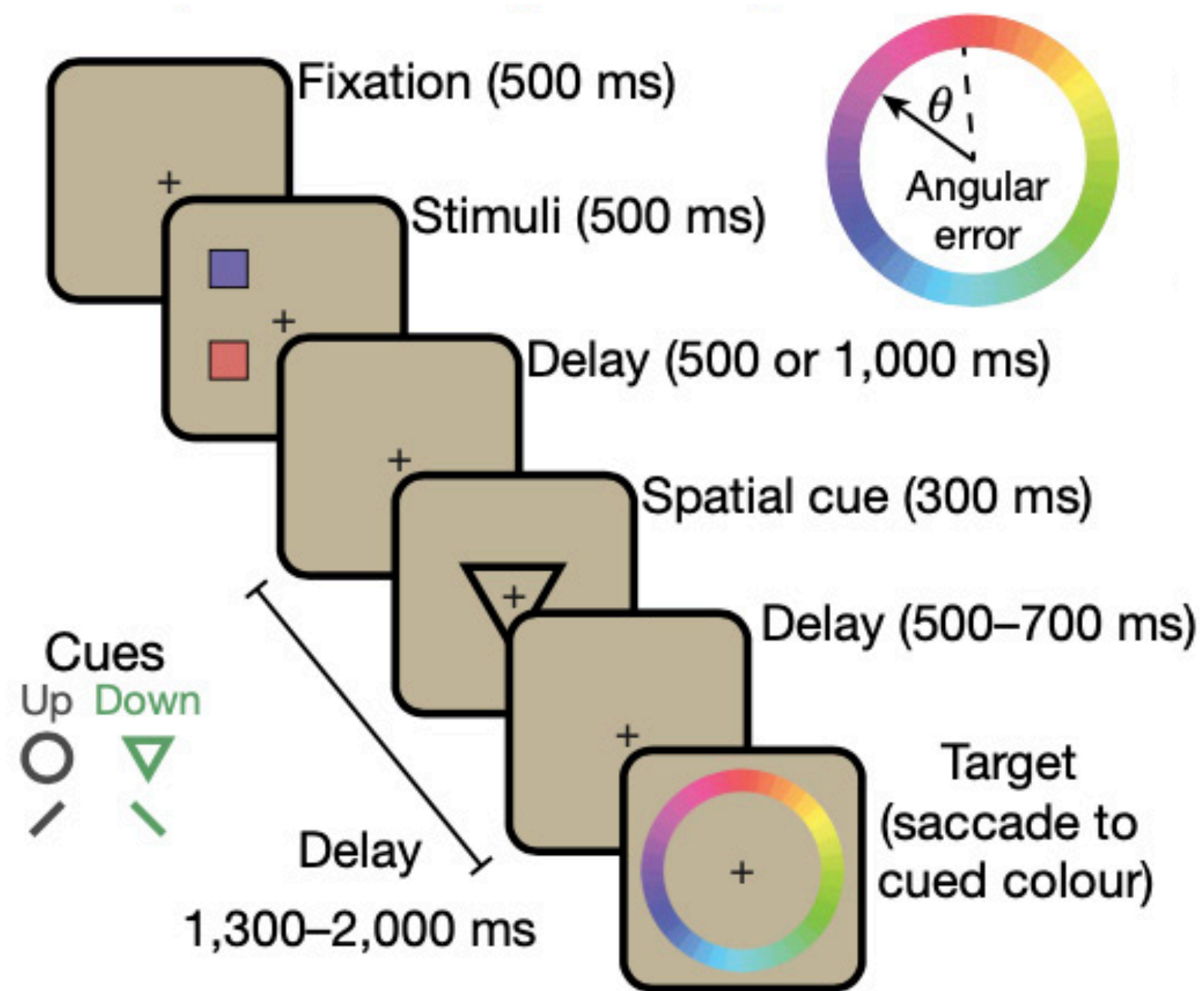
# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021



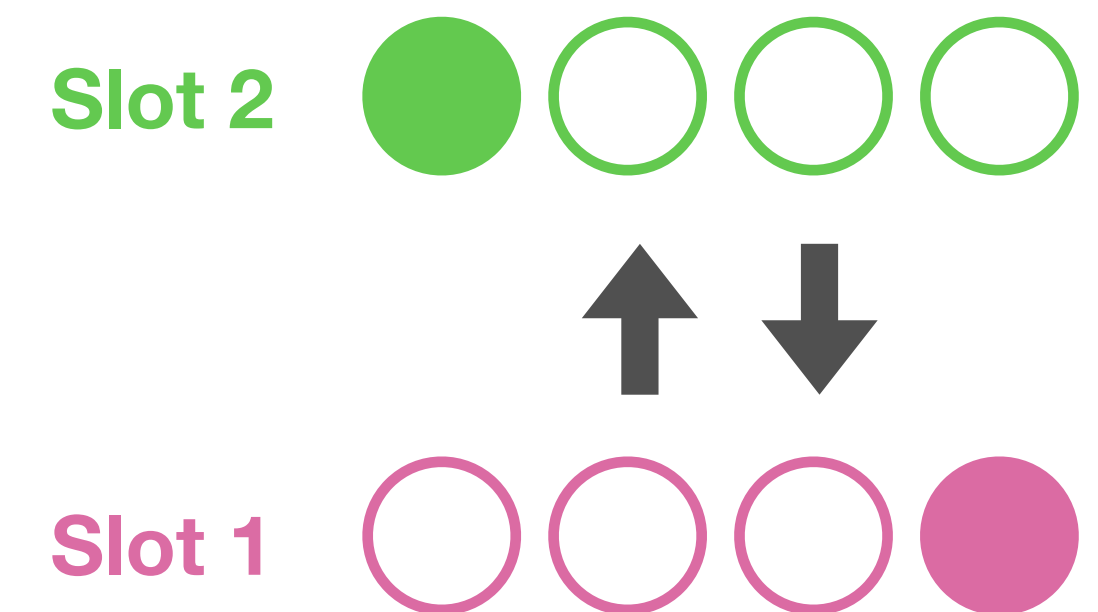
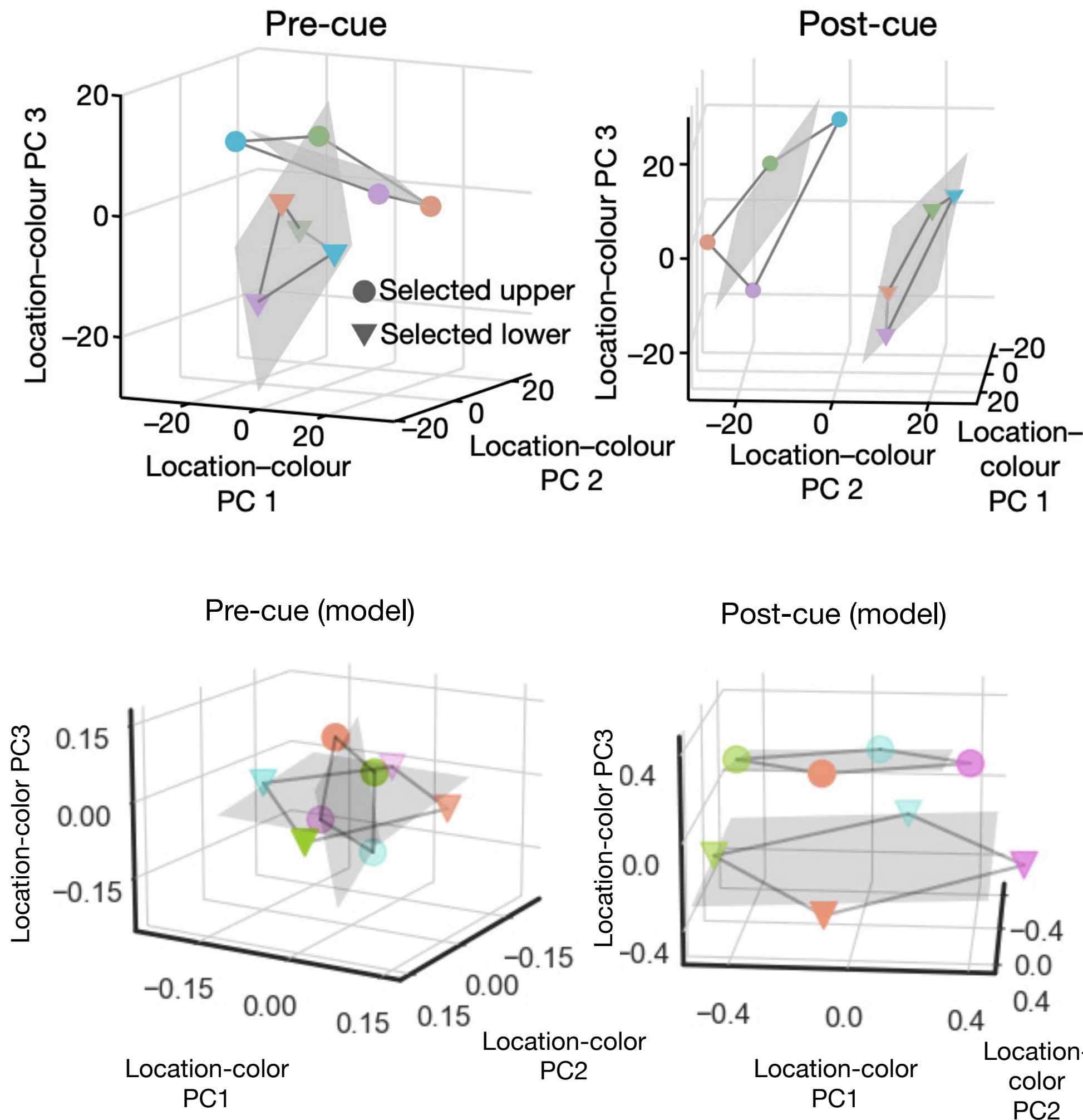
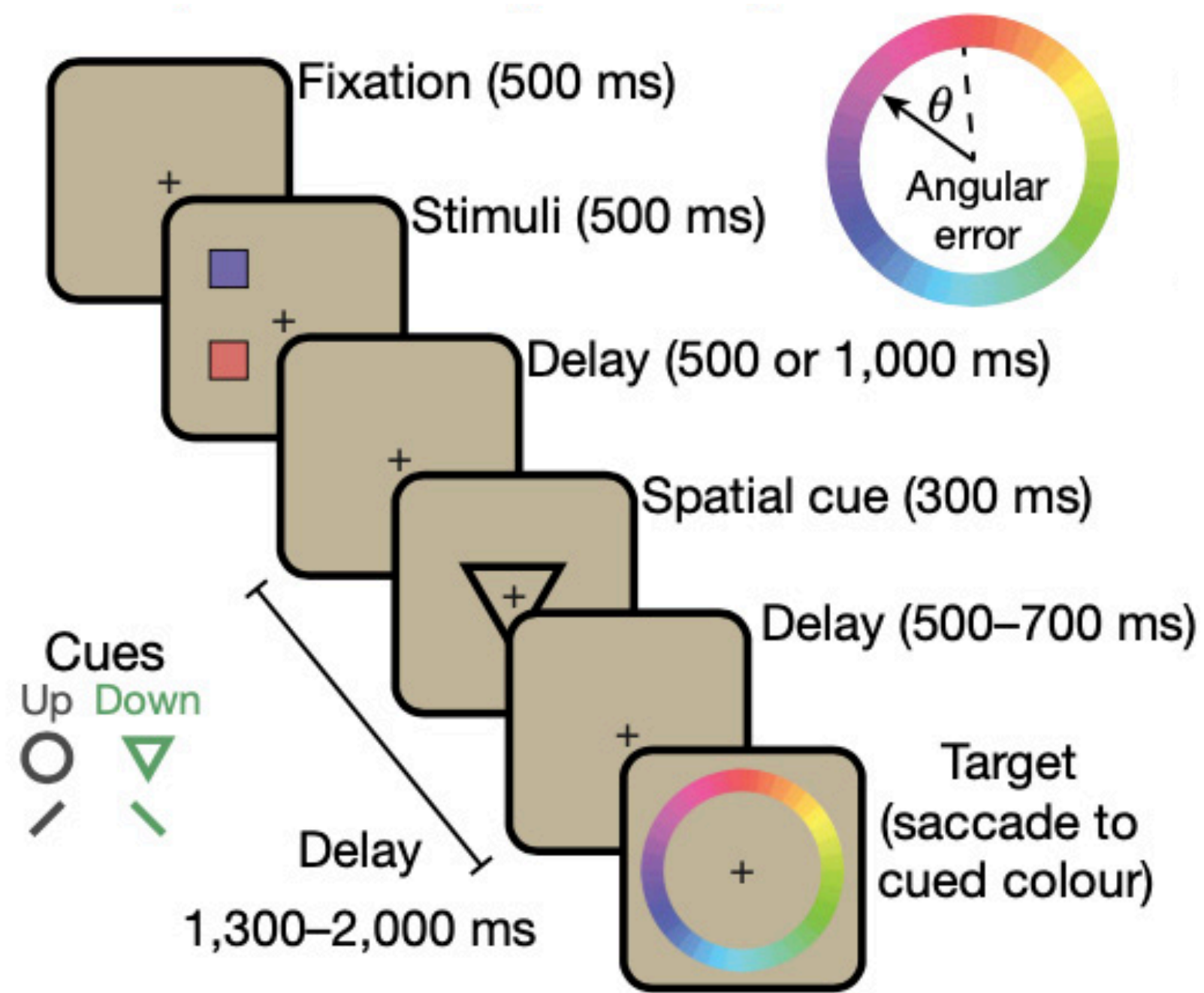
# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021



# A slot based understanding unifies many representations in prefrontal cortex

Panichello & Buschman 2021

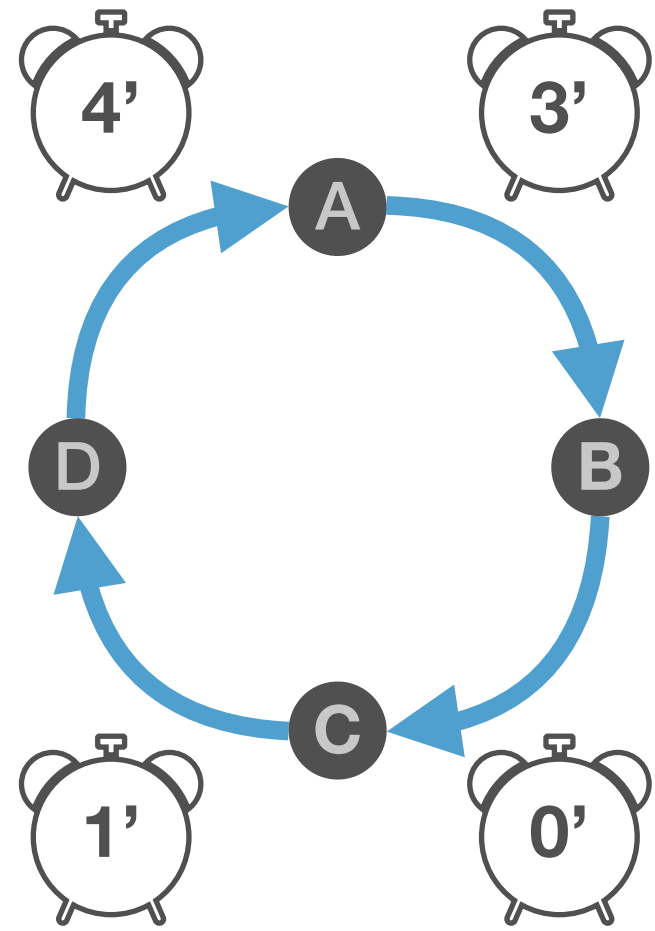


The cue behaves like an action signal!

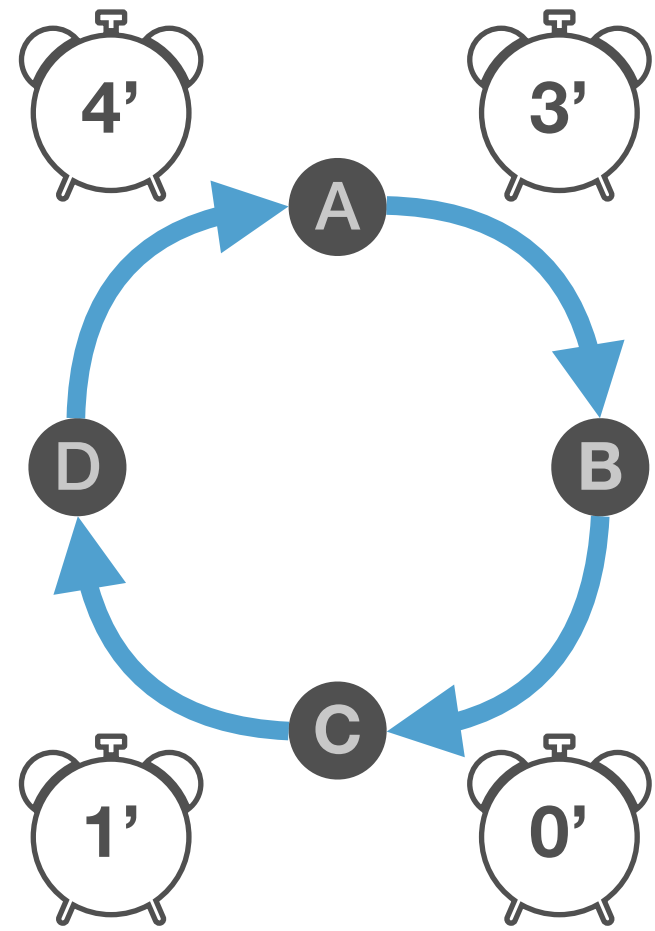


**Makes novel predictions we tested in prefrontal data**

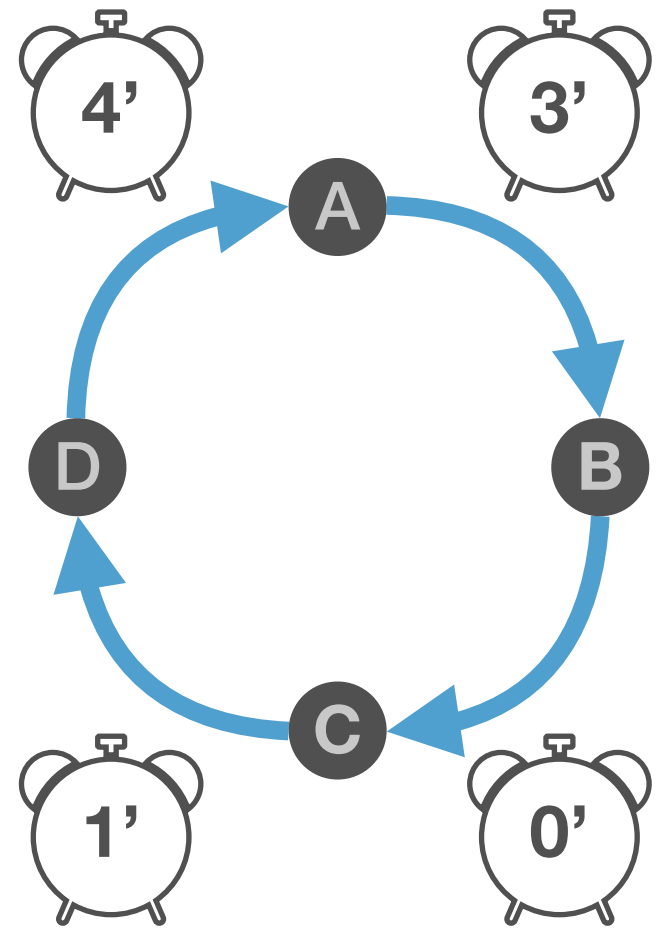
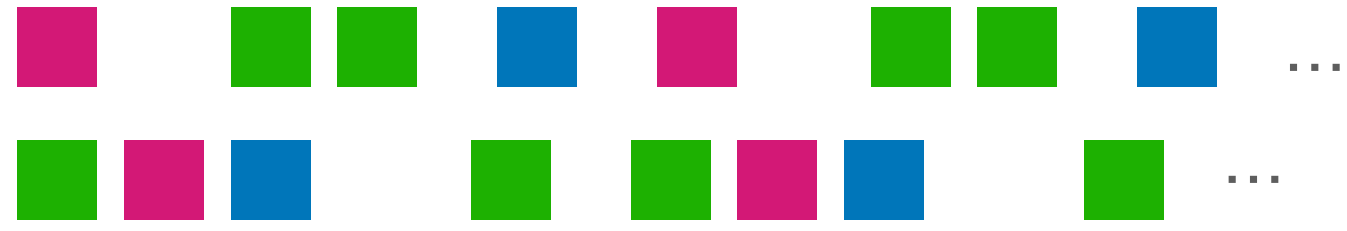
# Makes novel predictions we tested in prefrontal data



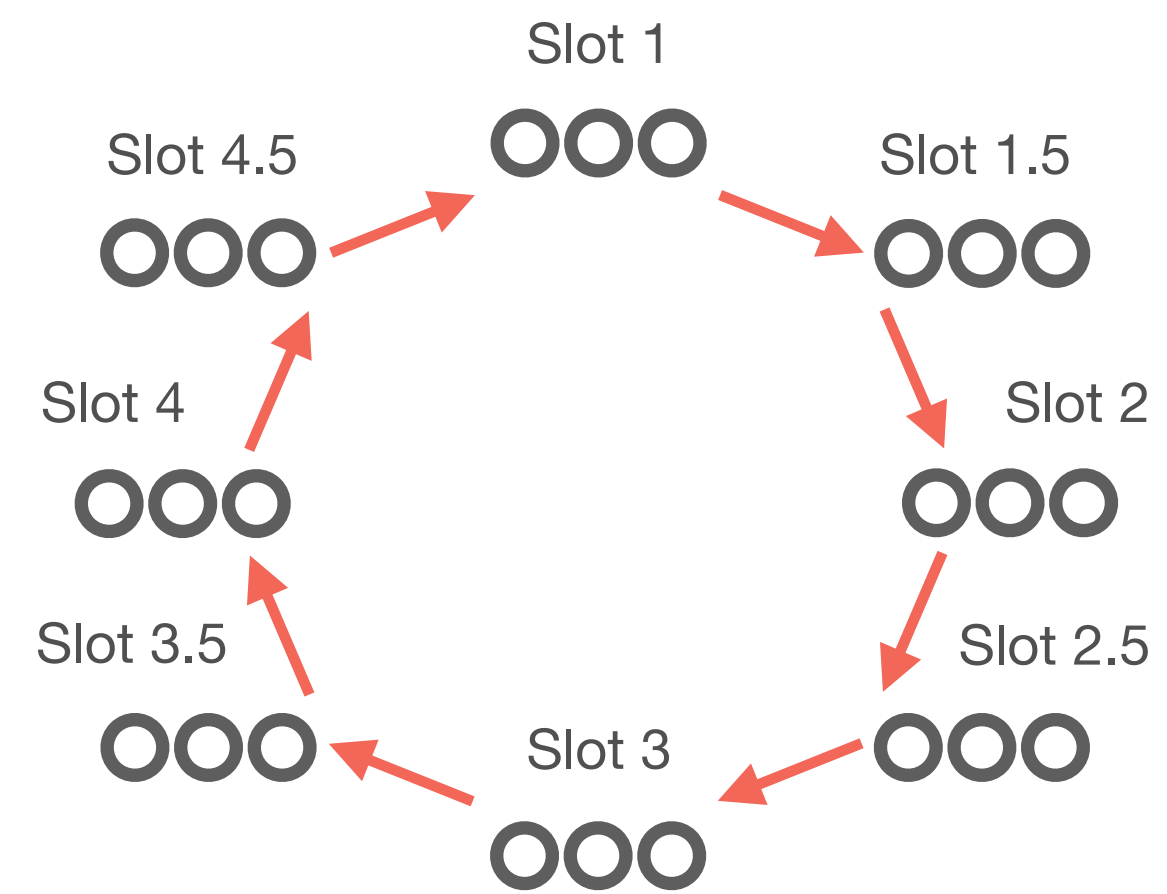
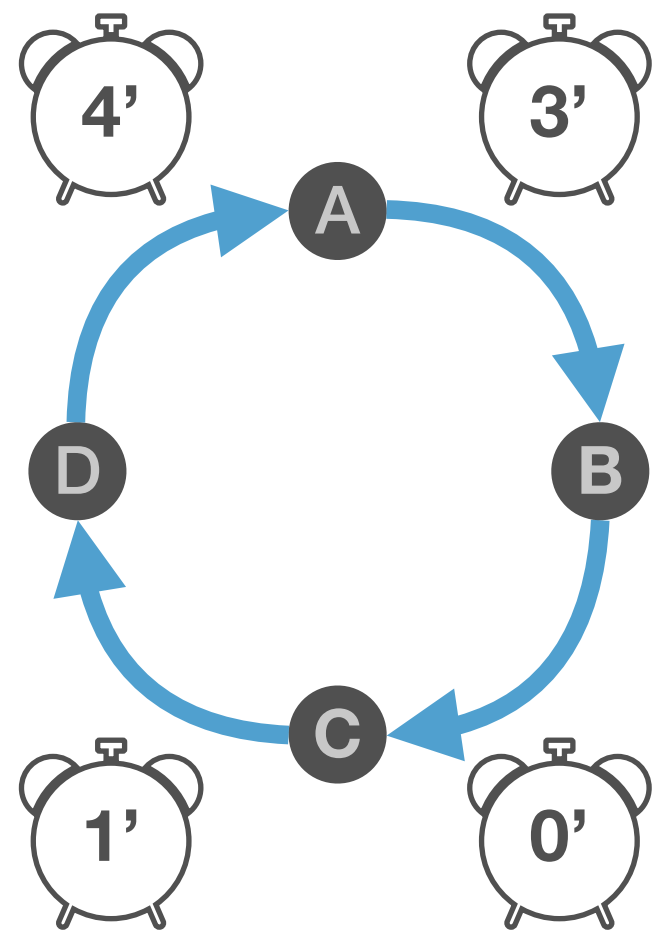
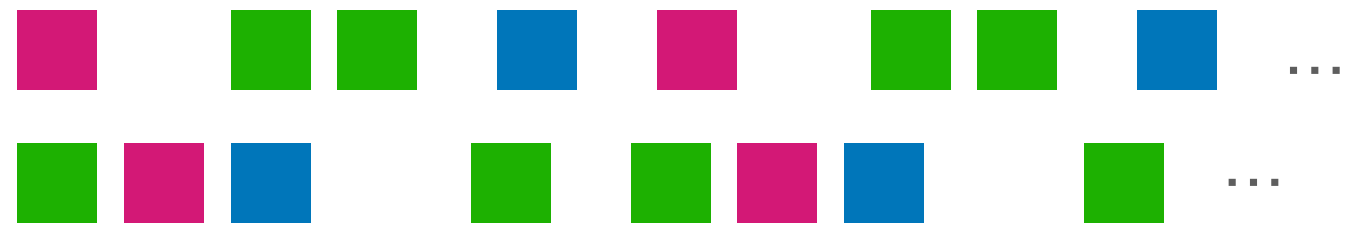
# Makes novel predictions we tested in prefrontal data



# Makes novel predictions we tested in prefrontal data



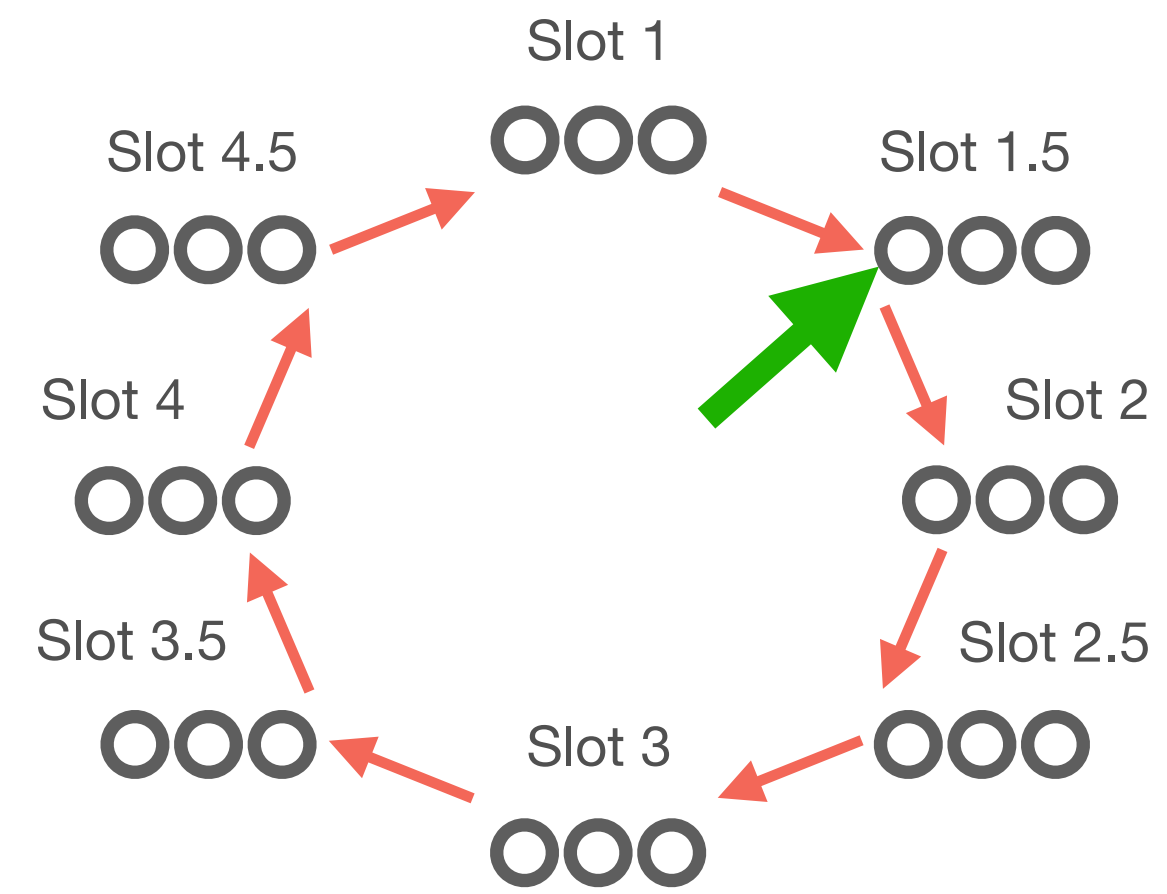
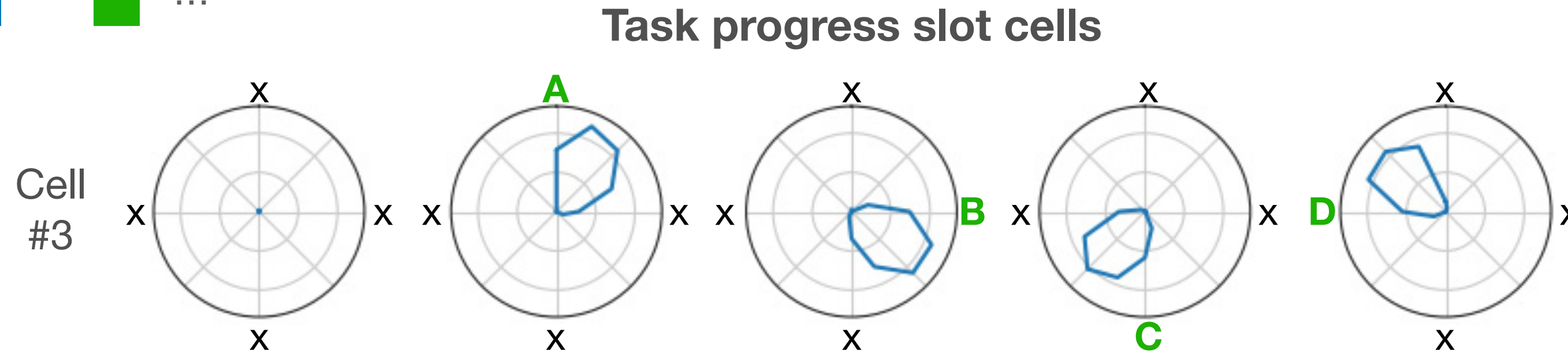
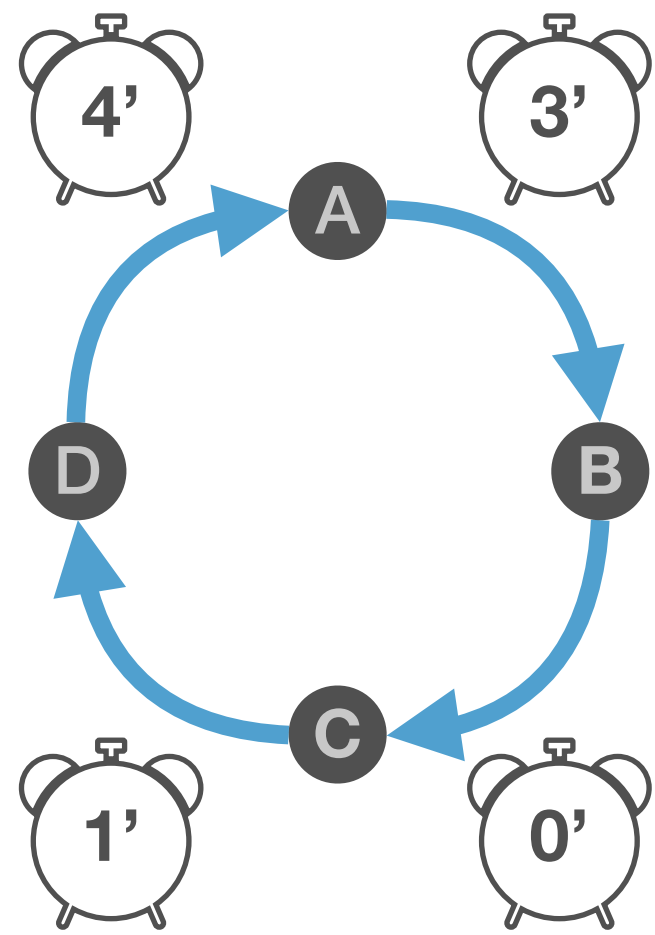
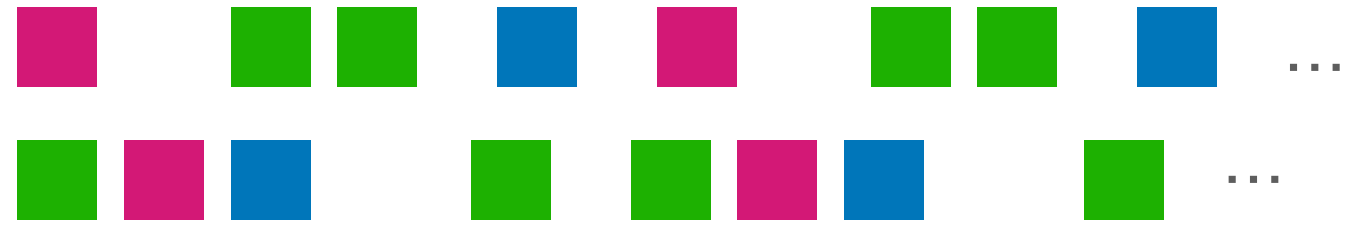
# Makes novel predictions we tested in prefrontal data



Whittington et al., 2023, *bioRxiv*

El-Gaby et al., 2023, *bioRxiv*

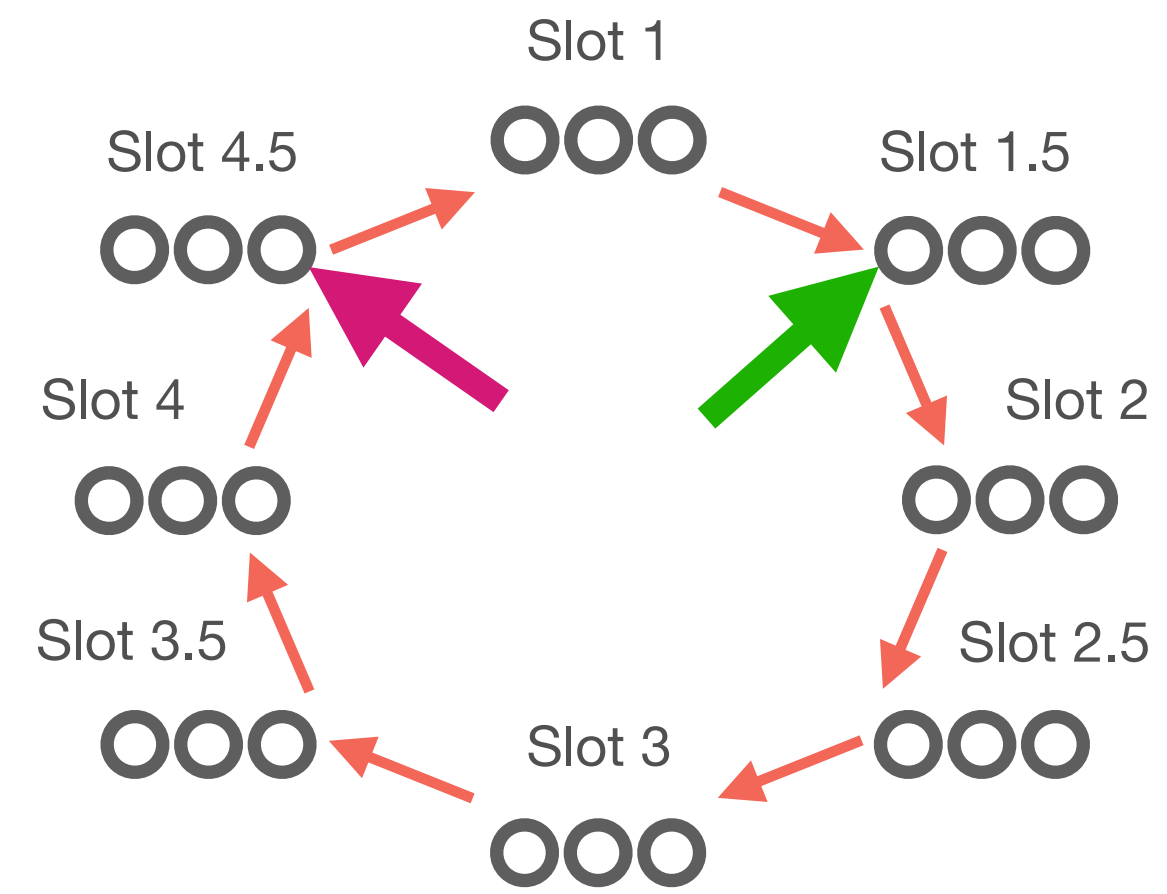
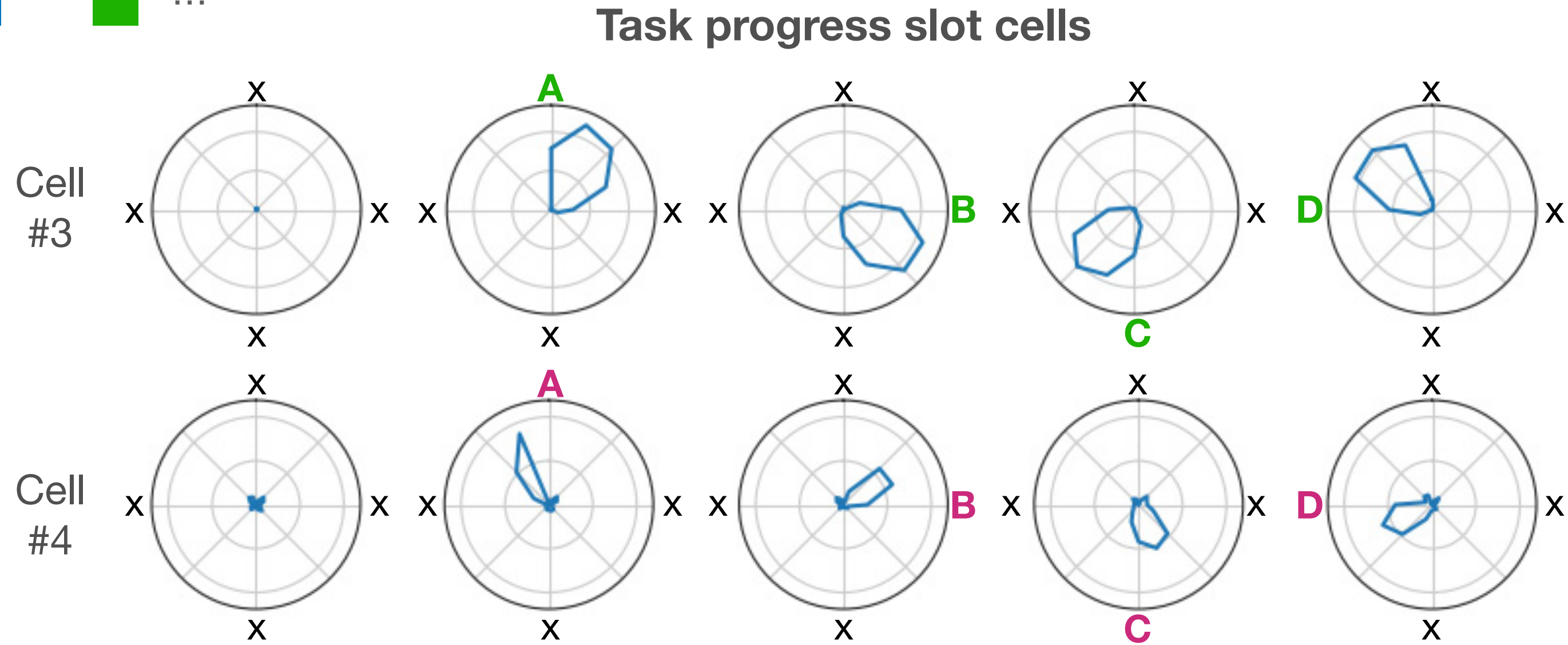
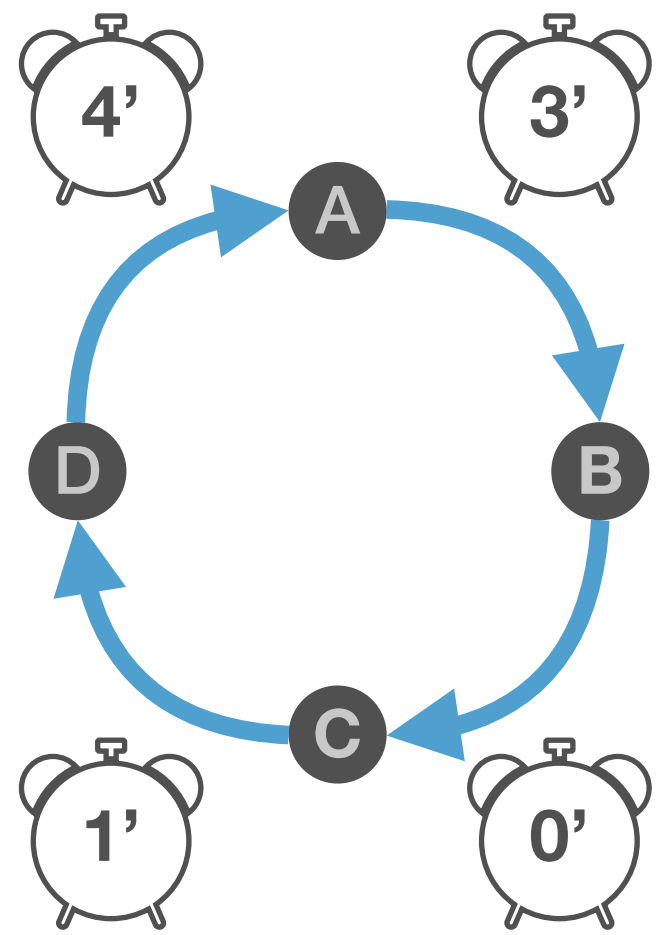
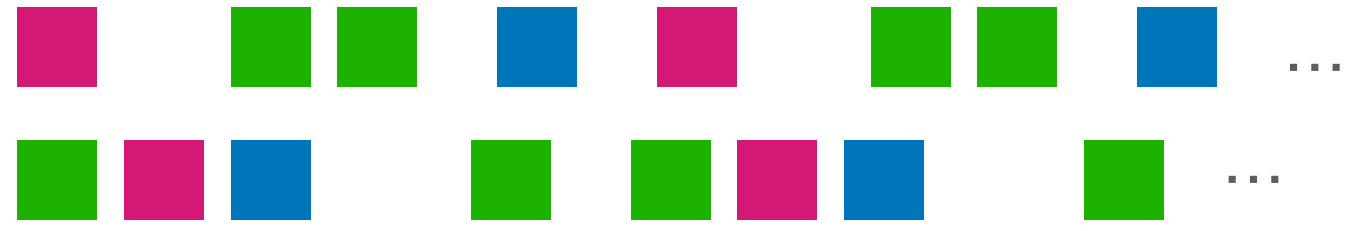
# Makes novel predictions we tested in prefrontal data



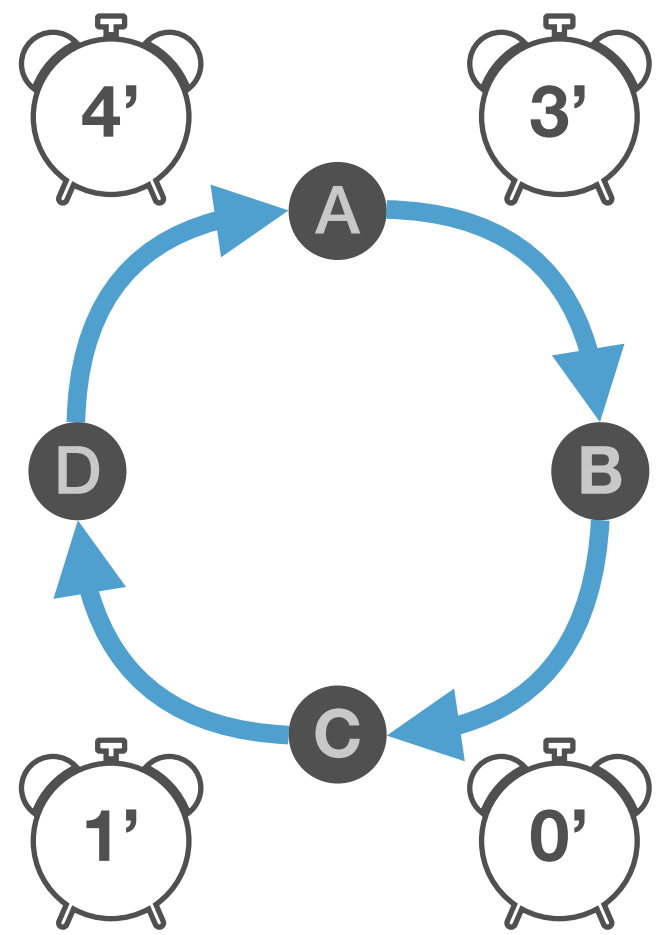
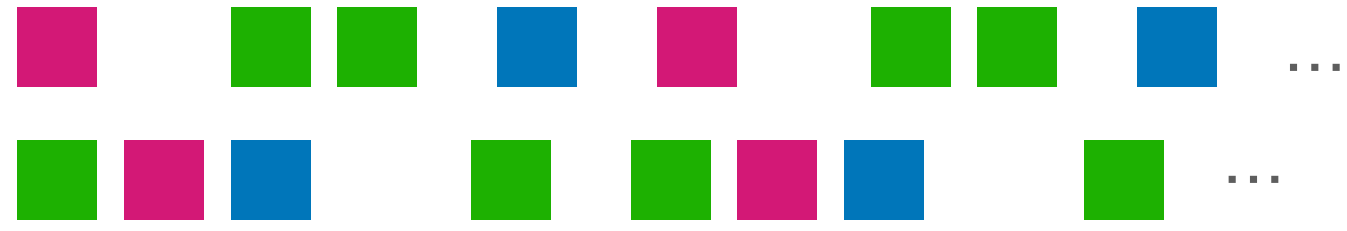
Whittington et al., 2023, *bioRxiv*

El-Gaby et al., 2023, *bioRxiv*

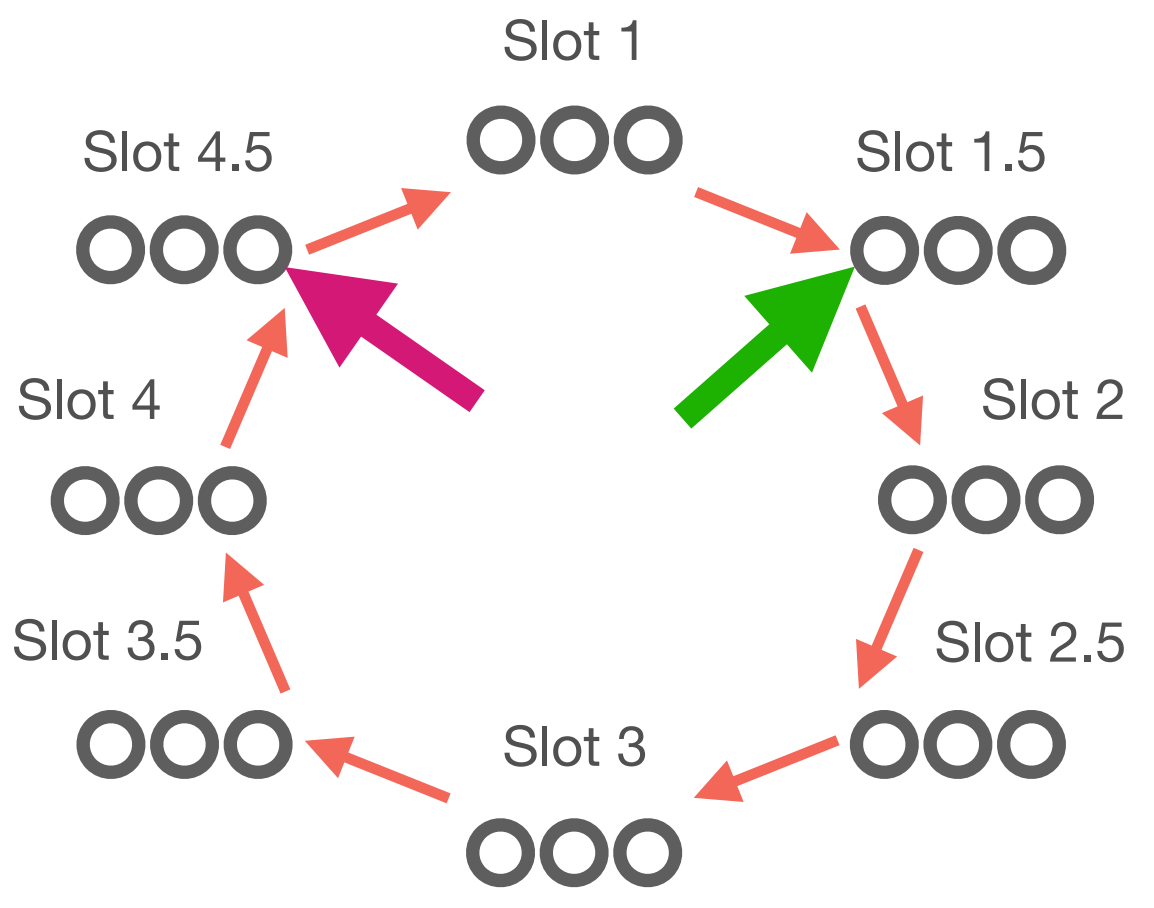
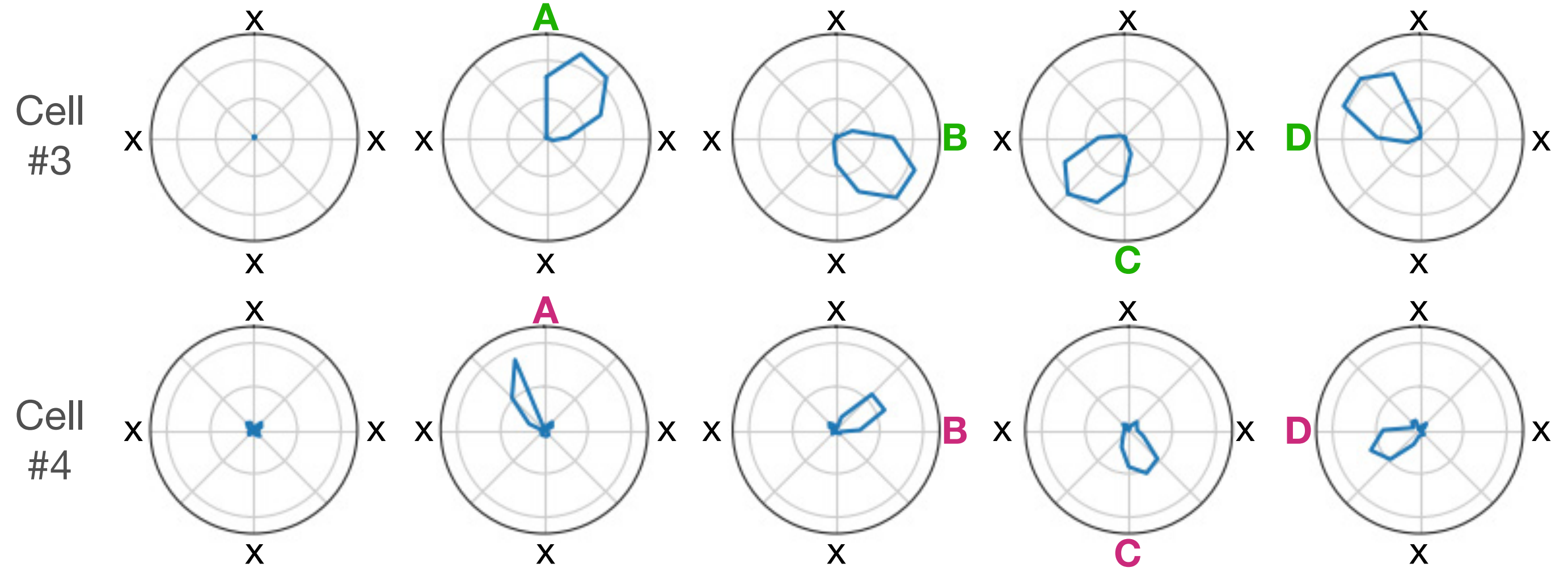
# Makes novel predictions we tested in prefrontal data



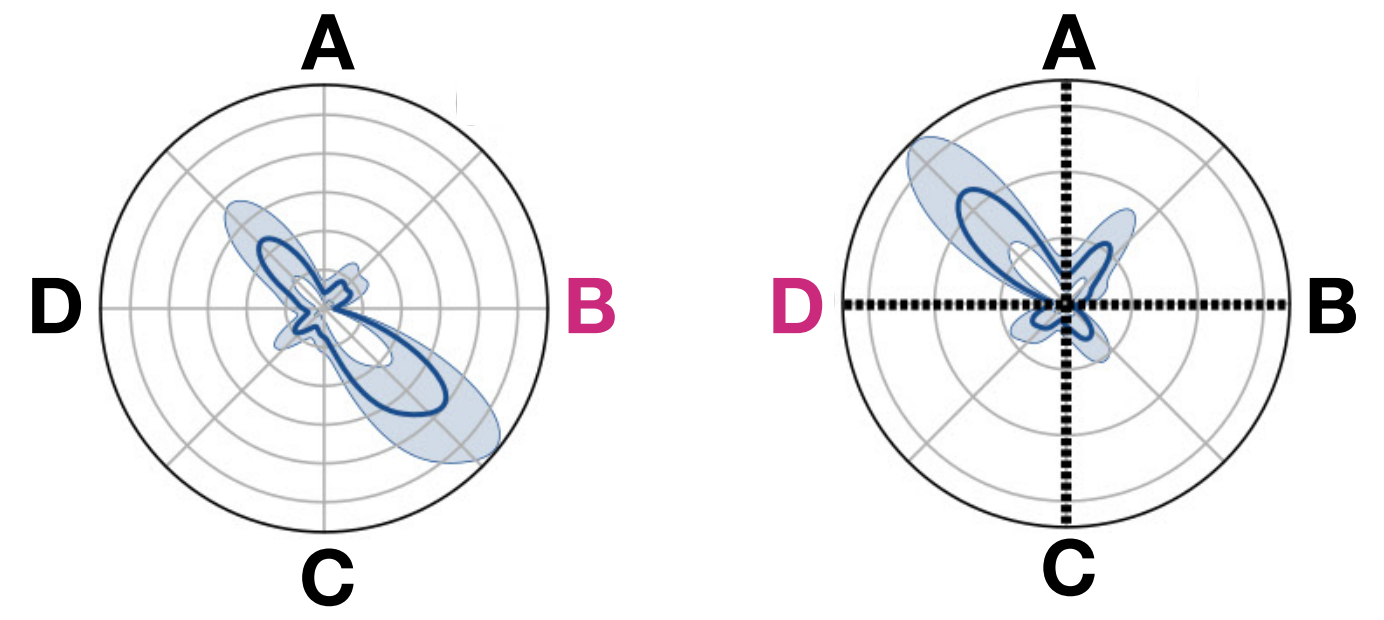
# Makes novel predictions we tested in prefrontal data



Task progress slot cells



PFC real 'purple 50% lag' neuron

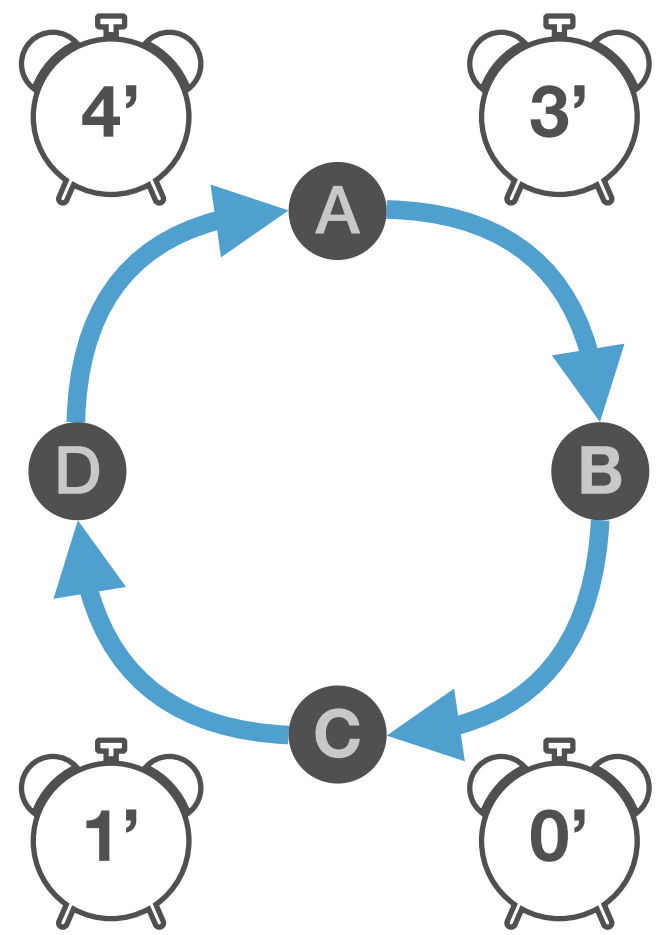
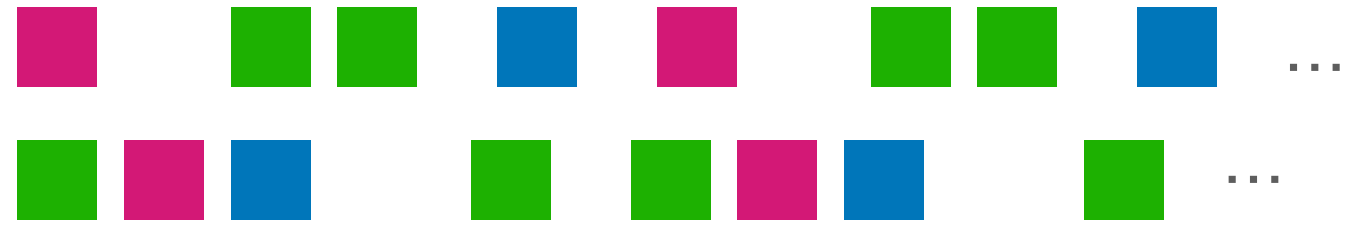


Whittington et al., 2023, *bioRxiv*

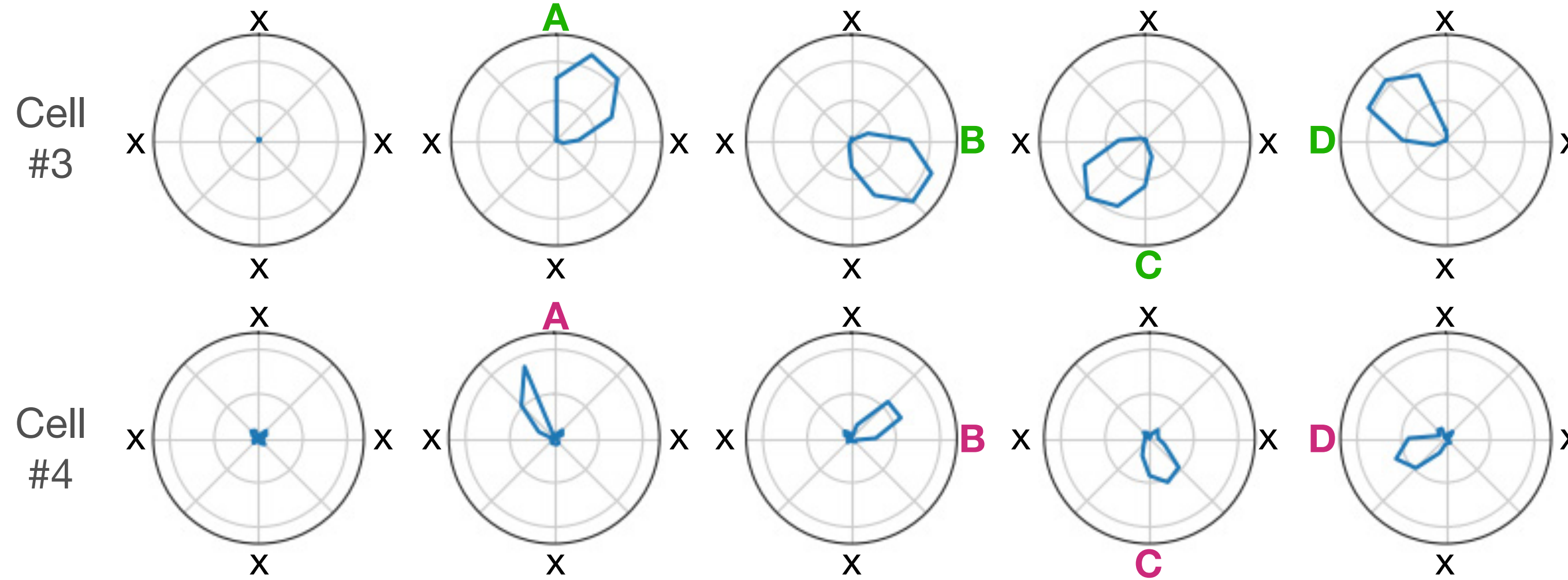
El-Gaby et al., 2023, *bioRxiv*



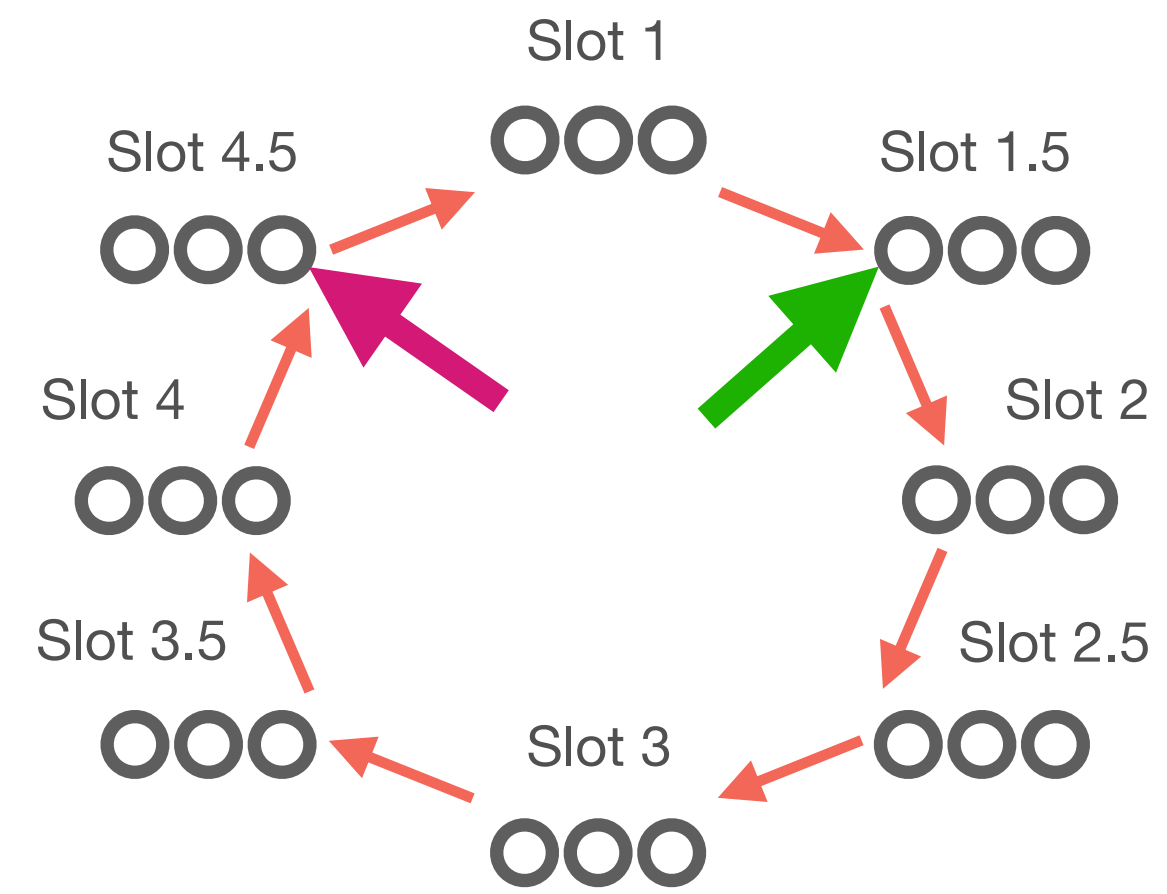
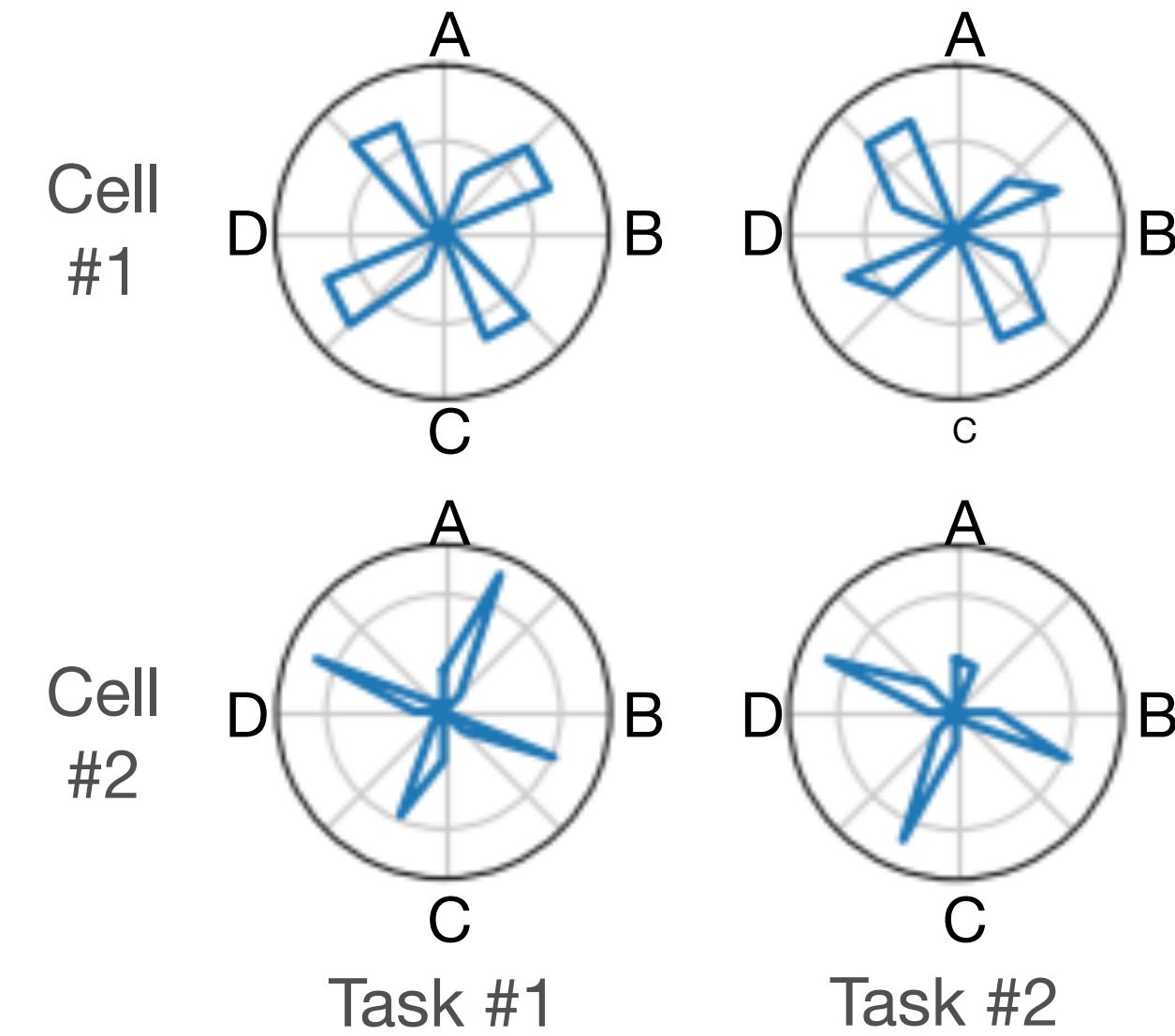
# Makes novel predictions we tested in prefrontal data



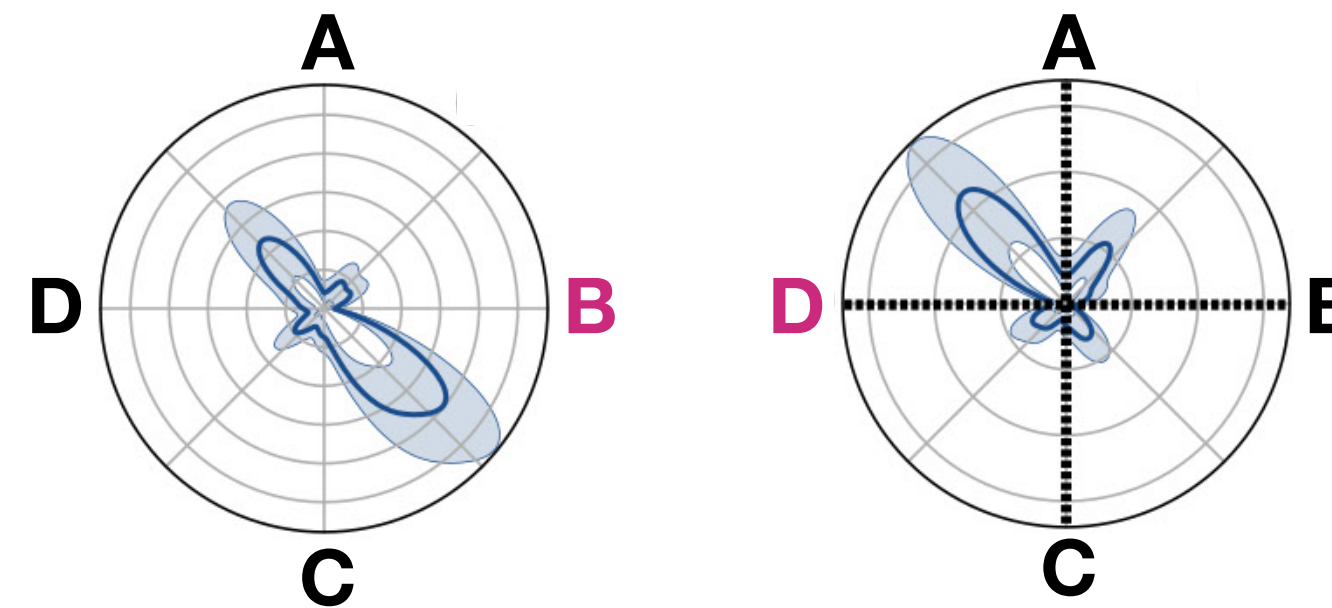
Task progress slot cells



Progress cells



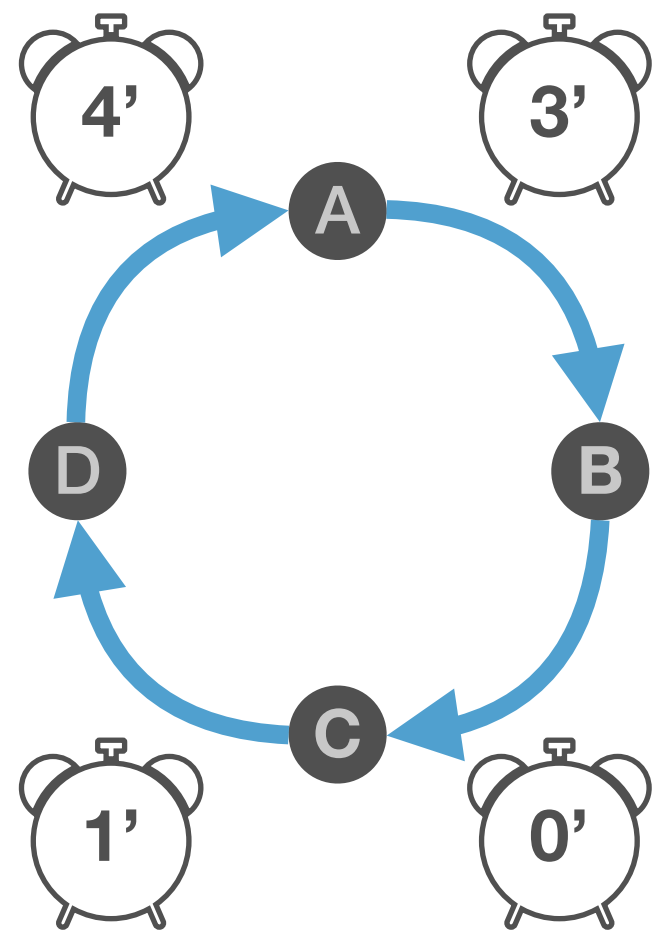
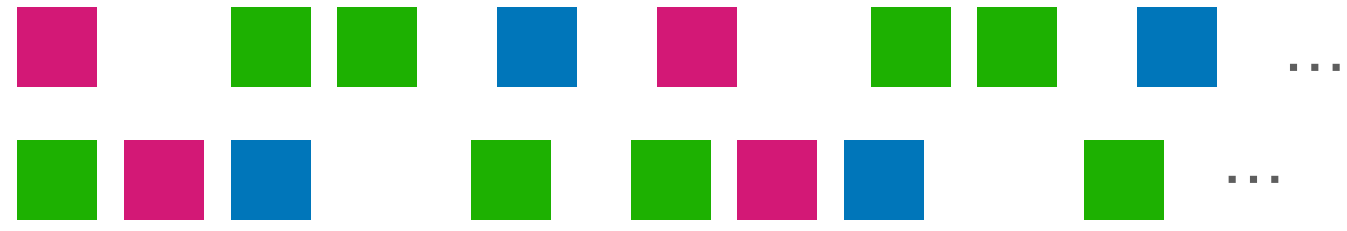
PFC real 'purple 50% lag' neuron



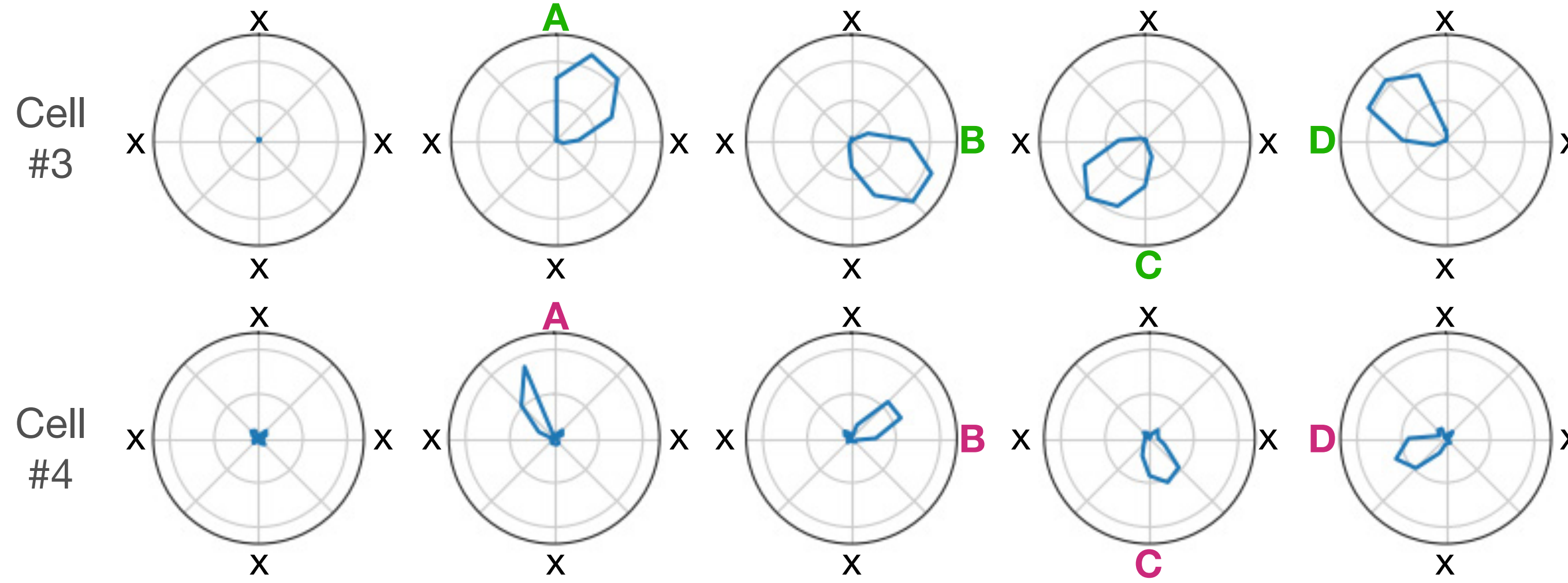
Whittington et al., 2023, *bioRxiv*

El-Gaby et al., 2023, *bioRxiv*

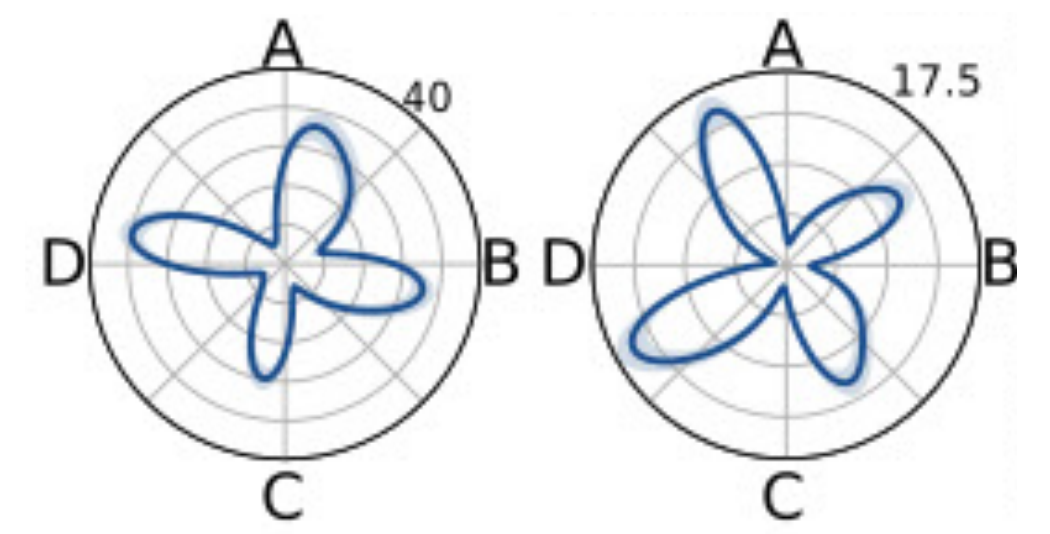
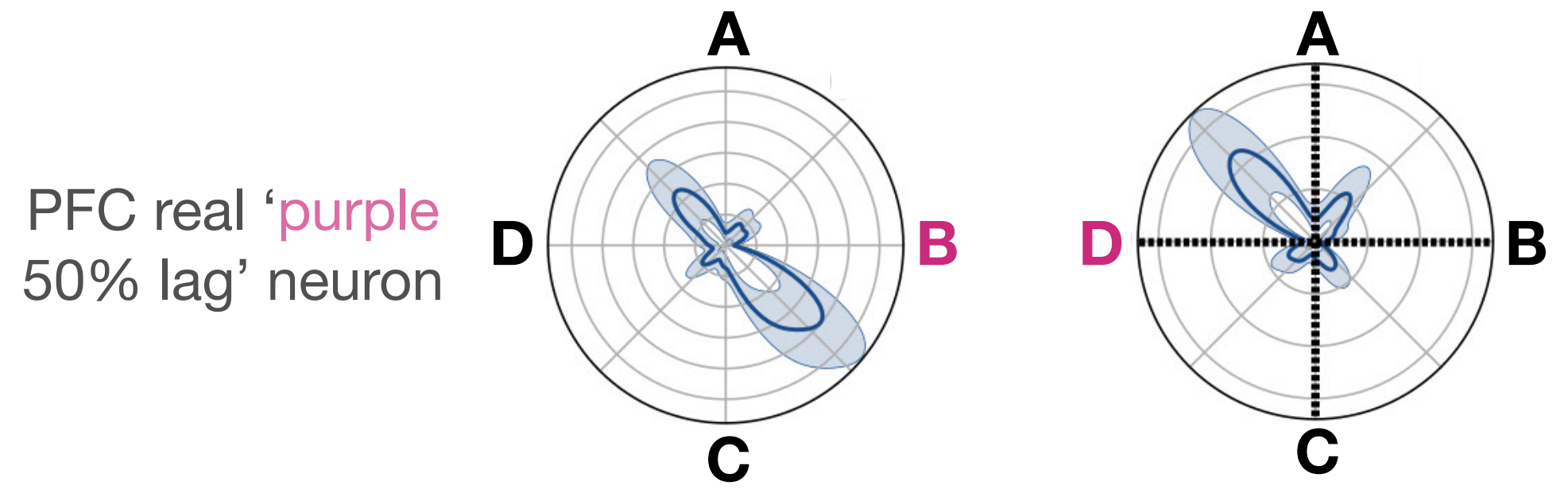
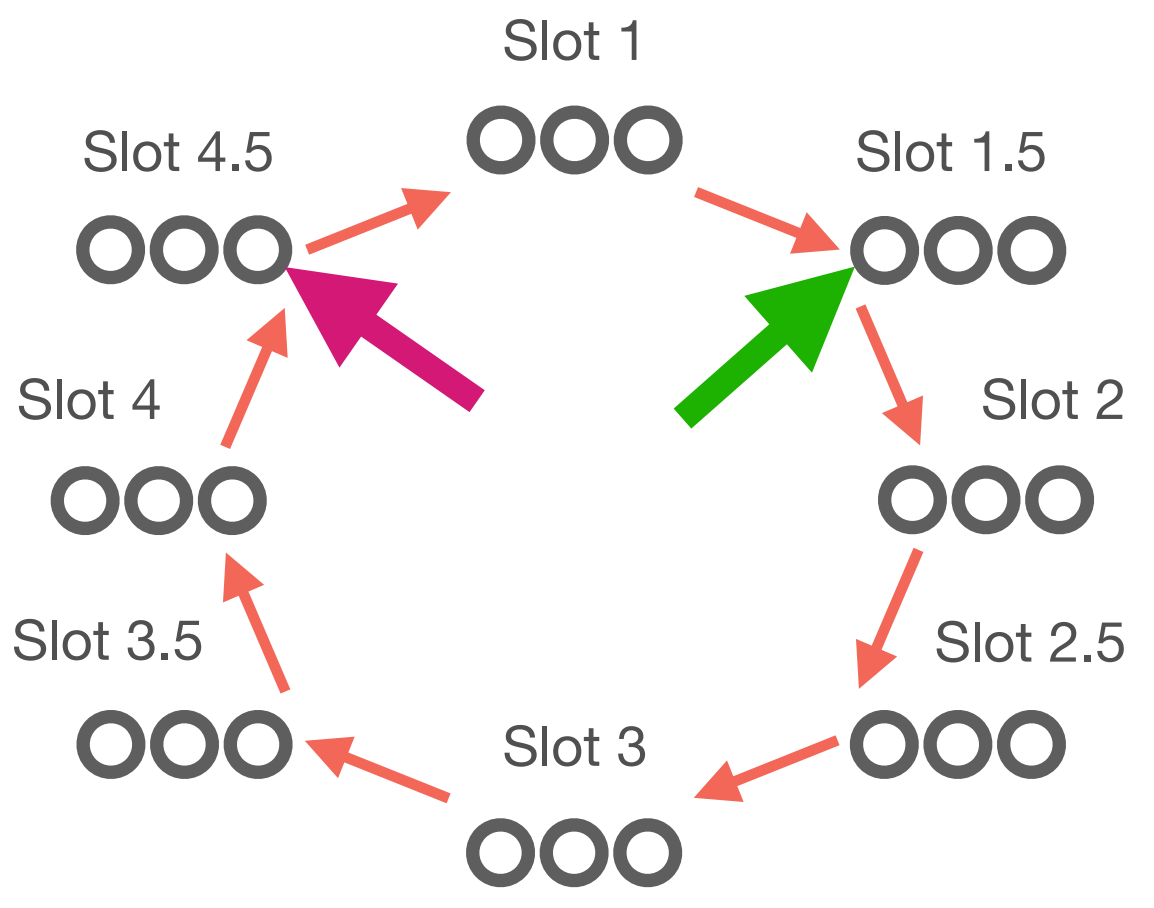
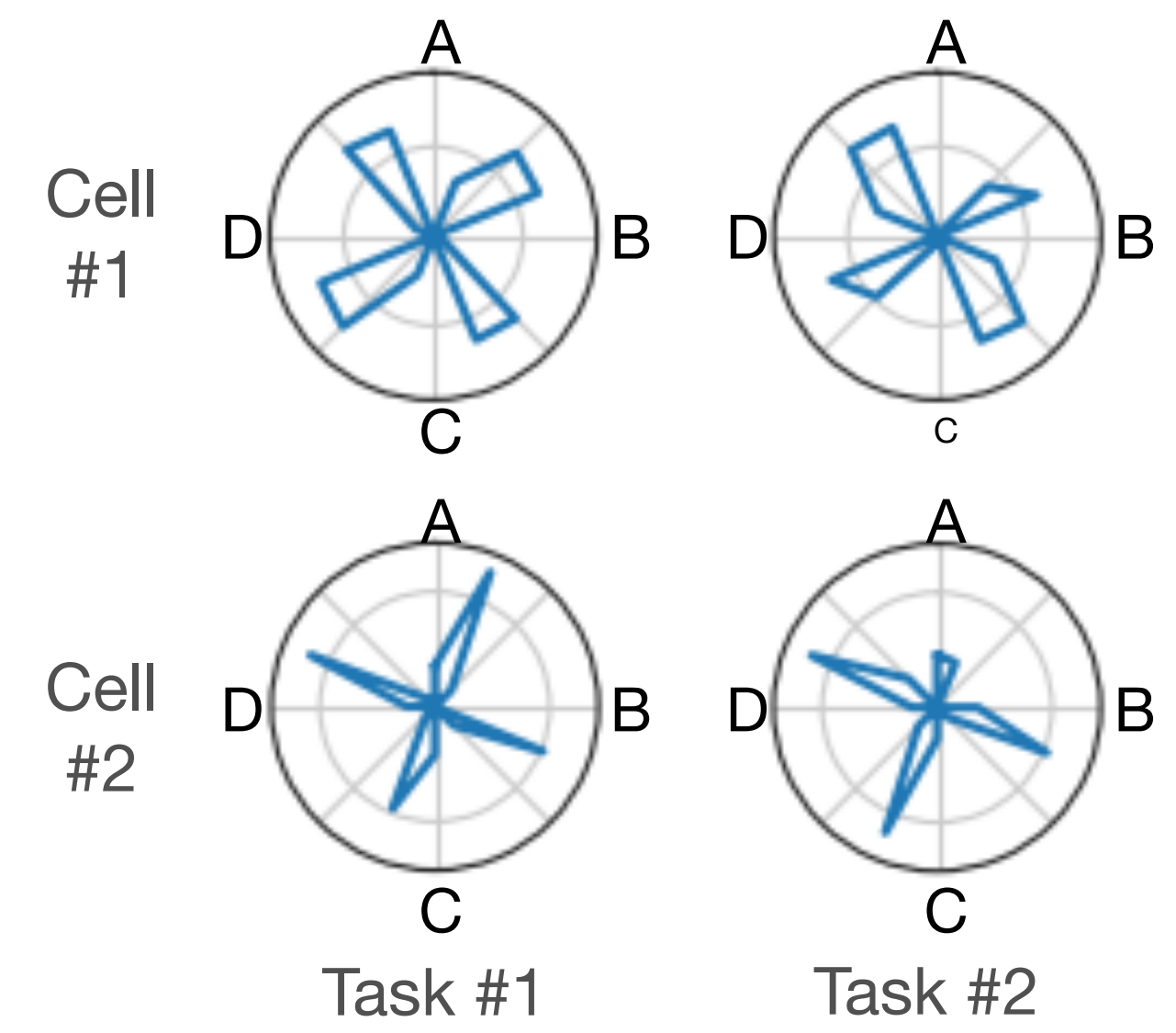
# Makes novel predictions we tested in prefrontal data



Task progress slot cells



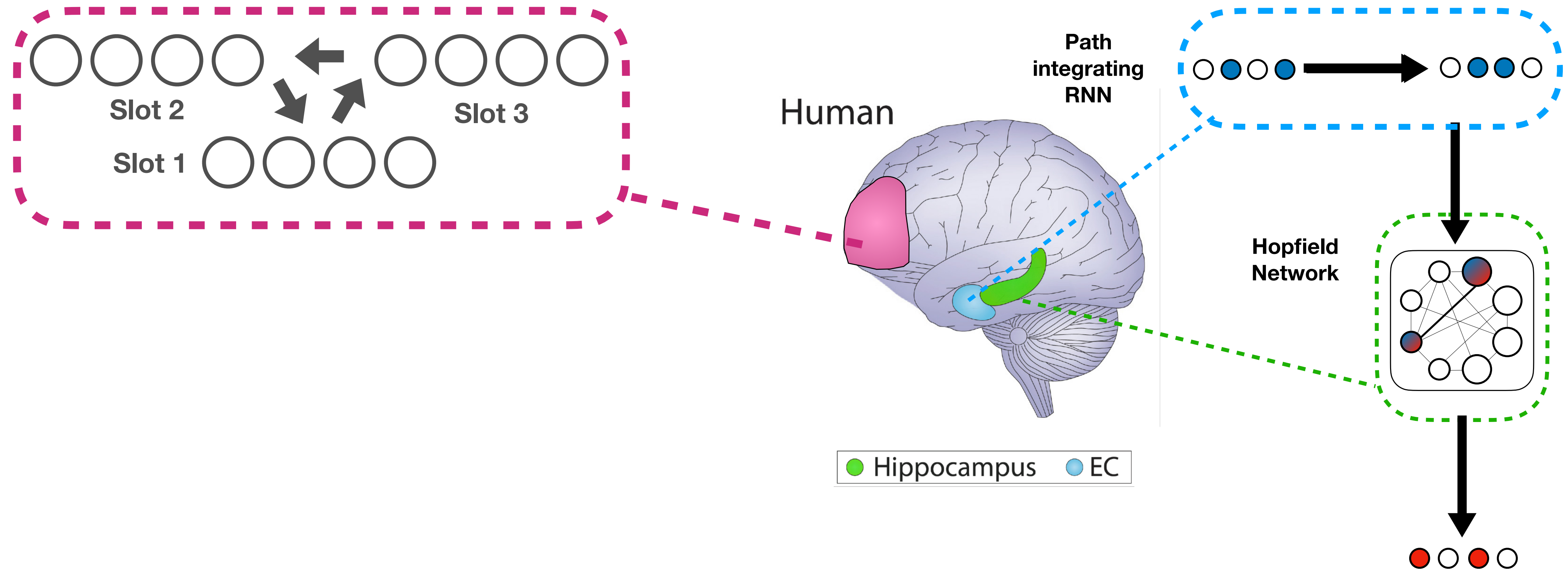
Progress cells



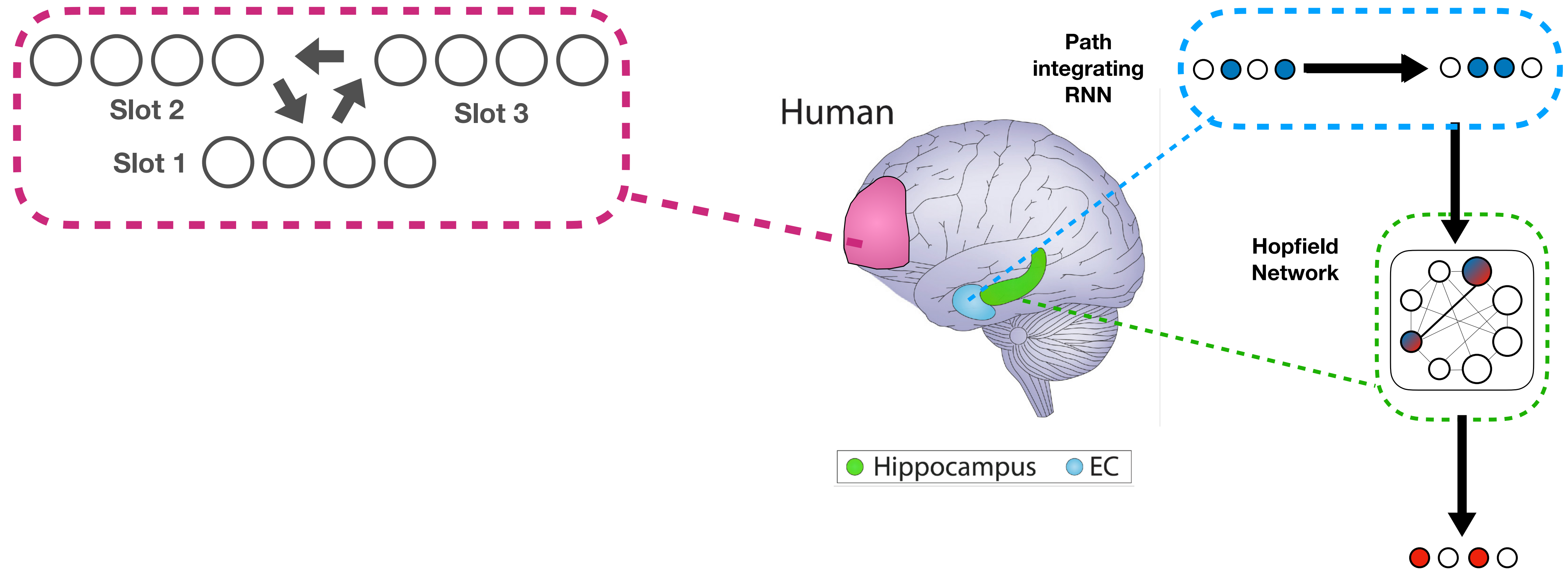
Whittington et al., 2023, *bioRxiv*

El-Gaby et al., 2023, *bioRxiv*

# Lastly, the two algorithms are mathematically equivalent



# Lastly, the two algorithms are mathematically equivalent



One stores memories in weights, and the other stores memories in neural activities

# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?



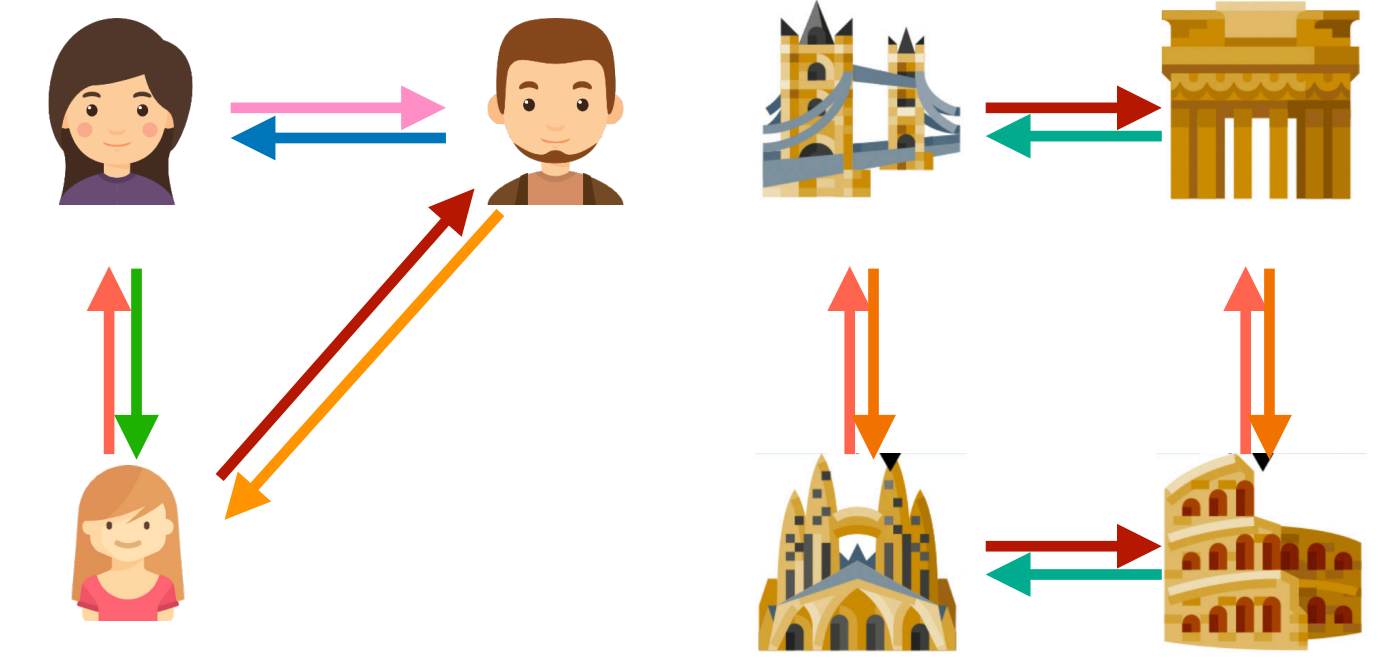
How can brains learn these maps?

Why do the neuron look the ways they do?

How do different brain regions solve the same problem in different ways?

# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?



How can brains learn these maps?

Why do the neuron look the ways they do?

How do different brain regions solve the same problem in different ways?

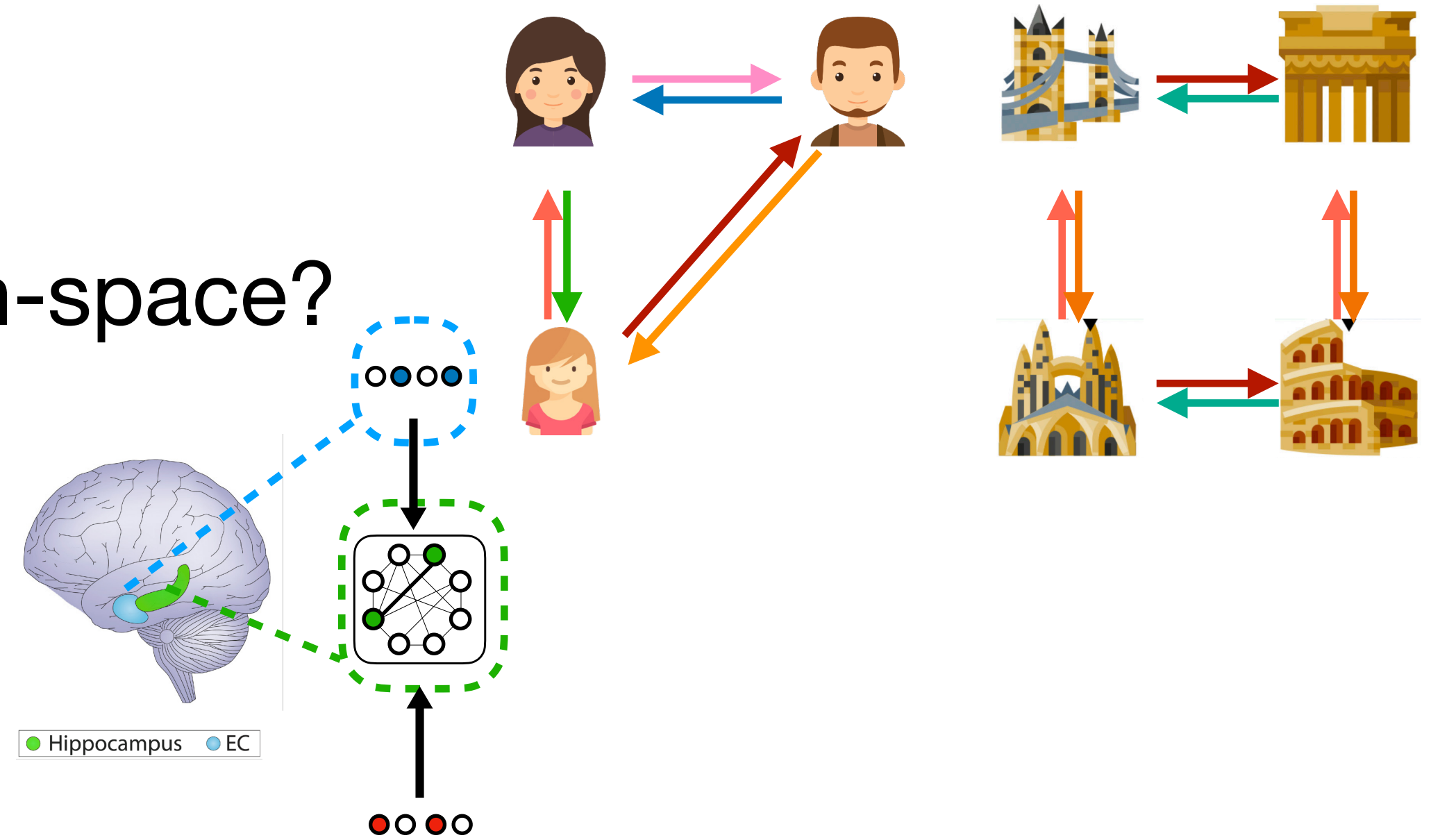
# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?

How can brains learn these maps?

Why do the neurons look the ways they do?

How do different brain regions solve the same problem in different ways?



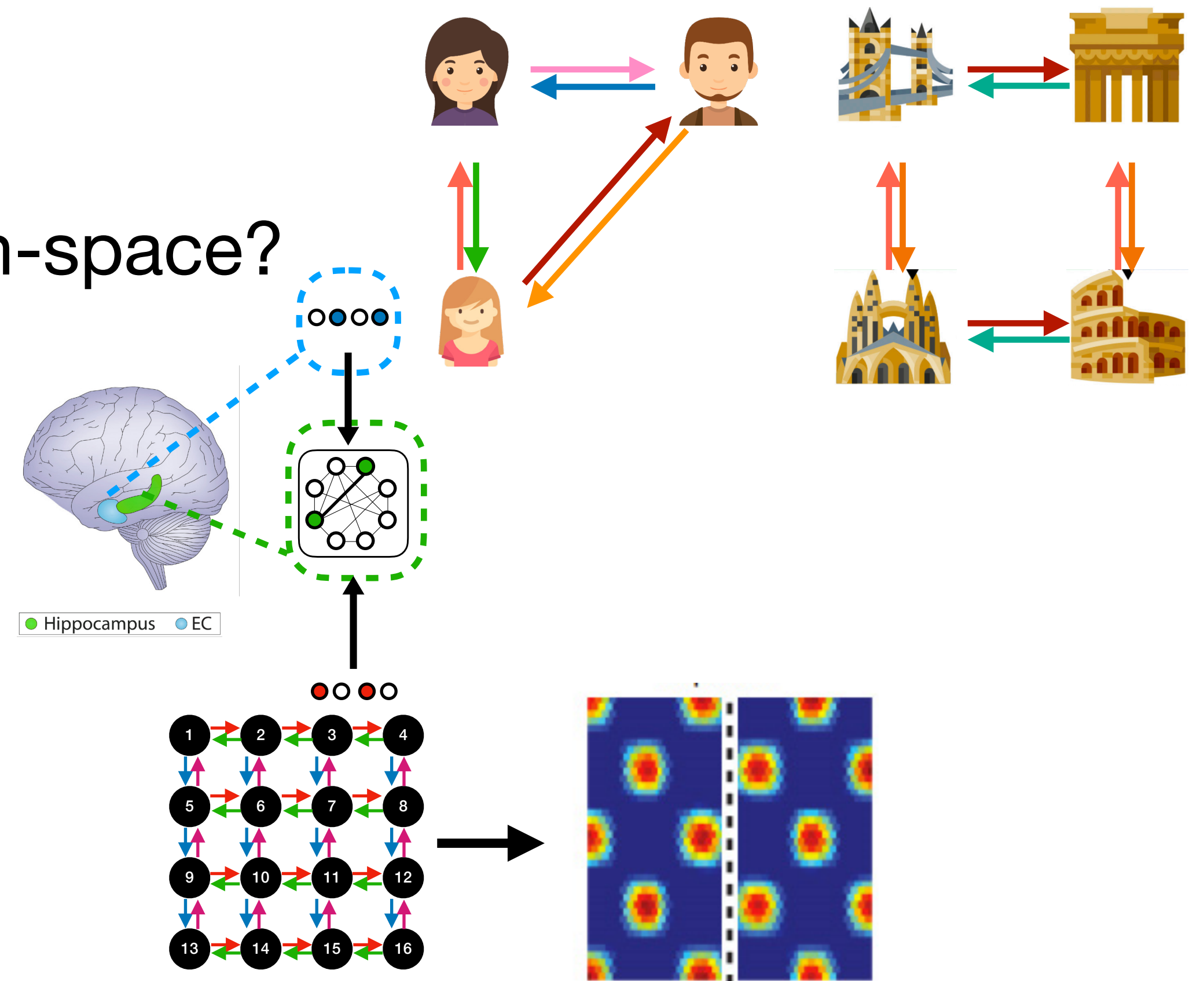
# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?

How can brains learn these maps?

Why do the neurons look the ways they do?

How do different brain regions solve the same problem in different ways?





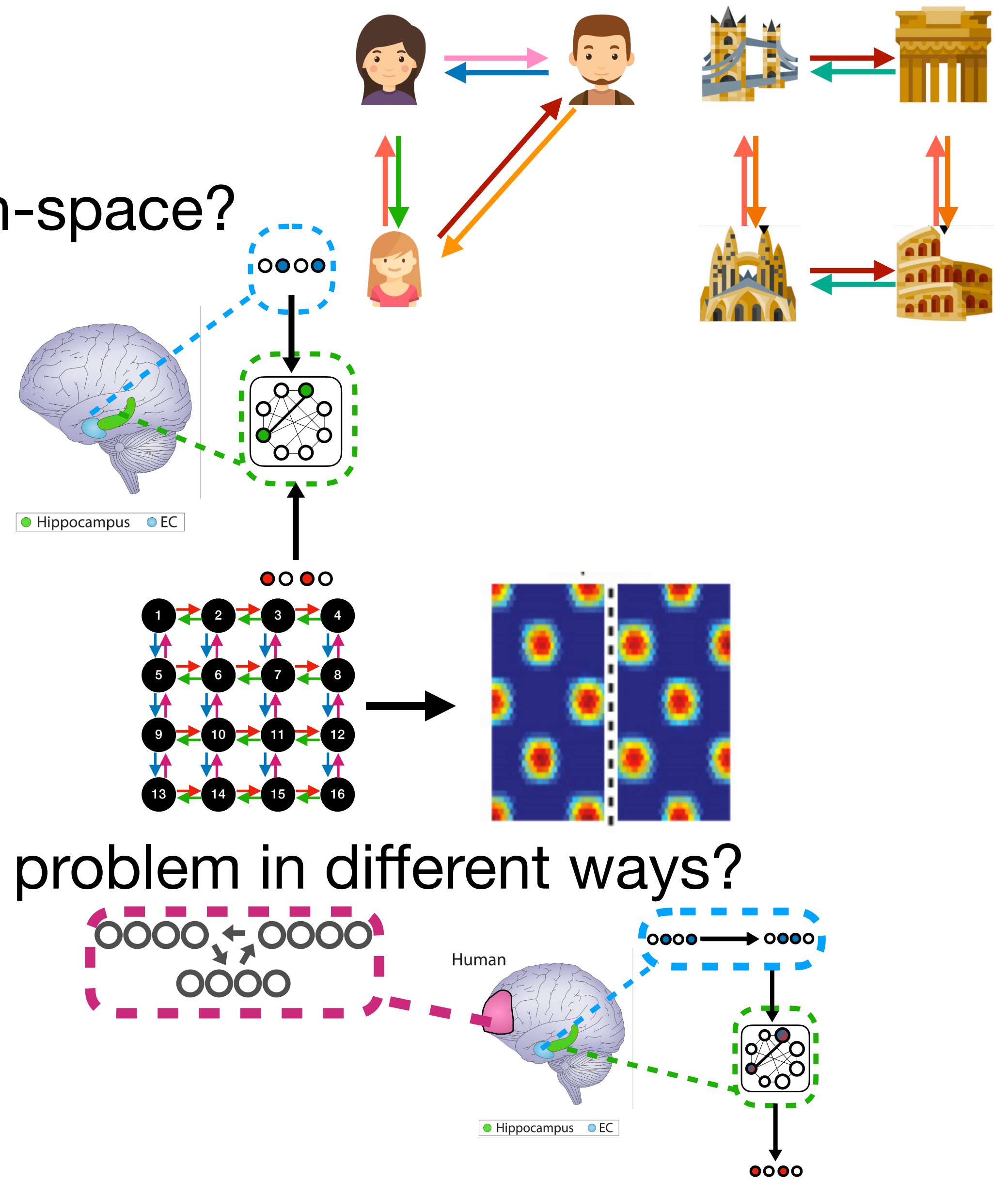
# Puzzles of cognitive maps in the brain

How does the same system do space and non-space?

How can brains learn these maps?

Why do the neurons look the ways they do?

How do different brain regions solve the same problem in different ways?



# Thanks!

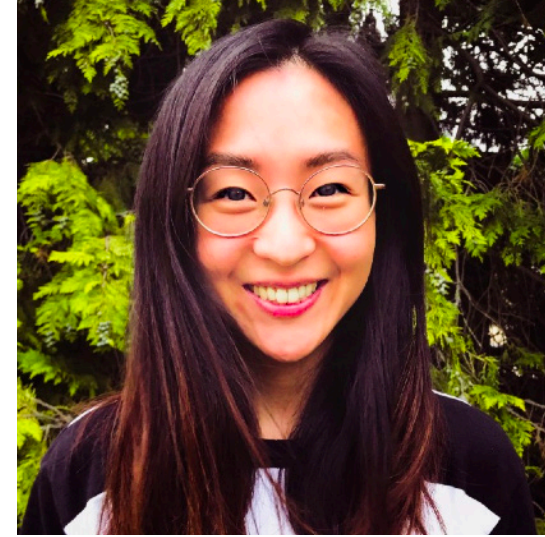
Jacob Bakermans



Tim Muller



Chongyu Qin



Will Dorrell



Shirley Mark



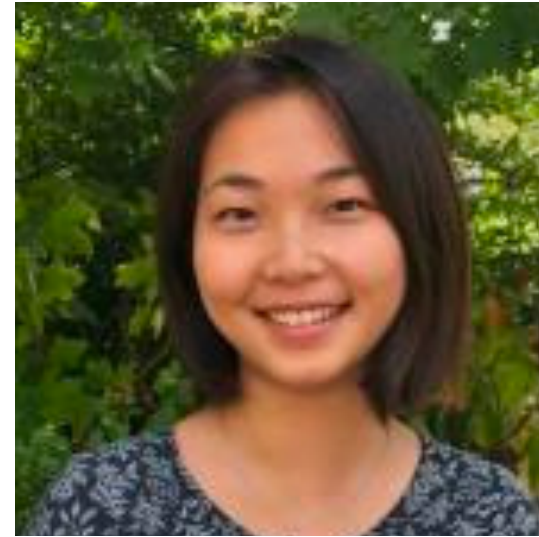
Guifen Chen



Jo Warren



Jiali Zhang



Mohamady El-Gaby



Tim Behrens



Surya Ganguli

