

# Collective behavior from surprise minimization

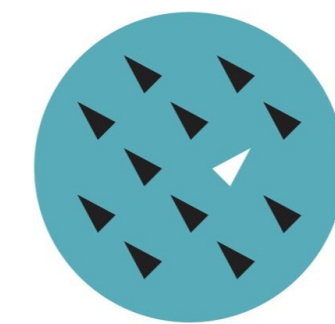


Conor Heins

 @conorheins

**VERSES**

VERSES Machine Learning Lab  
Max Planck Institute of Animal Behavior



Modelling Multi-Scale Collective Intelligences Workshop  
19 November 2024

# Overview

## **Part I: Background on Bayesian cognitive science & active inference**

- Perception as unconscious inference and Bayesian inference
- Minimizing prediction error as an algorithm for inference
- Active inference

## **Part II: Applying the concepts from Part I to model collective motion**

- Collective motion overview
- Phenomenological models vs 'cognitive' models
- Collective motion from multi-agent active inference
- Emergent information transfer & decision-making
- Learning one's model online, i.e., 'behavioral plasticity'

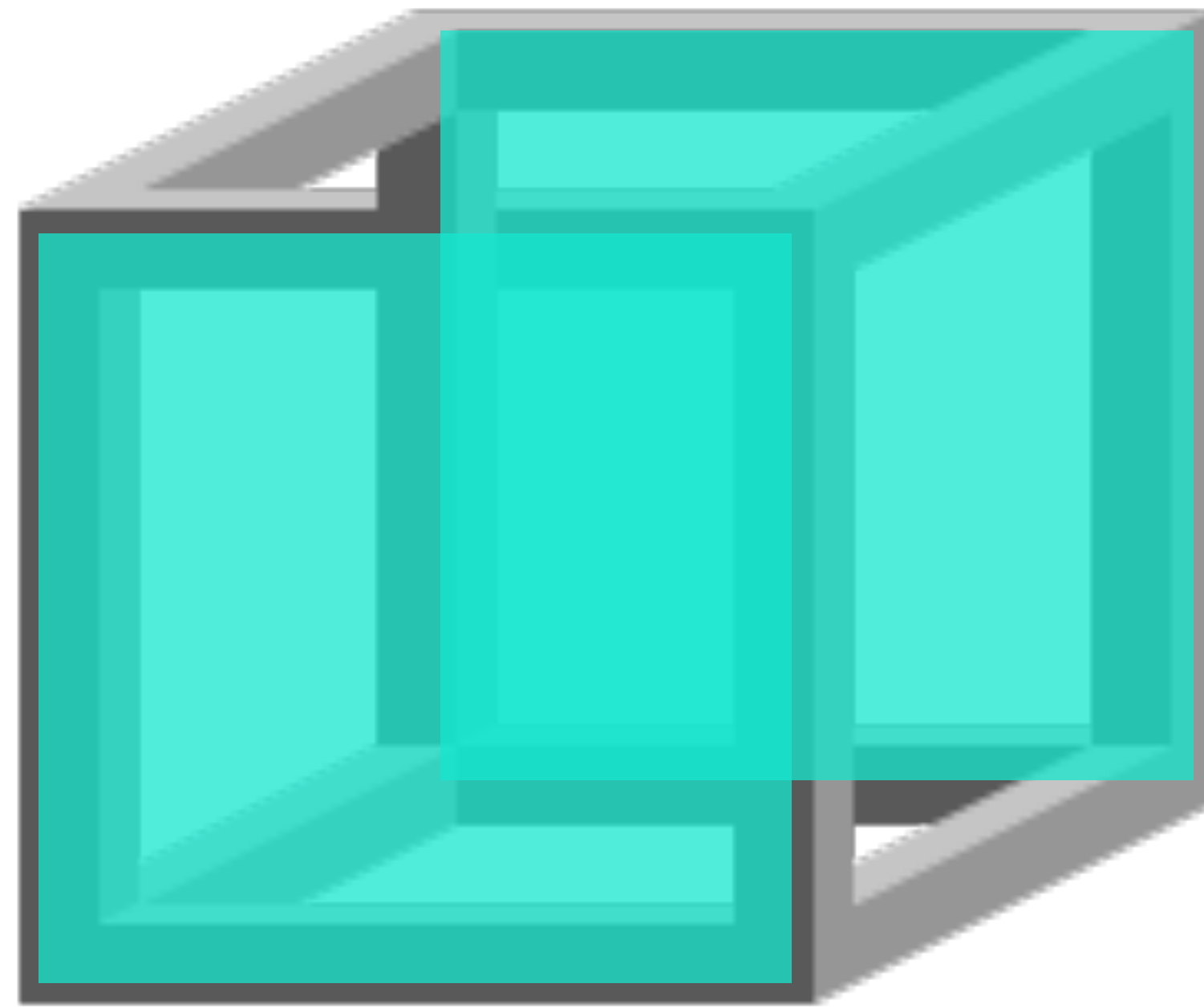
# The “Bayesian turn” in the cognitive sciences

- Hermann von Helmholtz, “Perception as unconscious inference” (*unbewusster Schluss*)

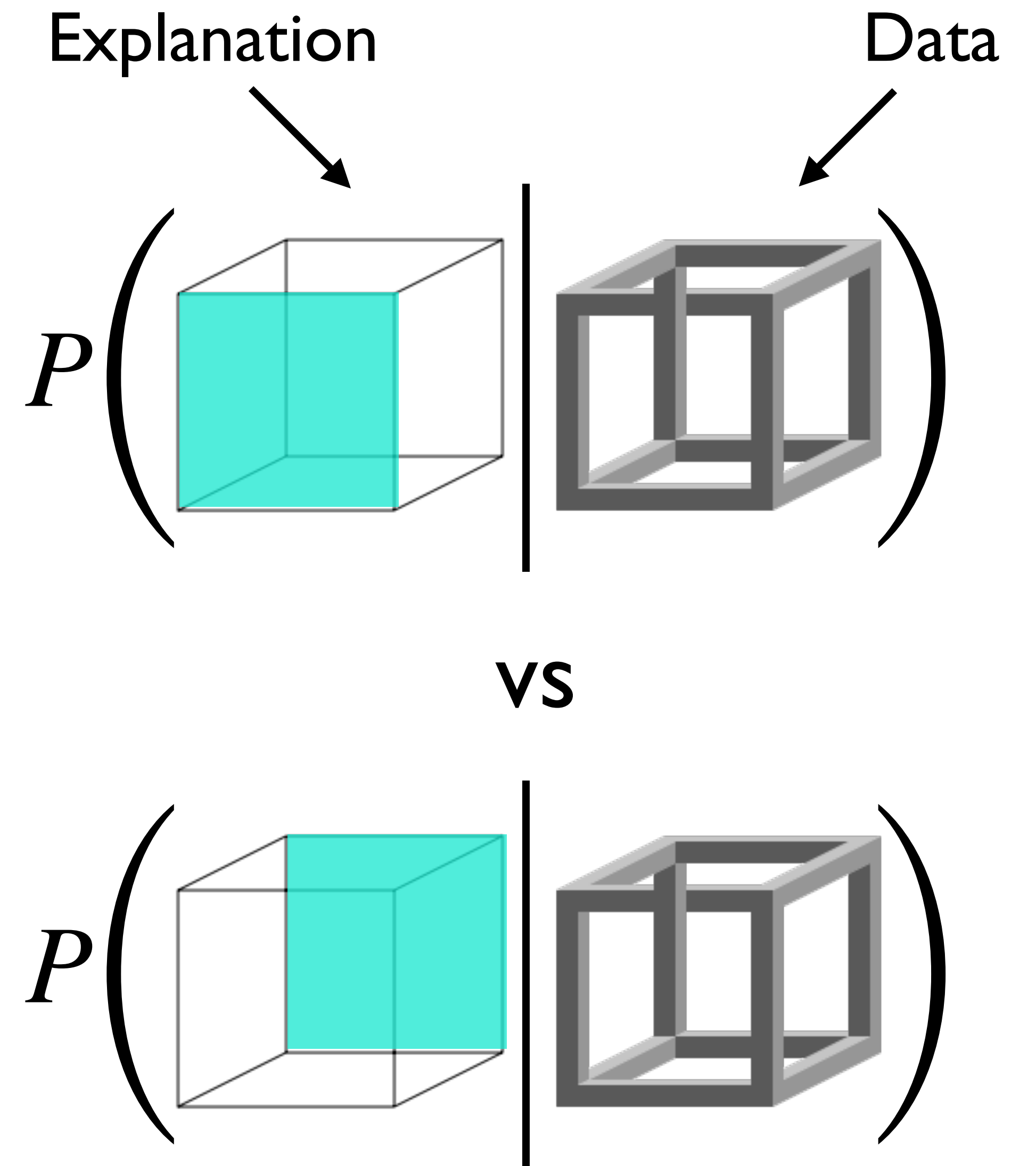


Hermann von Helmholtz

# Perception as explanation

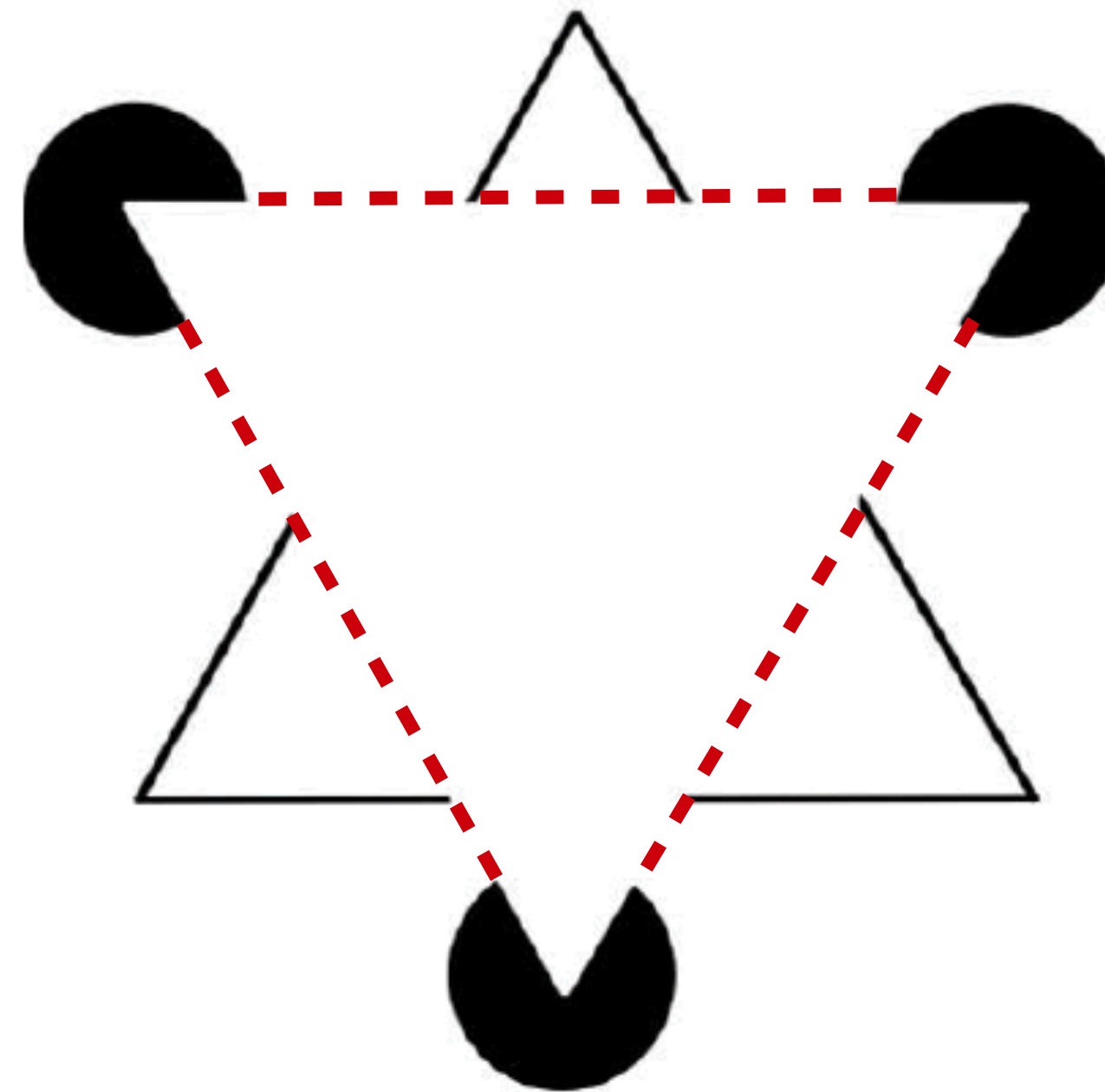


Neckar Cube



# Perception as explanation

Kanisza's triangle



Your perception is not the “raw data” (aka the pixel intensities on the screen), but an **inference** or **interpretation** of that data

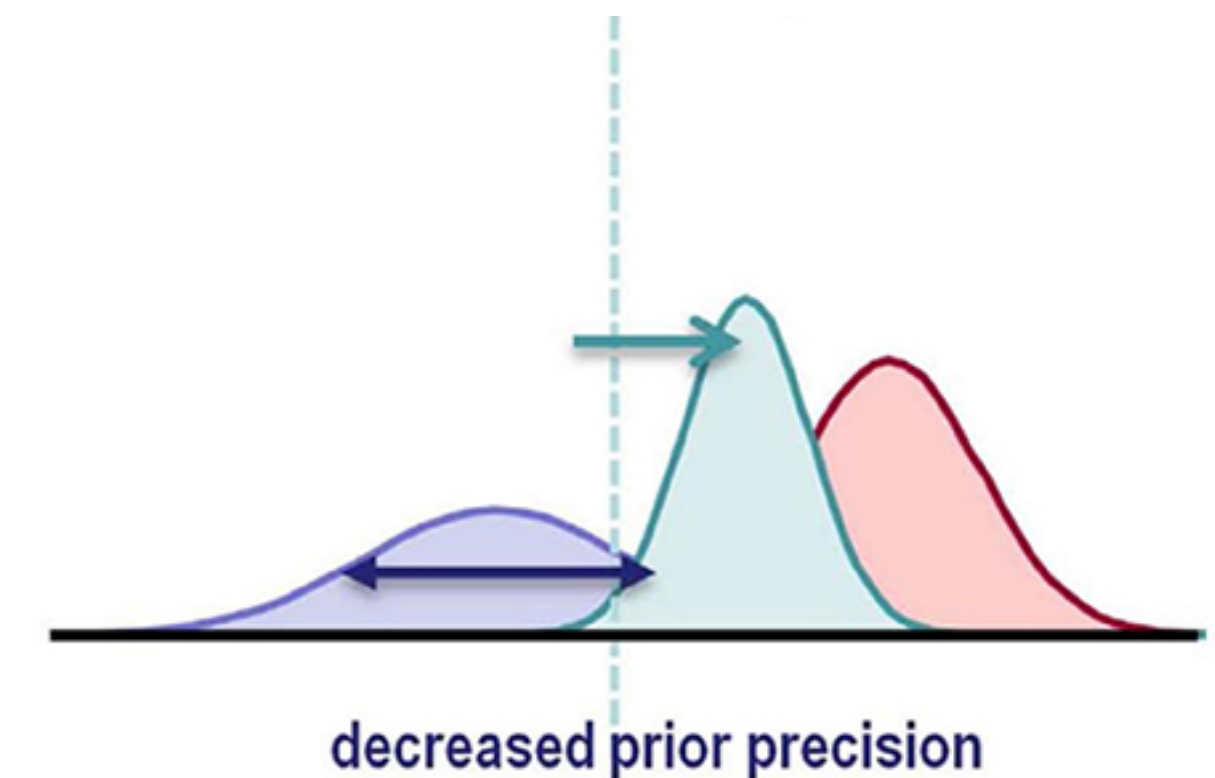
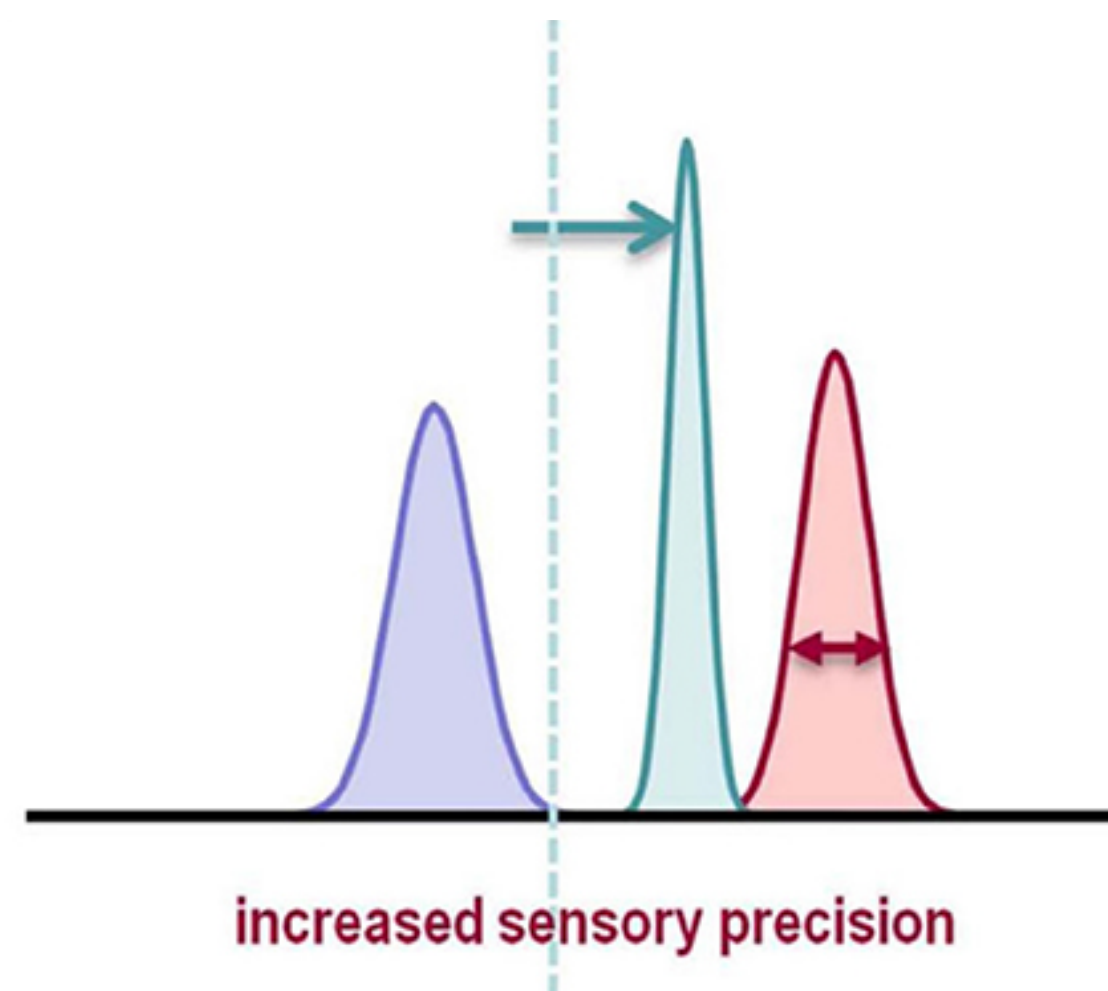
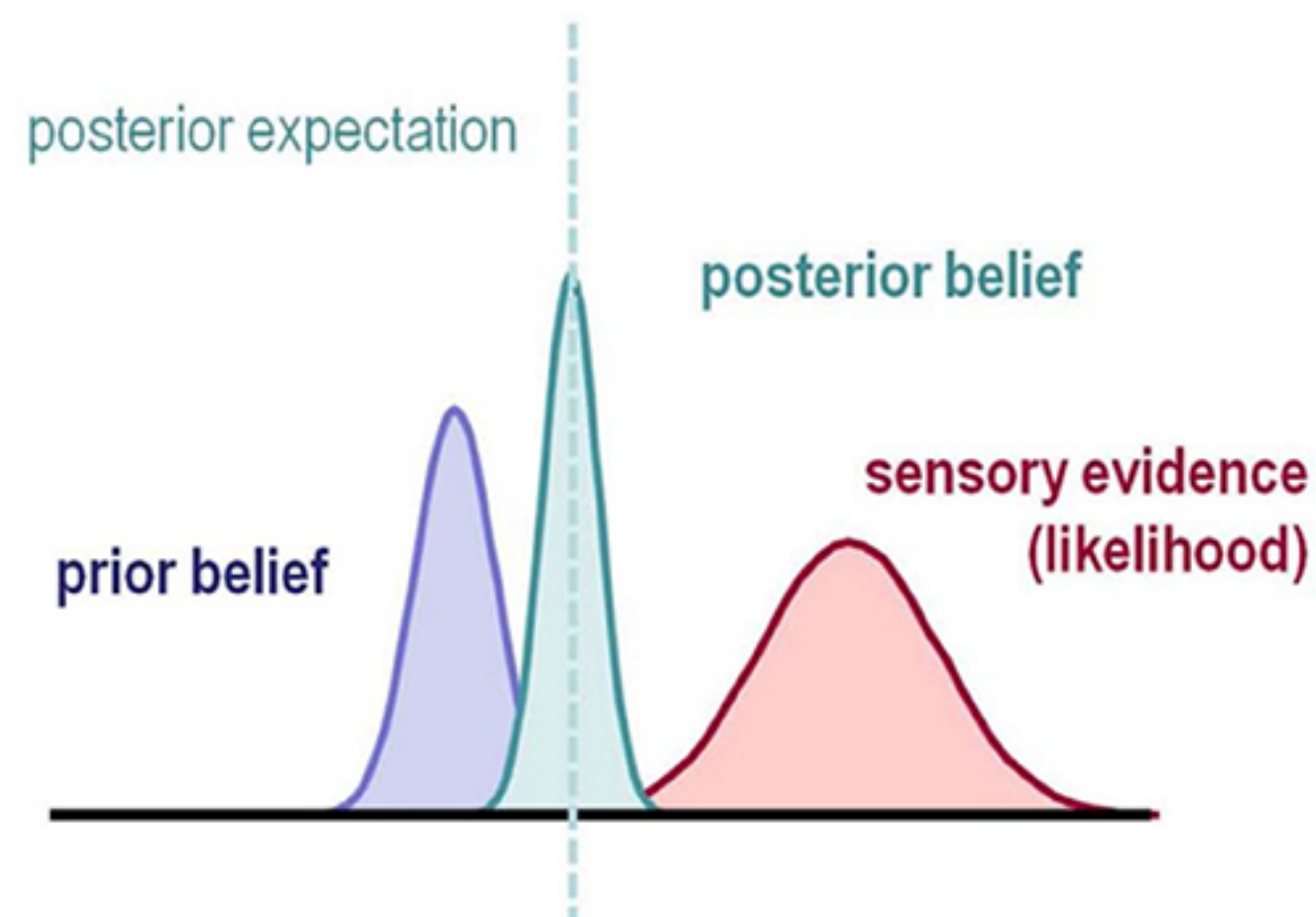
# The “Bayesian turn” in the cognitive sciences

- Hermann von Helmholtz, “Perception as unconscious inference” (*unbewusster Schluss*)
- Formalized later in the 20th century as probabilistic inference — use Bayes Rule to compute posterior probabilities

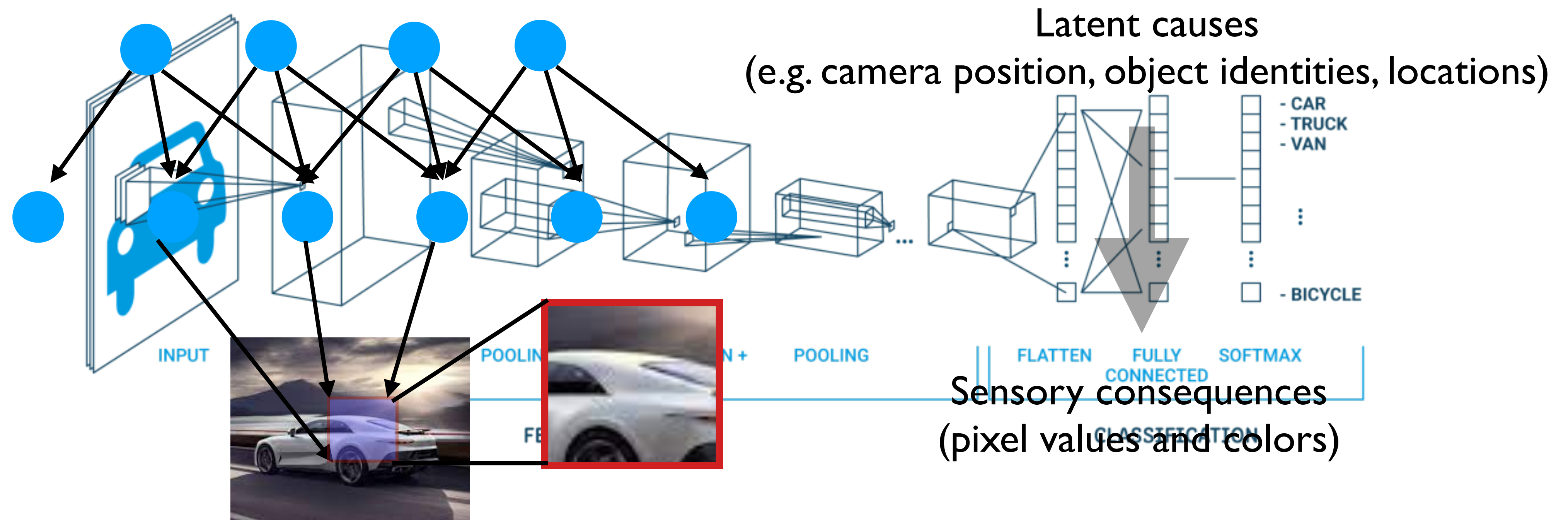


Hermann von Helmholtz

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

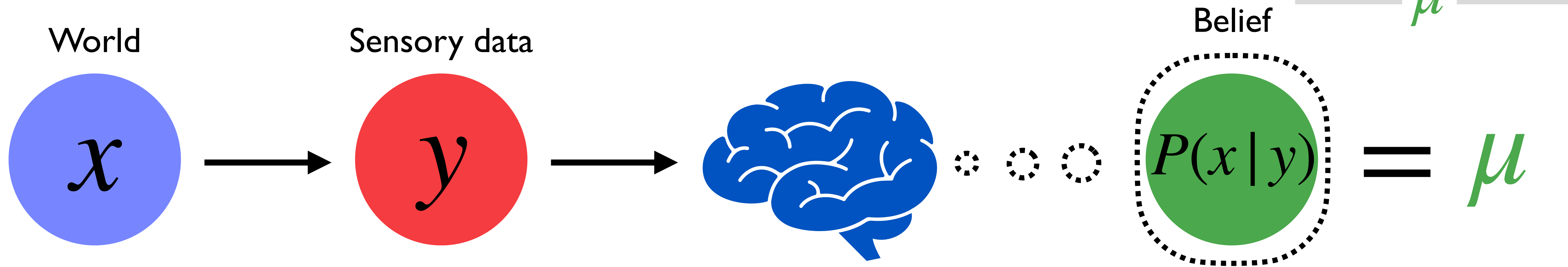


# Perception as ~~feature~~ ~~read~~ ~~detection~~ explanation



Feedforward cascade of feature detectors  
Inference using a generative model (inverse graphics)

# One way to infer: minimise **prediction error**



$$\mu^* = \arg \min_{\mu} (\text{Prediction error})$$

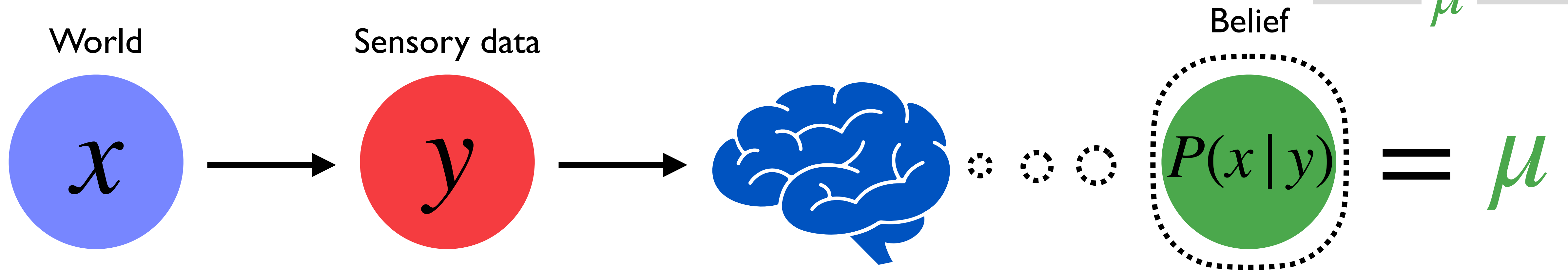
$\mu$   $\uparrow$   
What you see

$\uparrow$   
What you expect



# More formally:

use a gradient descent on free energy to update the posterior means



$$\frac{d\mu}{dt} = - \frac{\partial F(\mu, y)}{\partial \mu} \longrightarrow F \propto \epsilon_y^2 + \epsilon_\mu^2$$

Sensory prediction errors      Prior prediction errors

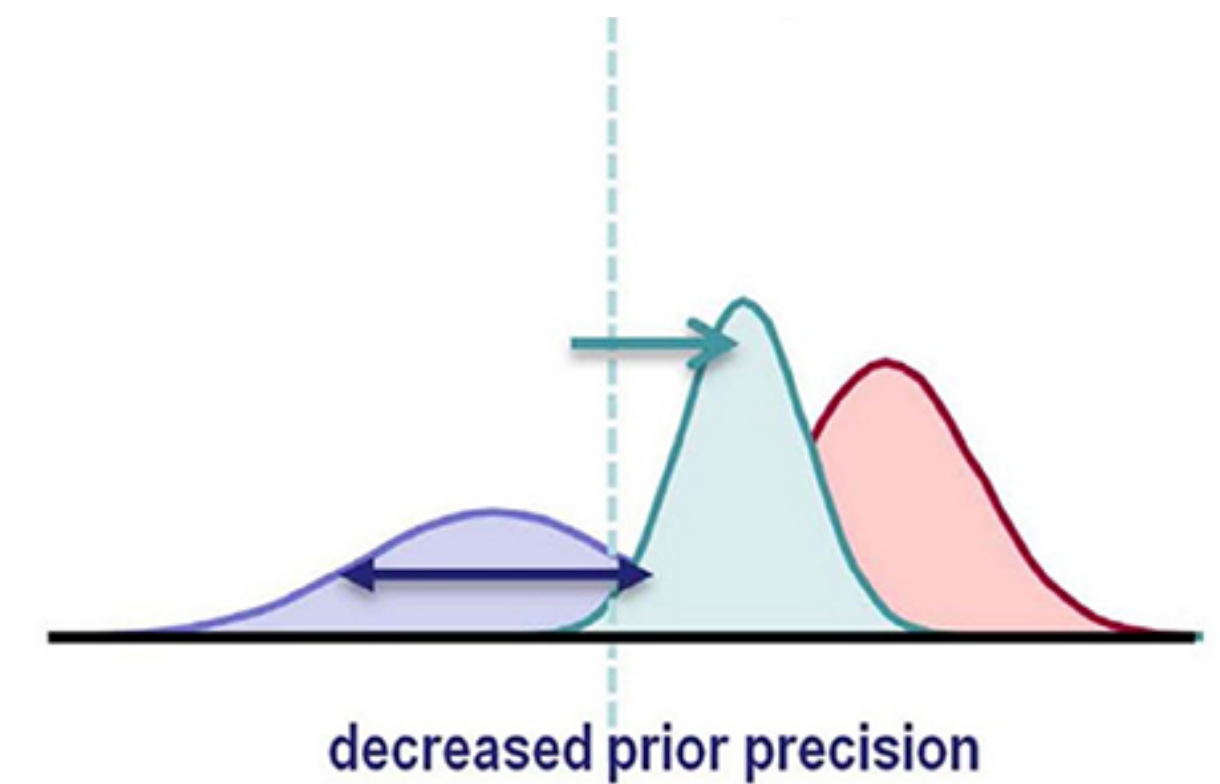
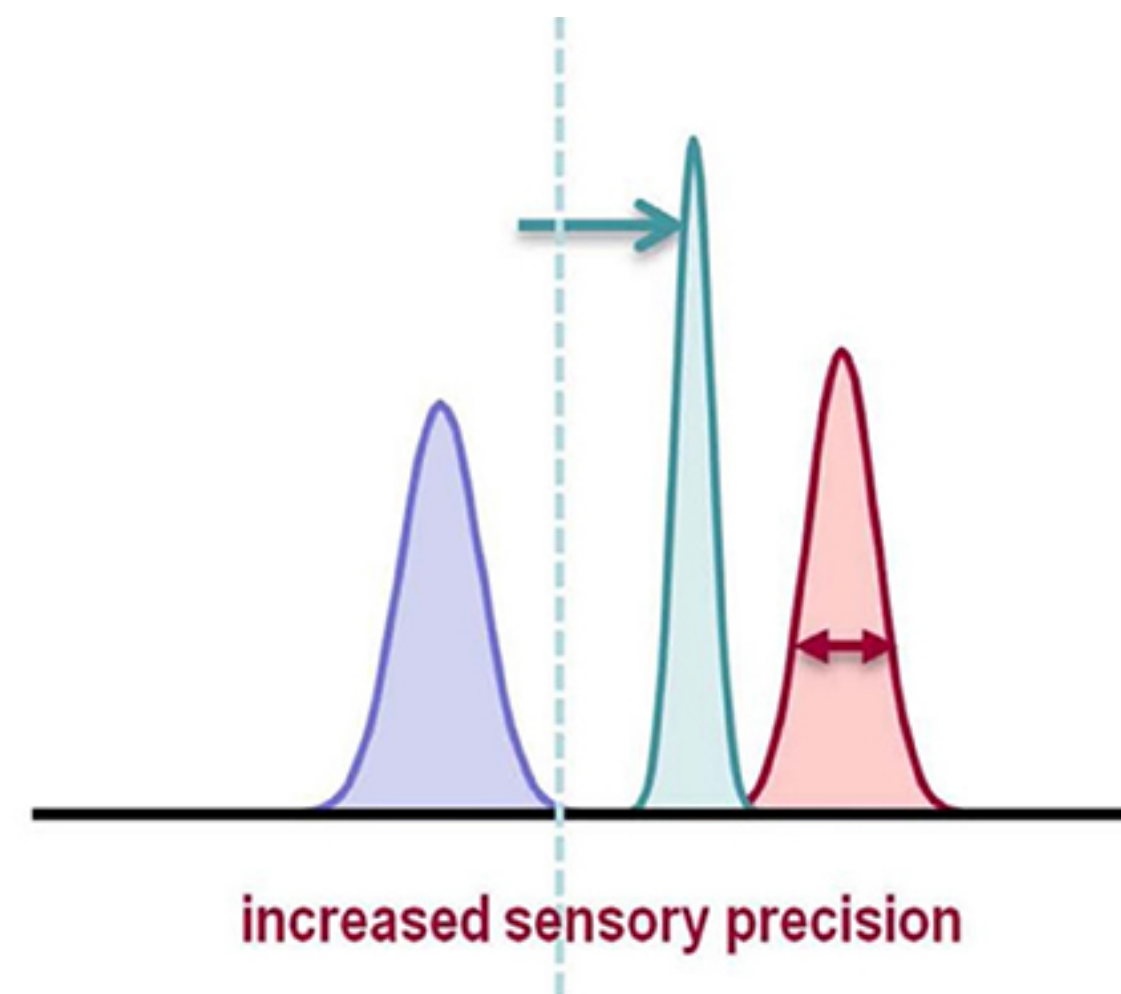
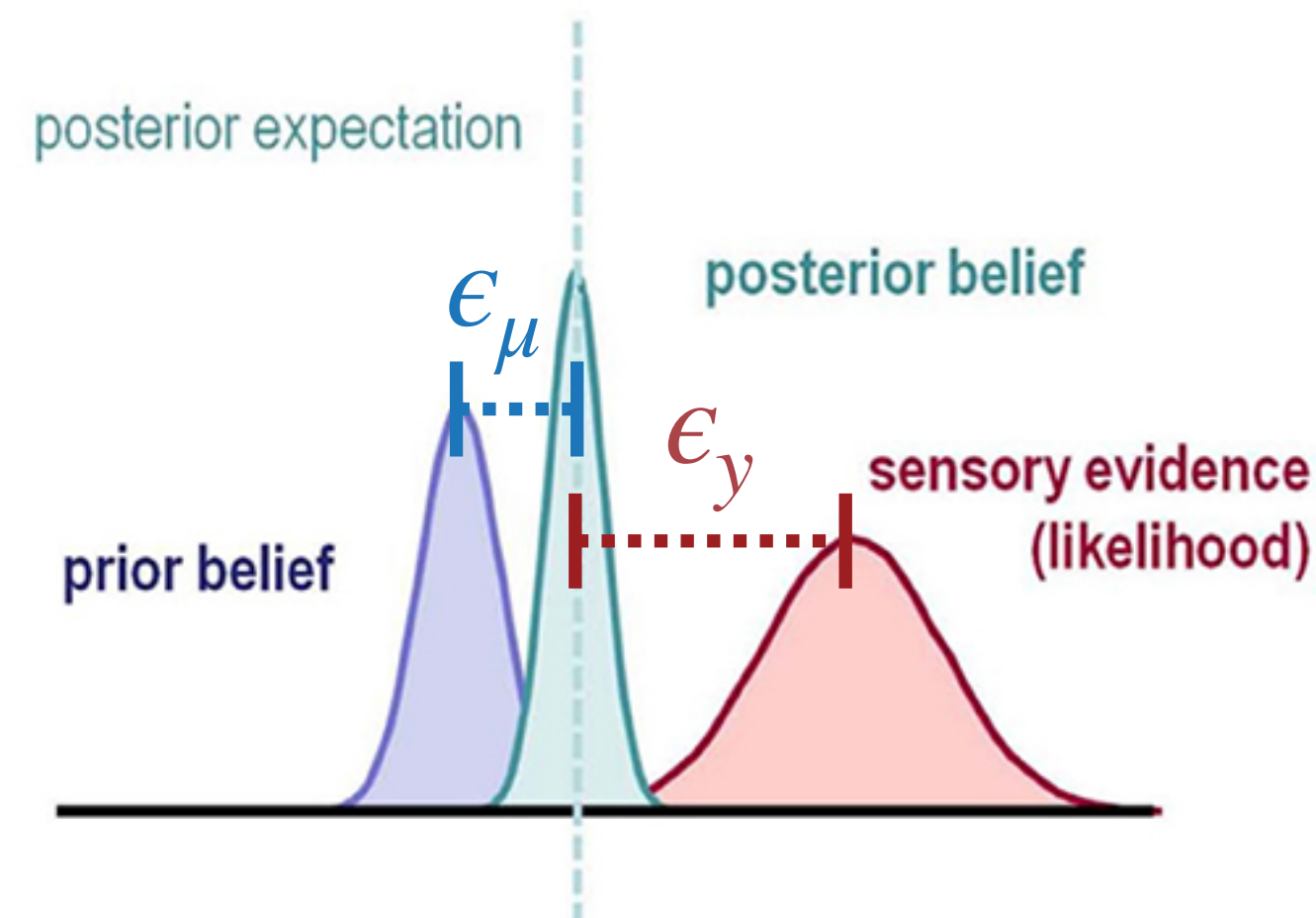
$$\epsilon_y = y - \mu \quad \epsilon_\mu = \mu - \eta$$

# Precision weighting is important!

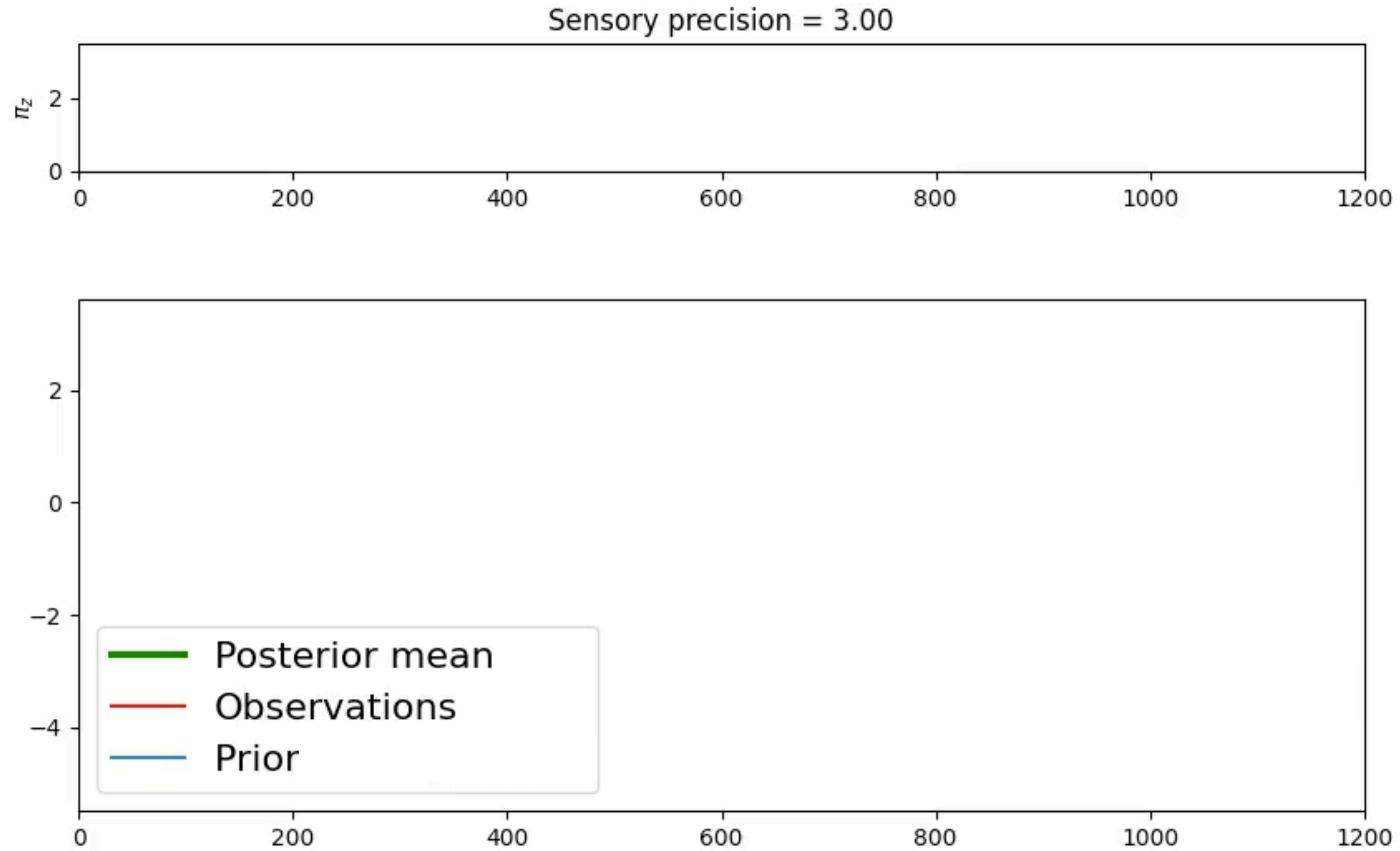
$$\frac{d\mu}{dt} = - \frac{\partial F(\mu, y)}{\partial \mu}$$

$$F \propto \pi_y \epsilon_y^2 + \pi_\mu \epsilon_\mu^2$$

Free energy  
 = surprise  
 = “sum of (squared) precision-weighted prediction errors”

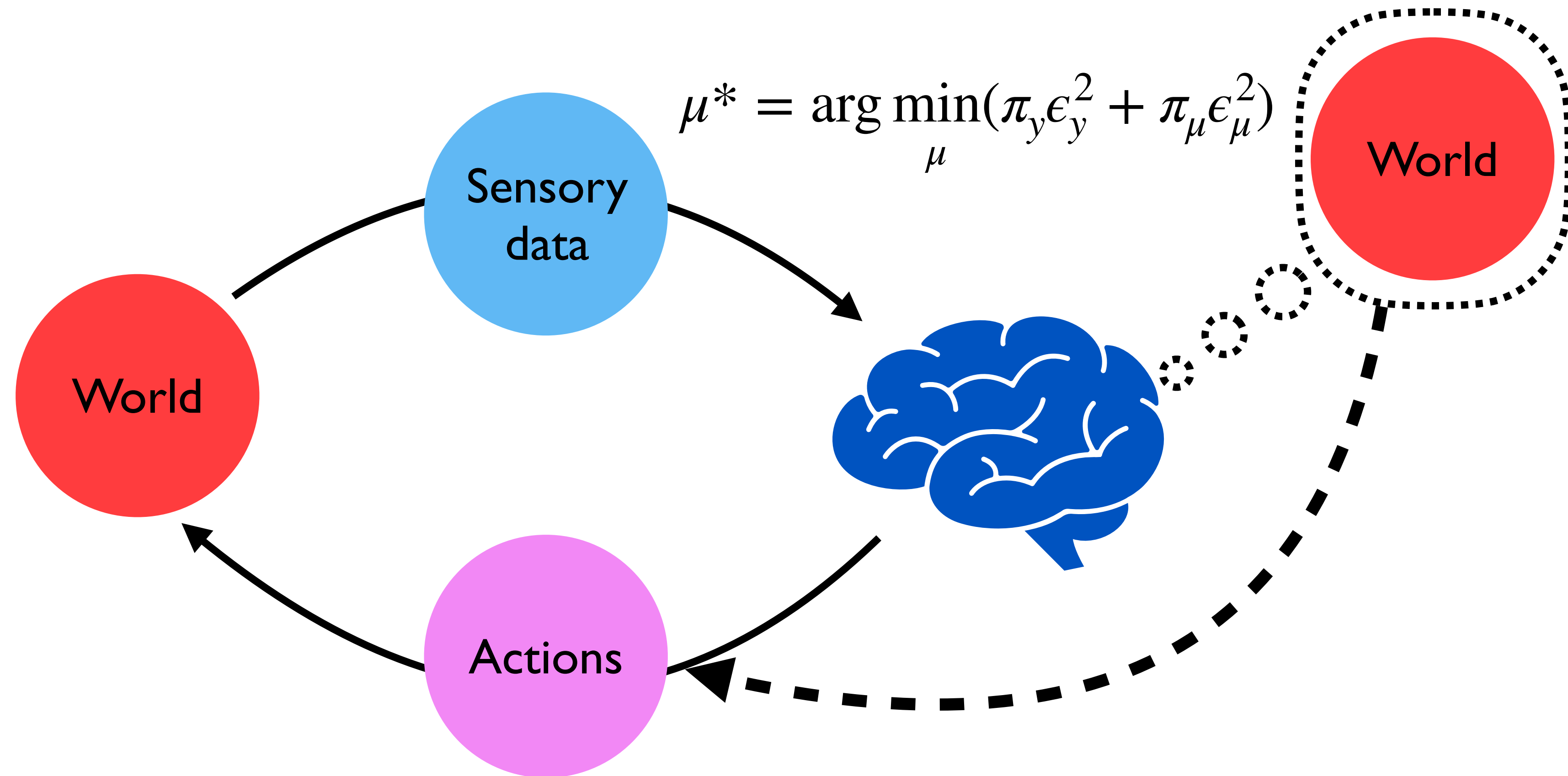


# Bayesian filtering example with attenuation in sensory precision



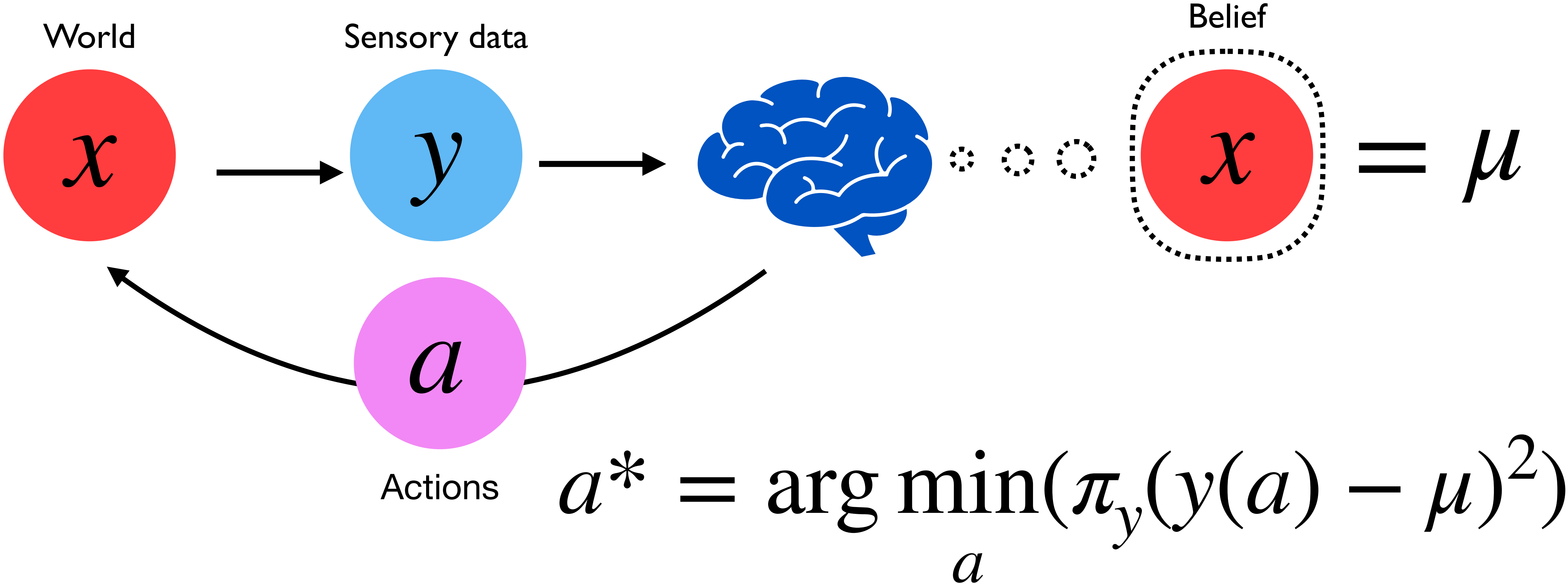
# What about action?

**Active** inference, aka belief-guided action



This is a standard reinforcement learning formulation  
— how is active inference any different?

# How to act? also minimise **prediction error**

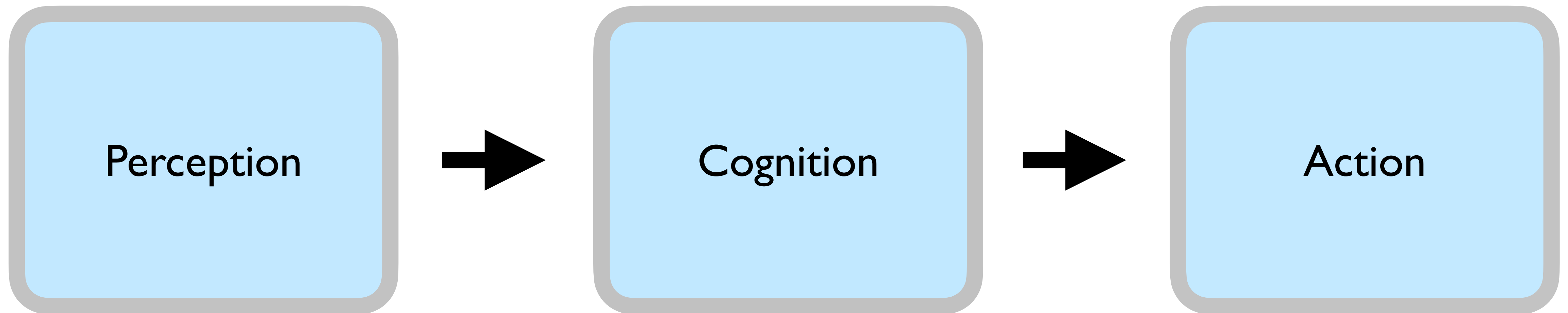


Under **active inference**, everything is about minimising **prediction error** aka “surprise”

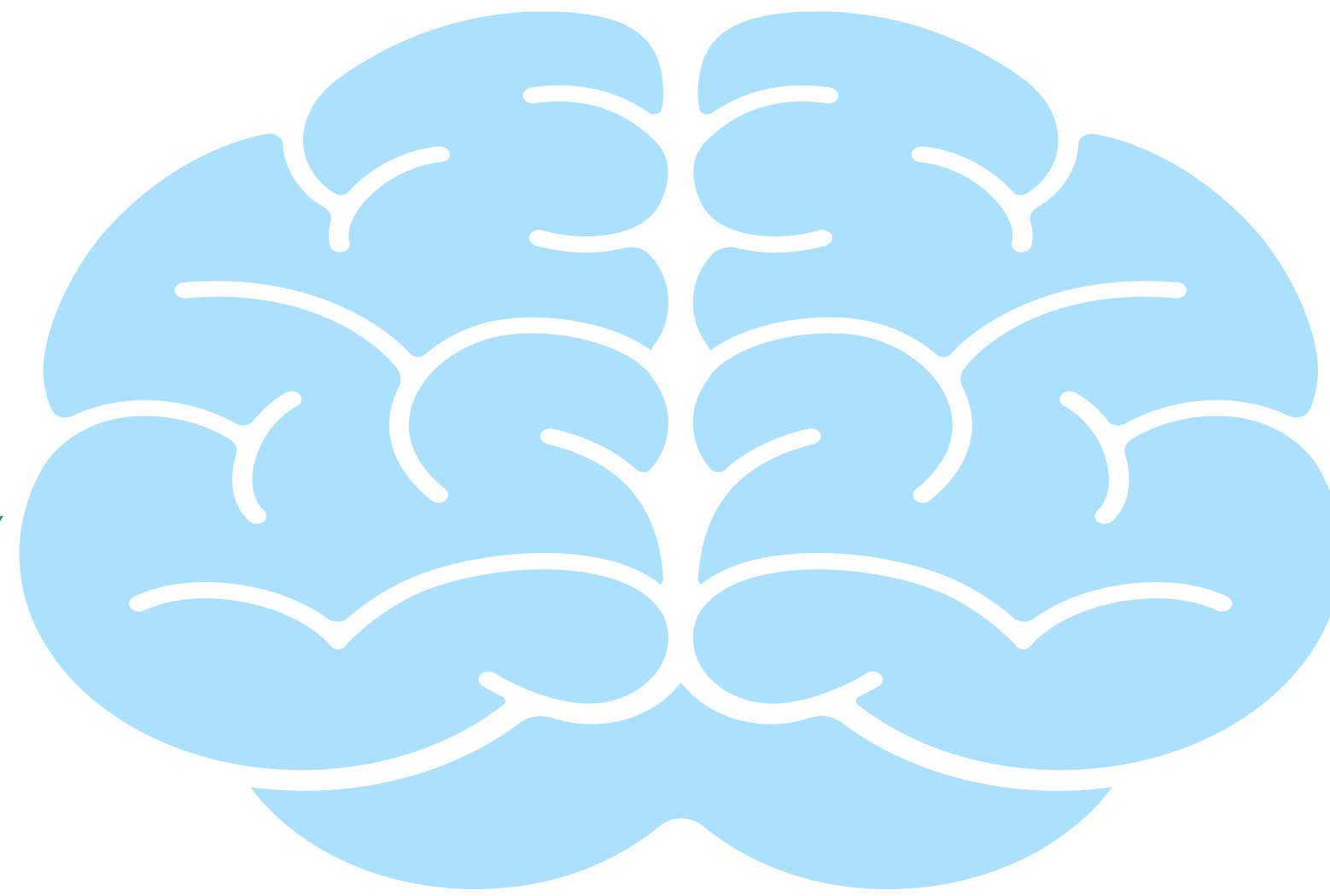
There are two ways to become less surprised

# Classical “Sandwich Model” of Cognition

- Unidirectional information flow from sensory states to motor effectors
- Distinct stages of processing
- Cognition plays role of generating “actionable representations”



Adams, Shipp, Friston (2013).  
*Predictions, not Commands:*  
*Active inference in the motor system*



Predictions

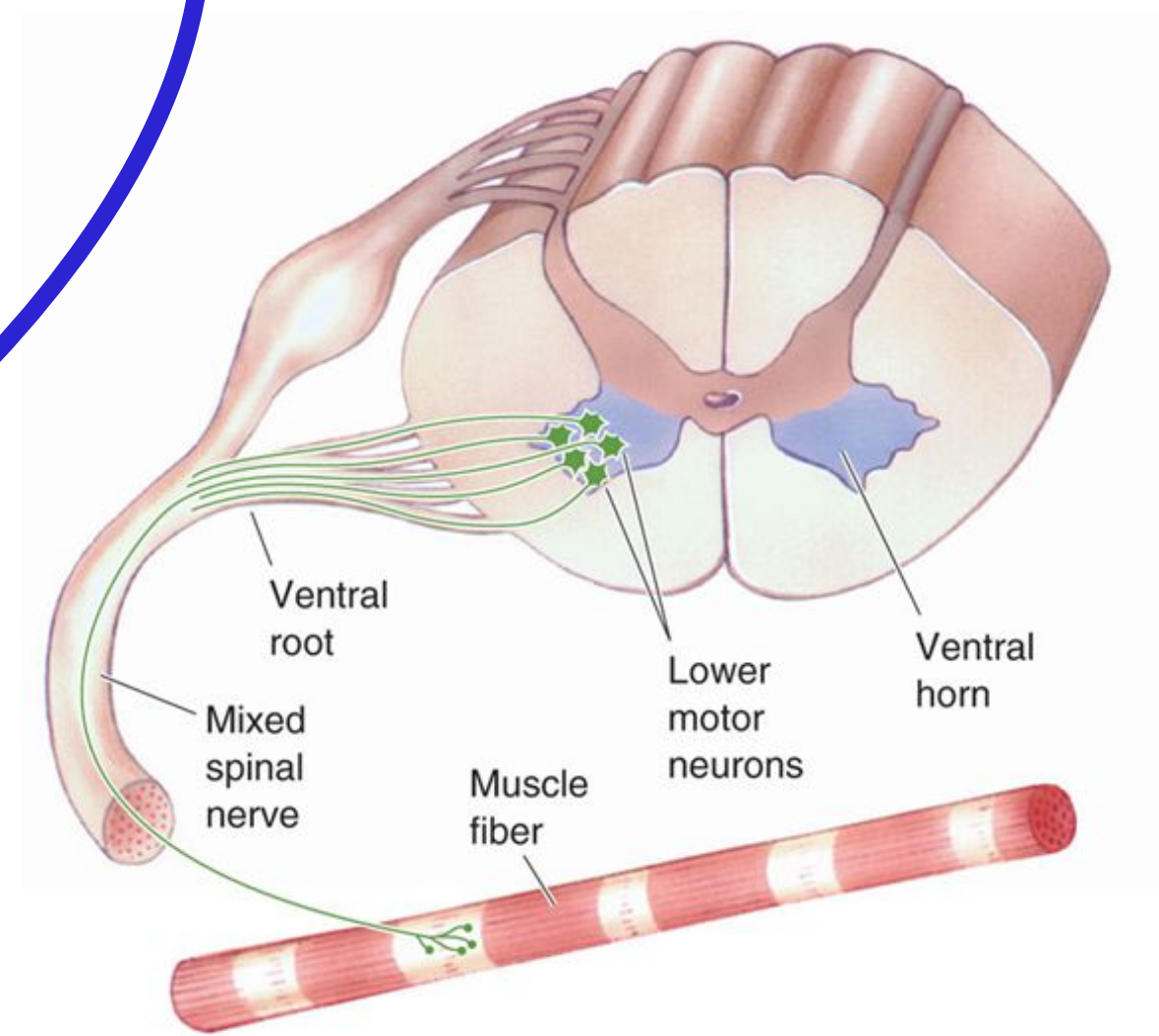
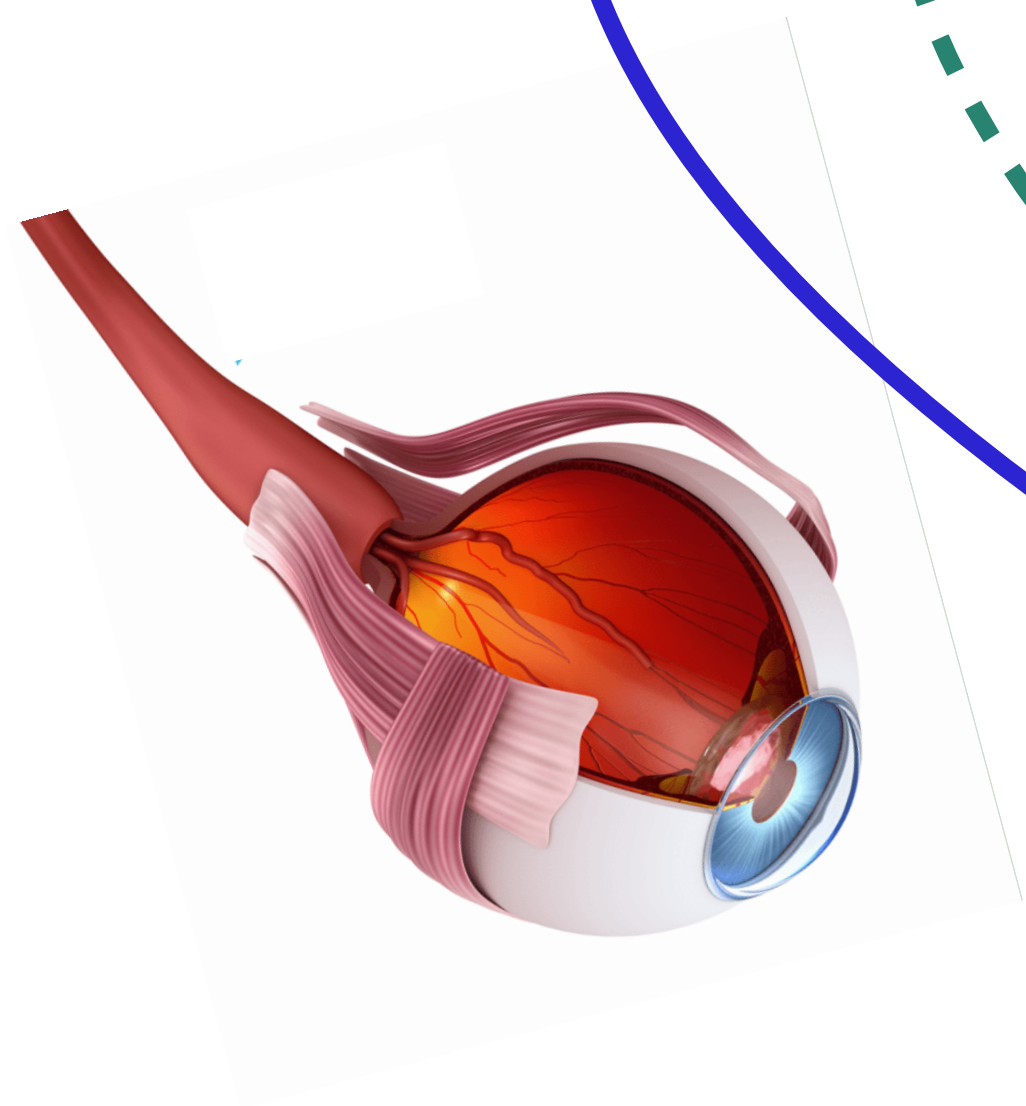
Prediction errors

Predictions

Prediction errors

Sensory data

$y$

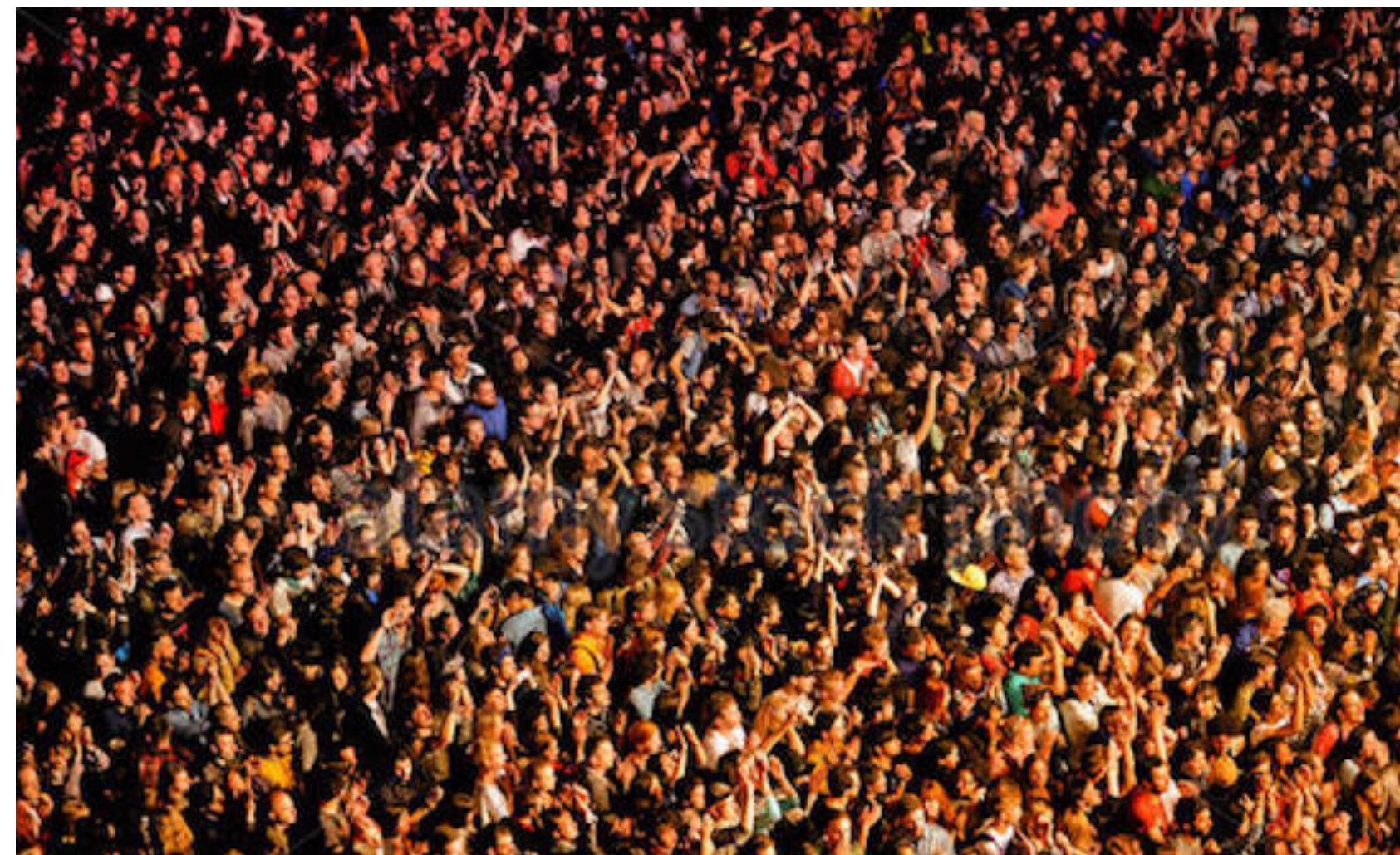
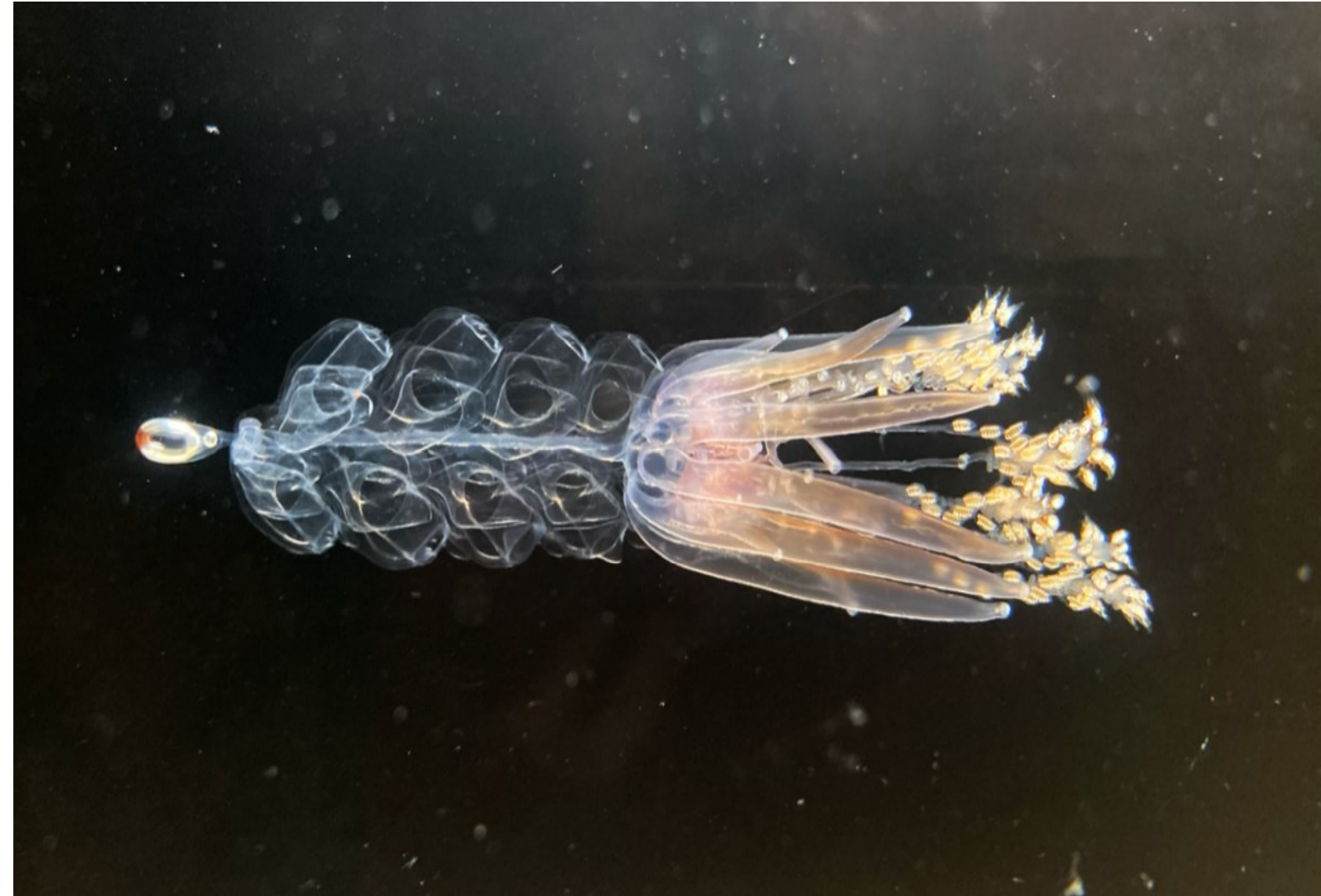




# Applications of active inference

- Eye movements and visual foraging (Friston et al. 2012, Mirza et al. 2016, Parr & Friston 2018, ...)
- Kinematic and postural control (Maselli et al. 2022, Priorelli et al. 2023)
- Embodied spatial decision-making (Priorelli et al. 2024)
- Emotion recognition (Smith et al. 2019, Hesp et al. 2021, Mirza et al. 2021)
- Economic decision-making under uncertainty (Smith et al. 2020, Markovic et al. 2021)
- Language understanding and speech (Parr & Pezzulo 2021, Friston et al. 2020)
- Active sensing — e.g., whisking in rodents (Mannella et al. 2021)

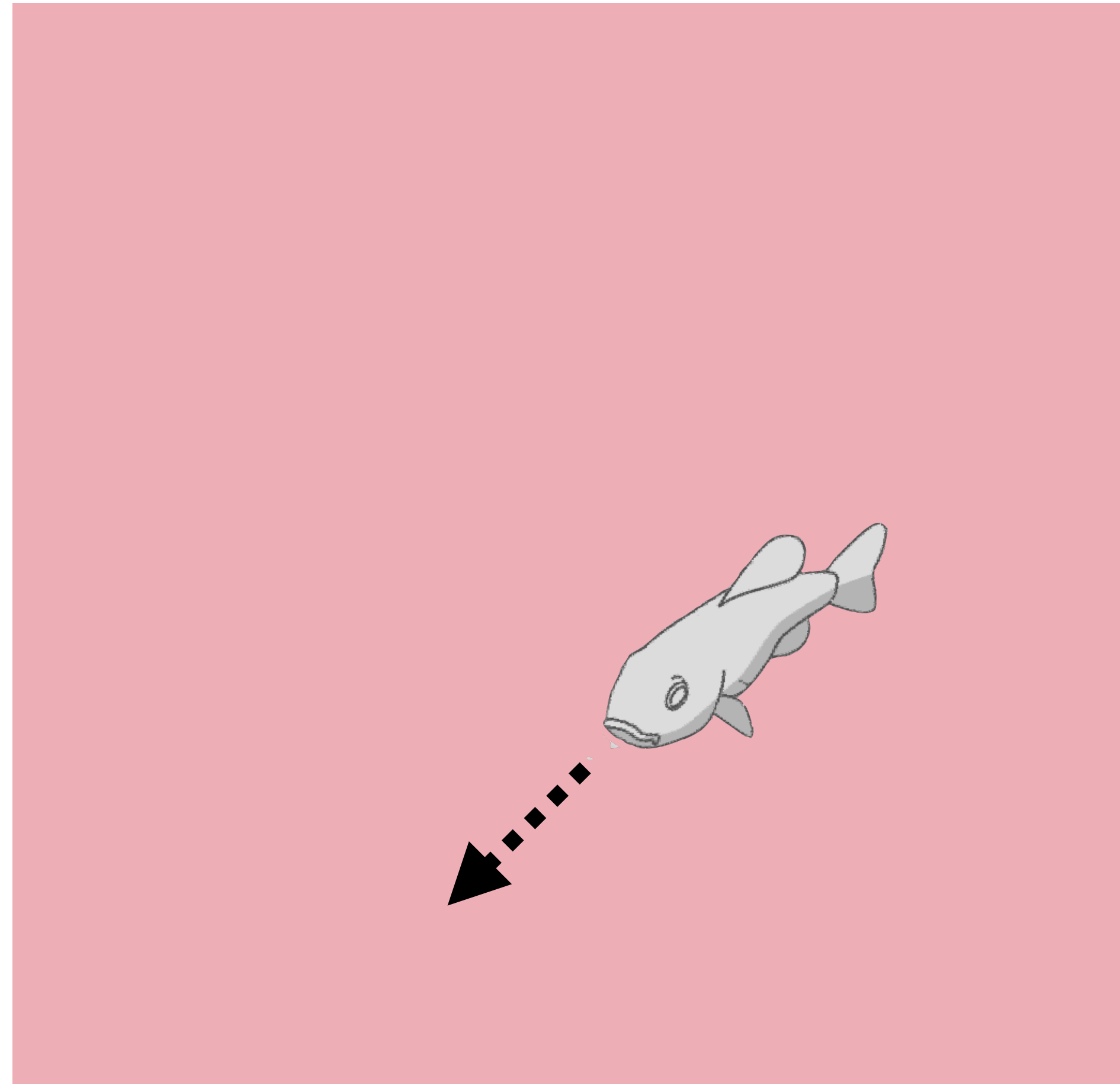
# Collective motion in nature



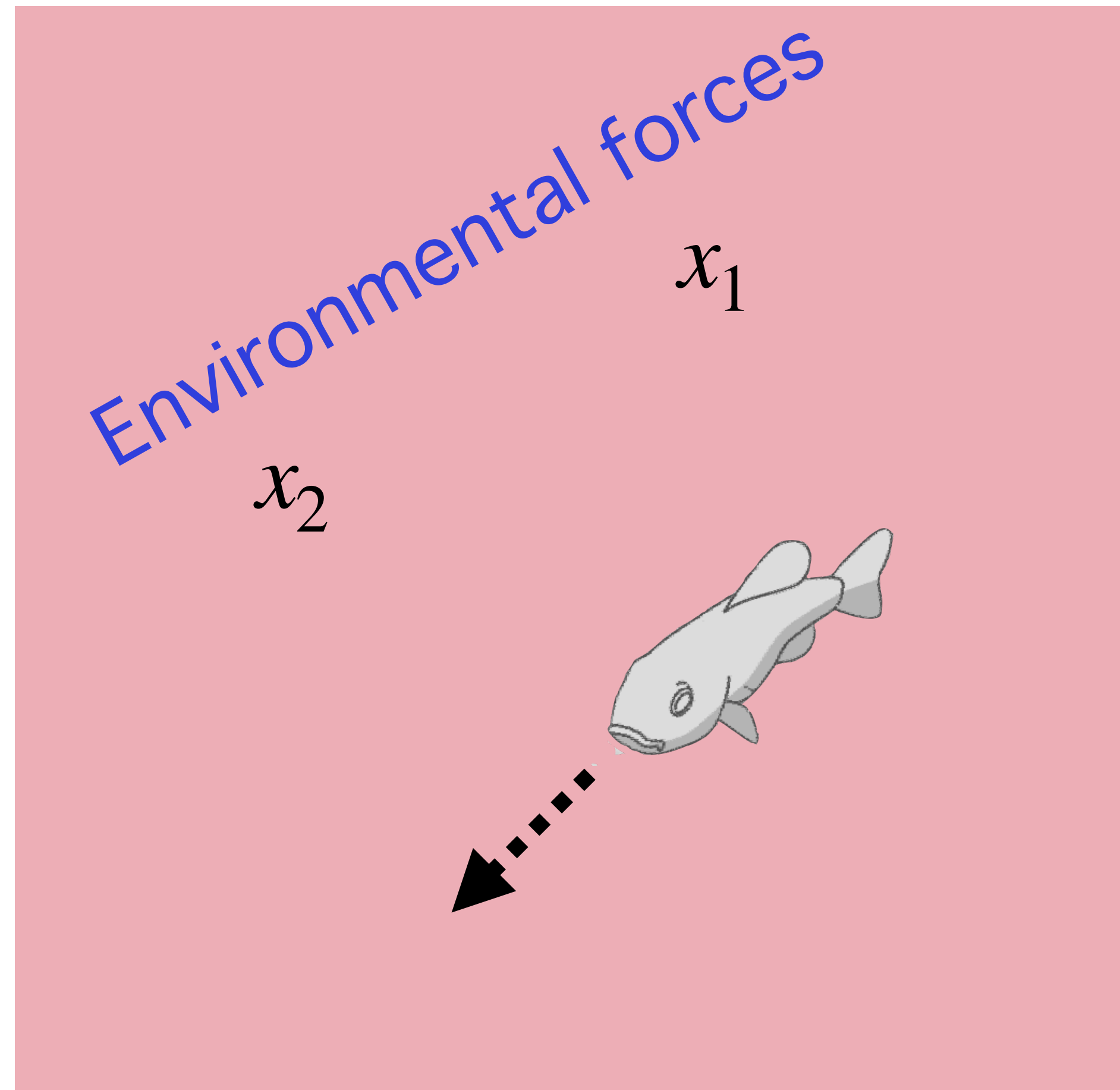
# Collective motion in nature

- Example of emergent order from simple, decentralised interactions
- Universality? Shows up across natural, technological, artificial disciplines
- Collective motion is relevant to biologists for many reasons (evolutionary, ecological, cognitive, neurobiological)

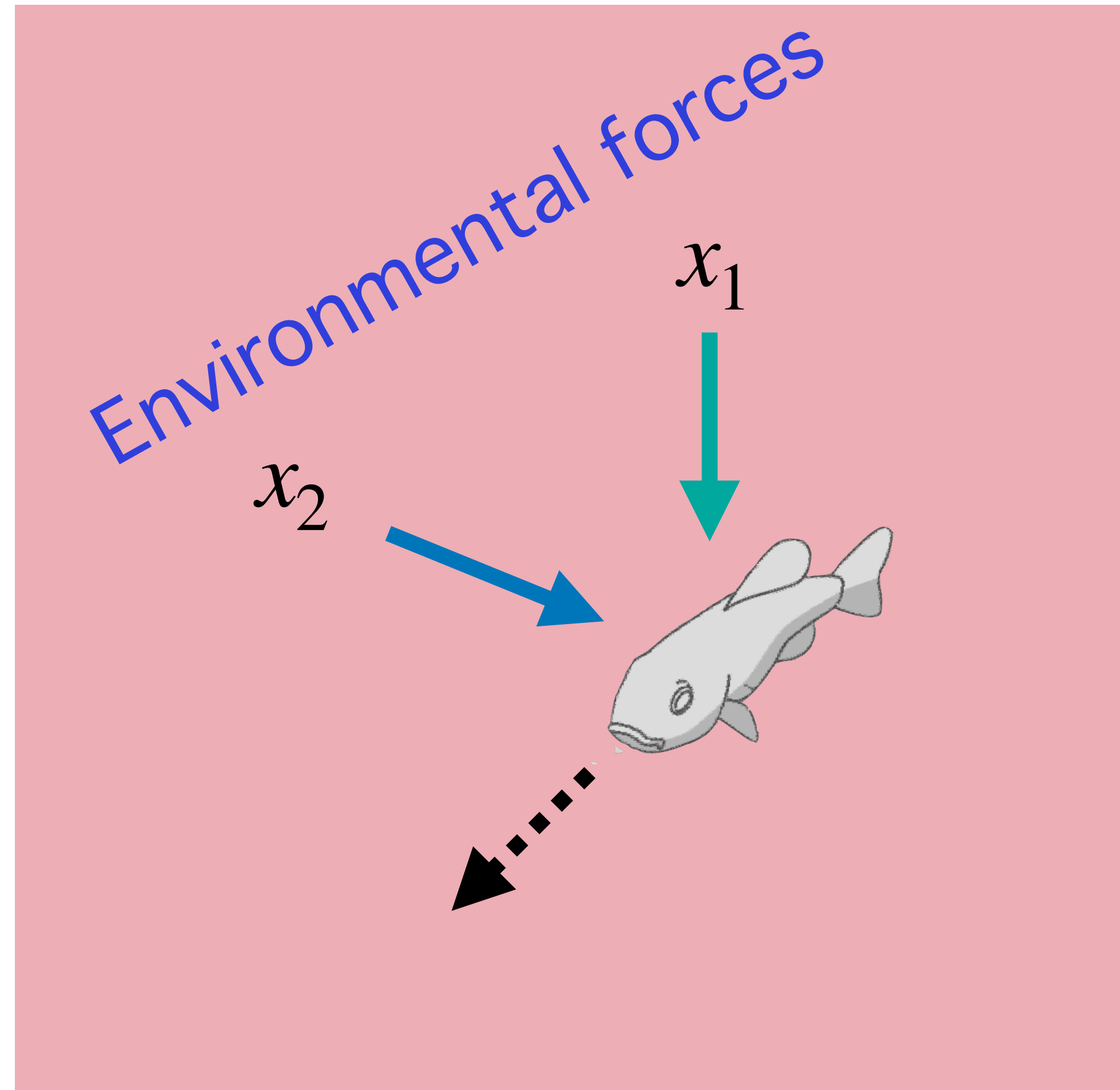
# Physics-based models (Vicsek model, etc.)



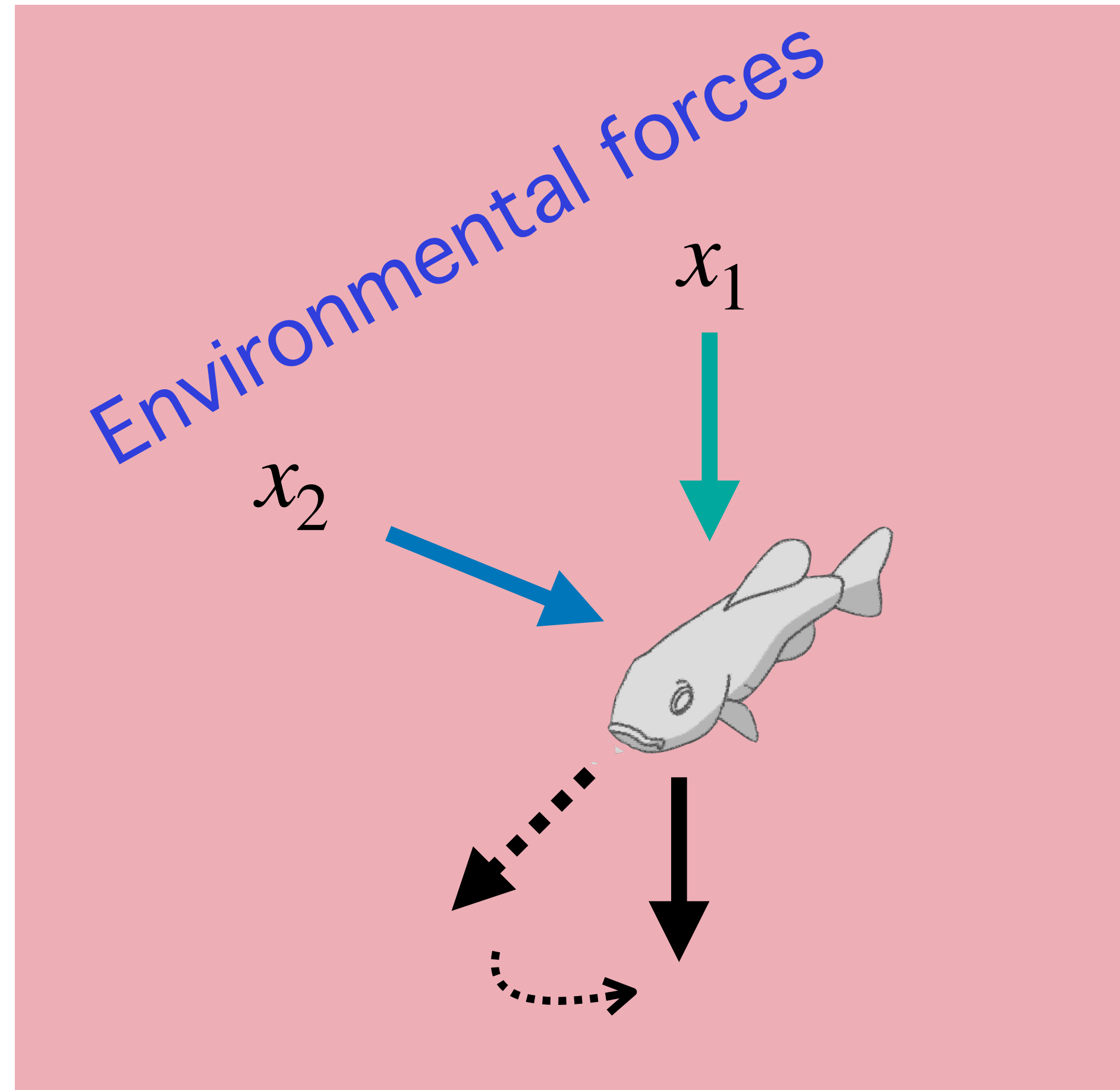
# Physics-based models (Vicsek model, etc.)



# Physics-based models (Vicsek model, etc.)

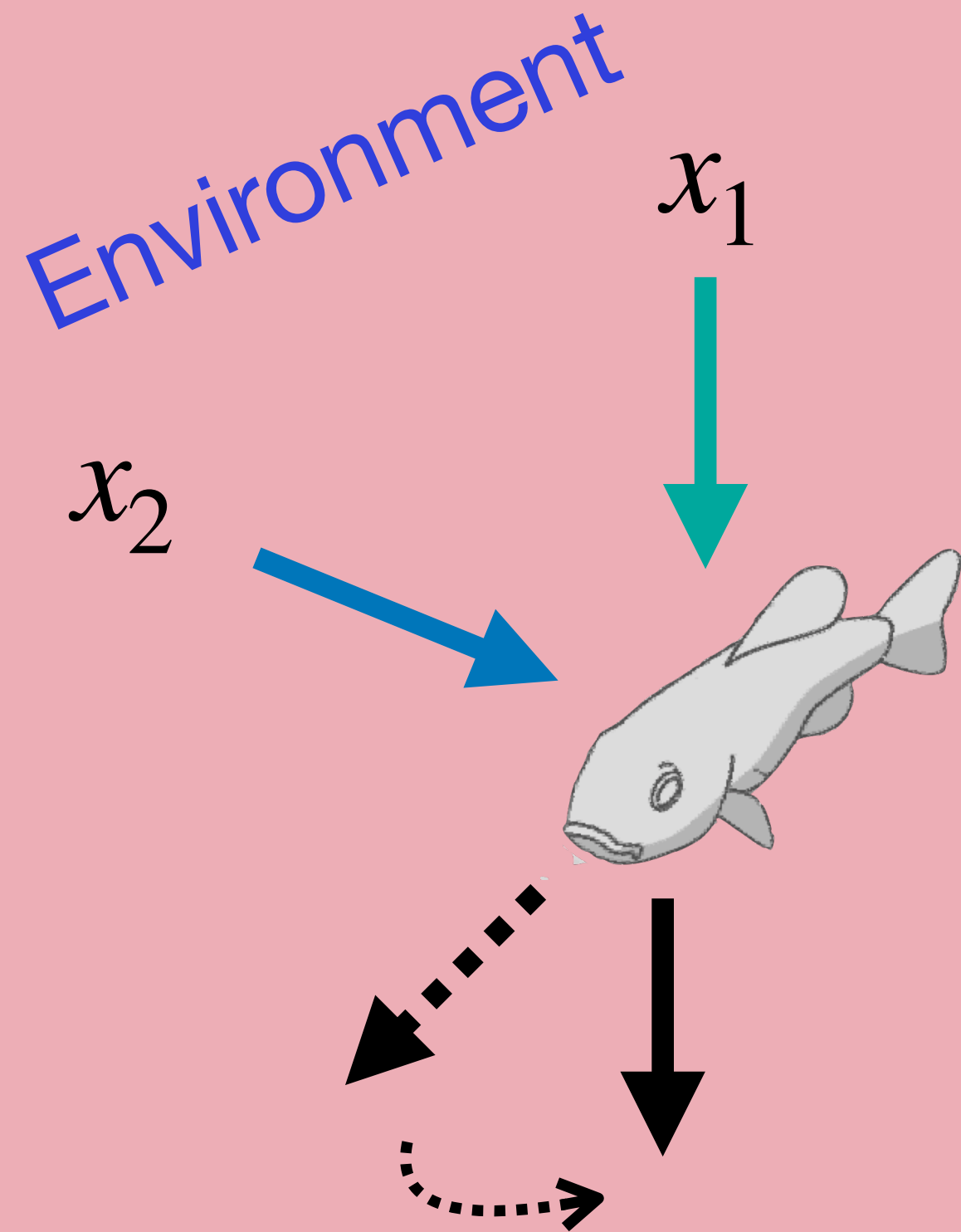


# Physics-based models (Vicsek model, etc.)



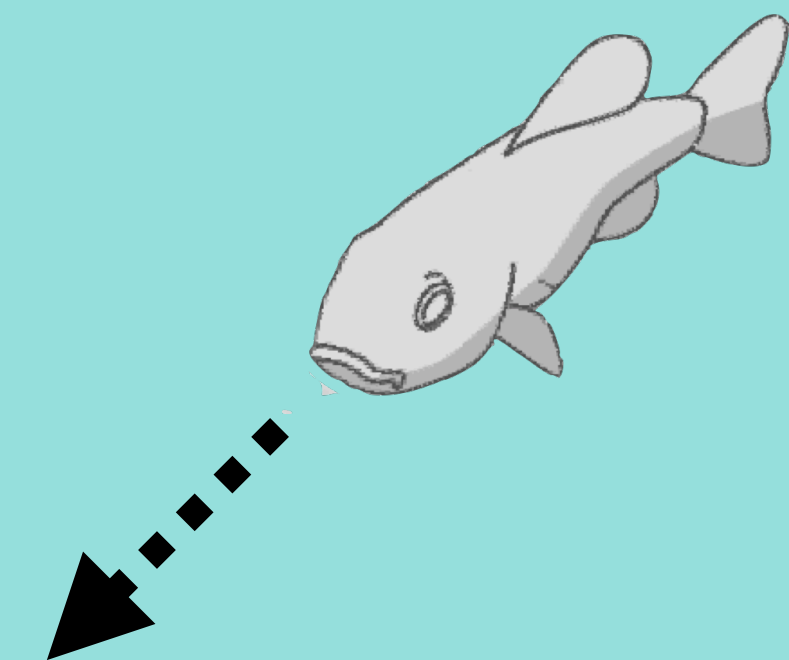
# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)



Explicitly *belief*-based approach

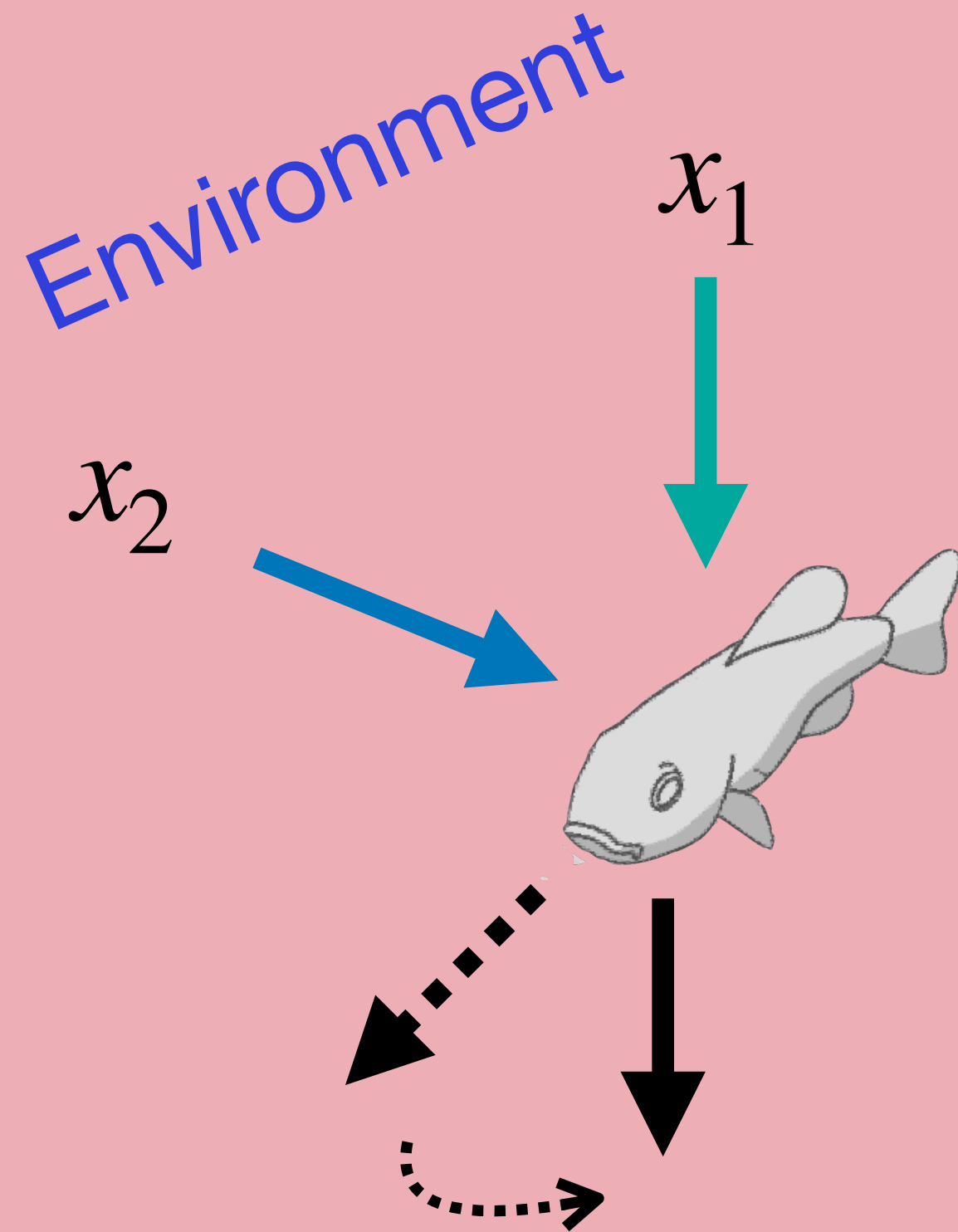
Environment  
 $x_1$   
 $x_2$



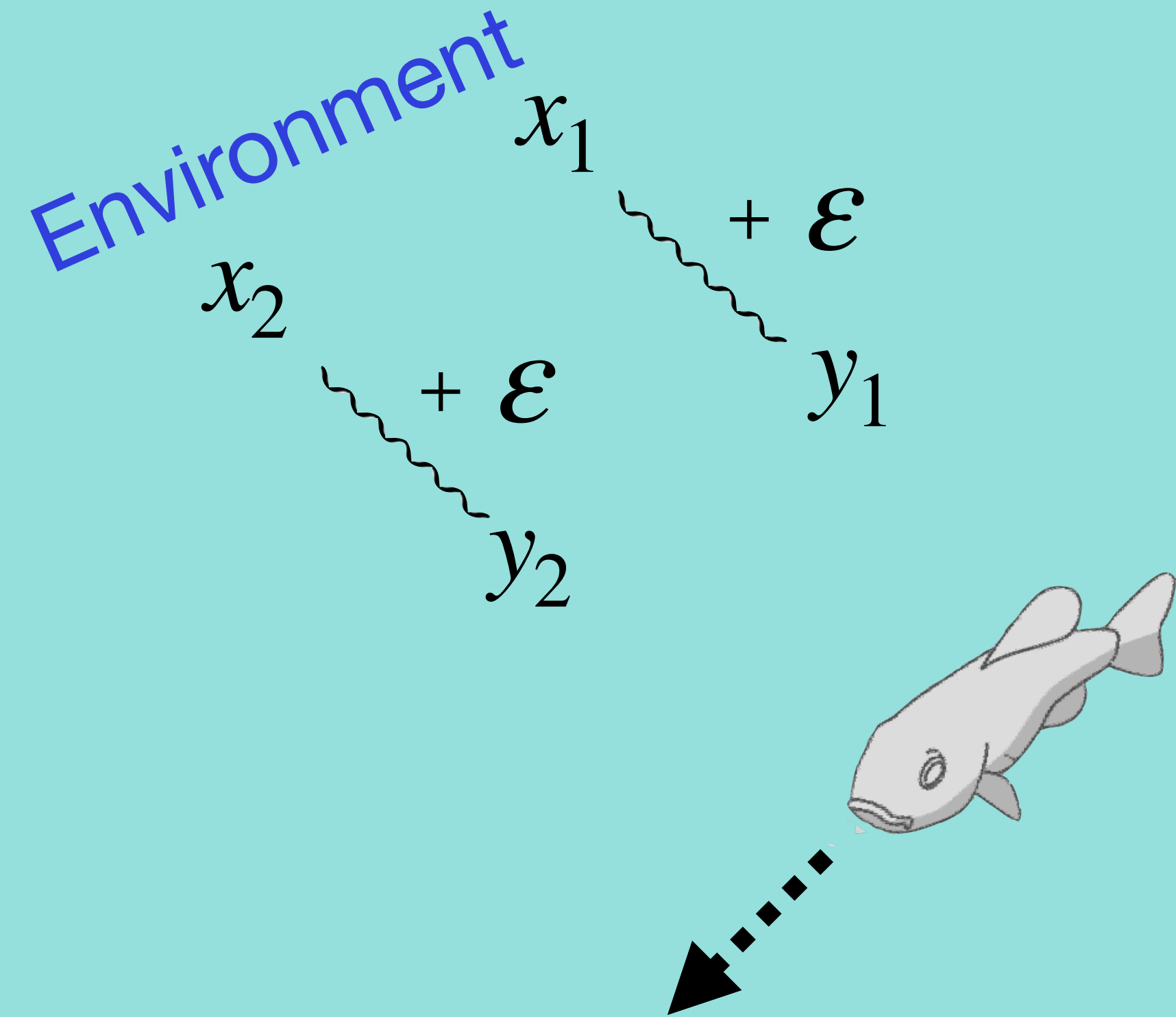


# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)

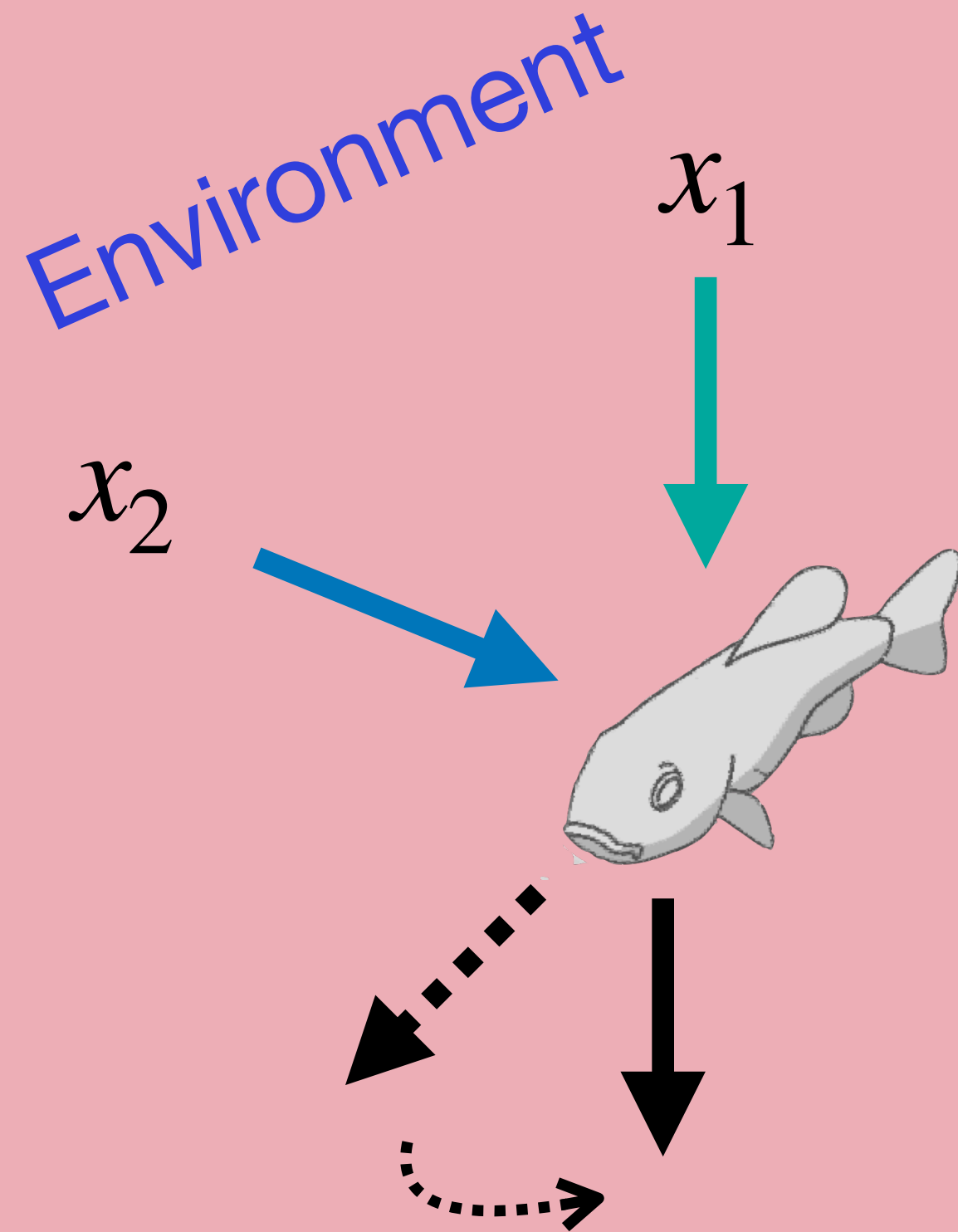


Explicitly *belief*-based approach

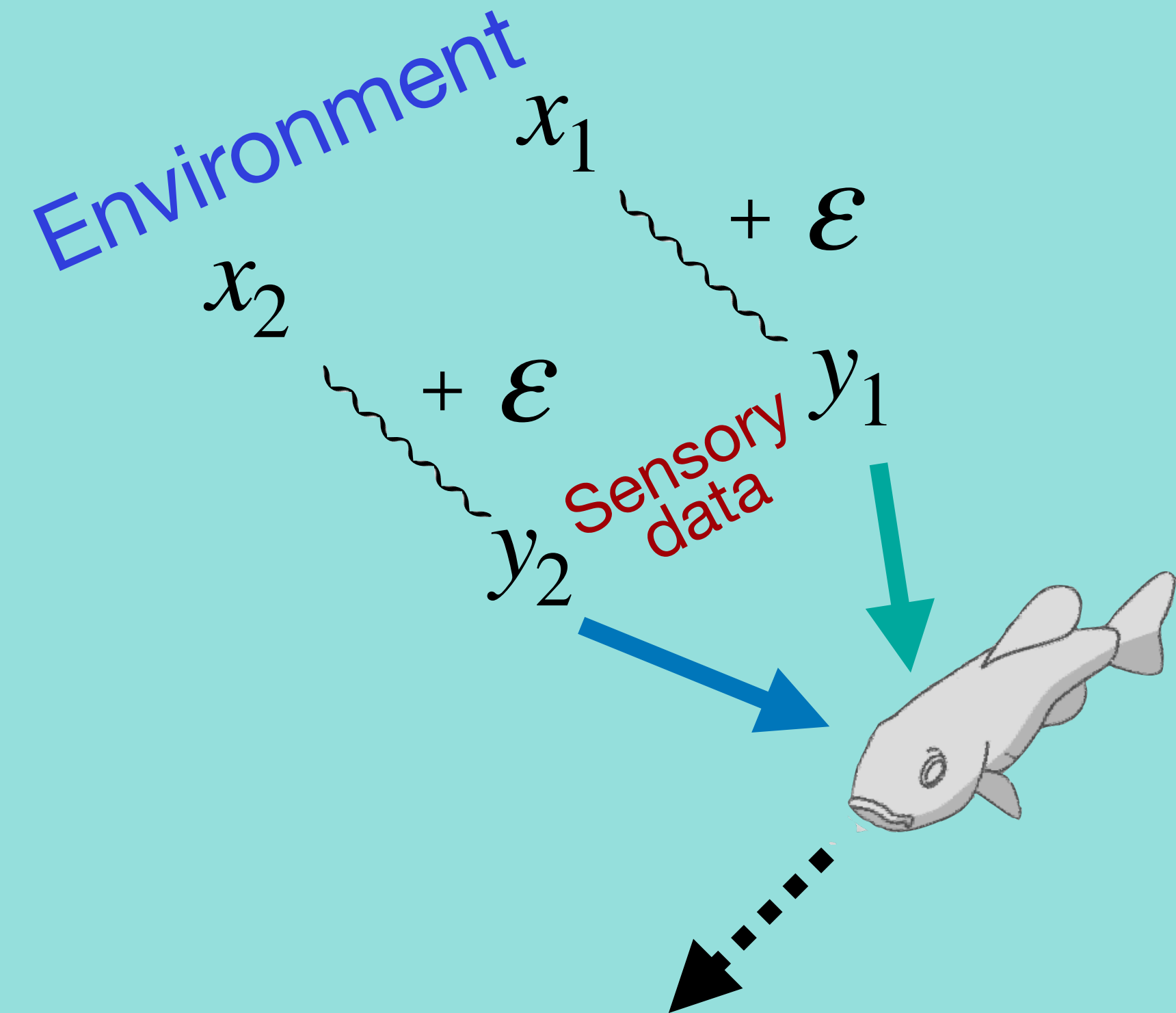


# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)

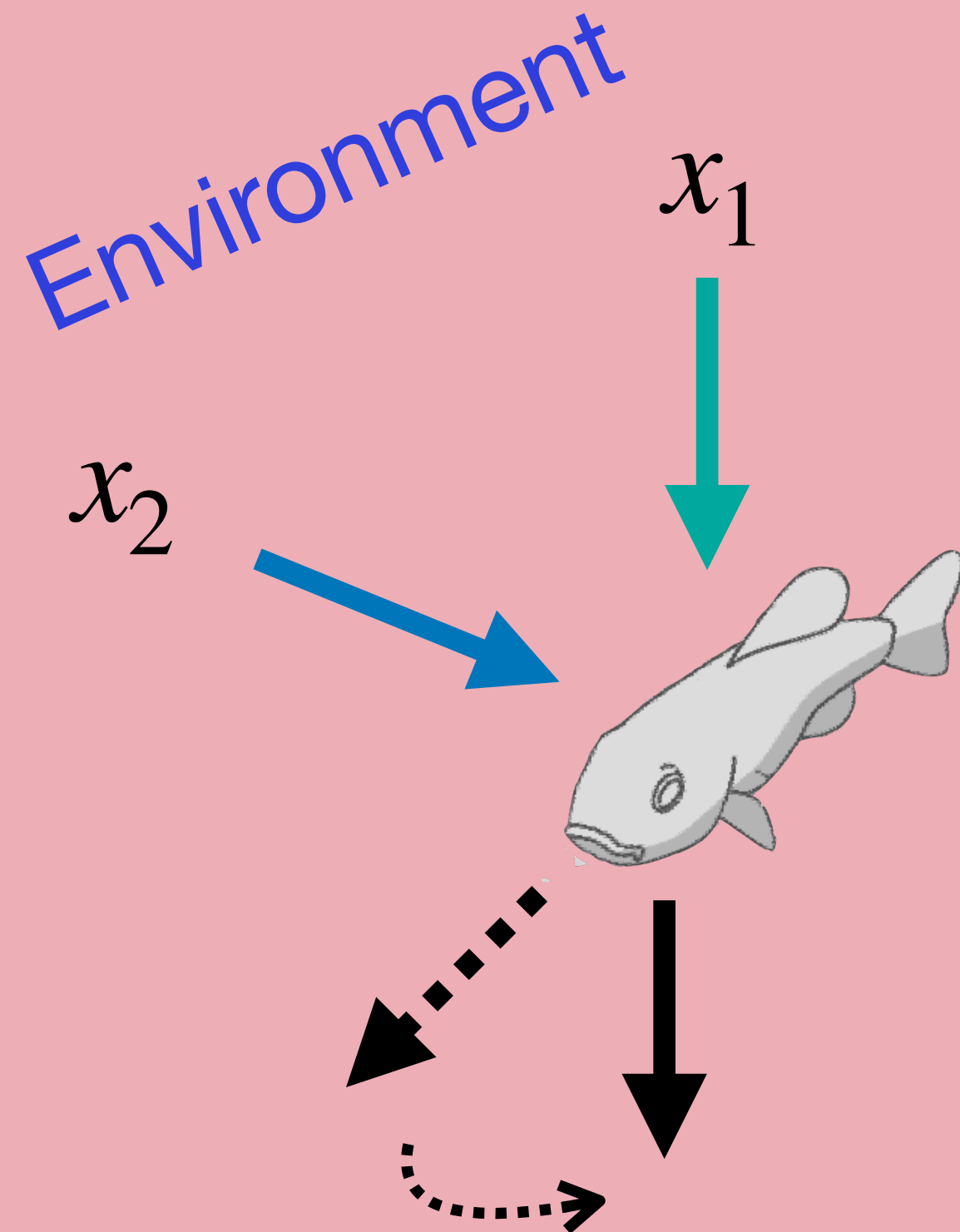


Explicitly *belief*-based approach

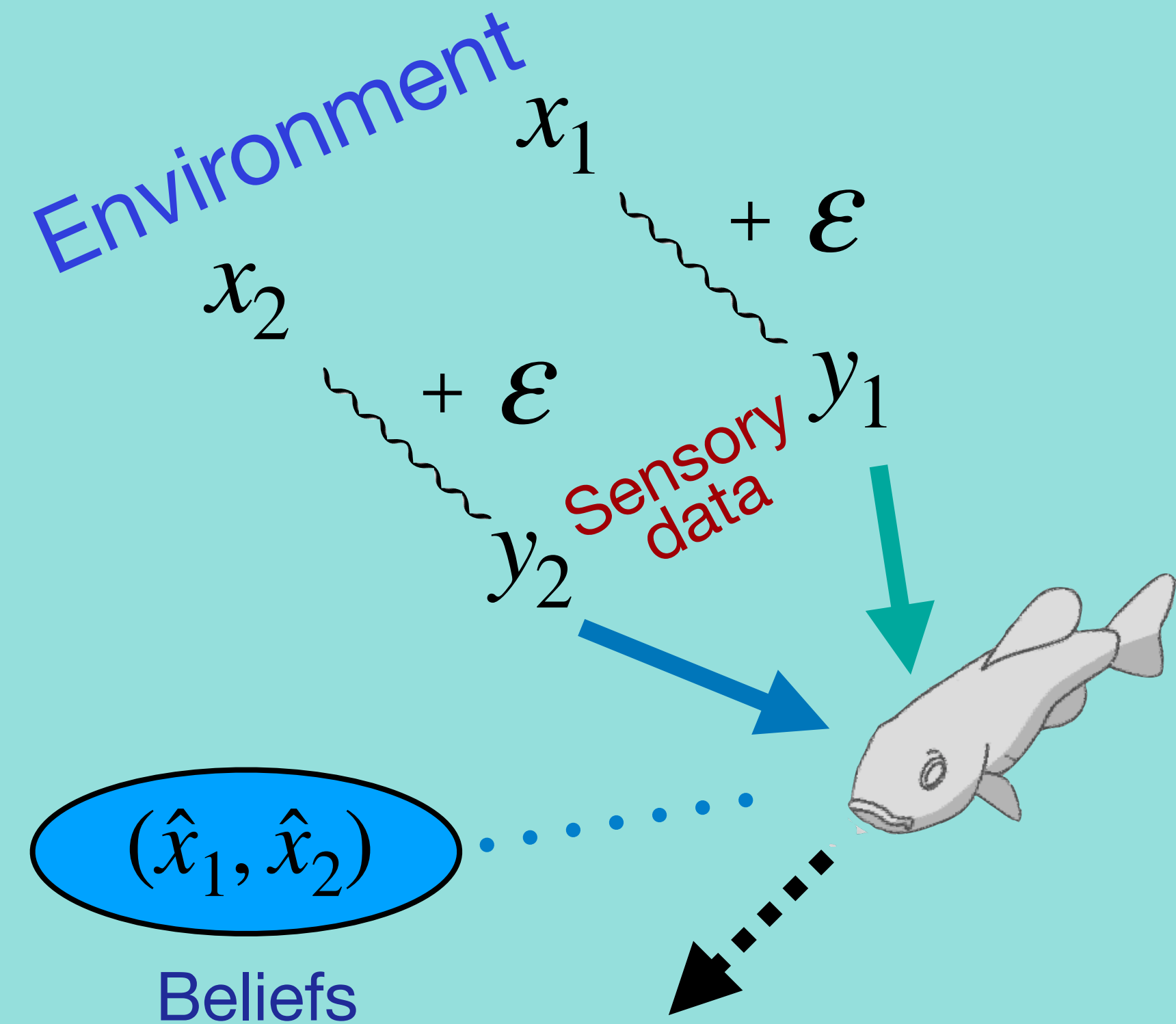


# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)

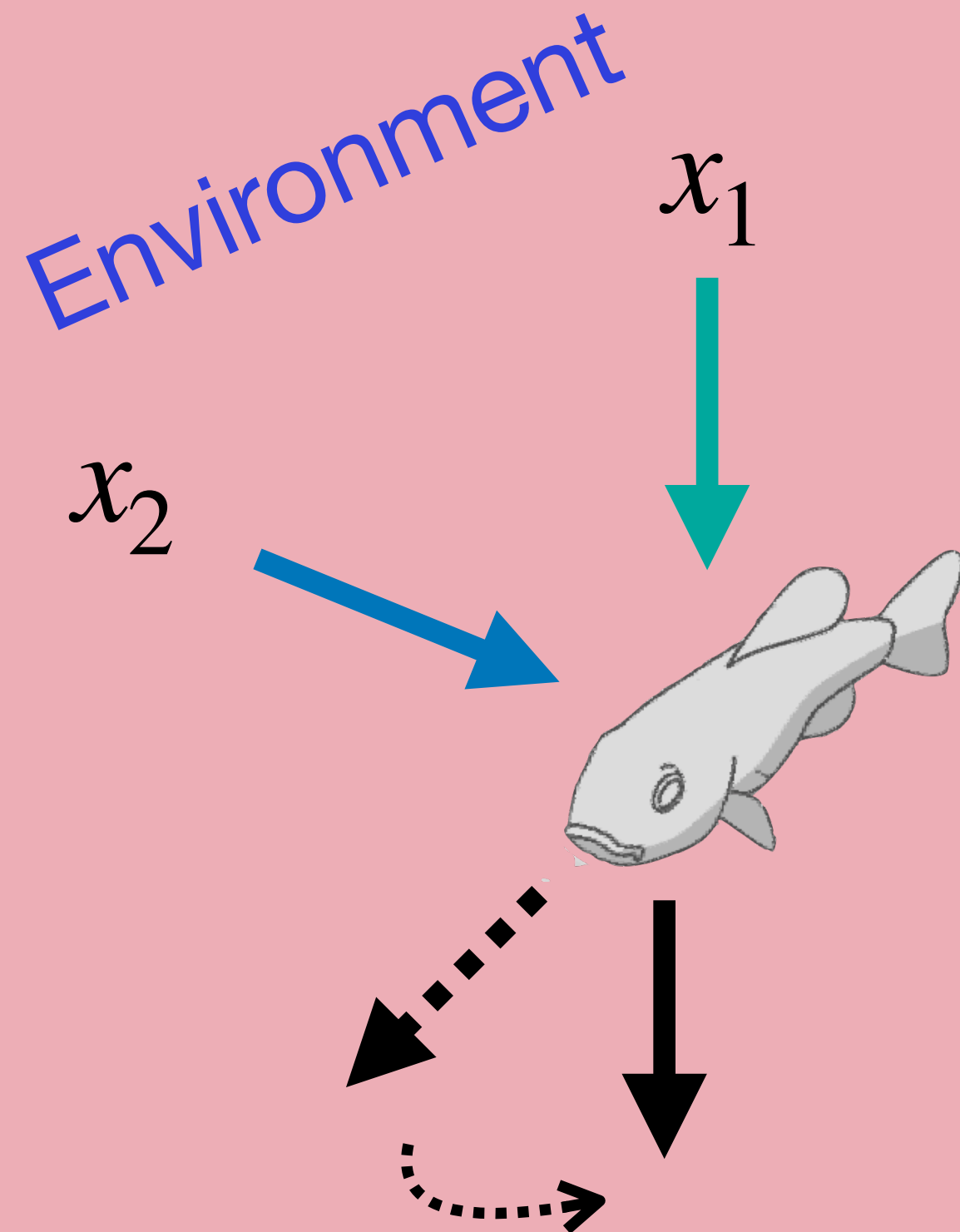


Explicitly *belief*-based approach

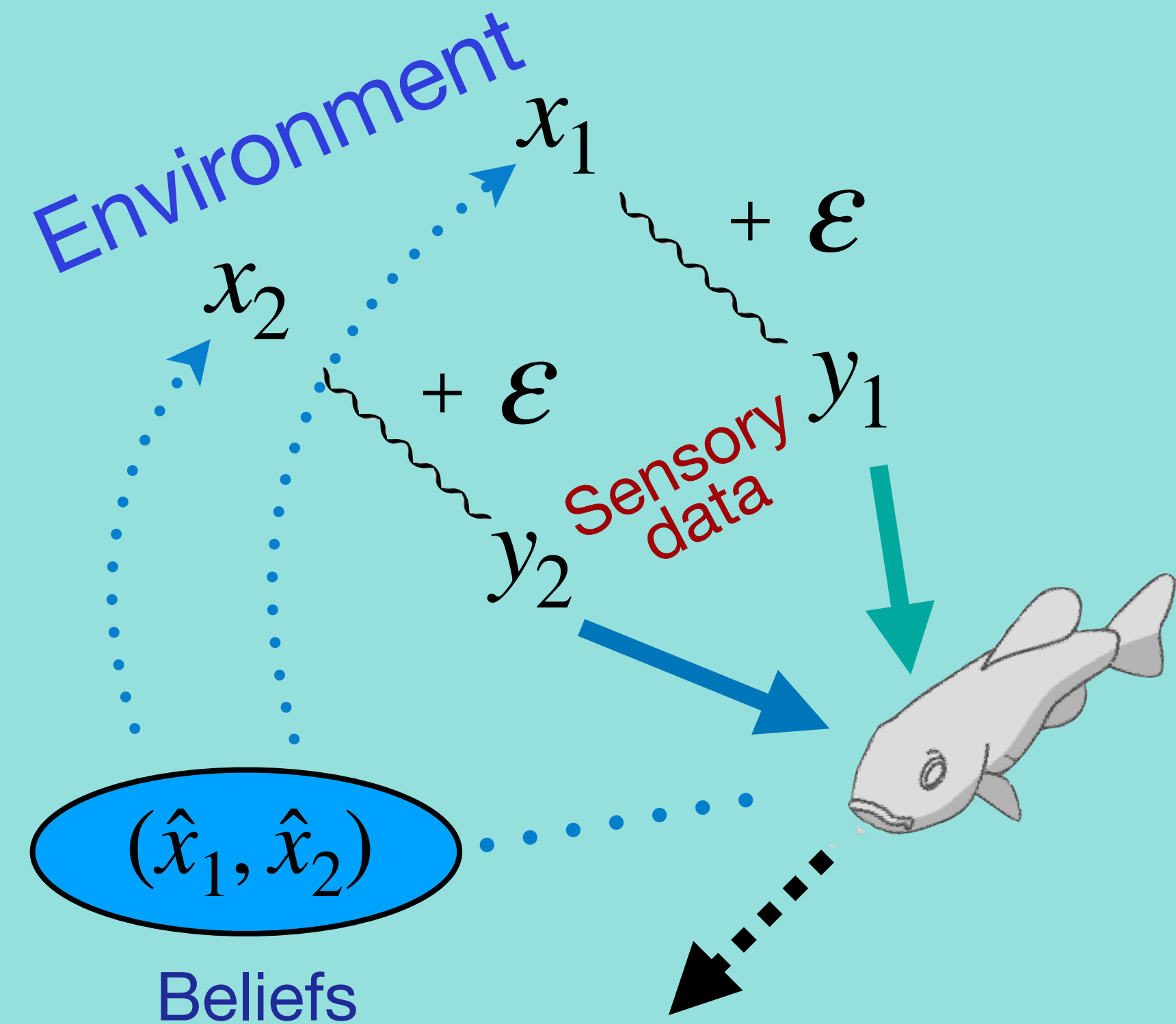


# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)

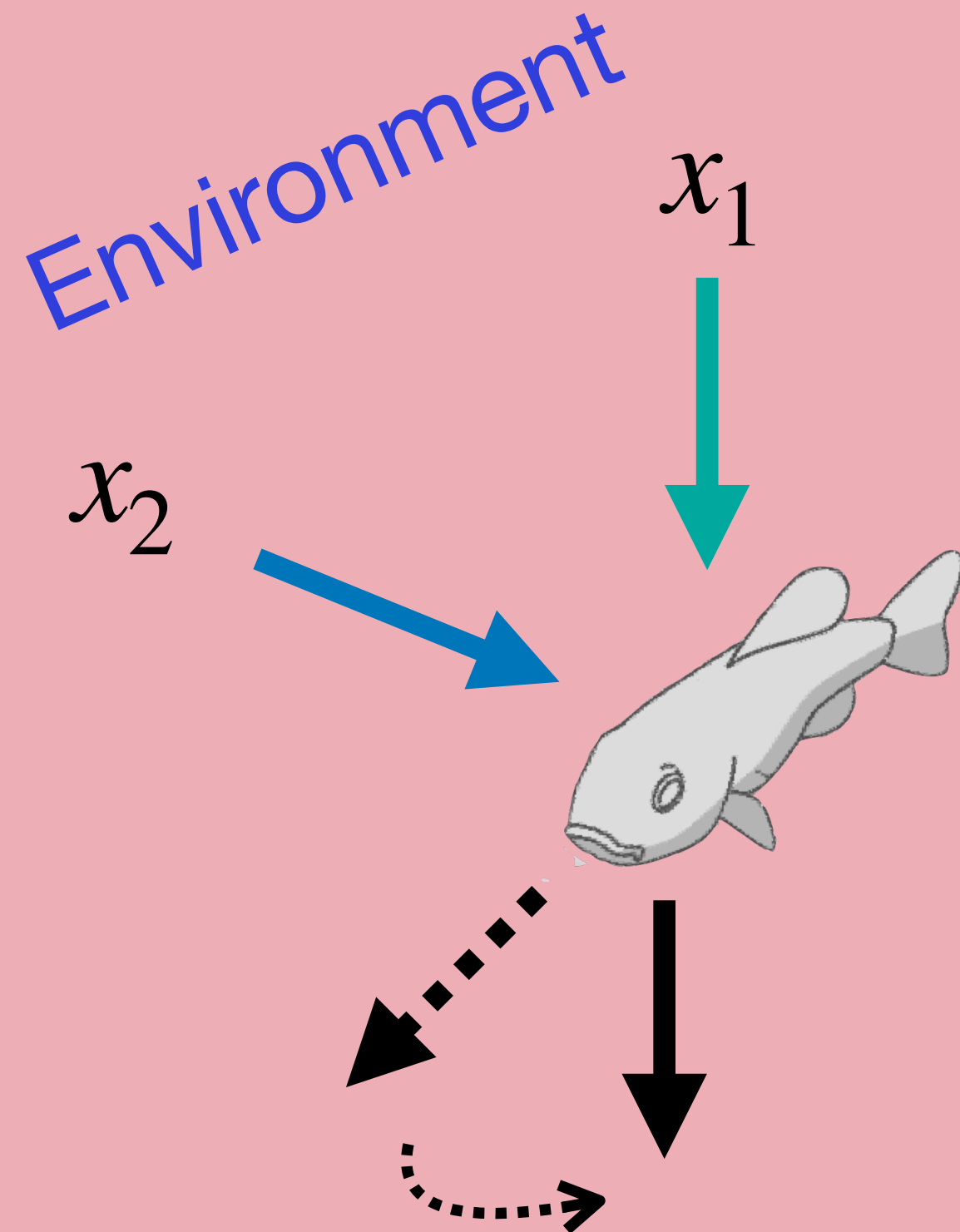


Explicitly *belief*-based approach

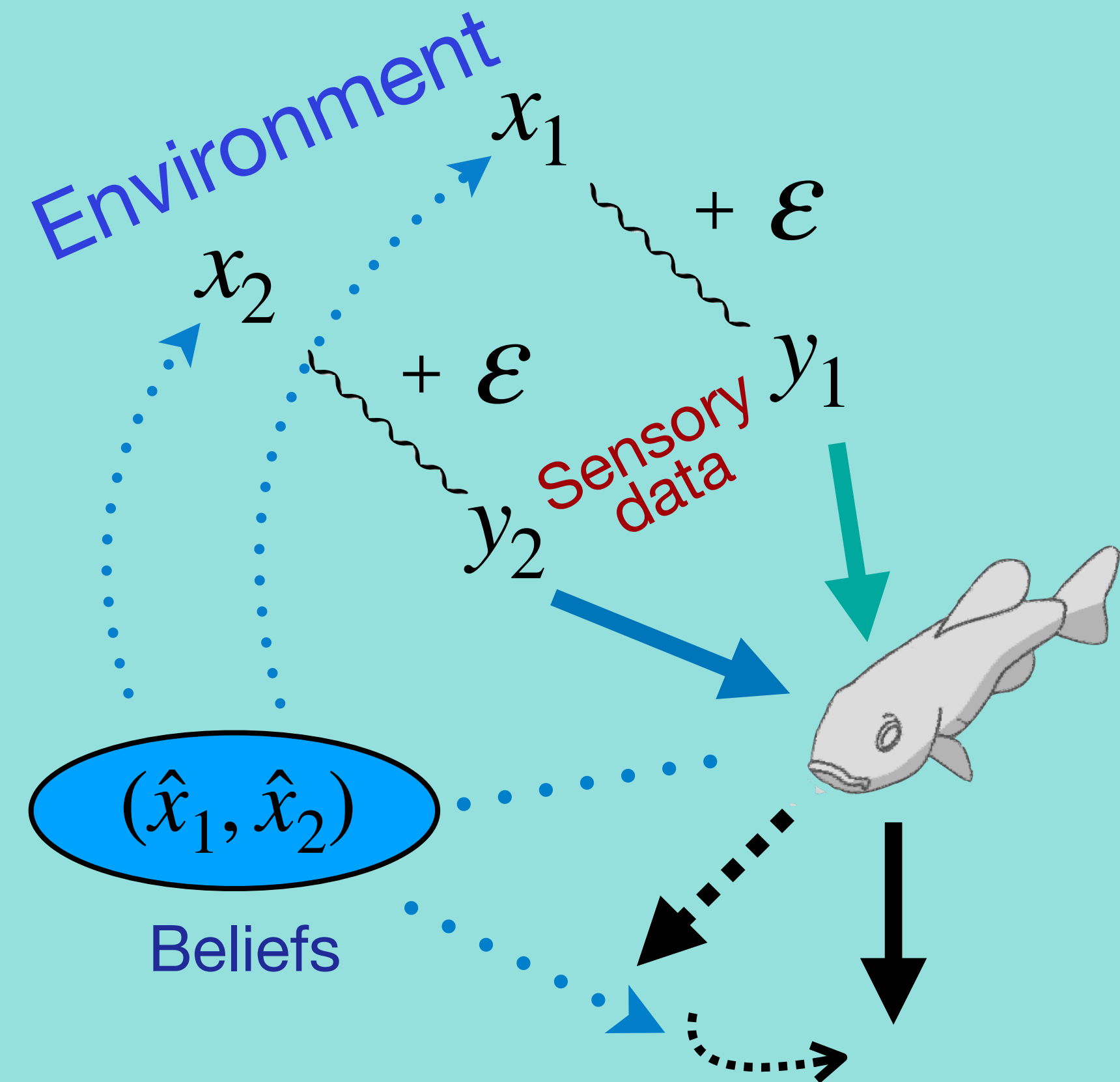


# Comparing classical and Bayesian approaches

“Classical” force-based approach  
(e.g. self-propelled particle models)



Explicitly *belief*-based approach

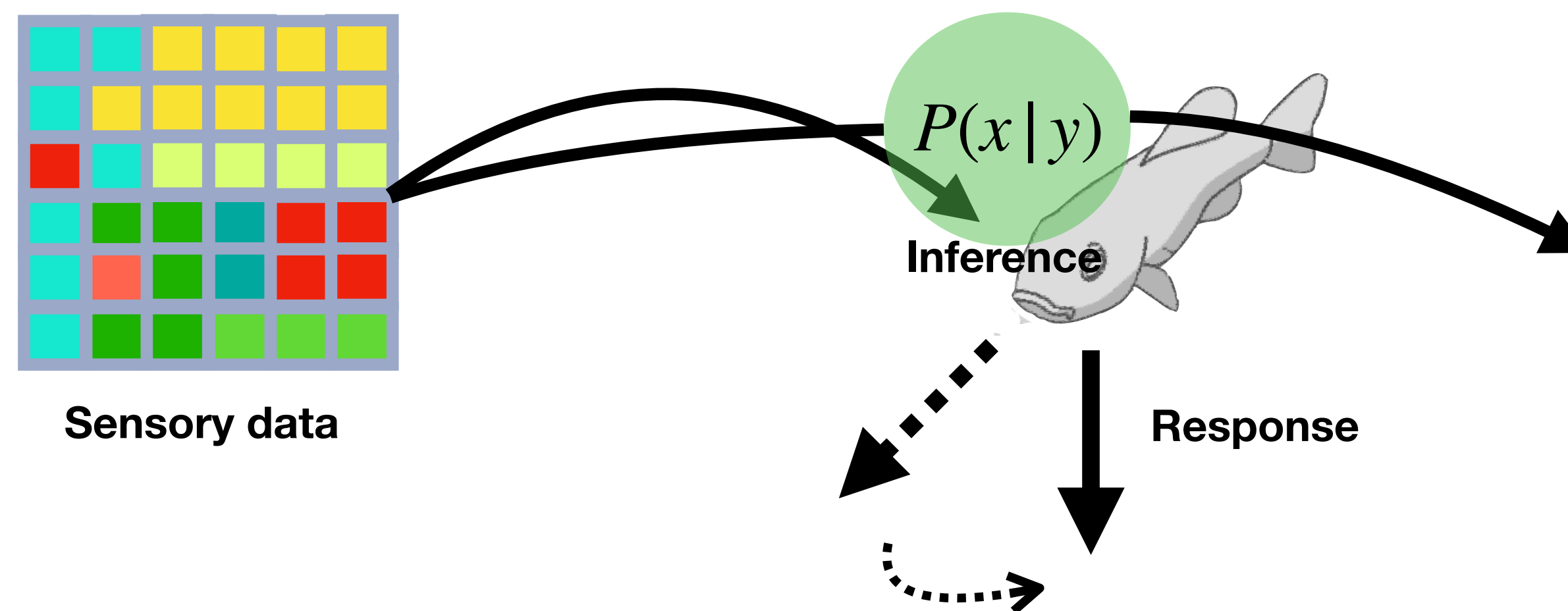


## PHYSICS

# A model of collective behavior based purely on vision

Renaud Bastien<sup>1,2,\*†</sup> and Pawel Romanczuk<sup>3,4†</sup>

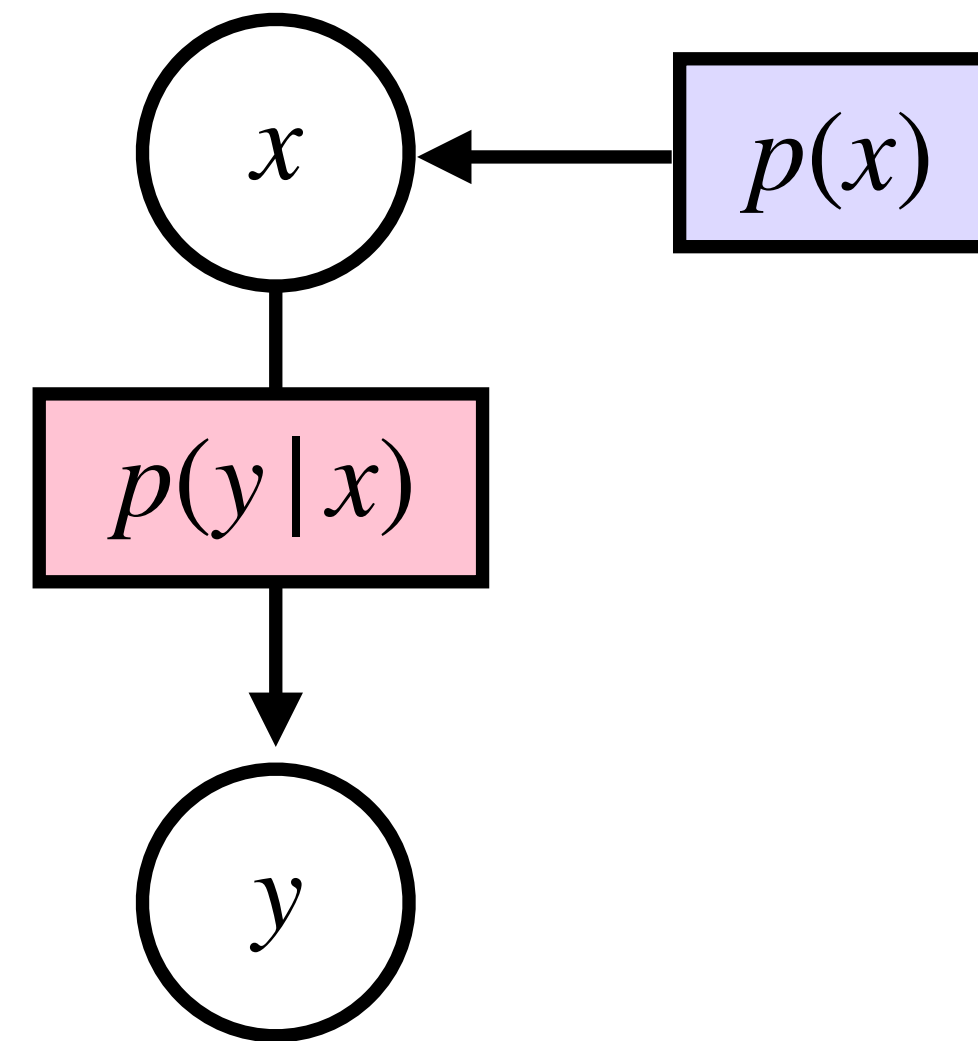
*“Collective behavior crucially depends on the sensory information available to individuals; thus, **ignoring perception by relying on ad hoc rules strongly limits our understanding of the underlying complexity of the problem.** Besides, it obstructs the interdisciplinary exchange between biology, neuroscience, engineering, and physics.”*



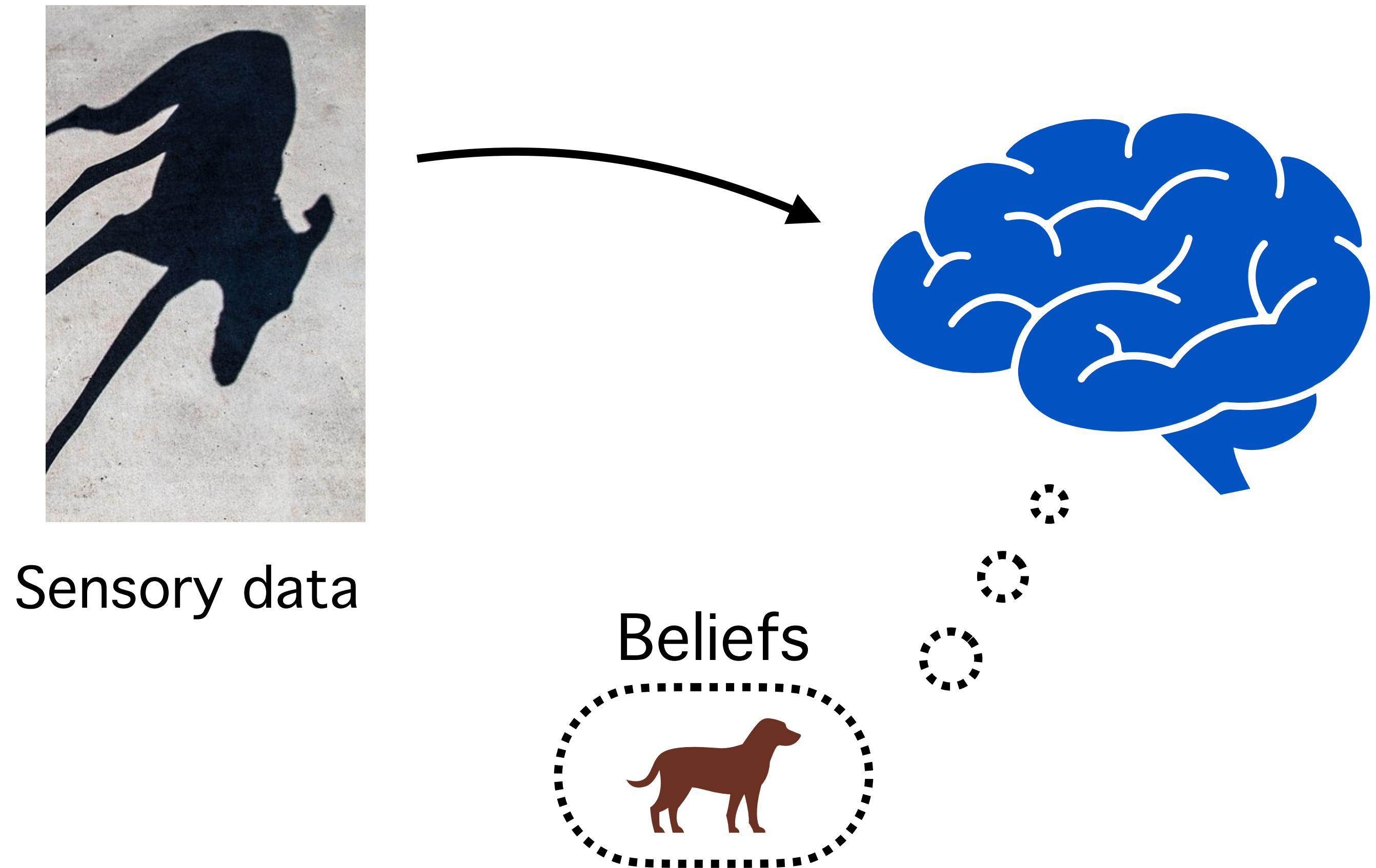
# Writing down an agent's world model a.k.a. the generative model

Generative model

$$p(y, x) = p(y | x)p(x)$$



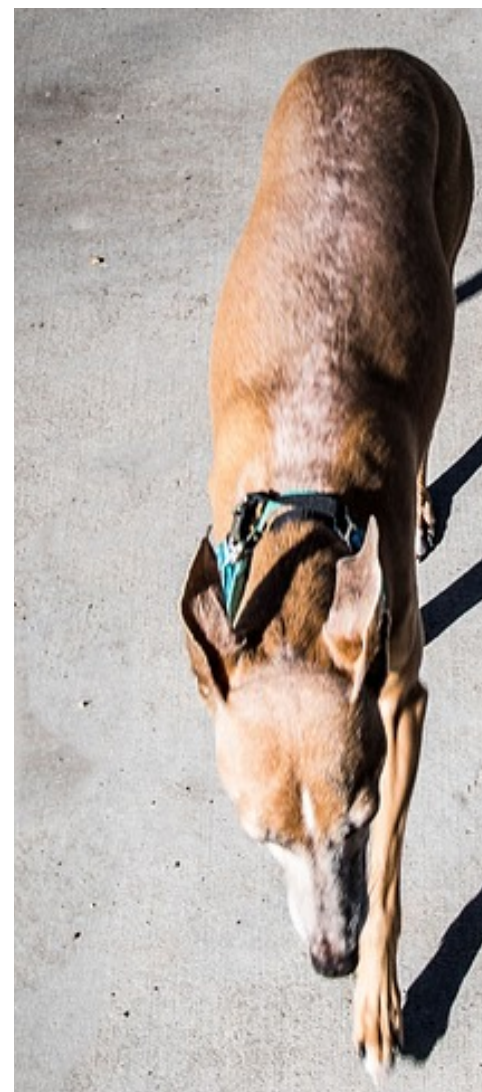
# Generative model needed for inference





# Generative model needed for inference

Generative model captures in-built assumptions about optics, light refraction, prevalence of objects, etc.



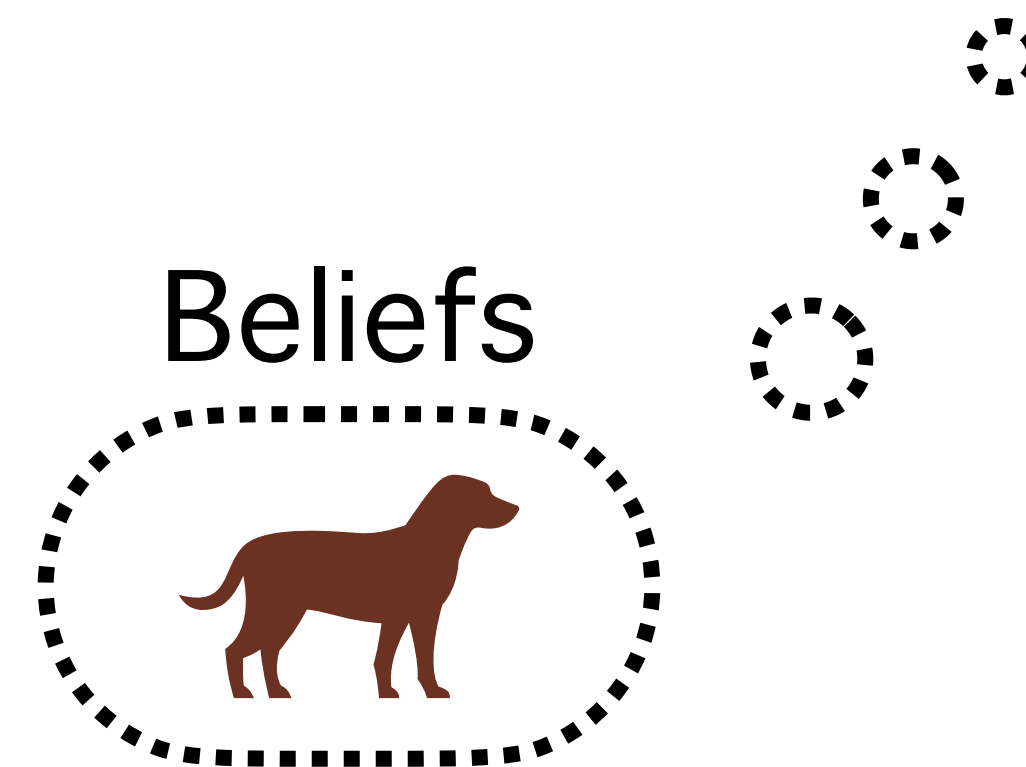
Hidden states



etc.



Sensory data



Writing down an agent's internal model  
a.k.a. the generative model

Generative model

$$p(y, x) = p(y | x)p(x)$$



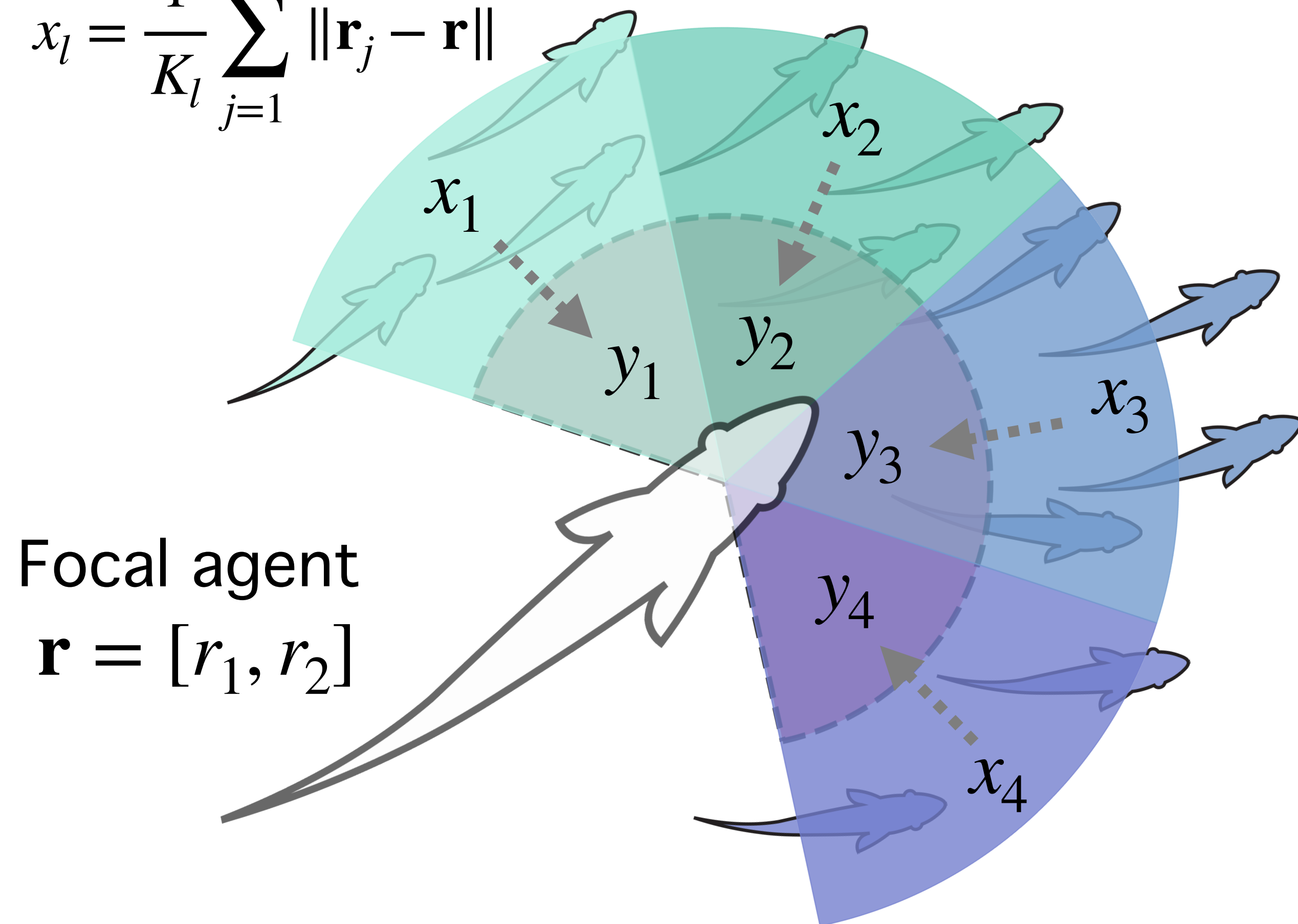
What sort of generative model  
might an individual in a mobile group have?

# Generative model for an individual

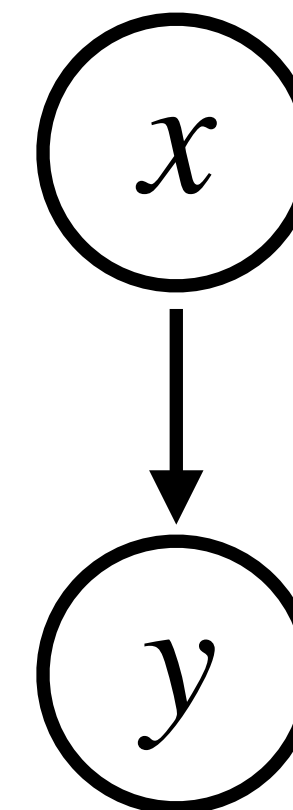
$$p(y, x) = p(y | x)p(x)$$

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



Hidden states  $x_t$  comprise the agent's environment

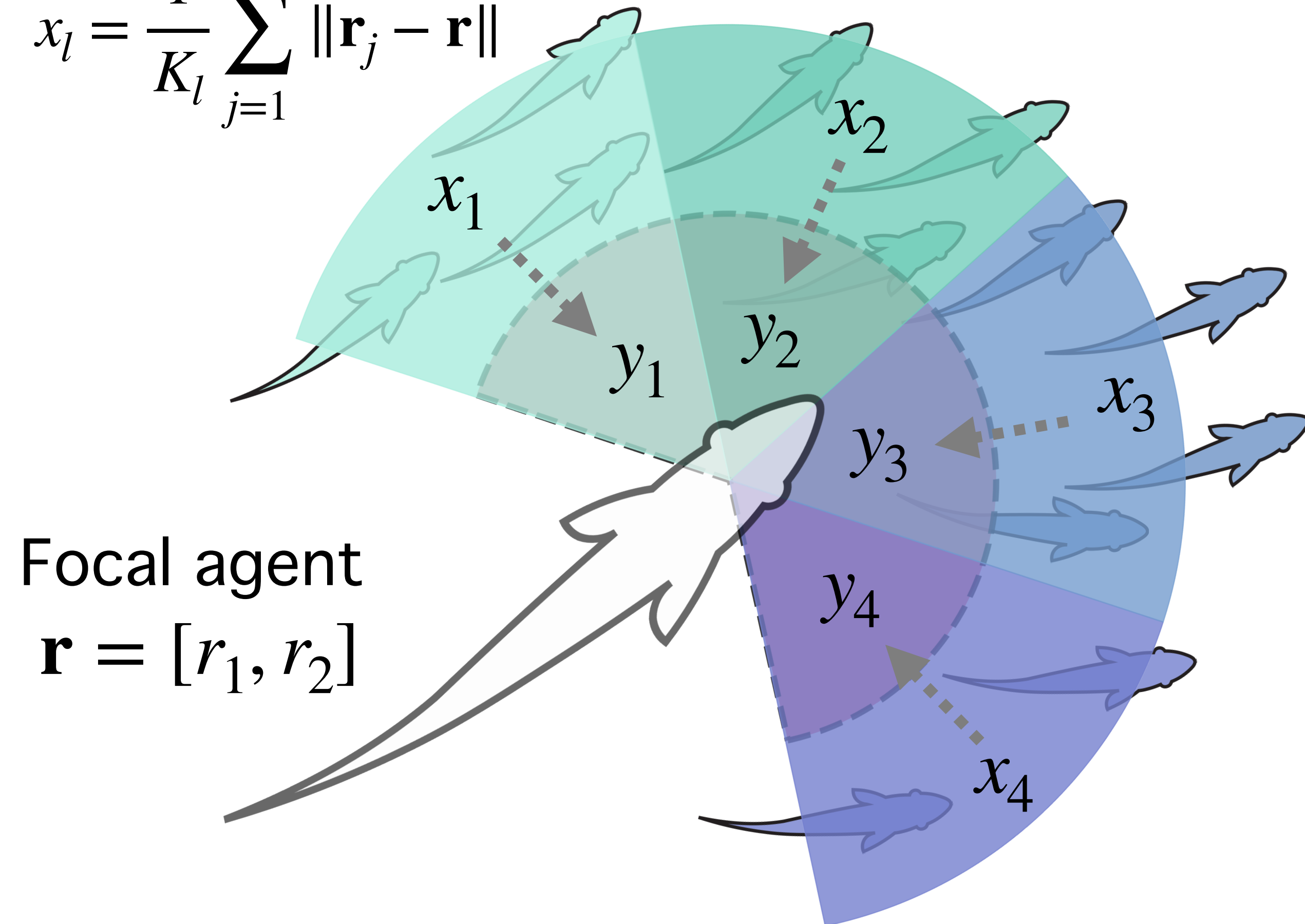


# Priors about social distance $x_l$

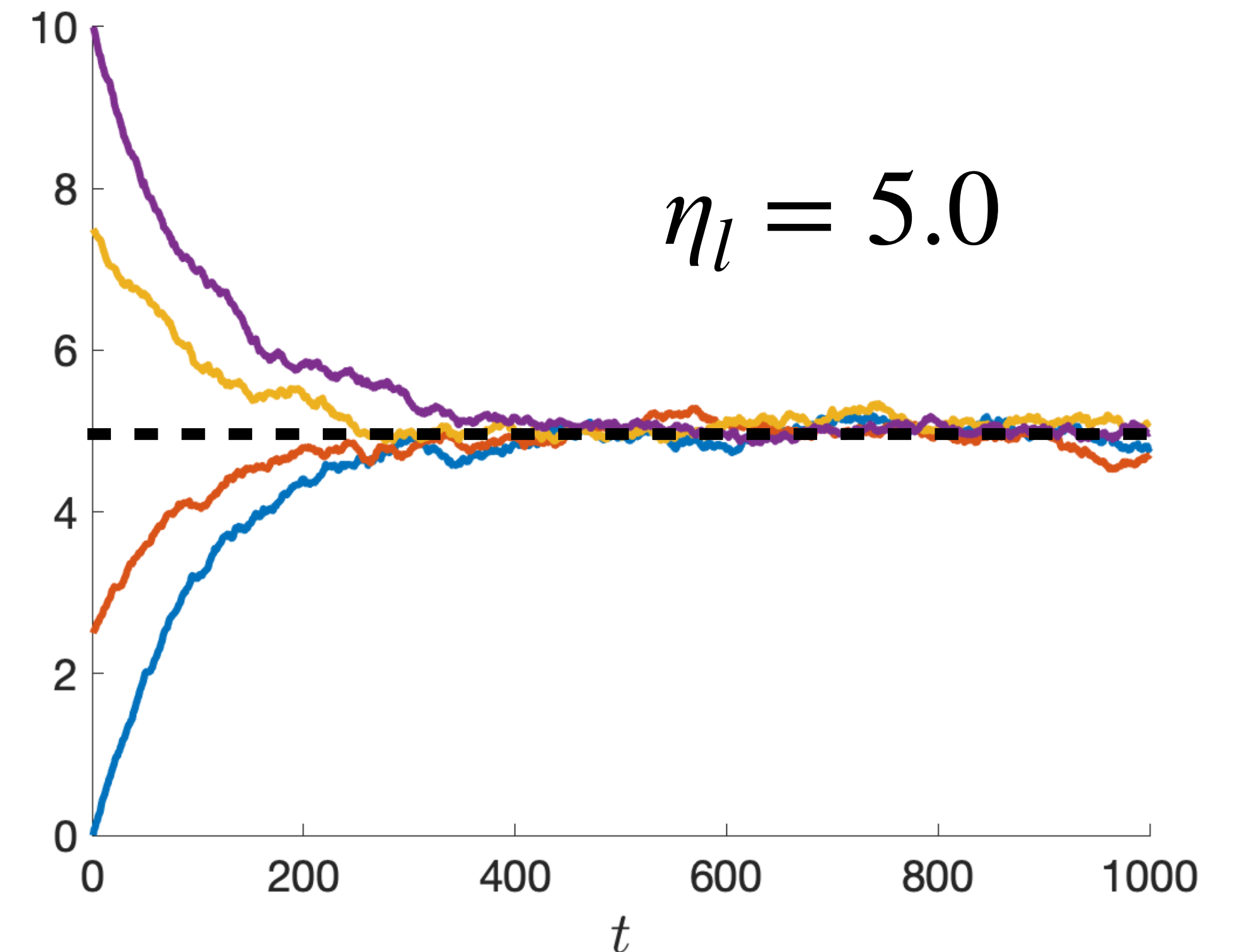
$$P(x_l) = N(\eta_l, \sigma_\omega)$$

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



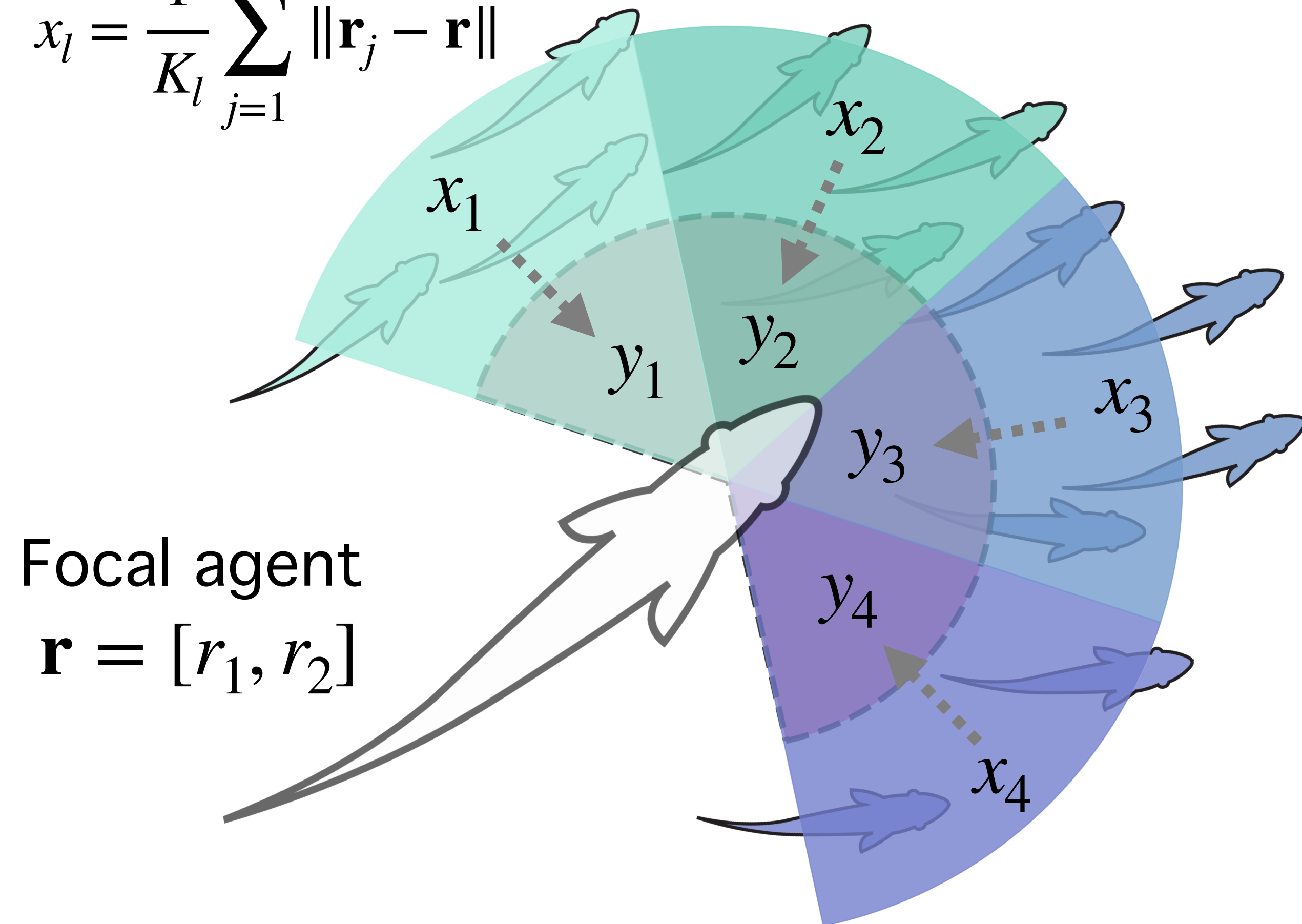
Prior belief about the social distance in a particular sector



# Social forces emerge from predictive control

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



$\mathbf{v}$  = Direction vector

Sensorimotor contingency      Precision-weighted PE

$$\frac{d\mathbf{v}}{dt} = - \begin{bmatrix} \frac{\partial y(\mathbf{v})}{\partial \mathbf{v}} & \frac{\partial F}{\partial y(\mathbf{v})} \end{bmatrix}$$

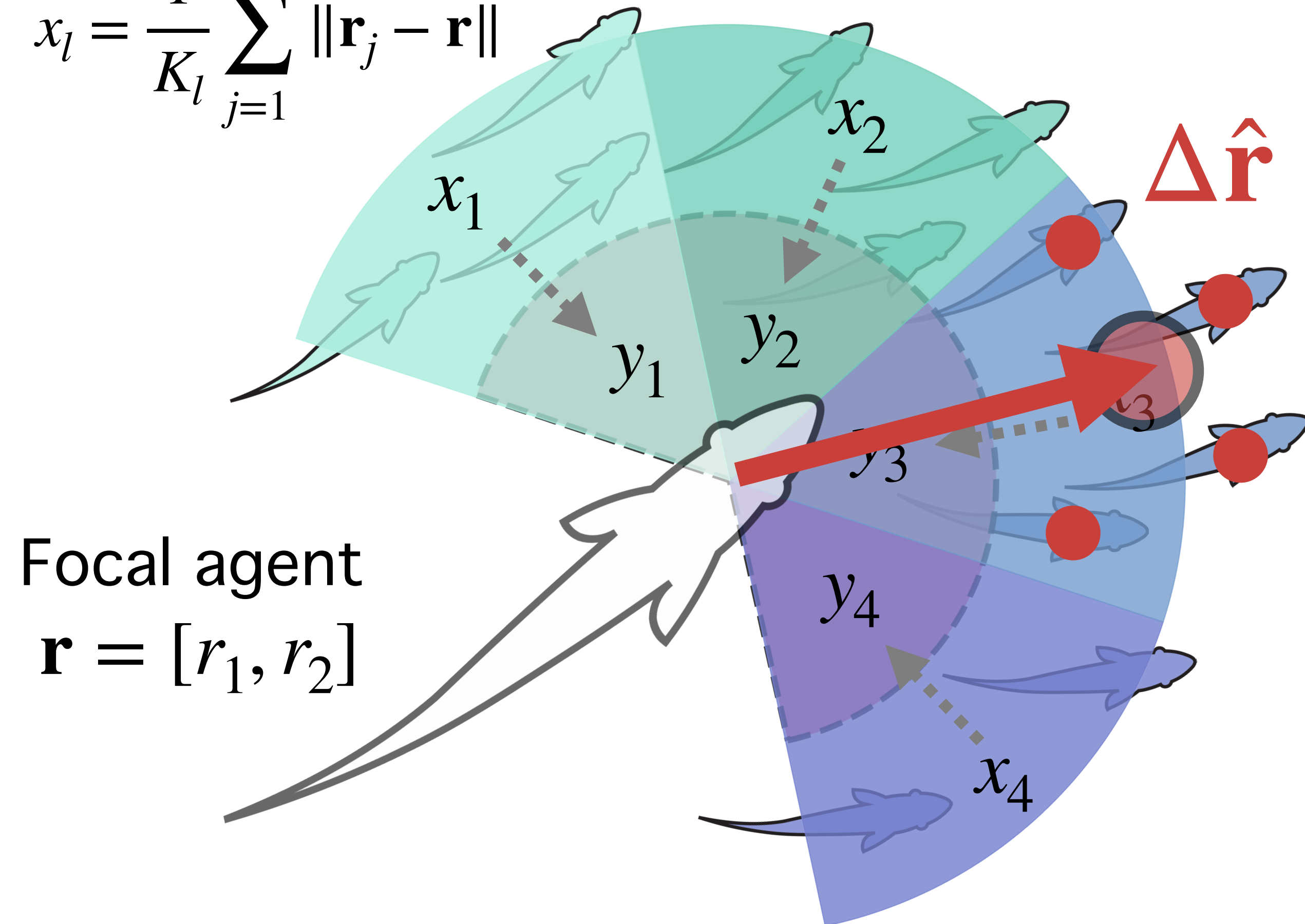
F doesn't directly depend on actions  
(chain rule)

# Social forces emerge from predictive control

$\mathbf{v}$  = Direction vector

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



Sensorimotor contingency      Precision-weighted PE

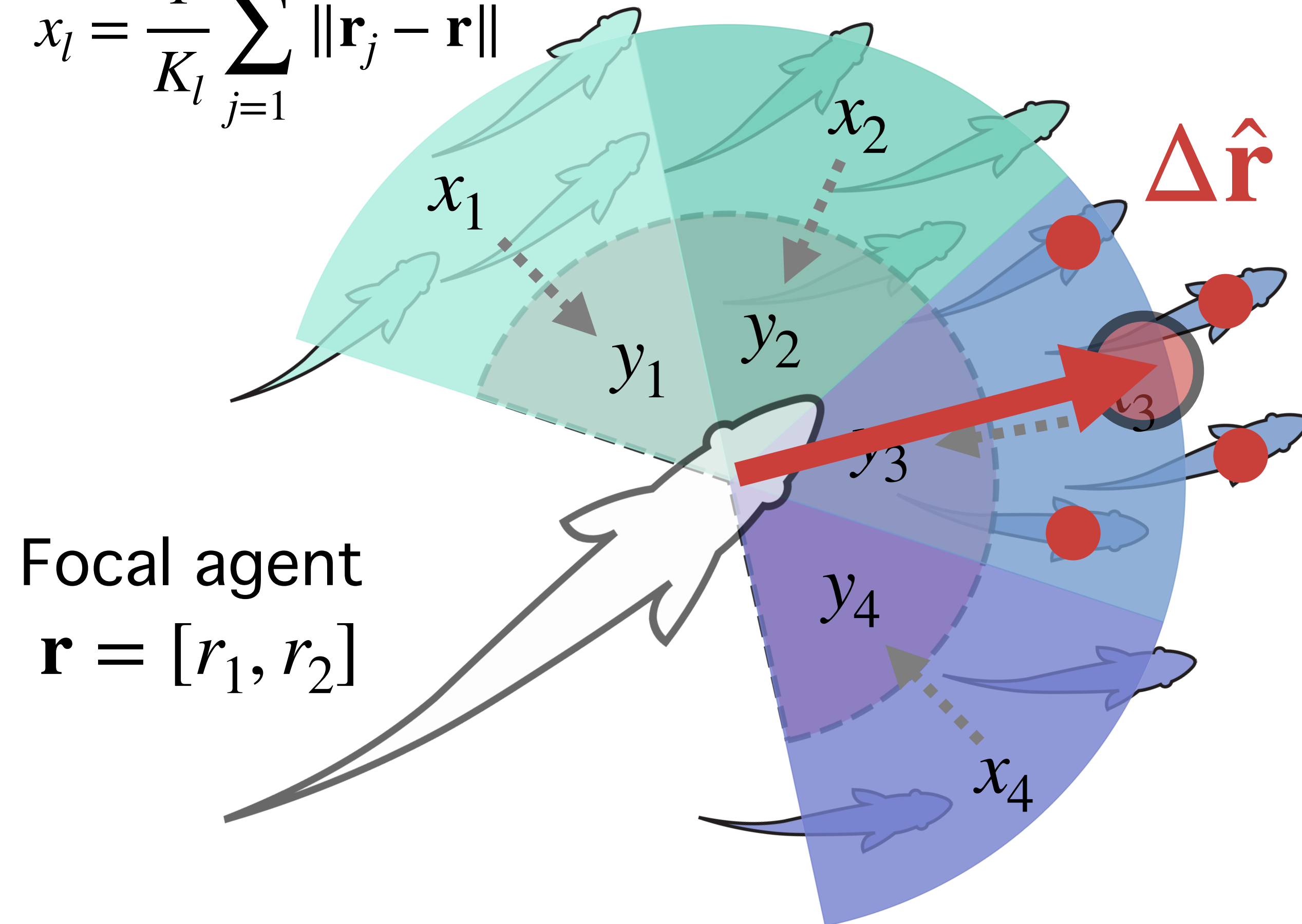
$$\frac{d\mathbf{v}}{dt} = \Delta \hat{\mathbf{r}}^\top \pi_z(y_l - \mu_l)$$

# Social forces emerge from predictive control

$\mathbf{v}$  = Direction vector

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



Sensorimotor contingency      Precision-weighted PE

$$\frac{d\mathbf{v}}{dt} = \Delta \hat{\mathbf{r}}^\top \pi_z(y_l - \mu_l)$$

$y_l - \mu_l > 0$       - - ->      Attraction

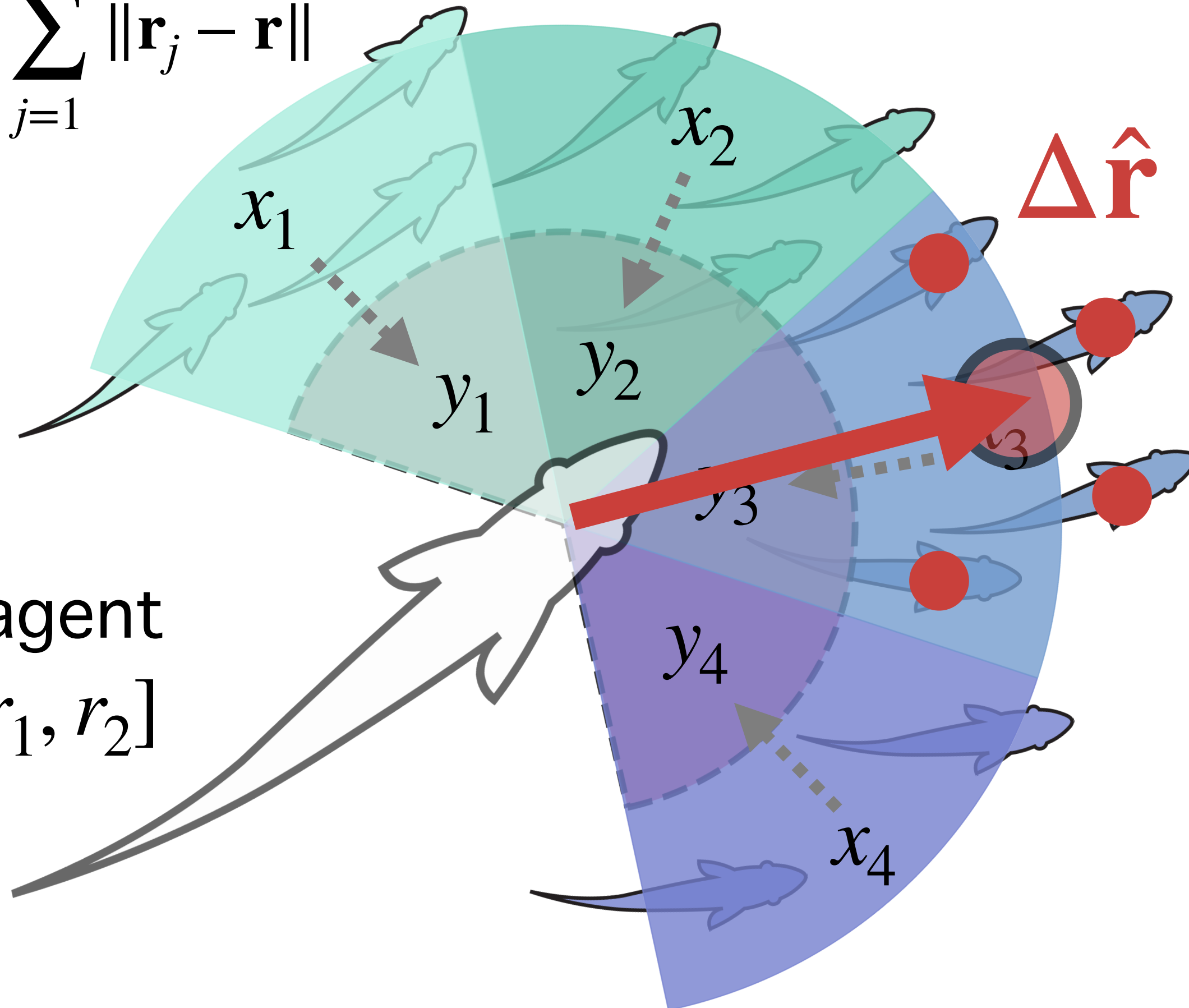
$y_l - \mu_l < 0$       - - ->      Repulsion

# Social forces emerge from predictive control

$\mathbf{v}$  = Direction vector

Sector-specific average distance

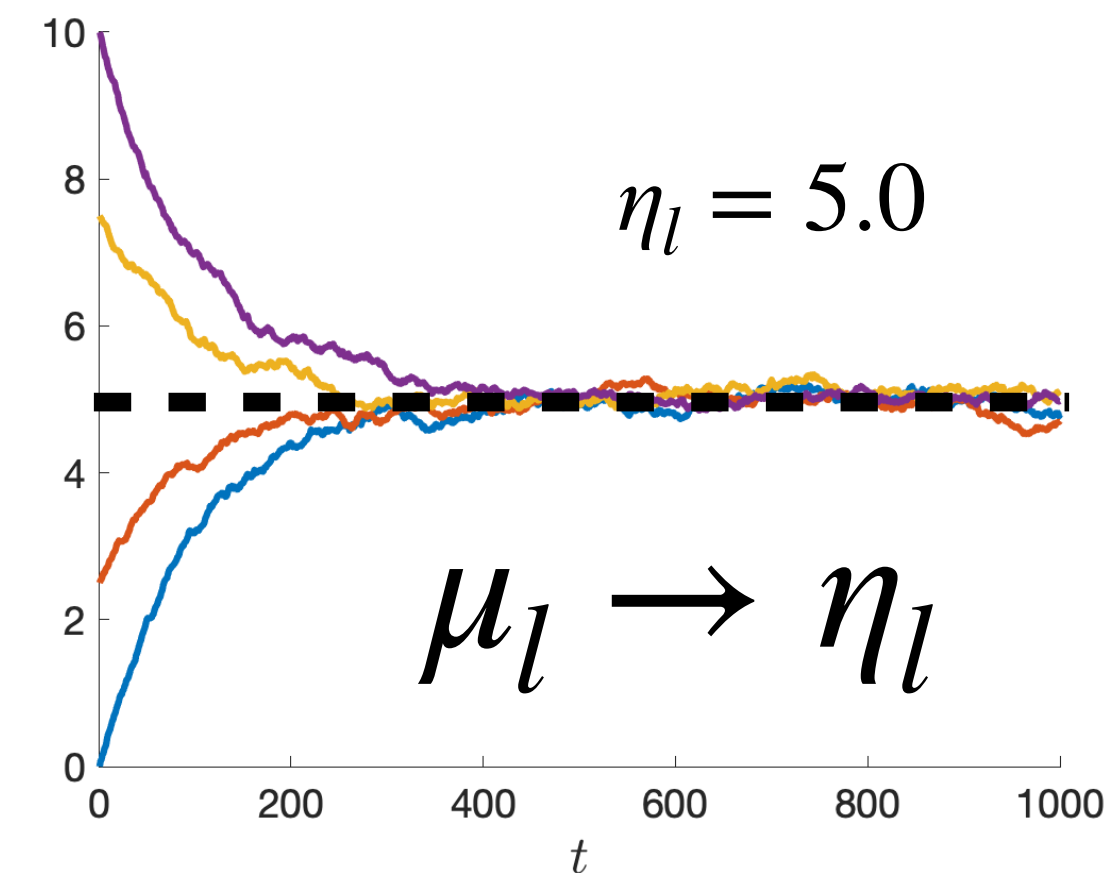
$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



Focal agent  
 $\mathbf{r} = [r_1, r_2]$

Sensorimotor contingency      Precision-weighted PE

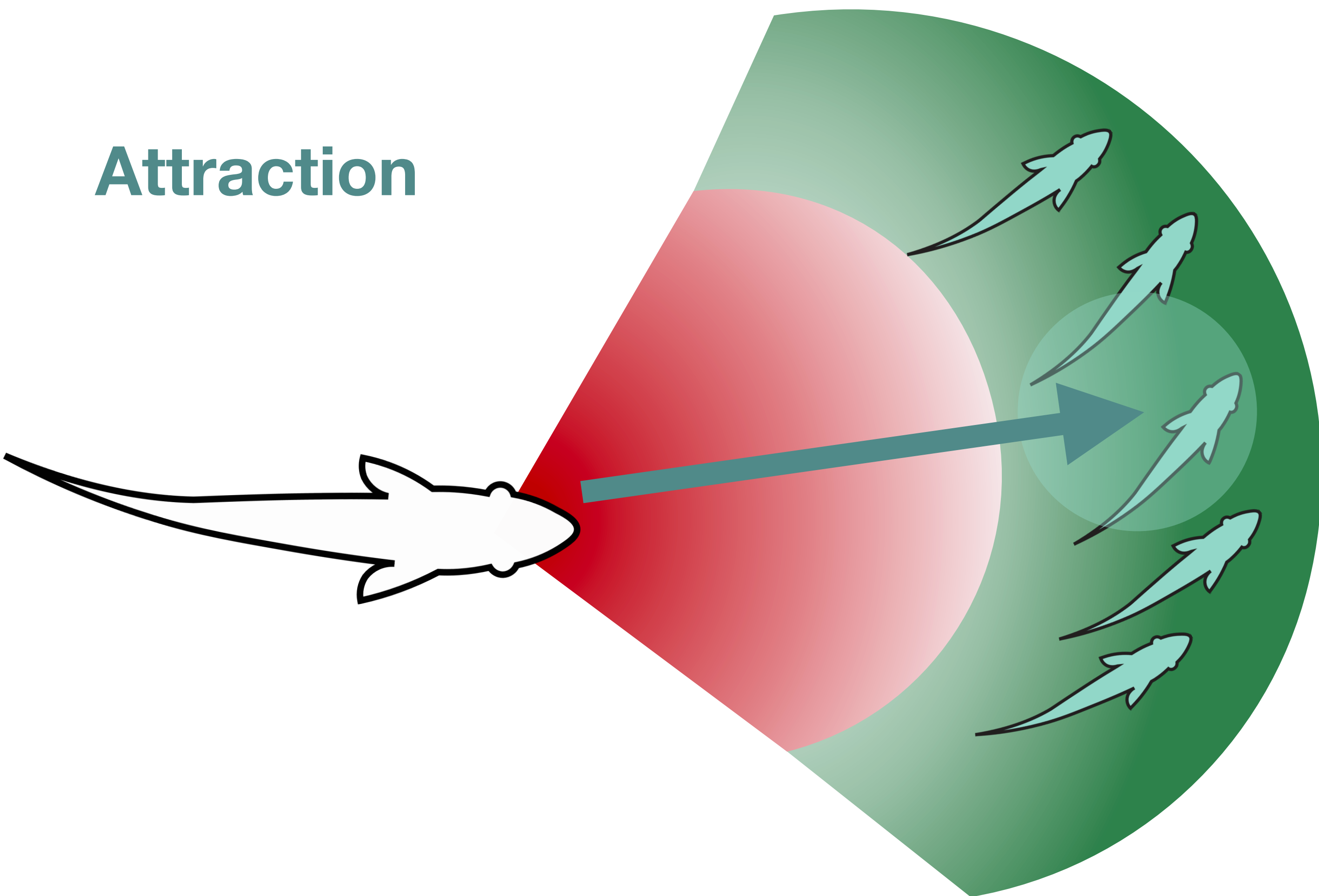
$$\frac{d\mathbf{v}}{dt} = \Delta \hat{\mathbf{r}}^\top \pi_z(y_l - \mu_l)$$



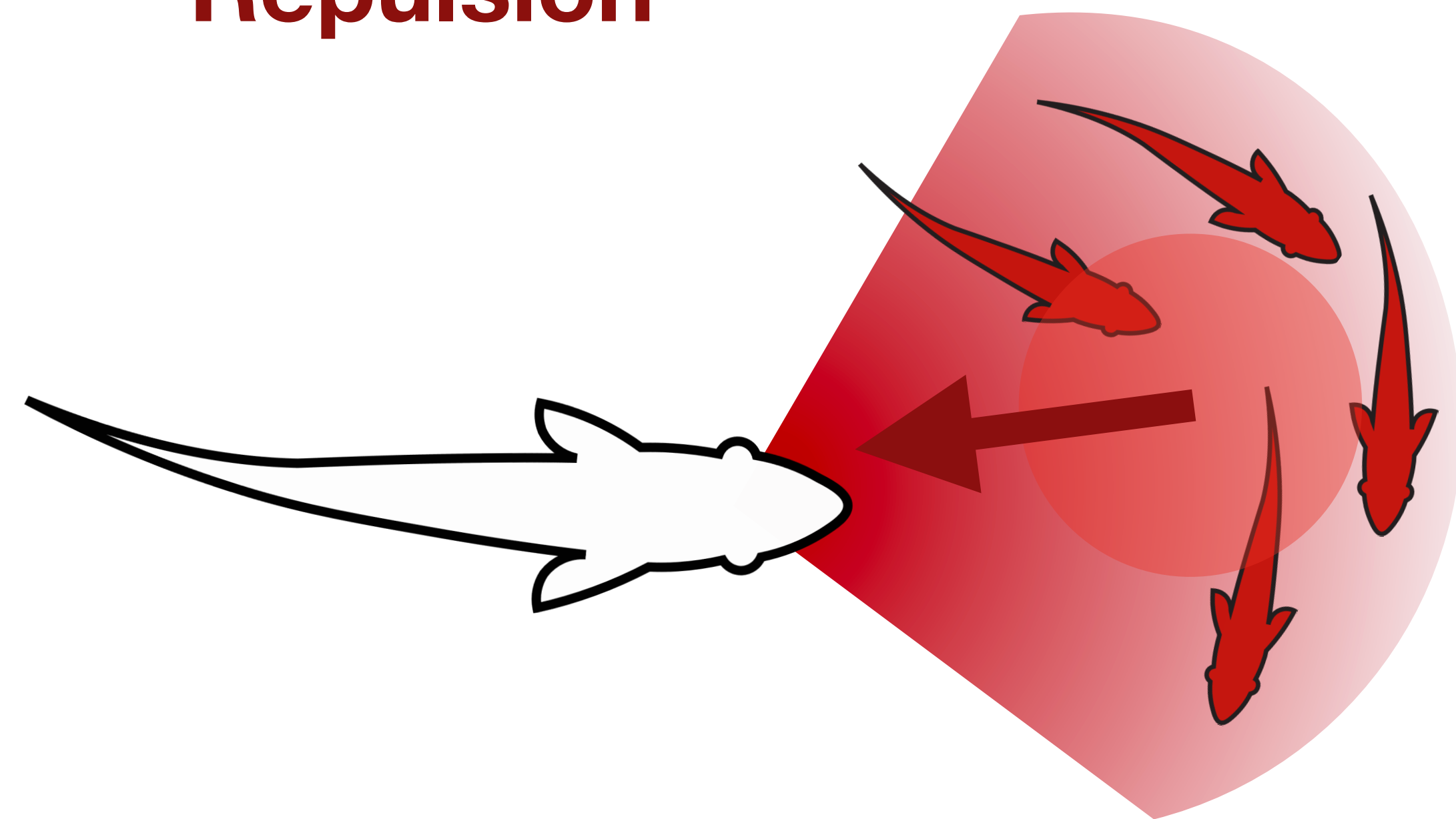


# Zonal social force models

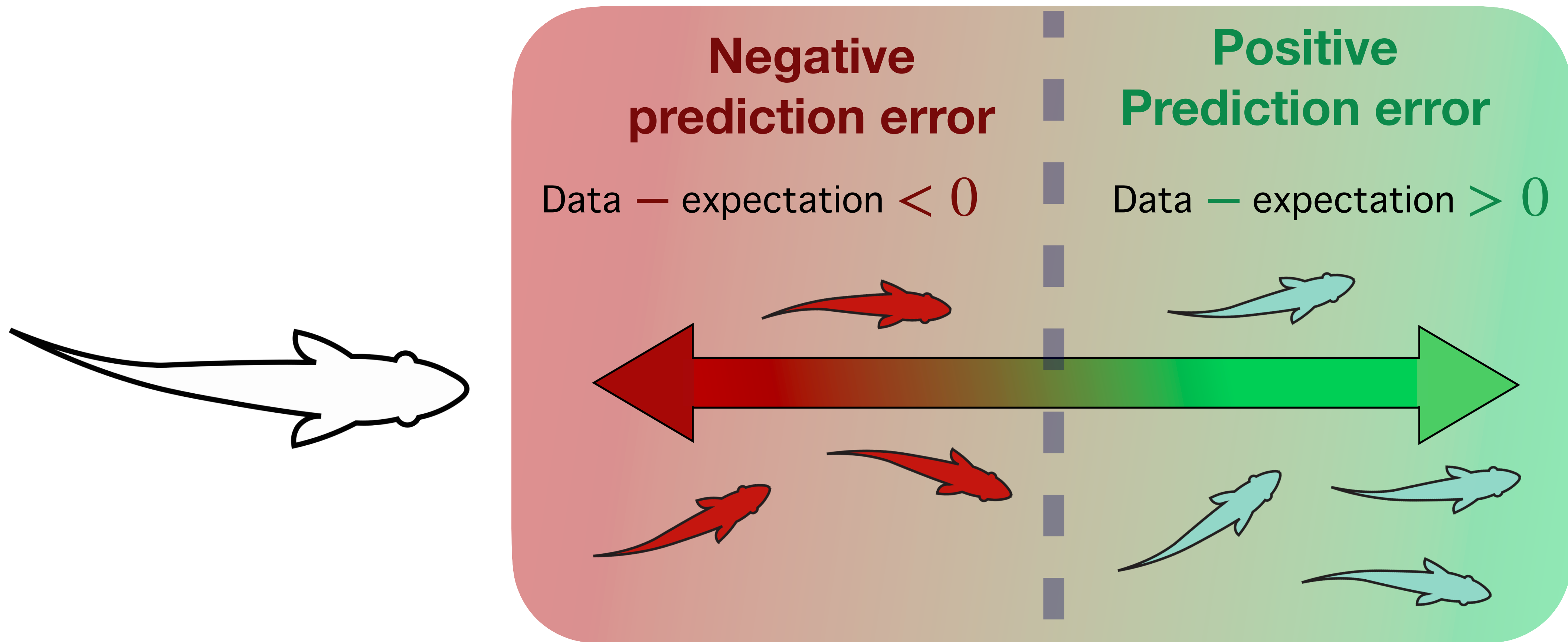
**Attraction**



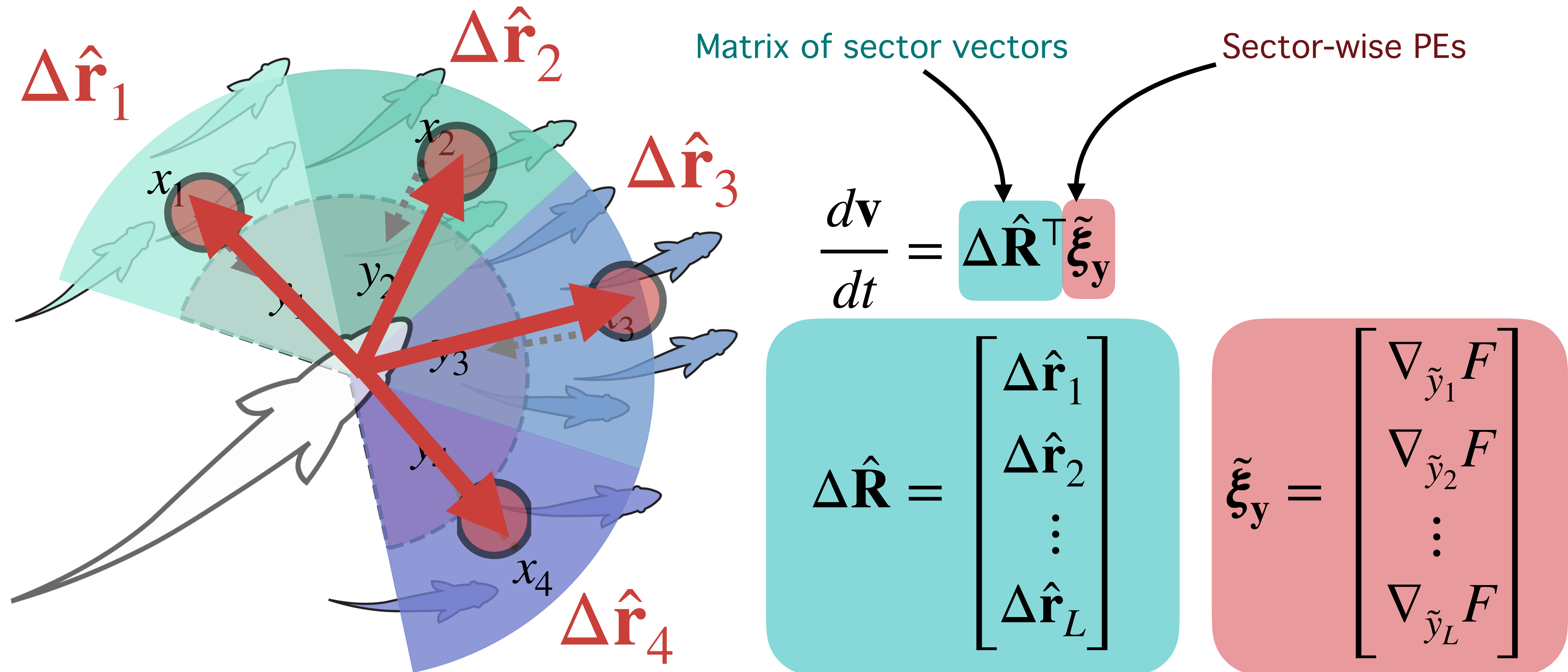
**Repulsion**



# Active control of prediction error



# Social forces emerge from (multivariate) predictive control



# Important addendum for collective motion theorists!

Sensorimotor contingency

Precision-weighted PE

$$\frac{d\mathbf{v}}{dt} = \Delta \hat{\mathbf{r}}^\top \pi_z(y_l - \mu_l)$$

$$y_l - \mu_l > 0 \quad \text{---} \rightarrow \text{Attraction}$$

$$y_l - \mu_l < 0 \quad \text{---} \rightarrow \text{Repulsion}$$

$$\tilde{y} = \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} = \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} \text{ Sensed distance}$$

$$\frac{d\mathbf{v}}{dt} = \nabla_{\tilde{y}} \mathbf{v}^\top \tilde{\Pi}_z(\tilde{y}_l - \tilde{\mu}_l)$$

$$y'_l - \mu'_l > 0 \quad \text{---} \rightarrow \text{Attraction}$$

$$y'_l - \mu'_l < 0 \quad \text{---} \rightarrow \text{Repulsion}$$

# Important addendum for collective motion theorists!

$$\tilde{\mathbf{y}} = \begin{bmatrix} y \\ y' \\ y'' \\ \vdots \end{bmatrix} = \begin{bmatrix} y \\ \partial_t y \\ \partial_t^2 y \\ \vdots \end{bmatrix} \begin{array}{l} \text{Sensed distance} \\ \text{Sensed "distance velocity"} \\ \text{Sensed "distance acceleration"} \end{array}$$

$$\frac{d\mathbf{v}}{dt} = \nabla_{\tilde{\mathbf{y}}} \mathbf{v}^\top \tilde{\Pi}_z (\tilde{\mathbf{y}}_l - \tilde{\mu}_l)$$

$$y'_l - \mu'_l > 0 \quad \text{---} \rightarrow \quad \text{Attraction}$$

$$y'_l - \mu'_l < 0 \quad \text{---} \rightarrow \quad \text{Repulsion}$$

Research article

## Swarming and pattern formation due to selective attraction and repulsion

Pawel Romanczuk  and Lutz Schimansky-Geier

Published: 26 September 2012 | <https://doi.org/10.1098/rsfs.2012.0030>

$y'_l = \frac{dy'_l}{dt}$  is equivalent to the “relative velocity”,  
or the rate at which neighbouring individuals are receding (positive) vs. looming (negative)

$$\frac{dy_l}{dt} = (\mathbf{r}_i - \mathbf{r}_j) \cdot \mathbf{v}_i + \sum_{j \in S_l} \left( (\mathbf{r}_j - \mathbf{r}_i) \cdot \mathbf{v}_j \right)$$

$\mathbf{r}_i$  = Position vector of focal agent  $i$

$\mathbf{r}_j$  = Position vector of neighbour  $j$  in sector  $l$

# Evidence for use of prediction errors (i.e., unpredicted changes in sensory input), rather than absolute values

nature communications



Article

<https://doi.org/10.1038/s41467-024-53361-8>

## Body orientation change of neighbors leads to scale-free correlation in collective motion

Zhicheng Zheng<sup>1</sup>, Yuan Tao<sup>1</sup>, Yalun Xiang<sup>1</sup>, Xiaokang Lei<sup>2</sup> & Xingguang Peng<sup>1</sup> ✉

<sup>1</sup>School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an, Shaanxi 710072, P. R. China. <sup>2</sup>School of Information and Control Engineering, Xi'an University of Architecture and Technology, Xi'an, Shaanxi 710055, P. R. China. ✉e-mail: [pxg@nwpu.edu.cn](mailto:pxg@nwpu.edu.cn)

PNAS

RESEARCH ARTICLE

ENGINEERING  
BIOLOGICAL SCIENCES

OPEN ACCESS

## Individual error correction drives responsive self-assembly of army ant scaffolds

Matthew J. Lutz<sup>a,b,c,2,1</sup>, Chris R. Reid<sup>d,2,1</sup>, Christopher J. Lustri<sup>e</sup>, Albert B. Kao<sup>f</sup>, Simon Garnier<sup>g</sup>, and Iain D. Couzin<sup>a,b,c</sup>

nature communications



## Zebrafish capable of generating future state prediction error show improved active avoidance behavior in virtual reality

Makio Torigoe<sup>1</sup>, Tanvir Islam<sup>1,2</sup>, Hisaya Kakinuma<sup>1,2</sup>, Chi Chung Alan Fung<sup>3</sup>, Takuya Isomura<sup>4</sup>, Hideaki Shimazaki<sup>5</sup>, Tazu Aoki<sup>1</sup>, Tomoki Fukai<sup>3</sup> & Hitoshi Okamoto<sup>1,2</sup> ✉

SCIENCE ADVANCES | RESEARCH ARTICLE

NEUROSCIENCE

## Predictive neural computations in the cerebellum contribute to motor planning and faster behavioral responses in larval zebrafish

Sriram Narayanan, Aalok Varma, Vatsala Thirumalai\*

# Collective simulation achieved by minimizing individual free energy functionals

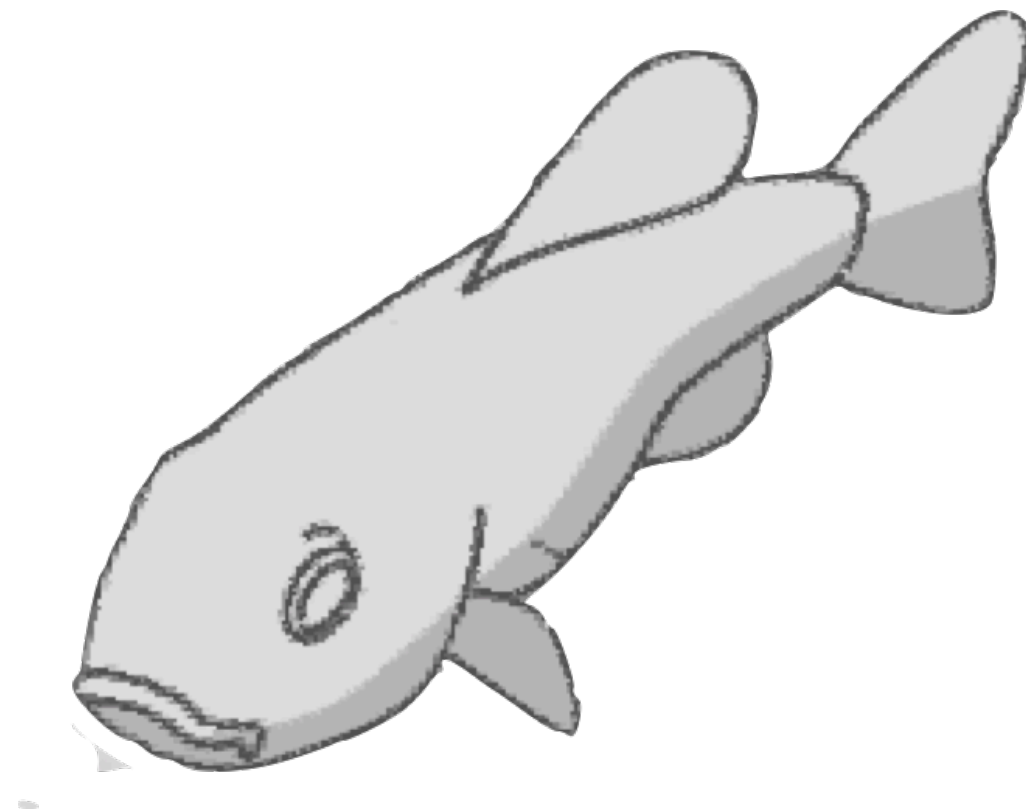
$F = \text{Surprise} = \text{Prediction error}$

$$\dot{\mu} = -\nabla_{\mu} F(\mu, y)$$

Perception

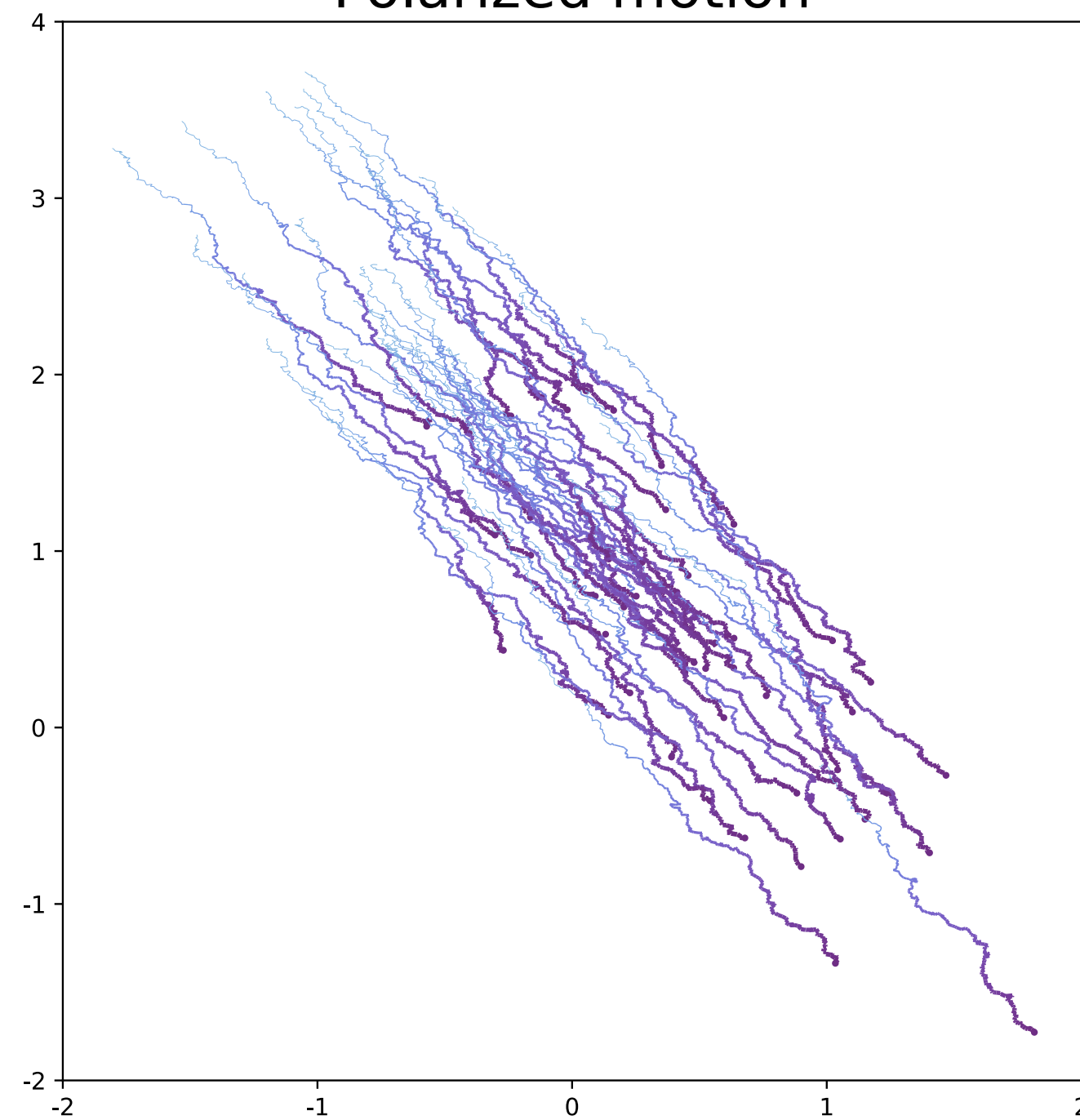
$$\dot{v} = -\nabla_{v} F(\mu, y(v))$$

Action

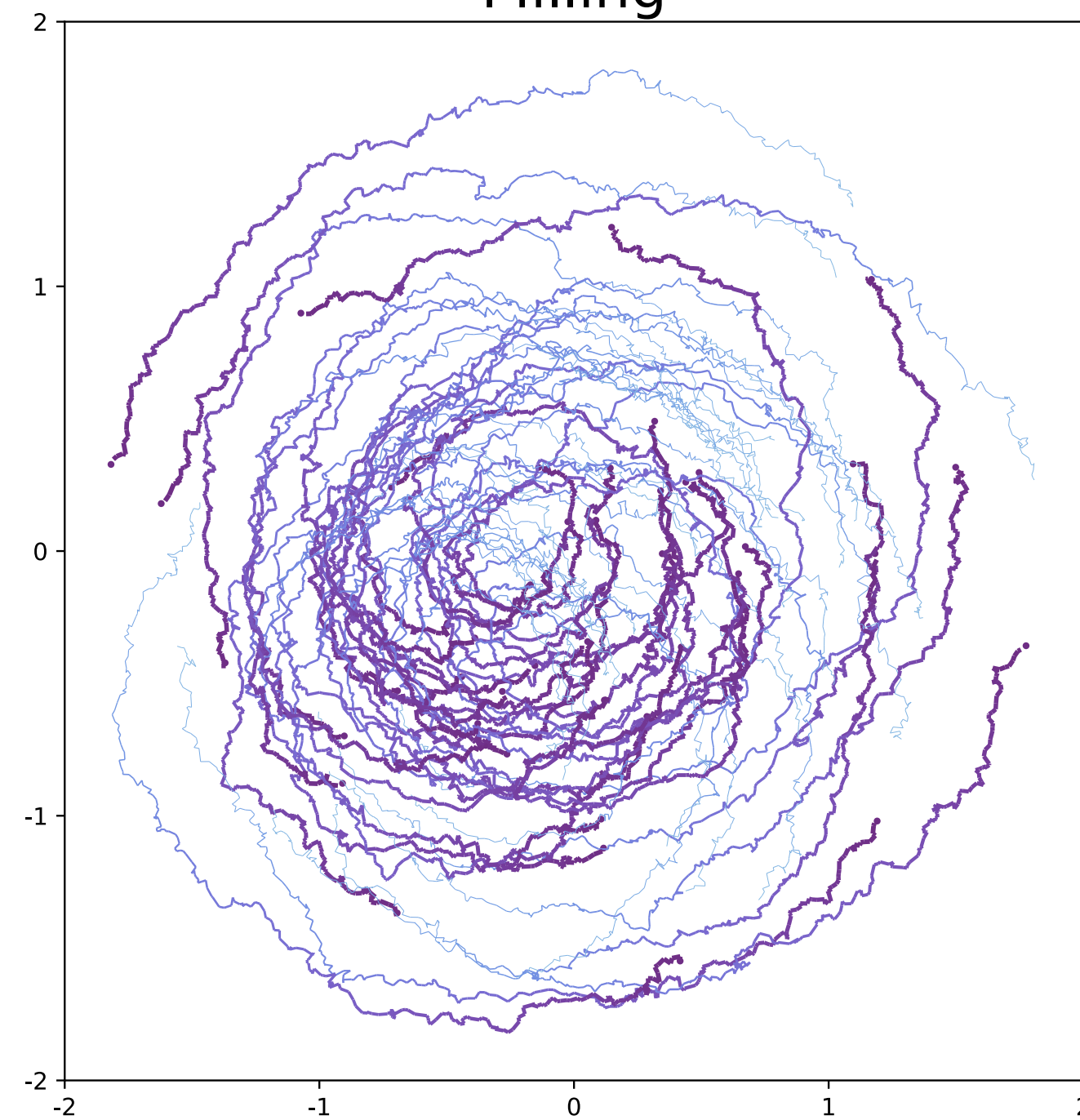


# Collective regimes

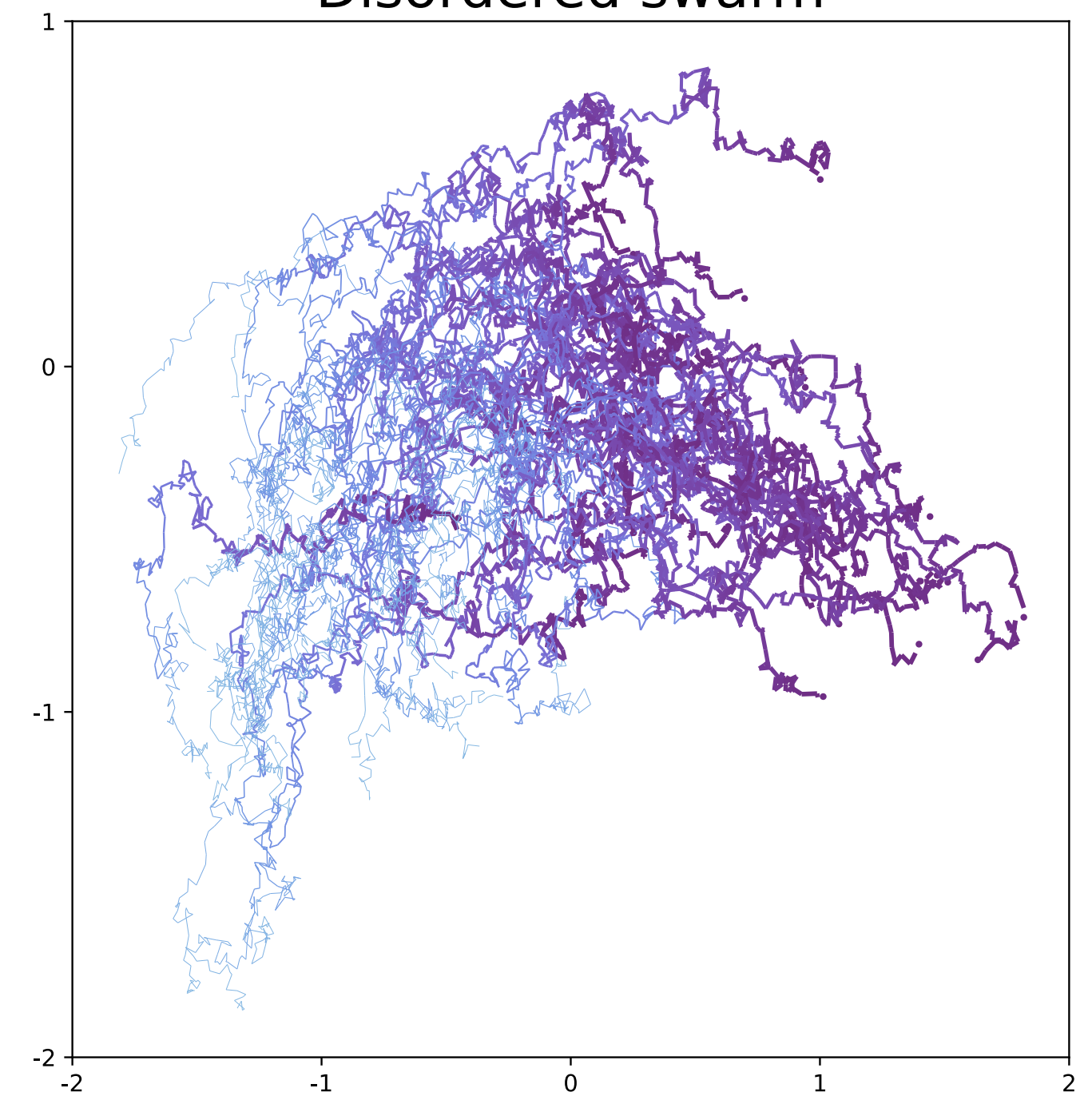
Polarized motion



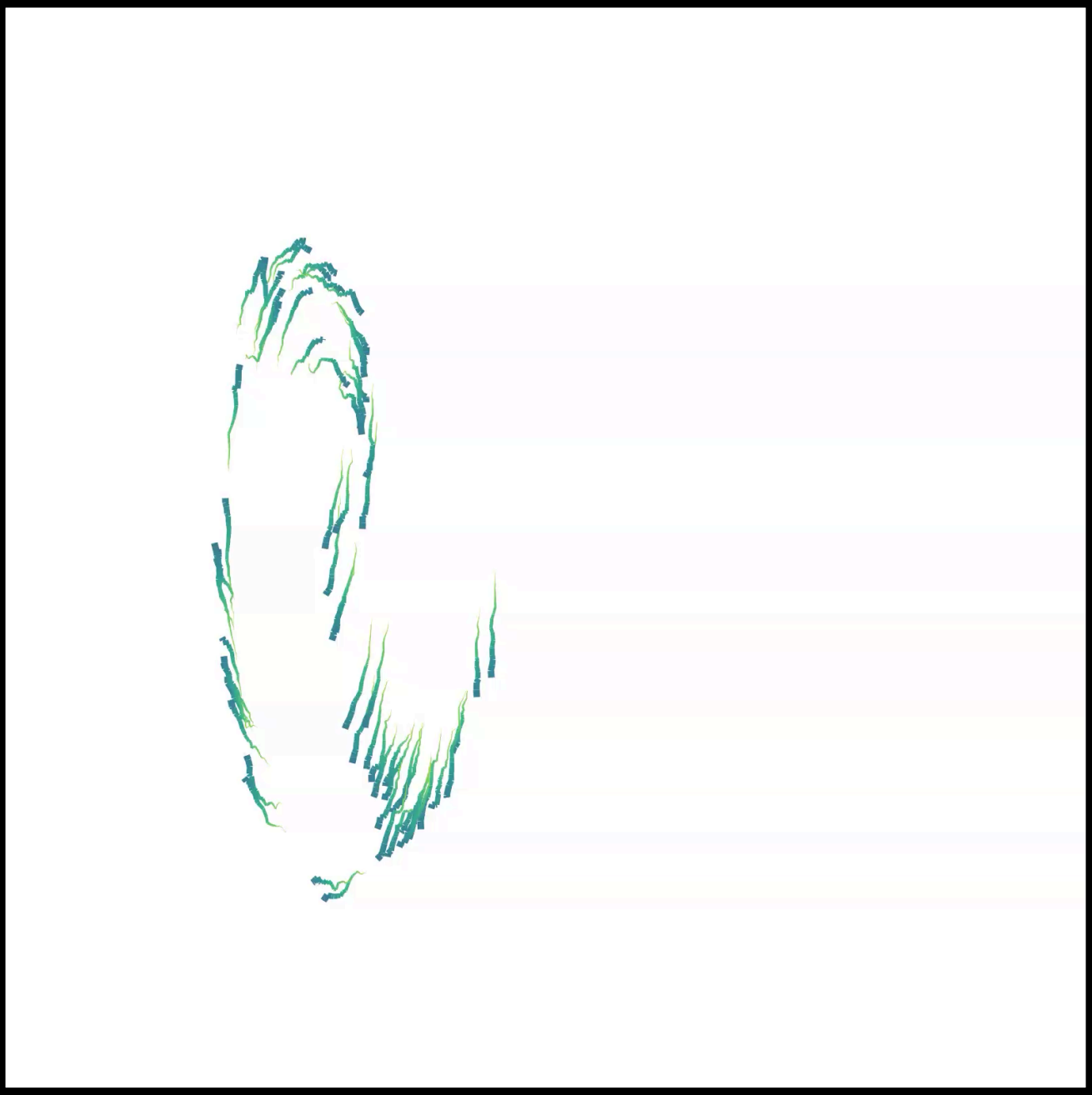
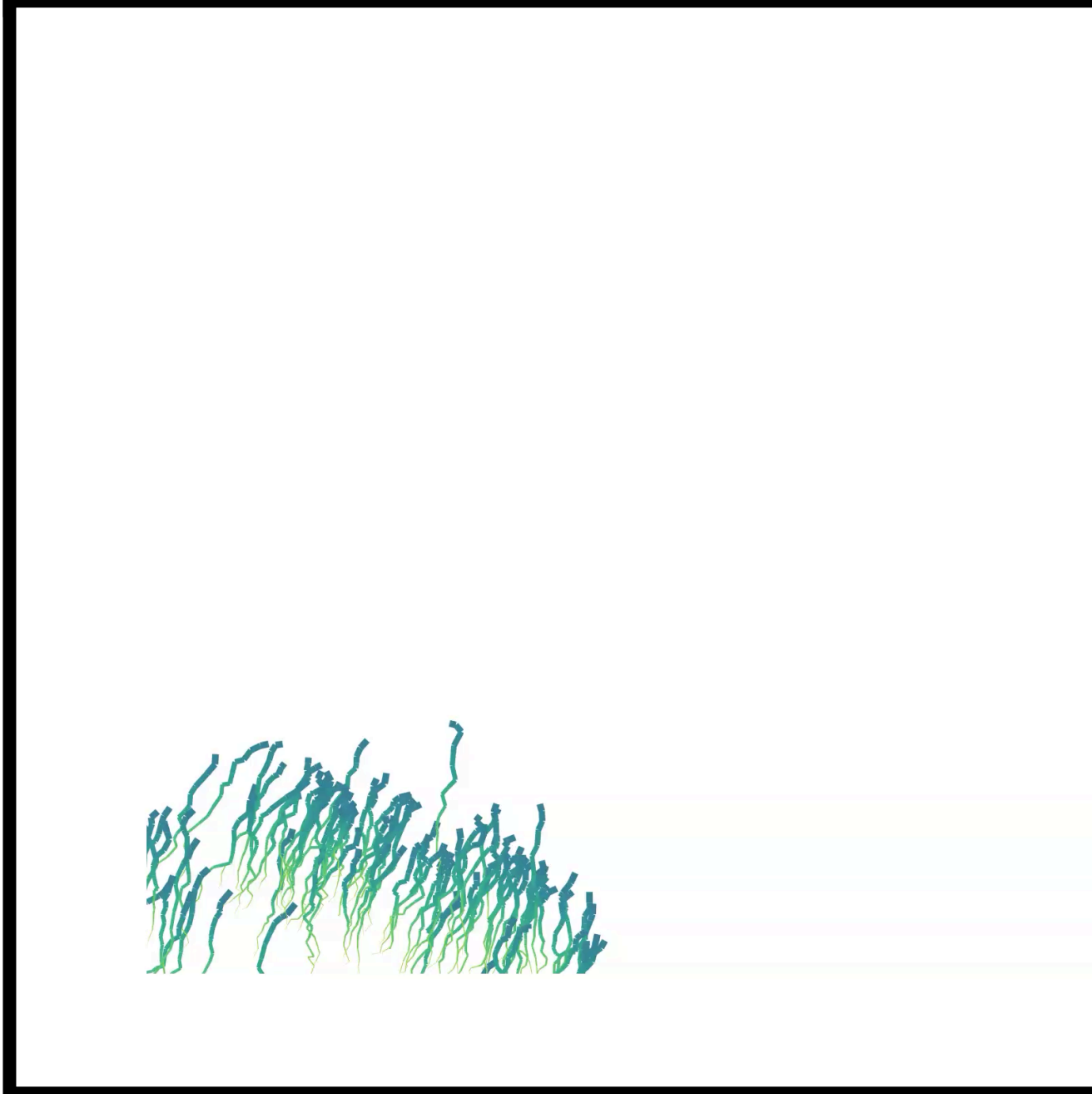
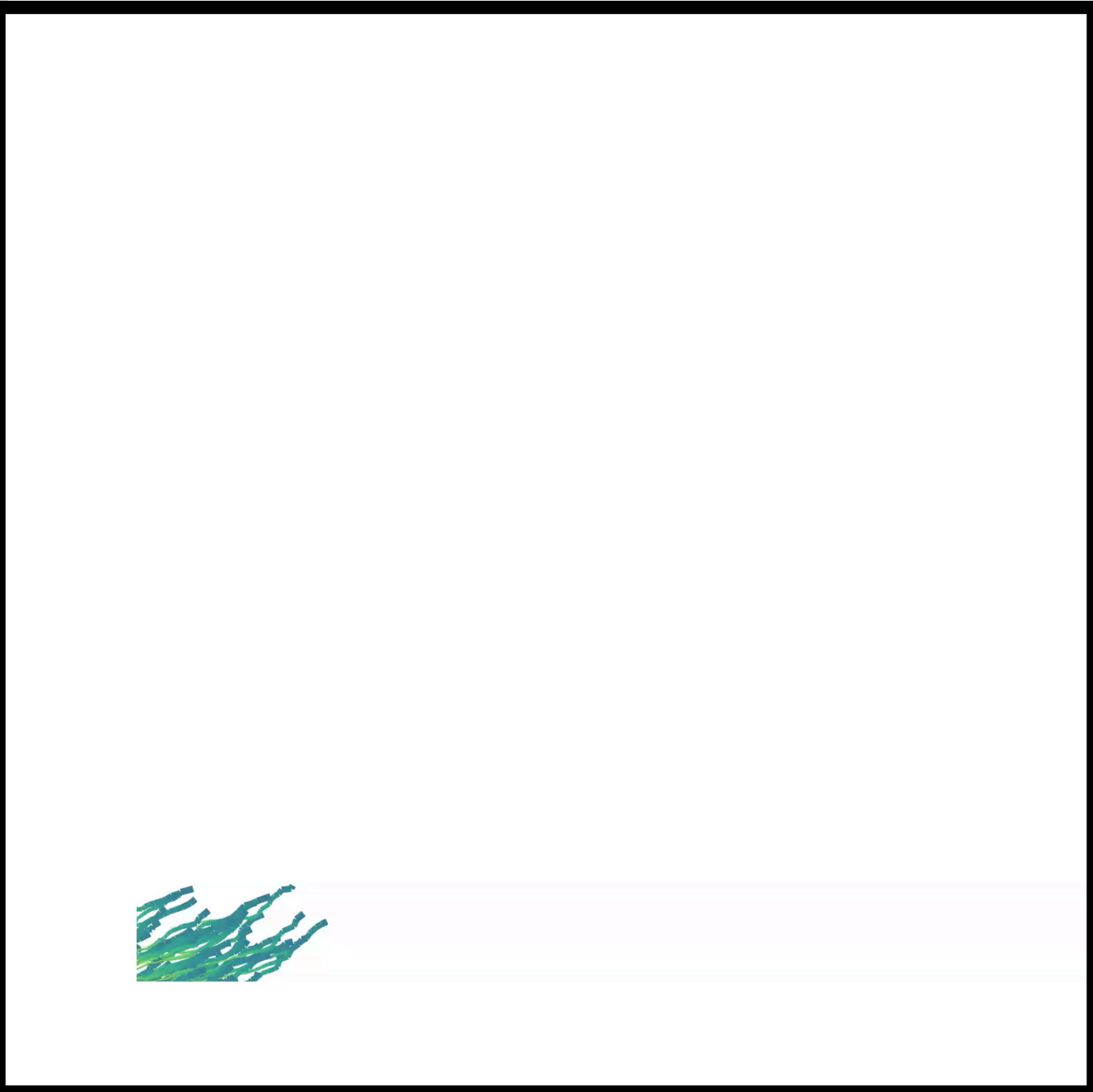
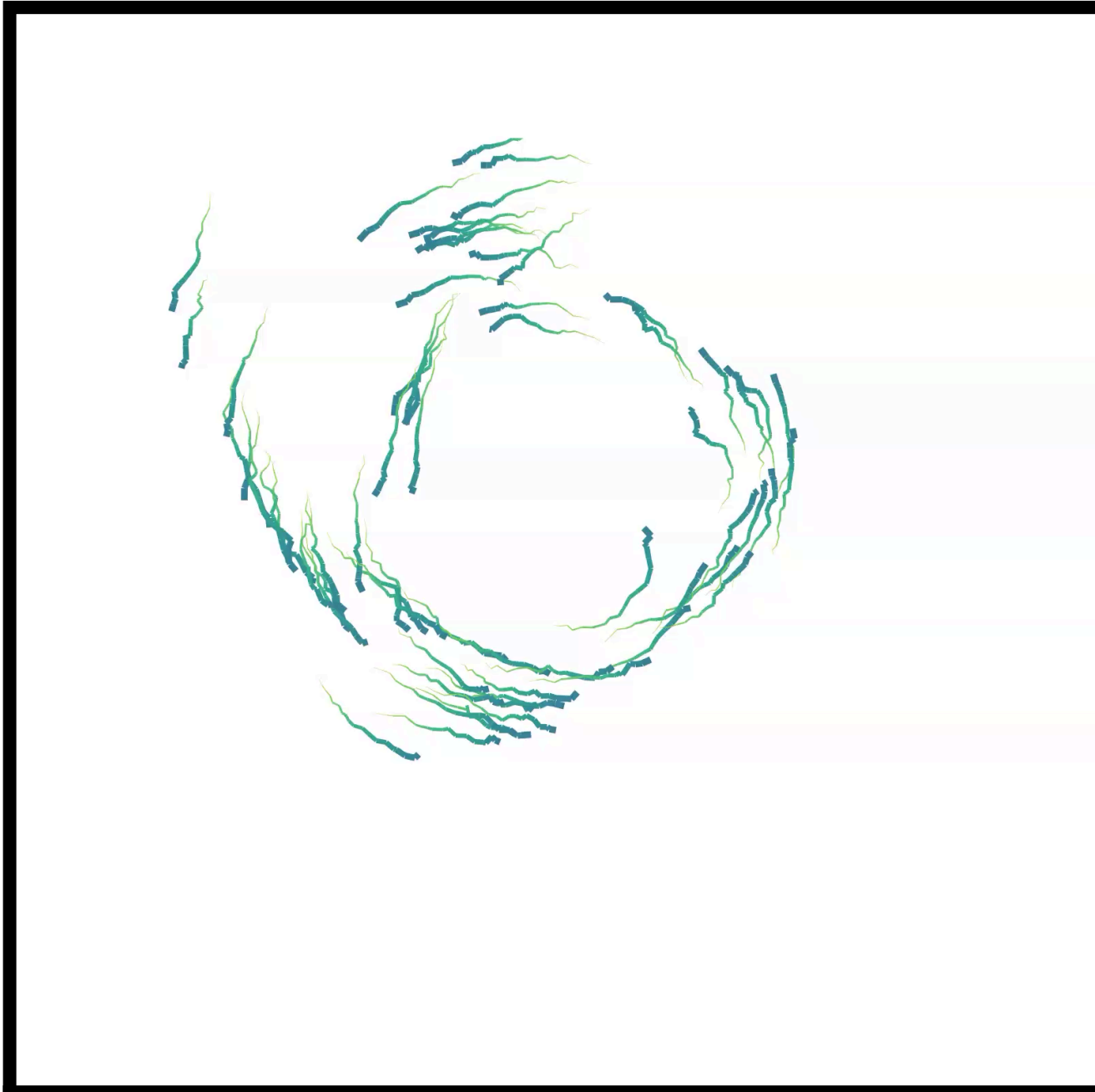
Milling



Disordered swarm





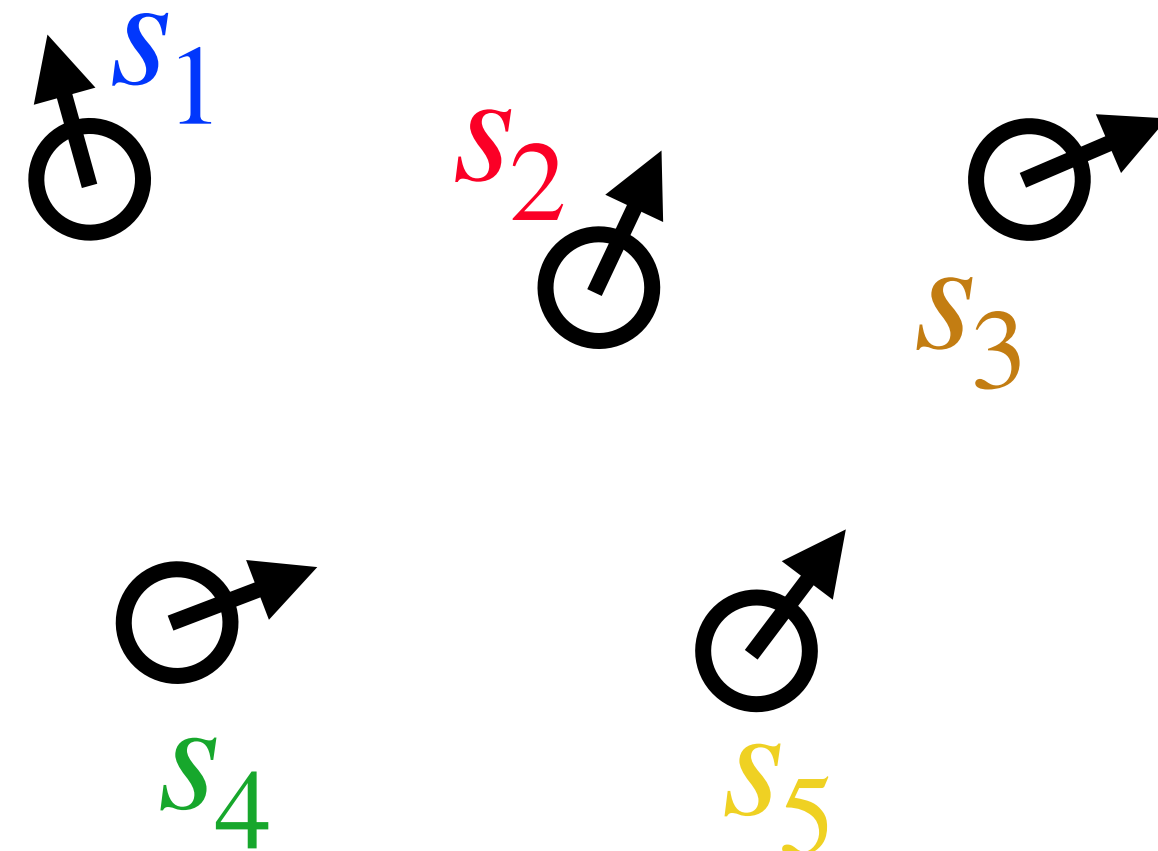


# Potential function vs. free energy function(al)

## Potential systems

- **Global function** of the configurational states of the system
- Individual dynamics move down gradients of the shared, global potential

$$E(s_1, s_2, s_3, \dots)$$



$$\dot{s}_1 \propto -\nabla_{s_1} E$$

$$\dot{s}_2 \propto -\nabla_{s_2} E$$

$$\dot{s}_3 \propto -\nabla_{s_3} E$$

⋮

# Potential function vs. free energy function(al)

## Collective Bayesian (active inference) systems

- **Local functional** of probabilistic beliefs about one's environment
- Individual dynamics driven by dual gradient flows (action and perception) on this moving functional

$$F_1(o_1, q_1)$$

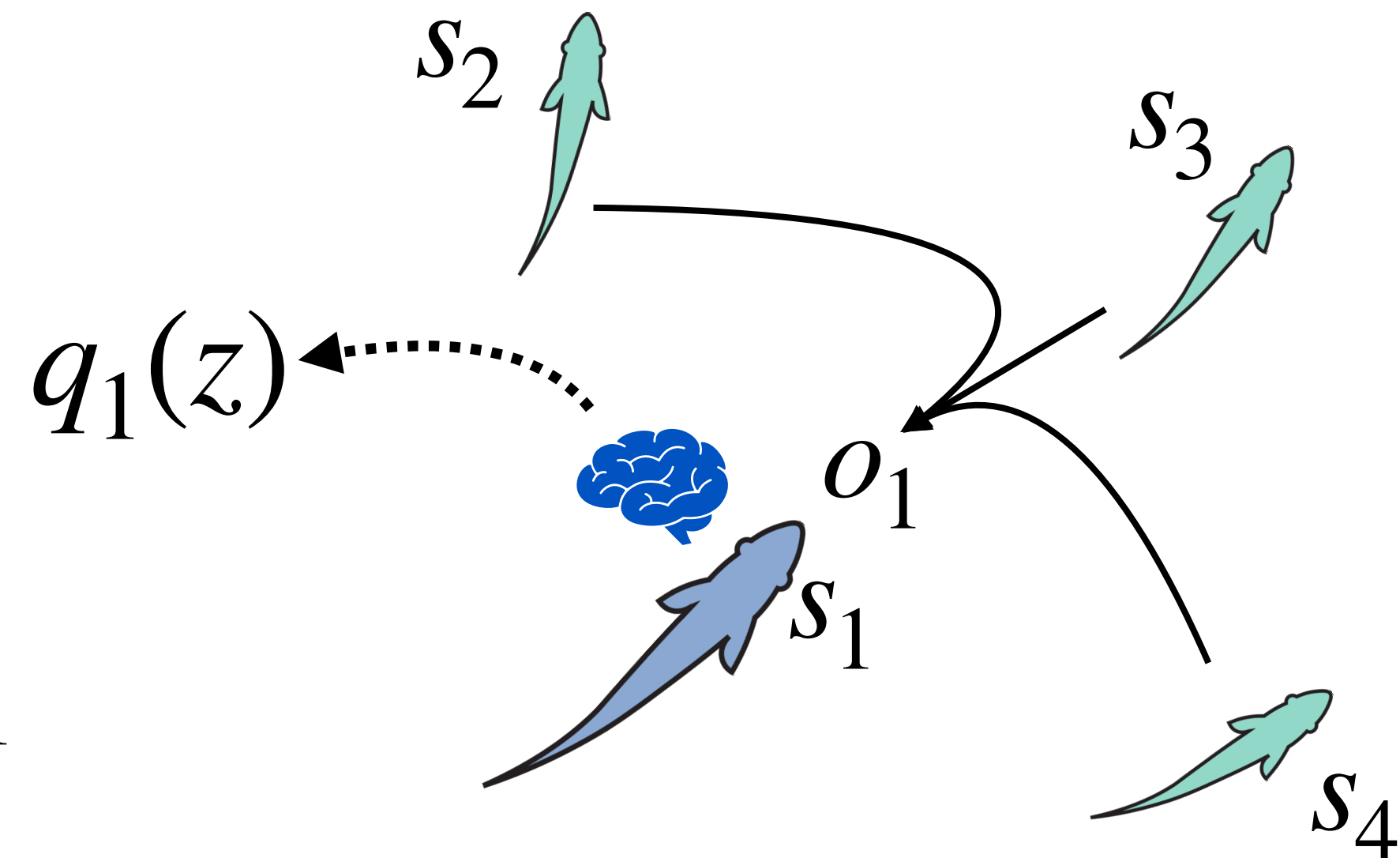
Free energy functional of  
probabilistic beliefs

$$\dot{q}_1 \propto -\nabla_{q_1} F_1$$

Perception

$$\dot{a}_1 \propto -\nabla_{a_1} F_1$$

Action



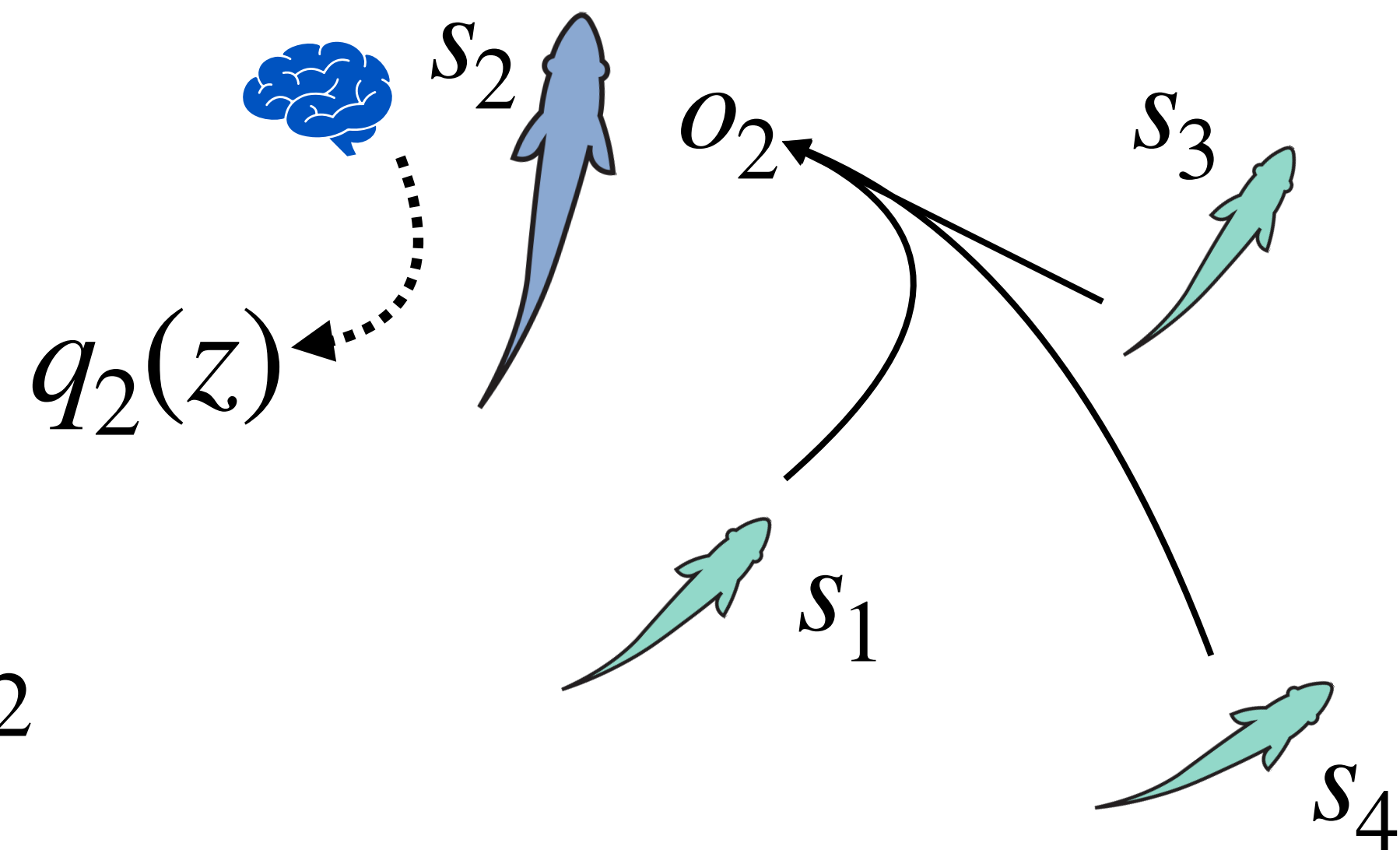
# Potential function vs. free energy function(al)

## Collective Bayesian (active inference) systems

- **Local functional** of probabilistic beliefs about one's environment
- Individual dynamics driven by dual gradient flows (action and perception) on this moving functional

$$F_2(o_2, q_2)$$

Free energy functional of  
probabilistic beliefs



$$\dot{q}_2 \propto -\nabla_{q_2} F_2$$

Perception

$$\dot{a}_2 \propto -\nabla_{a_2} F_2$$

Action

How do properties of individual models determine collective outcomes?

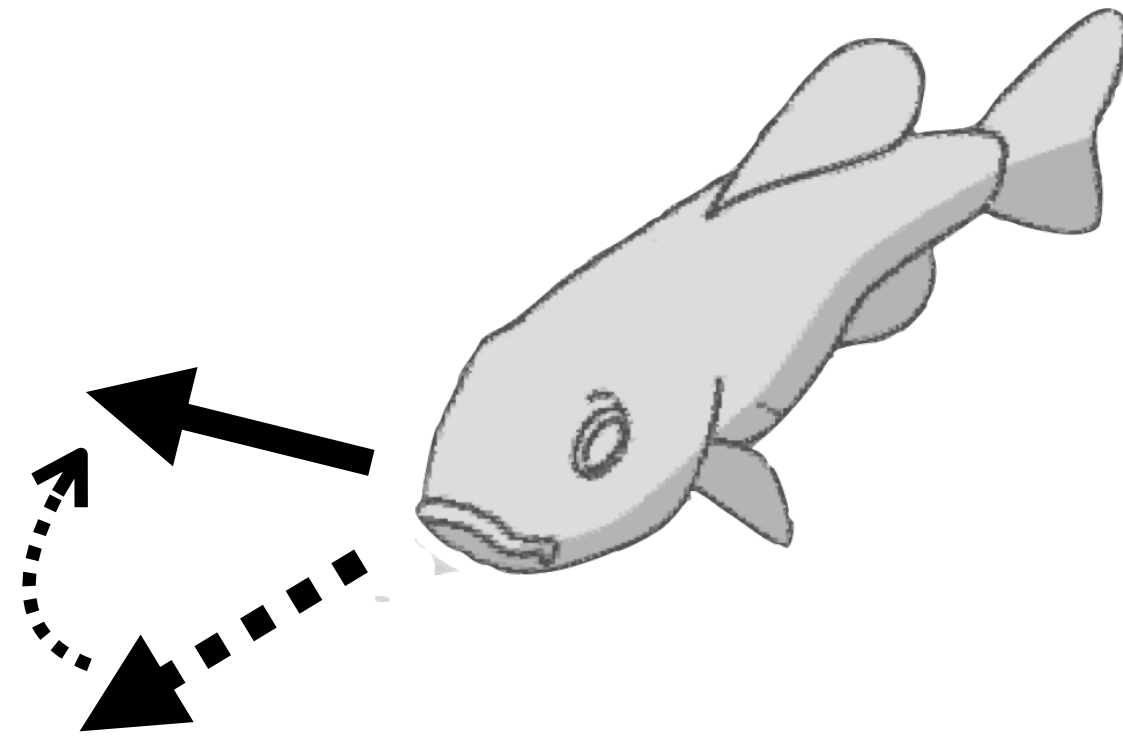
Modify the generative model of single agents and measure consequences

$$p(y, x) = p(y | x)p(x)$$



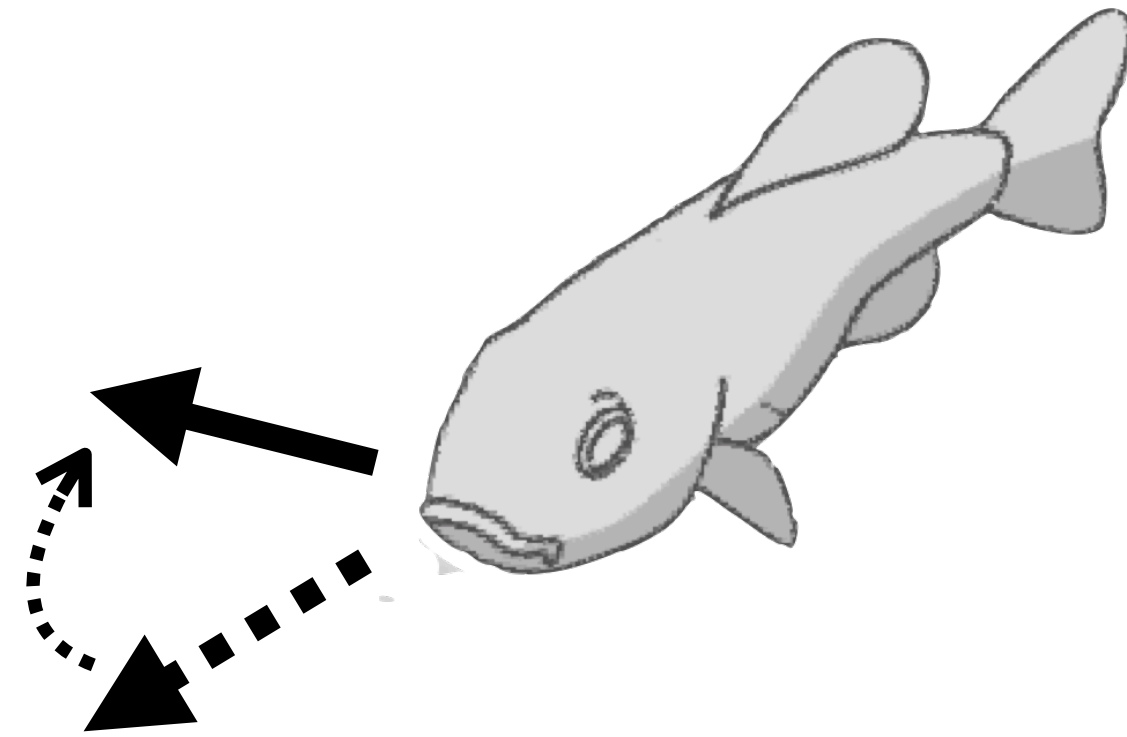
# Weights vs beliefs

Behavior  $= f_{social}(x) + f_{env}(z) + \epsilon$



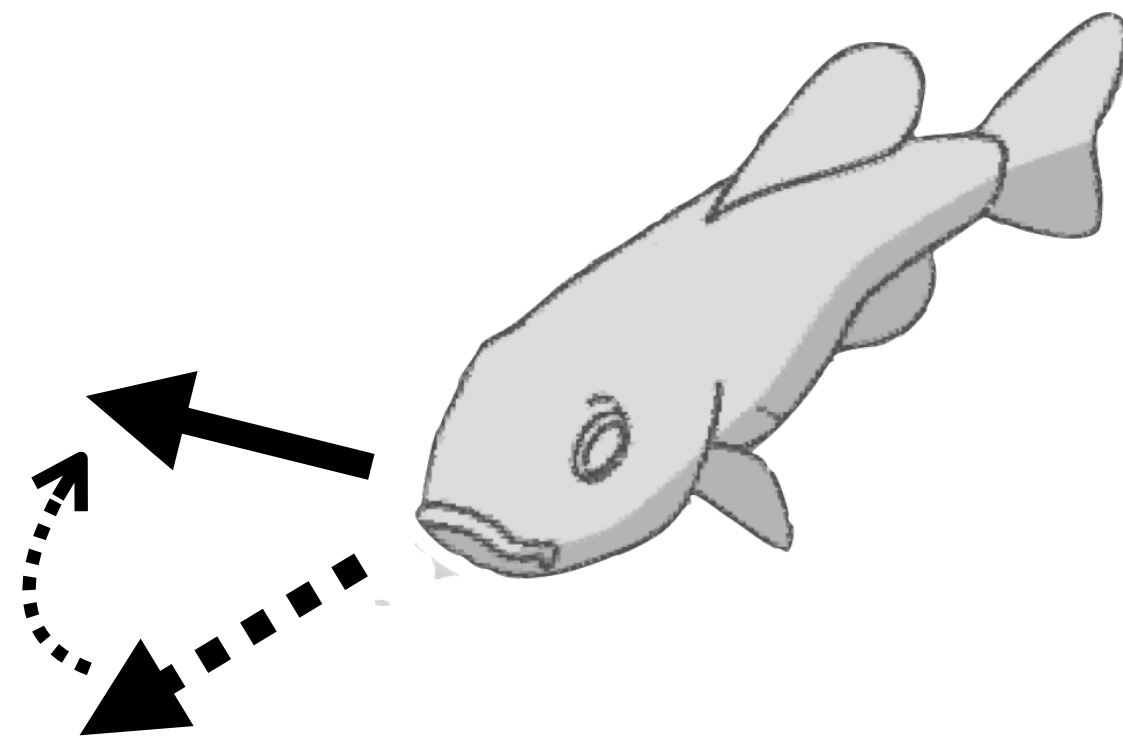
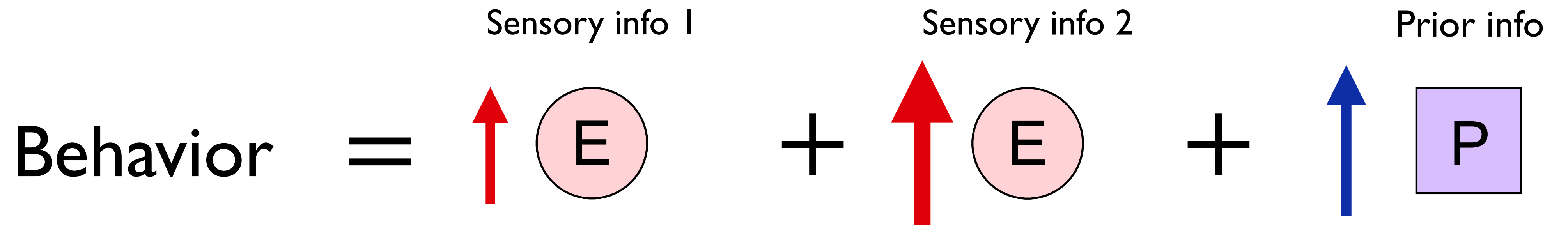
# Weights vs beliefs

$$\text{Behavior} = \omega_1 f_{\text{social}}(x) + \omega_2 f_{\text{env}}(z) + \epsilon$$



$\omega_1, \omega_2$  ?

# Weights vs beliefs

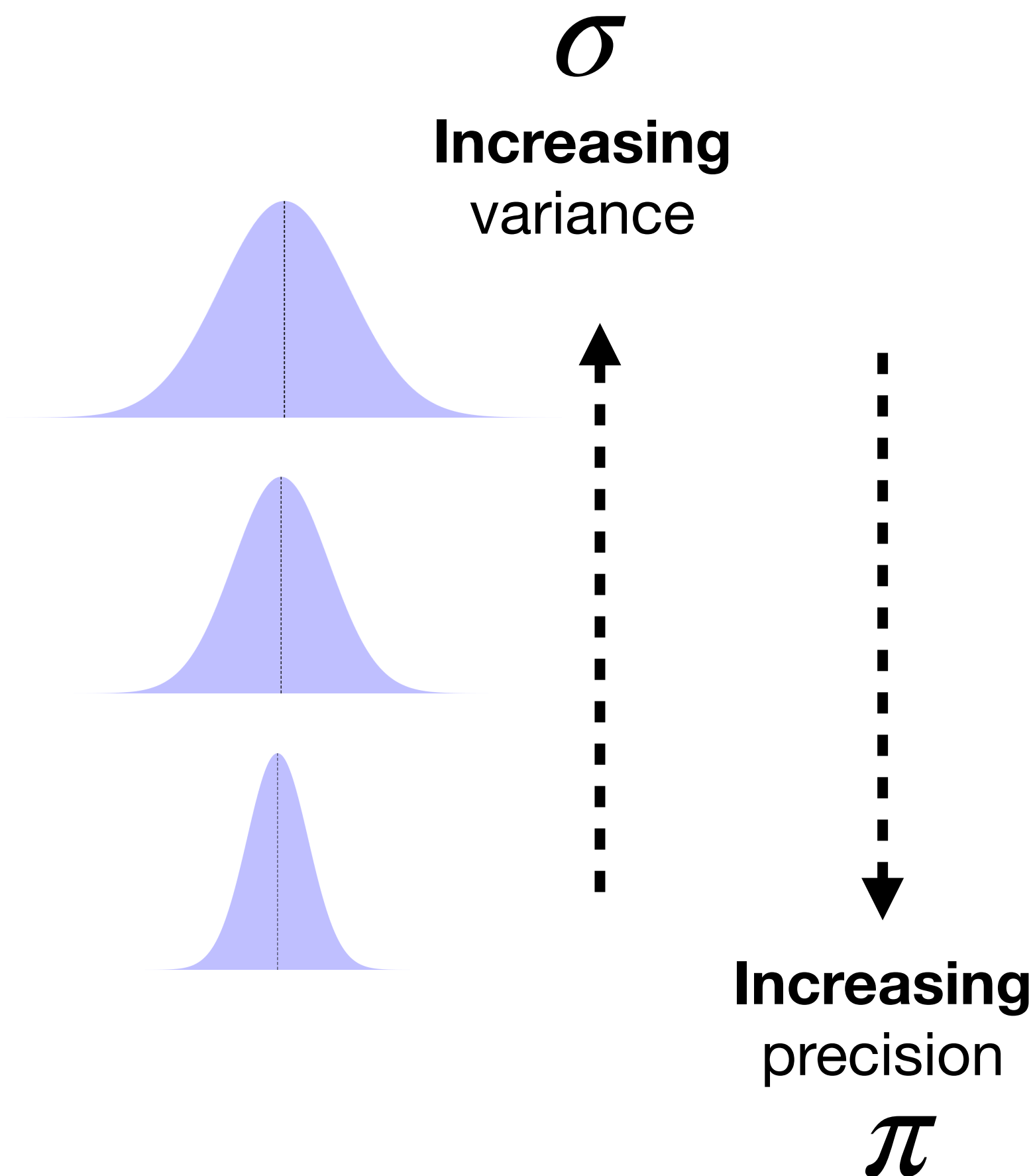
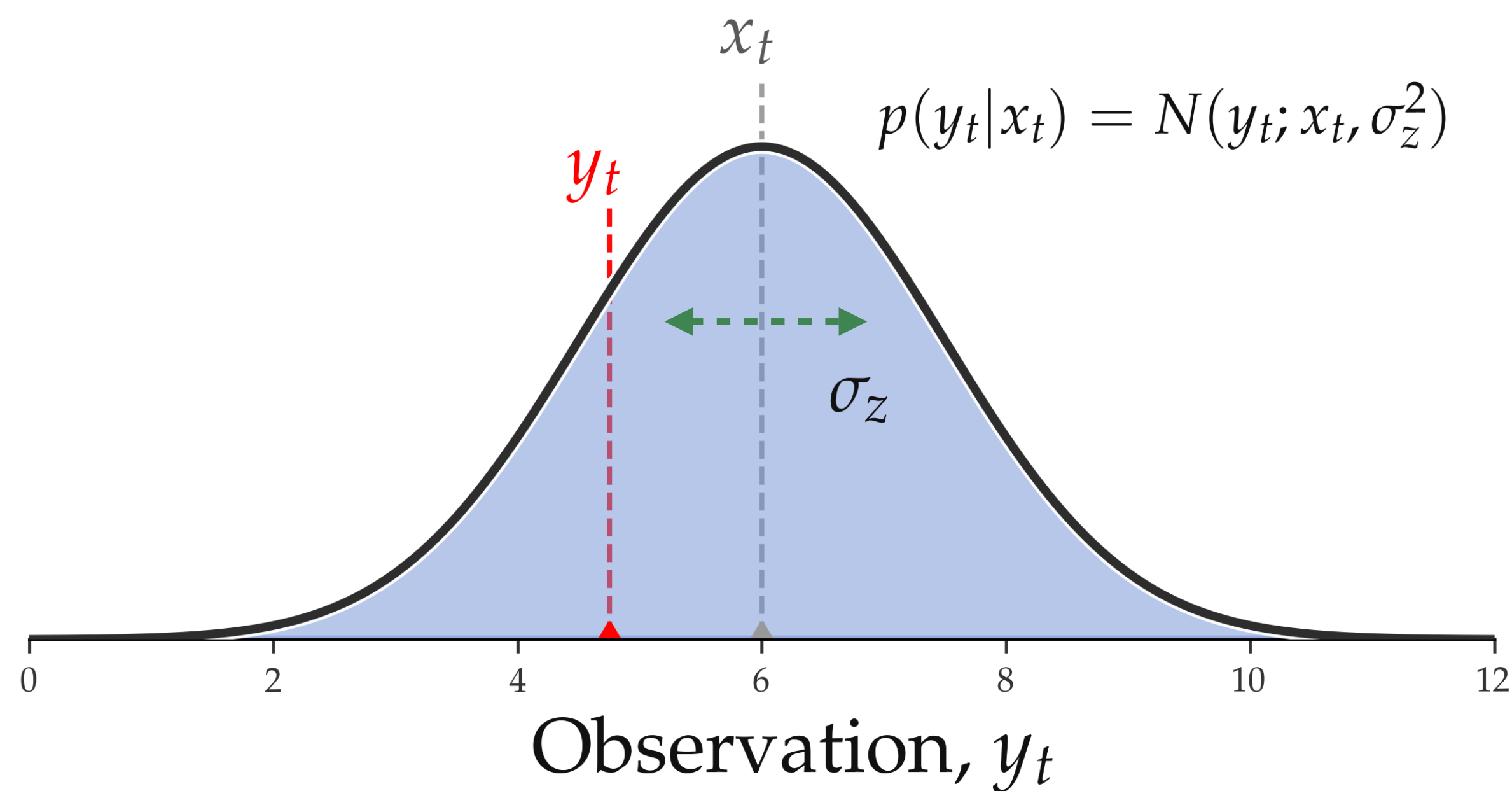
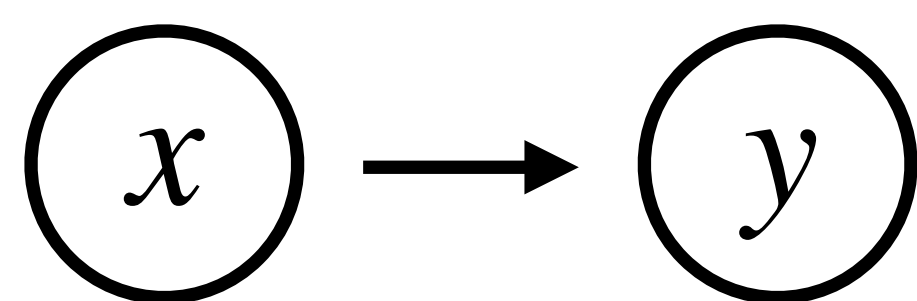


Beliefs about reliability of different types of **sensory** and **prior** information



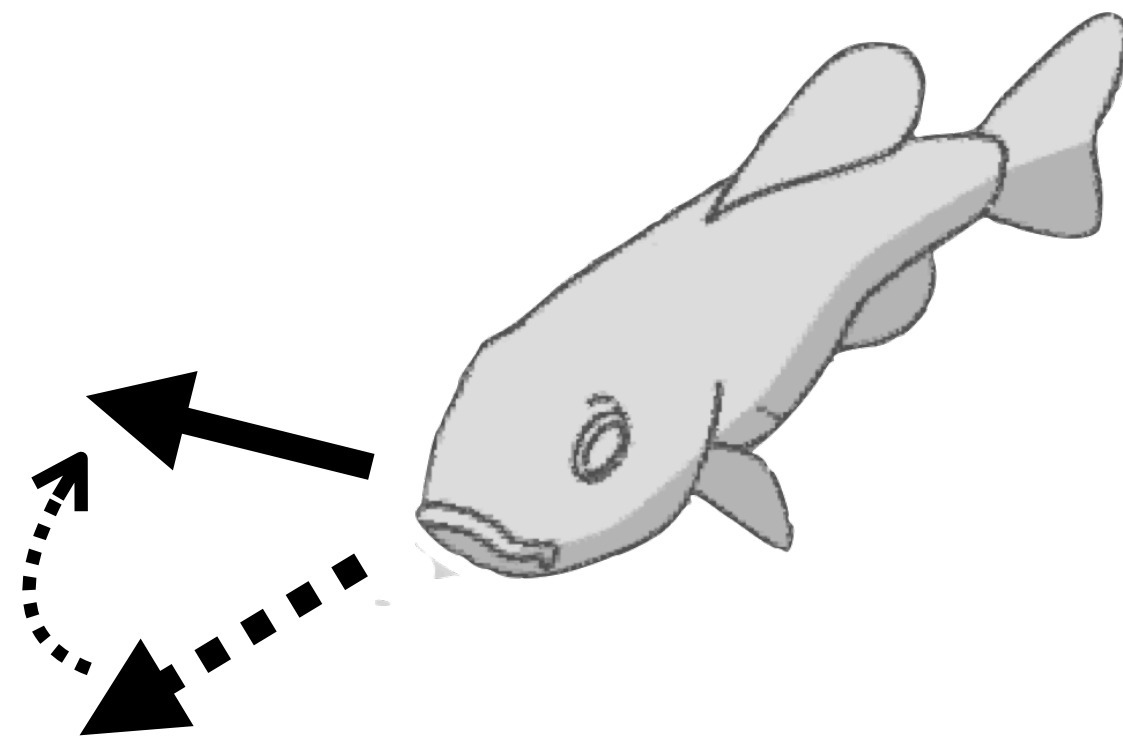
# Re-interpreting force-weights as beliefs about information reliability

Information from different sources can be more or less “trustworthy”



# Re-interpreting force-weights as beliefs about sensory reliability

$$\text{Movement} = \omega_1 f_{\text{social}}(x) + \omega_2 f_{\text{env}}(z) + \epsilon$$

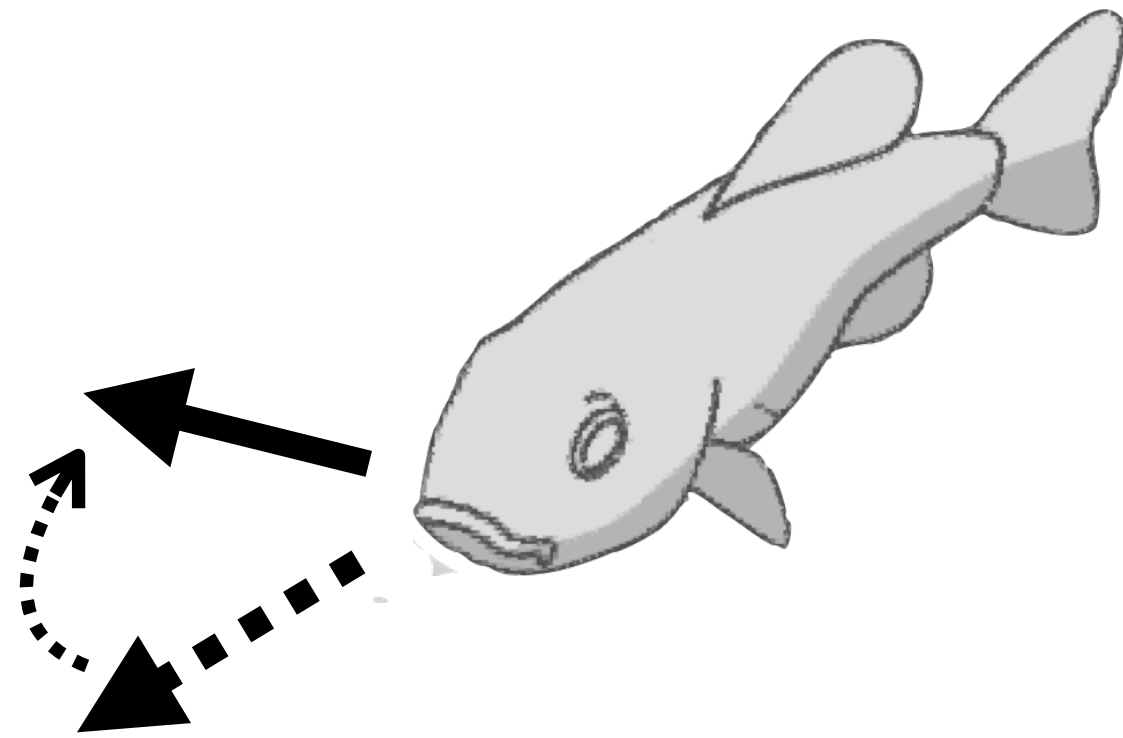


$\omega_1, \omega_2$  ?

# Re-interpreting force-weights as beliefs about sensory reliability

Prediction errors

$$\text{Movement} = \pi_1 \epsilon_1 + \pi_2 \epsilon_2$$



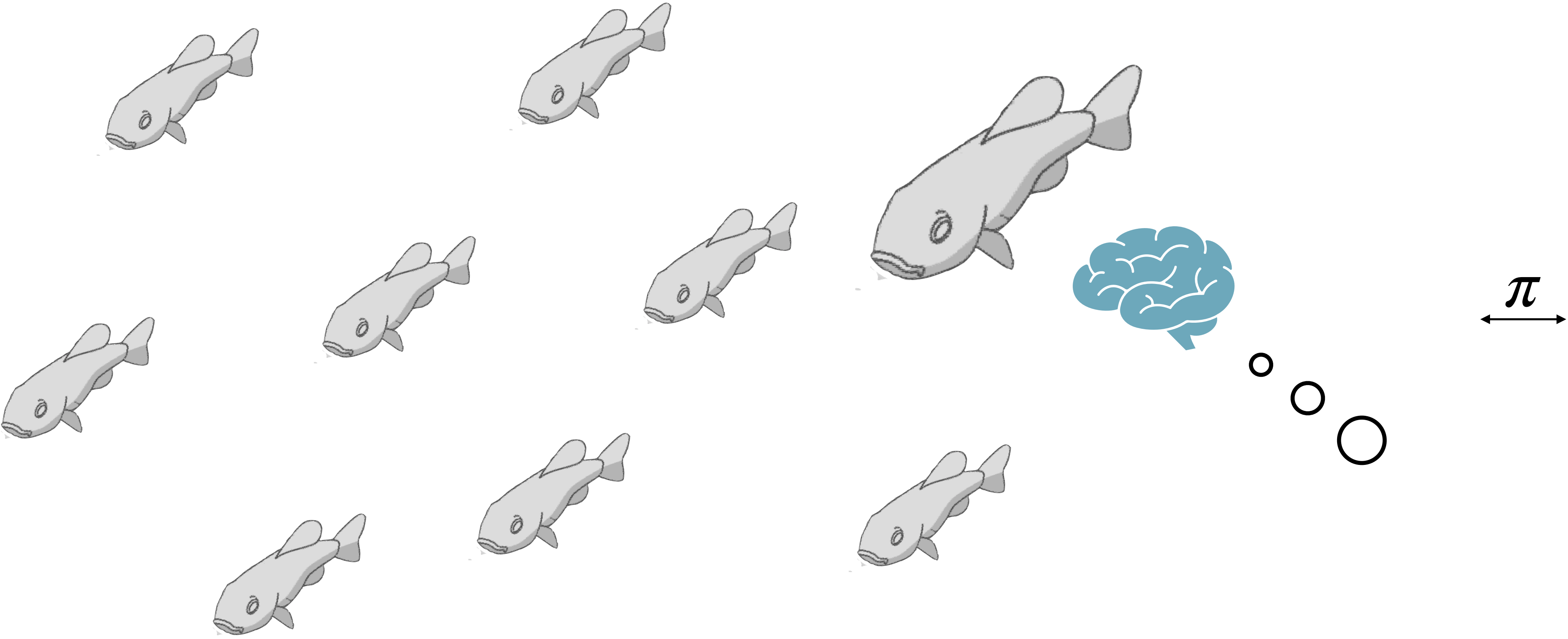
$\pi_1$

Beliefs about reliability of signal 1

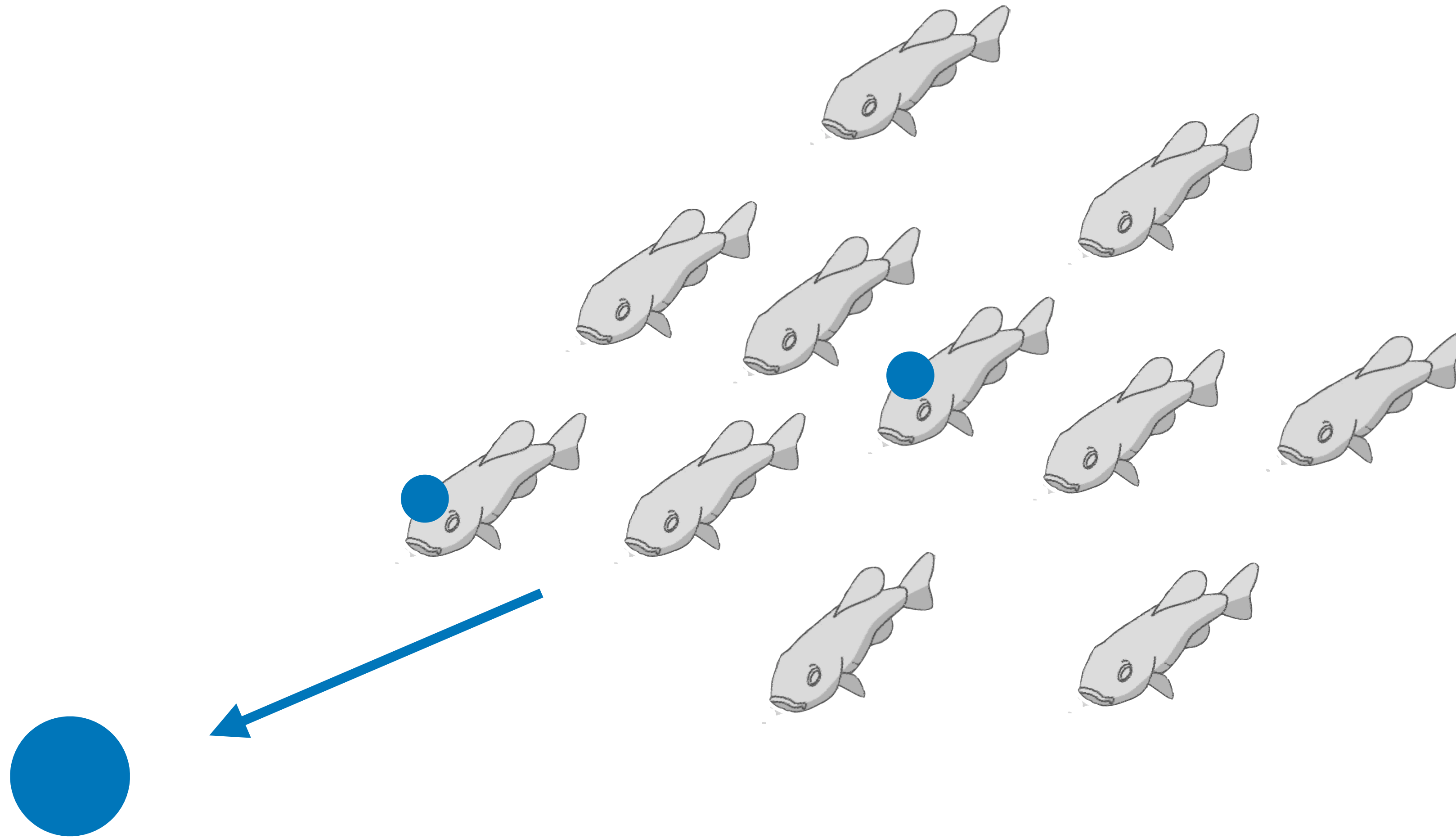
$\pi_2$

Beliefs about reliability of signal 2

# How do individual beliefs determine collective information processing?



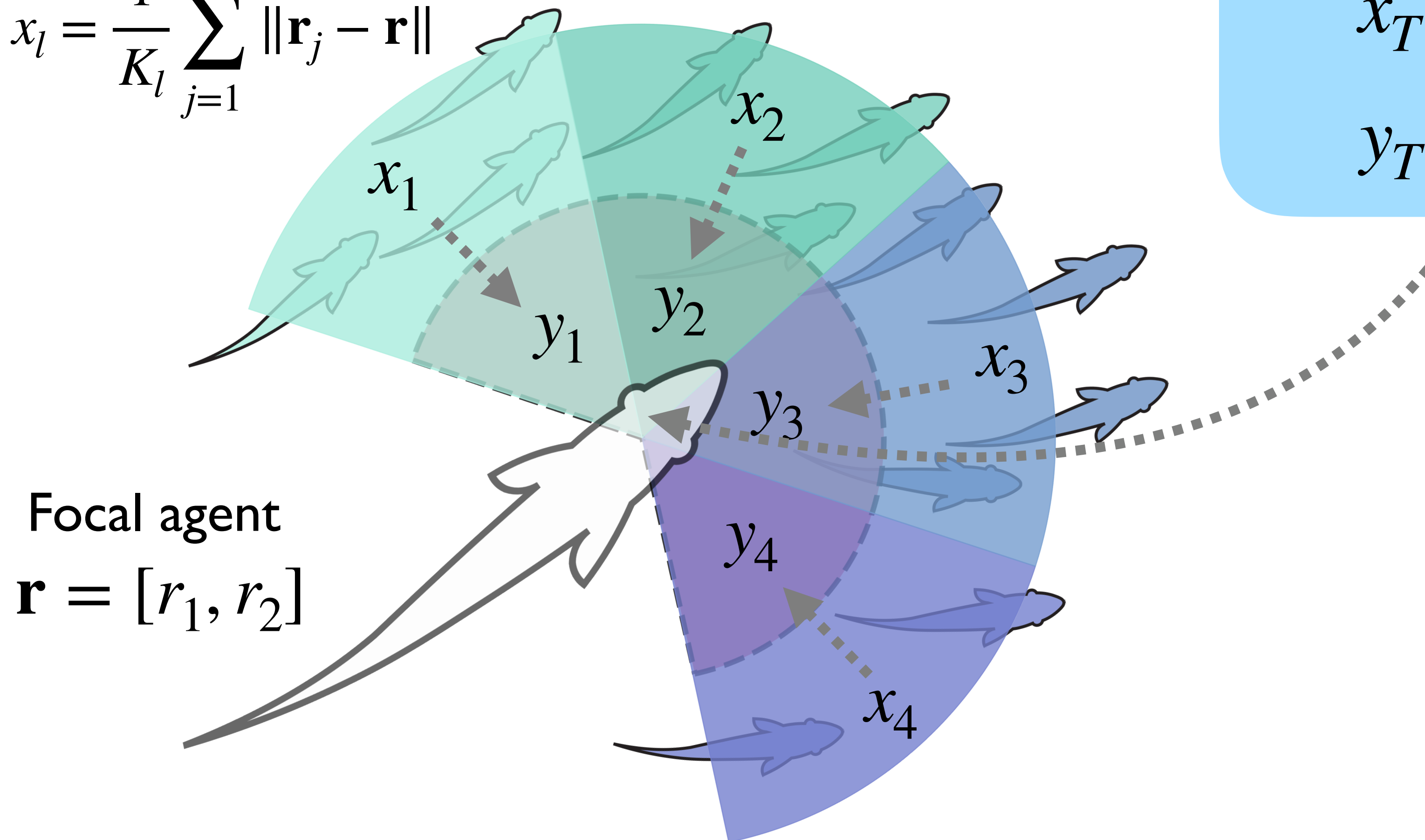
# Collective information transfer



# Some agents have an extra source of sensory information

Sector-specific average distance

$$x_l = \frac{1}{K_l} \sum_{j=1}^{K_l} \|\mathbf{r}_j - \mathbf{r}\|$$



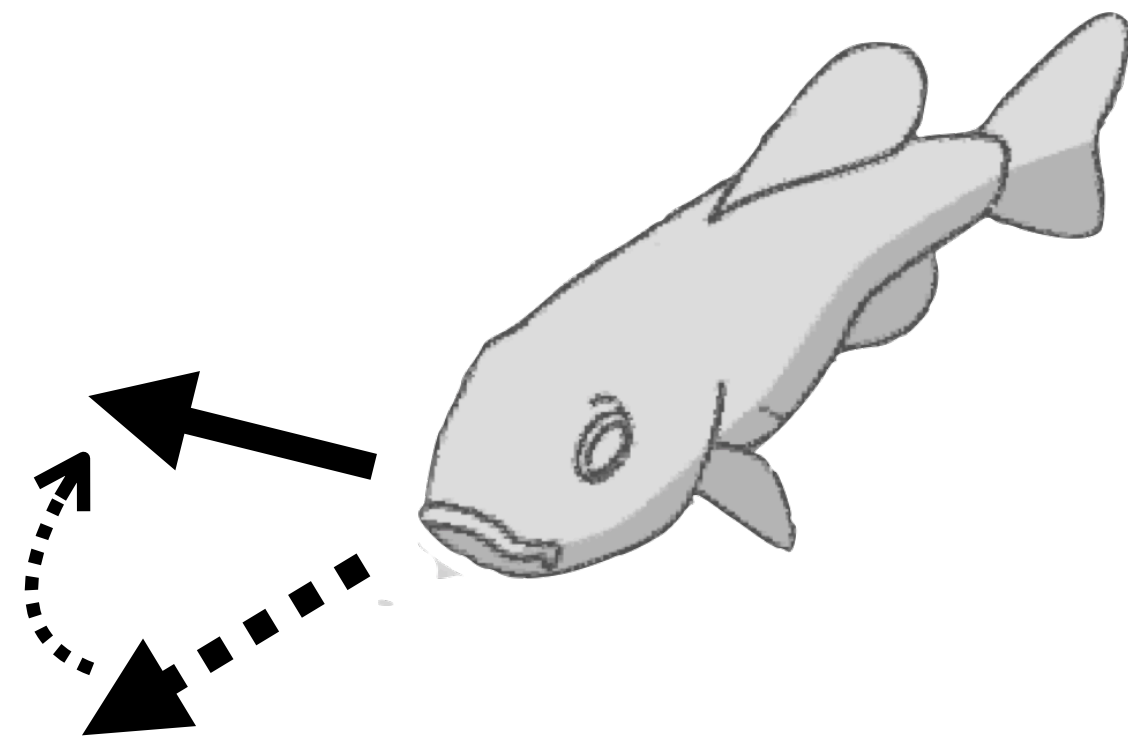
New sensory channel:  
Distance-to-target

$$x_T = \|\mathbf{r}_j - \mathbf{T}\|$$

$$y_T = x_T + z$$

# Action becomes a precision-weighted sum of vectors

$$\text{Heading direction} = \pi_{Soc} \times \uparrow + \pi_{Target} \times \uparrow$$



$\pi_{Soc}$

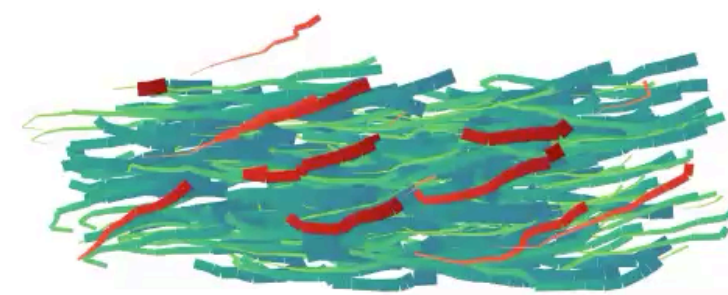
Beliefs about reliability of social info

$\pi_{Target}$

Beliefs about reliability of target info

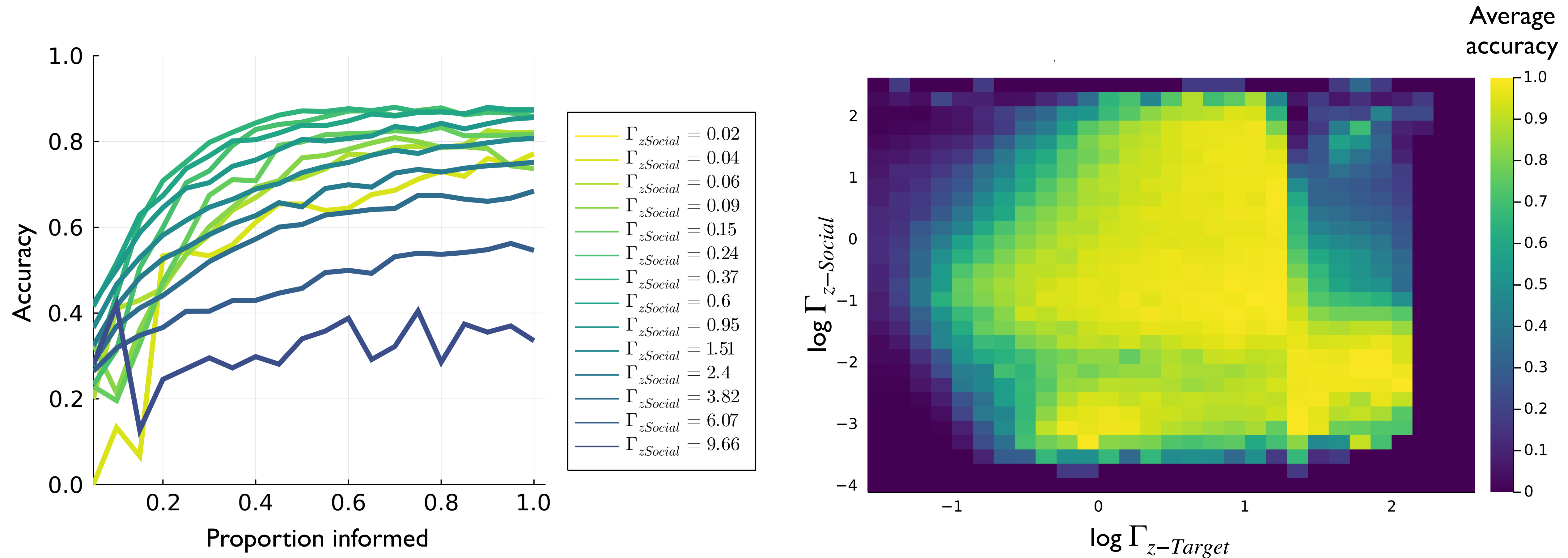
# Collective navigates to target despite uninformed individuals

● Target 1





# High accuracy within a large regime of social and target precisions

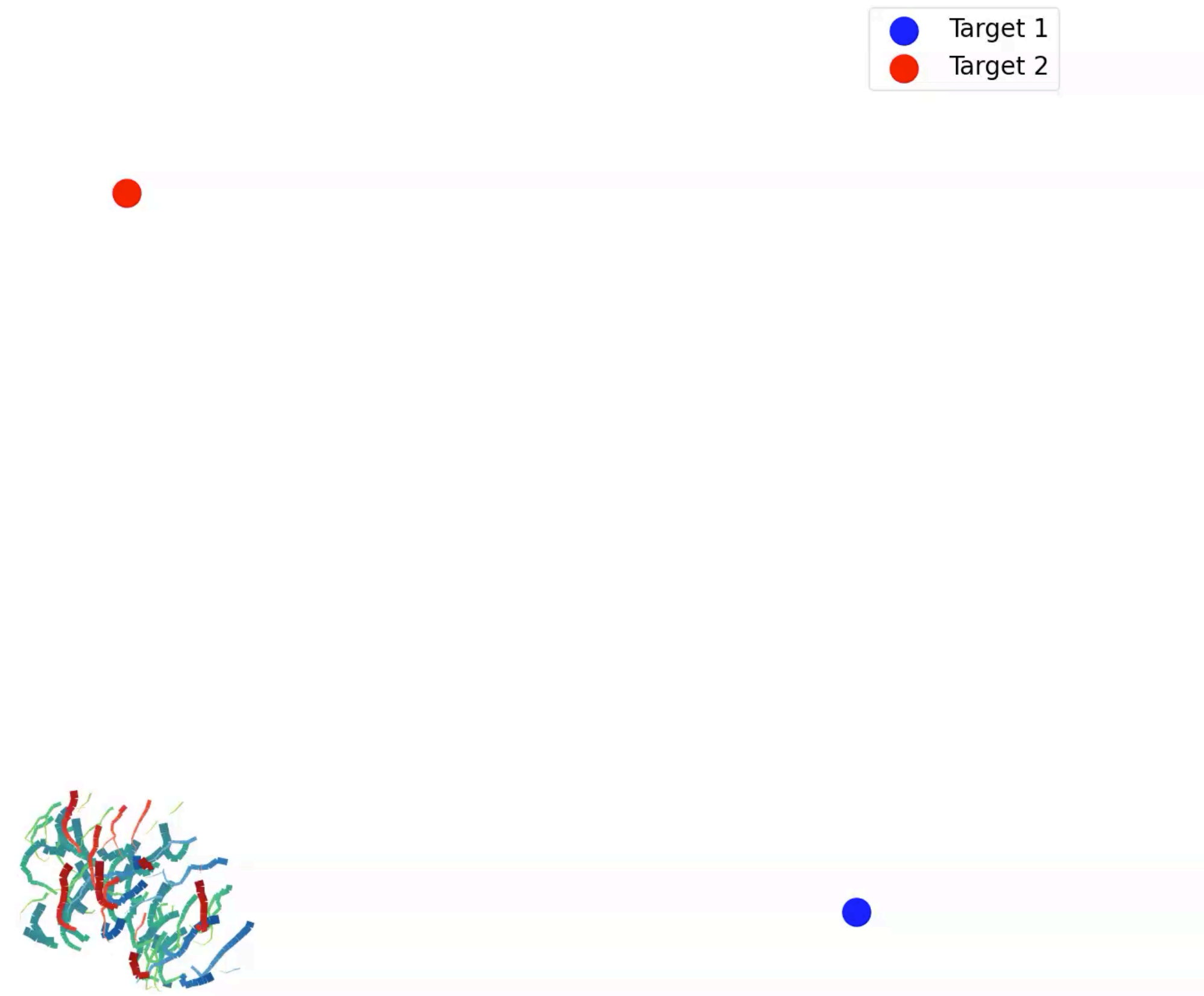


$$\Gamma_{z-Social} = \pi_{Social}$$

$$\Gamma_{z-Target} = \pi_{Target}$$

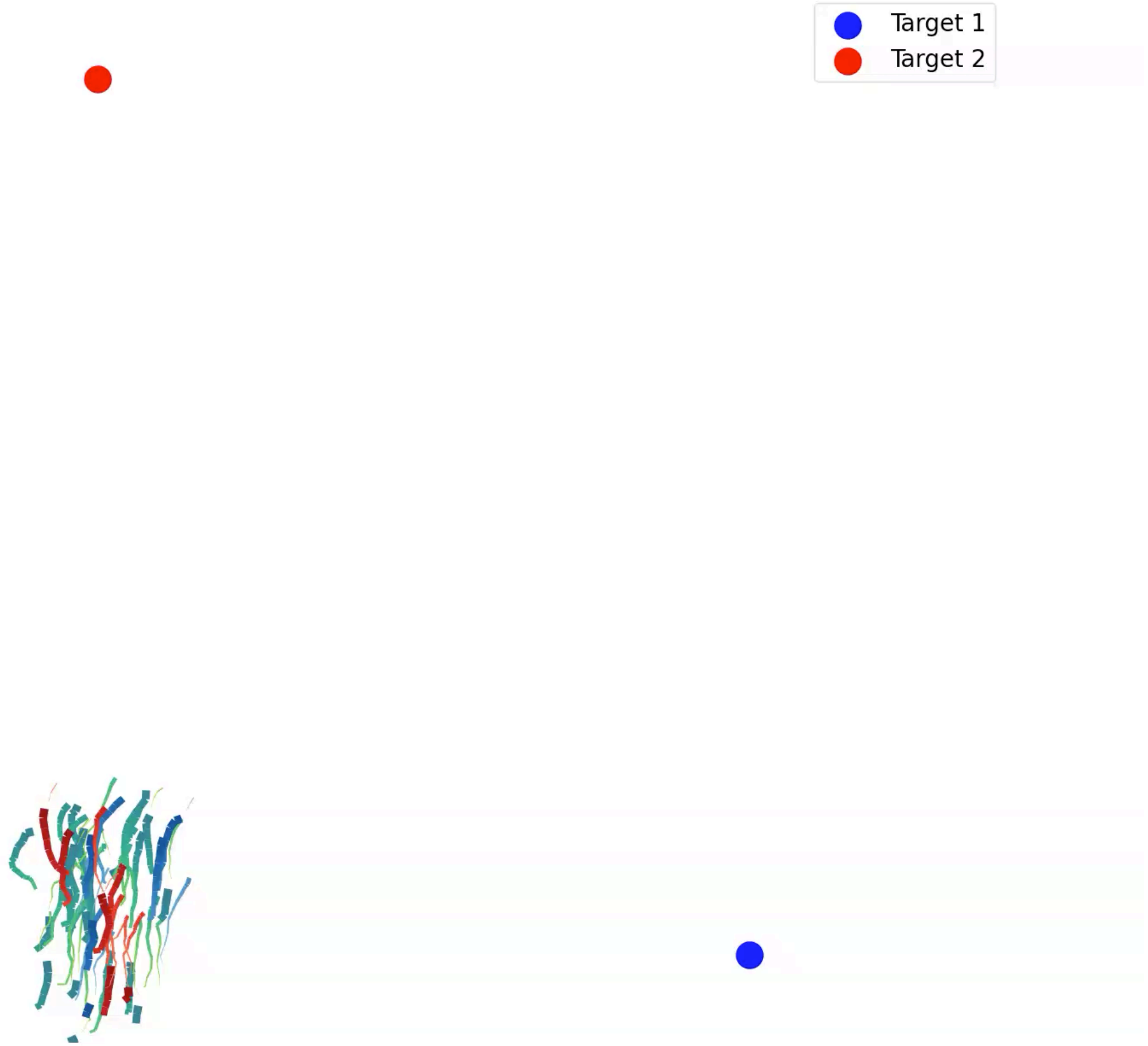
# Group fission in presence of multiple targets

$$\Gamma_{z-Social} = \pi_{Social} = 1$$



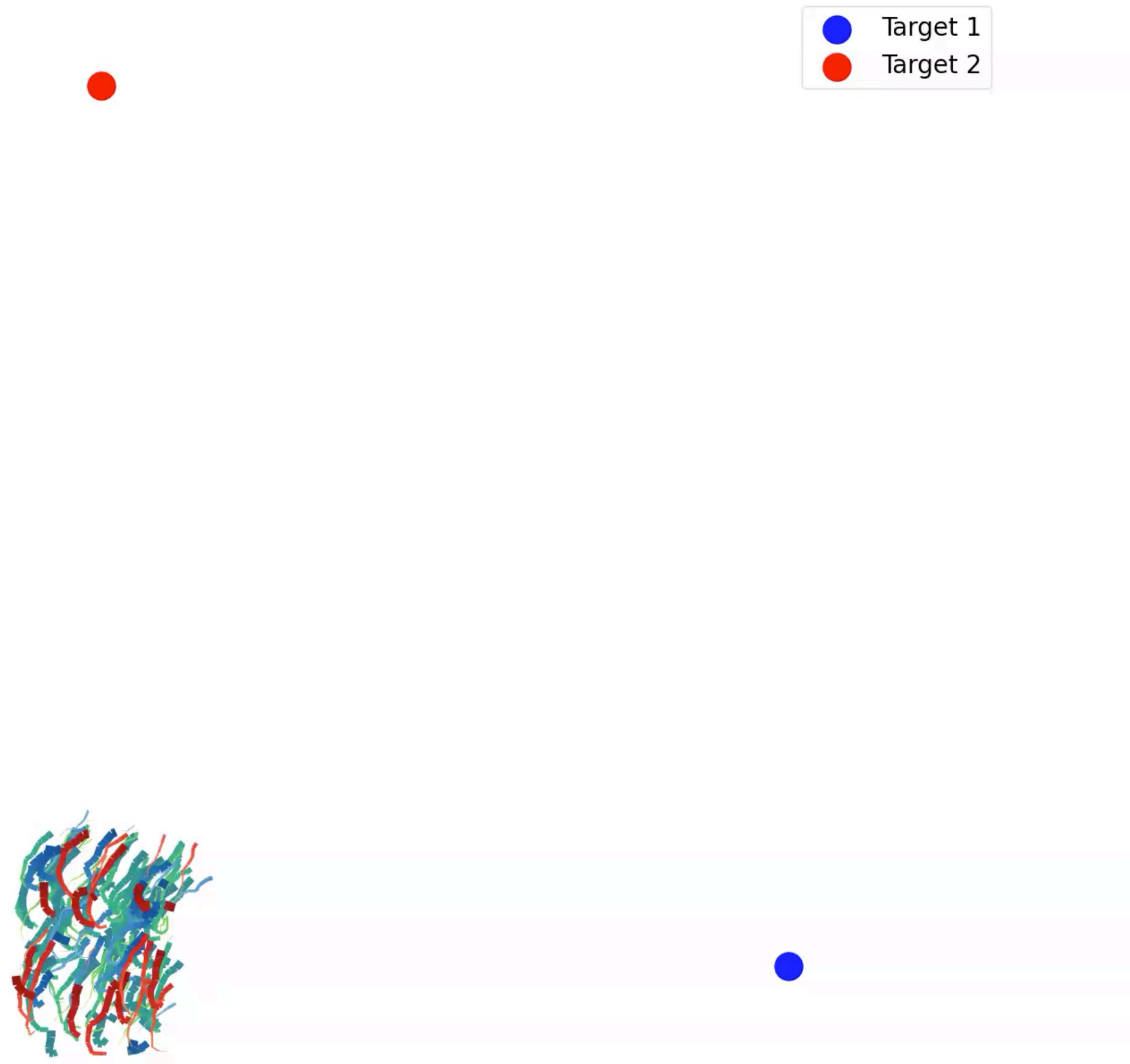
# Decreasing social precision leads to consensus

$$\Gamma_{z-Social} = \pi_{Social} = \frac{1}{4}$$



# Oscillations between the targets

$$\Gamma_{z-Social} = \pi_{Social} = \frac{1}{4}$$



# Adapting the generative model

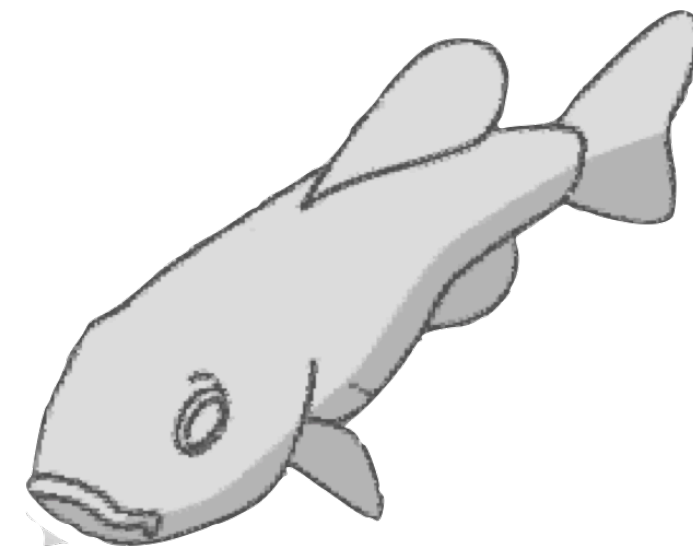
$F = \text{Surprise} = \text{Prediction error}$

$$\frac{d\mu}{dt} = -\nabla_{\mu} F(\mu, y)$$

Perception

$$\frac{d\mathbf{v}}{dt} = -\nabla_{\mathbf{v}} F(\mu, y(\mathbf{v}))$$

Action



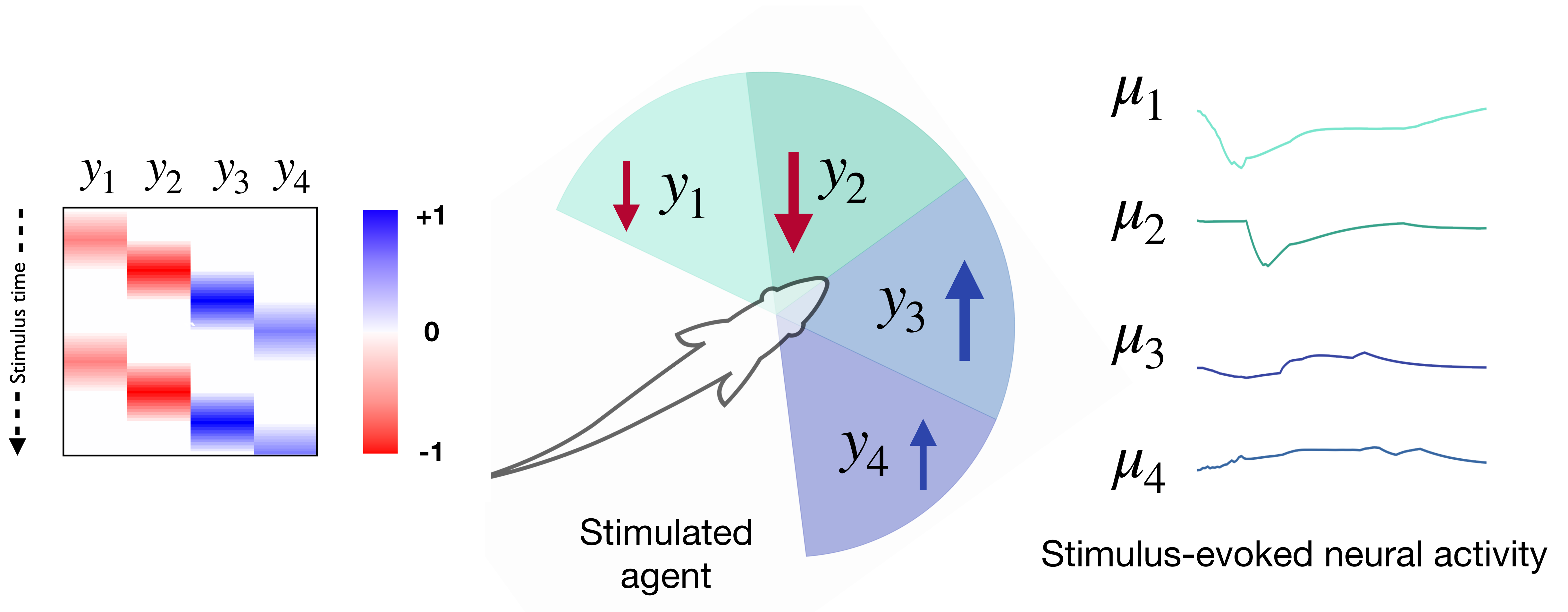
$$\frac{d\theta}{dt} = -\nabla_{\theta} F(\mu, y, \theta)$$

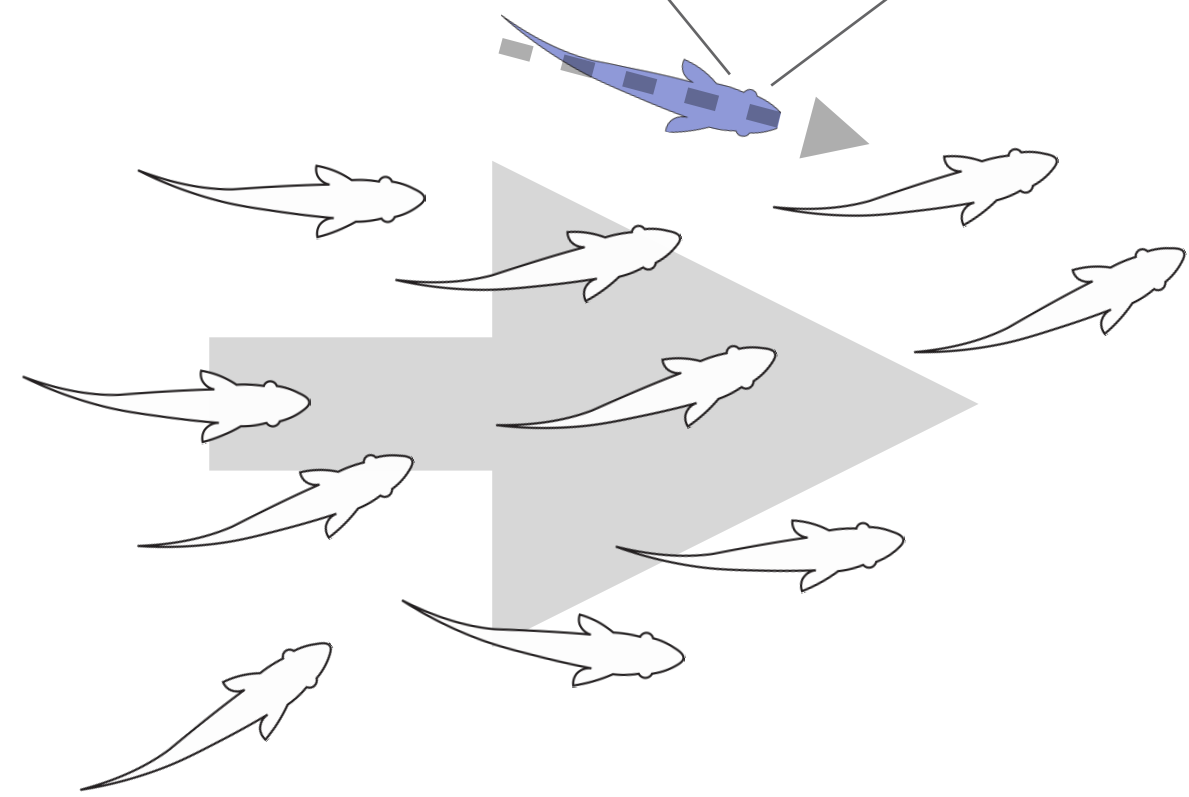
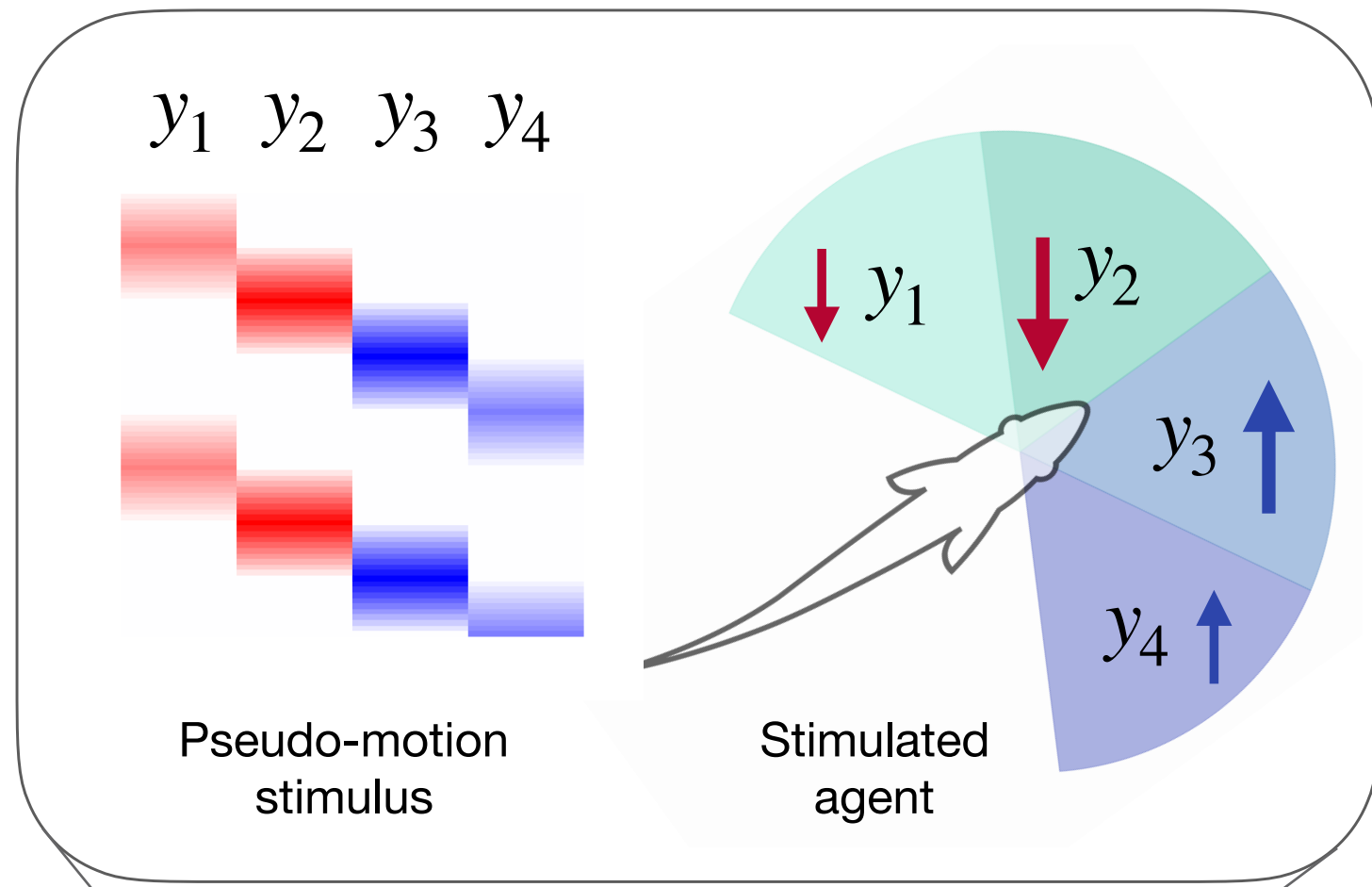
Plasticity

$\pi_z$

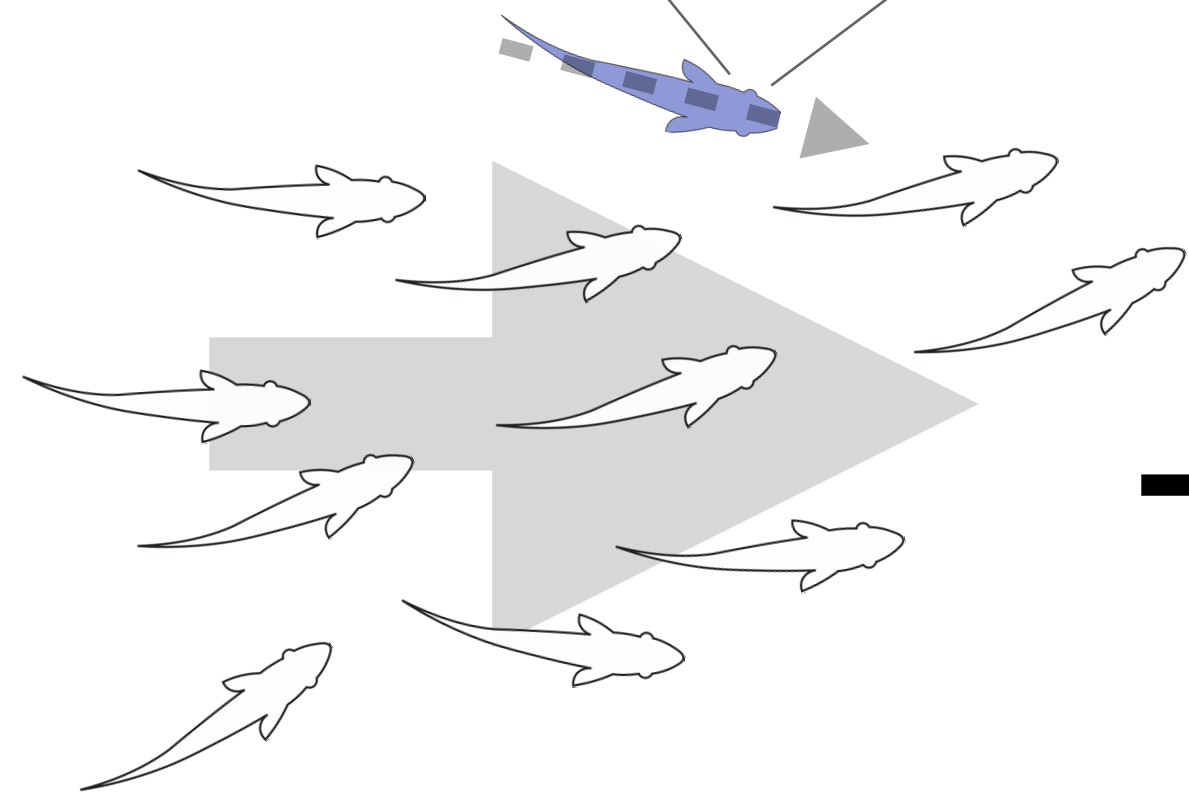
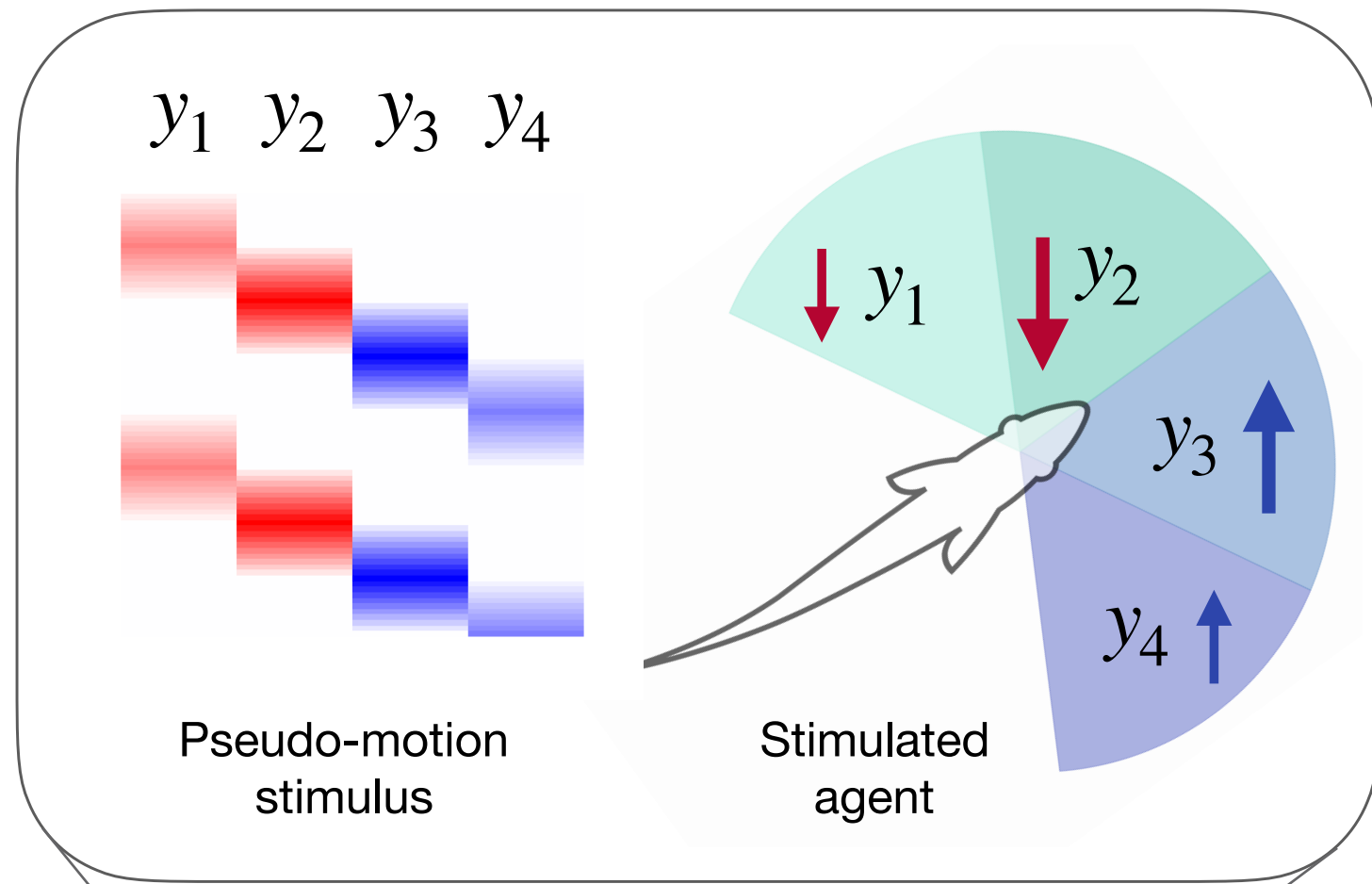
A horizontal double-headed arrow pointing both left and right, positioned below the symbol  $\pi_z$ .

# Phantom prediction errors in single agents

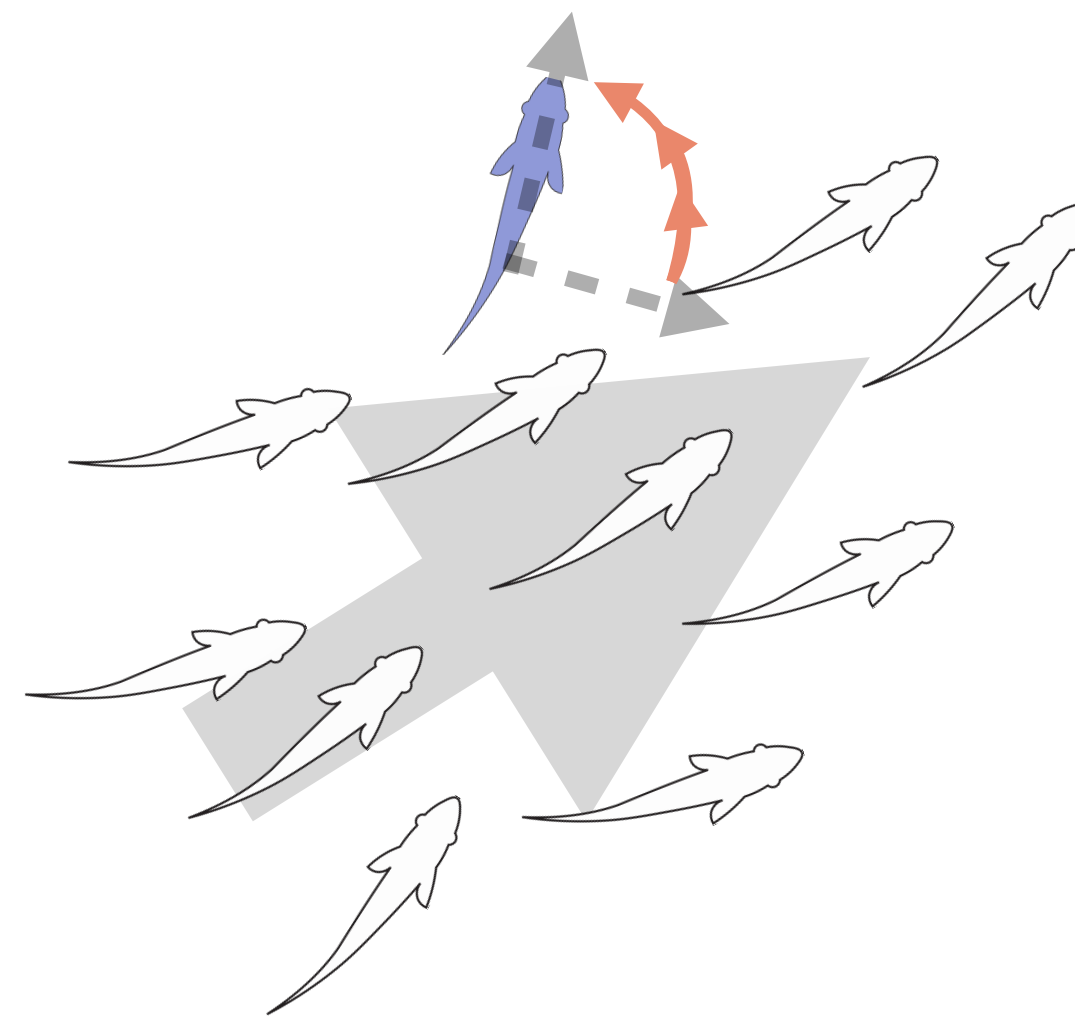




$T = 0$

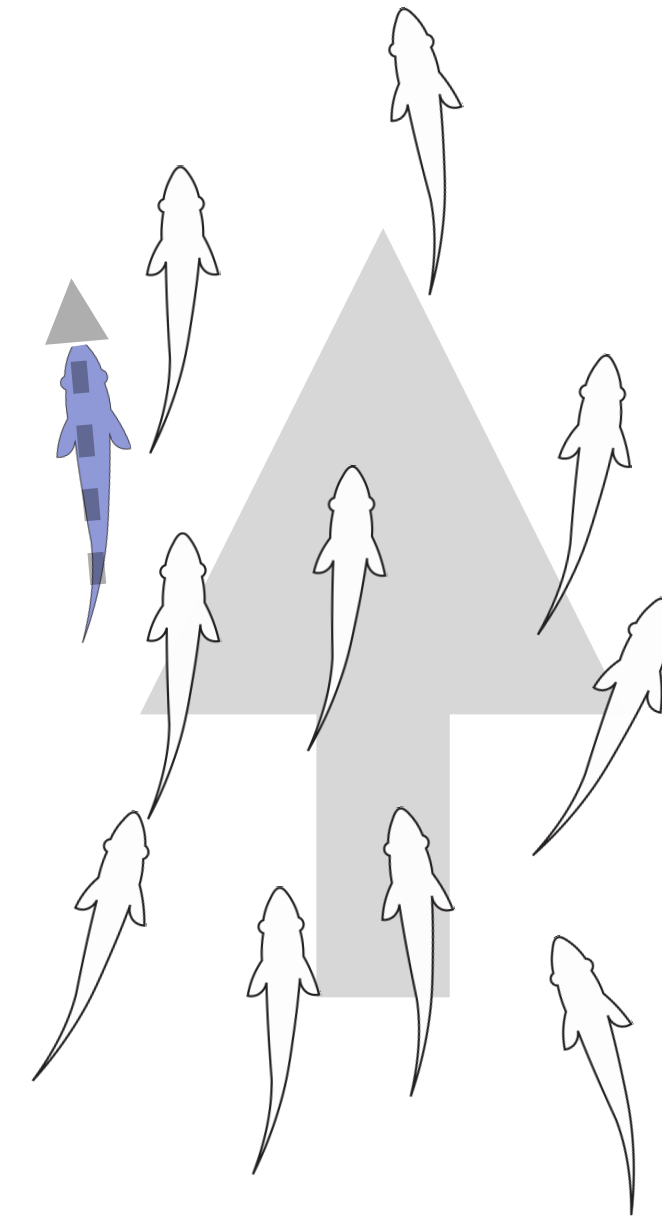
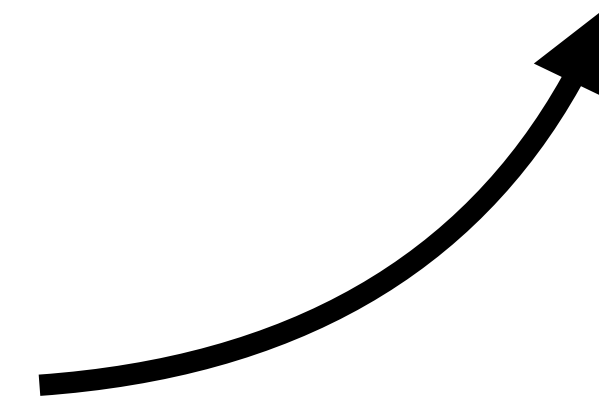
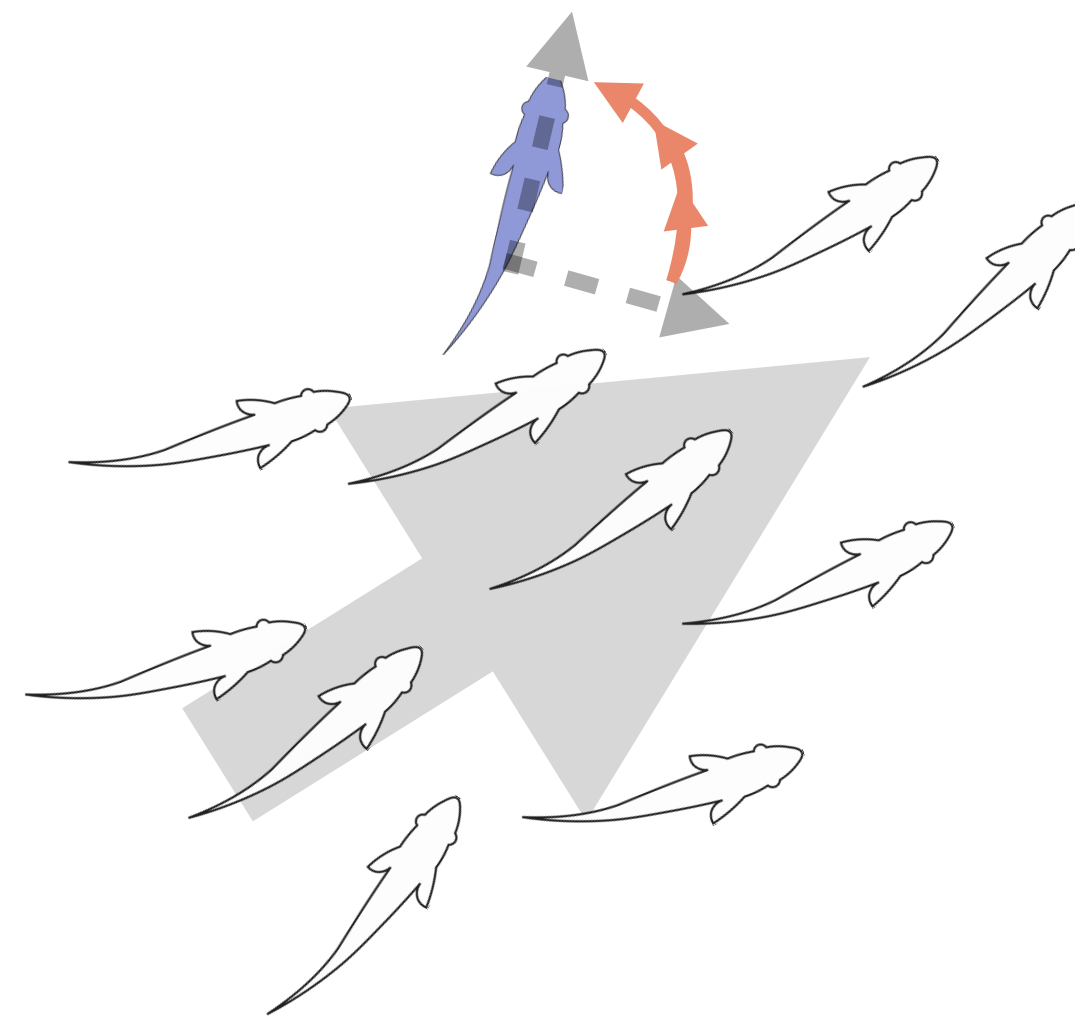
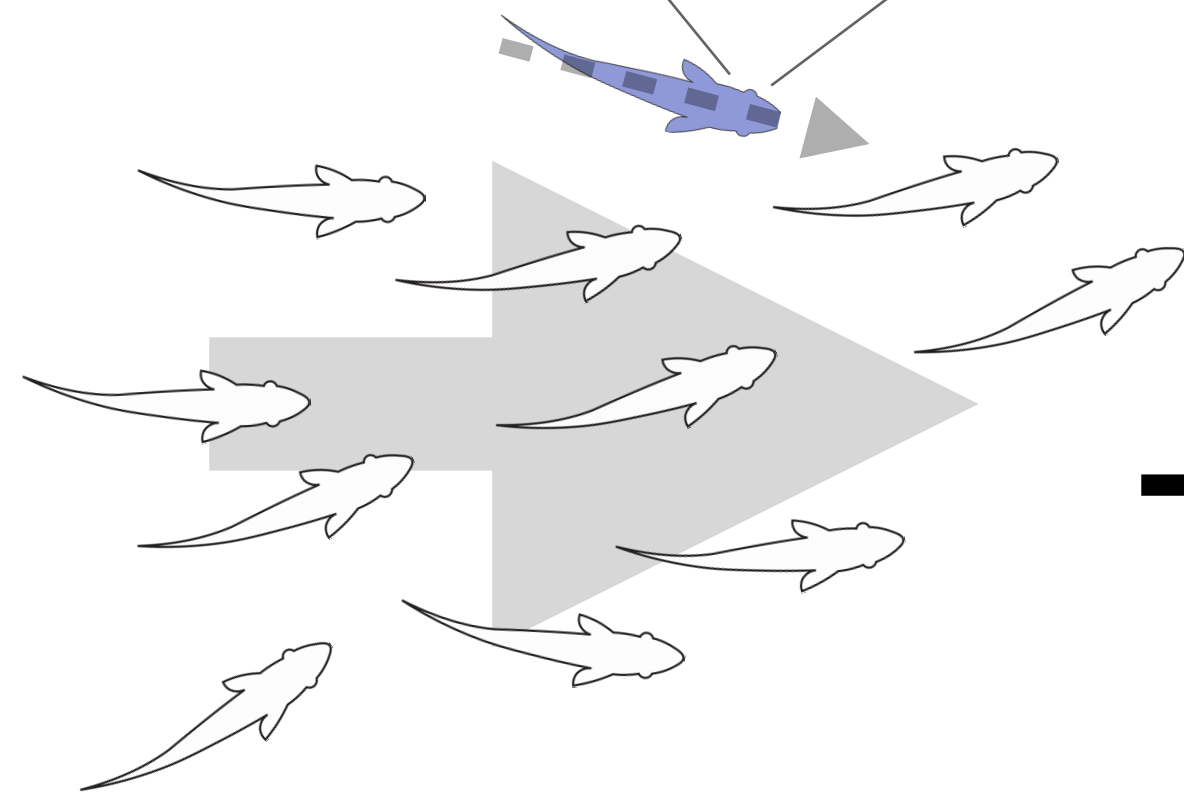
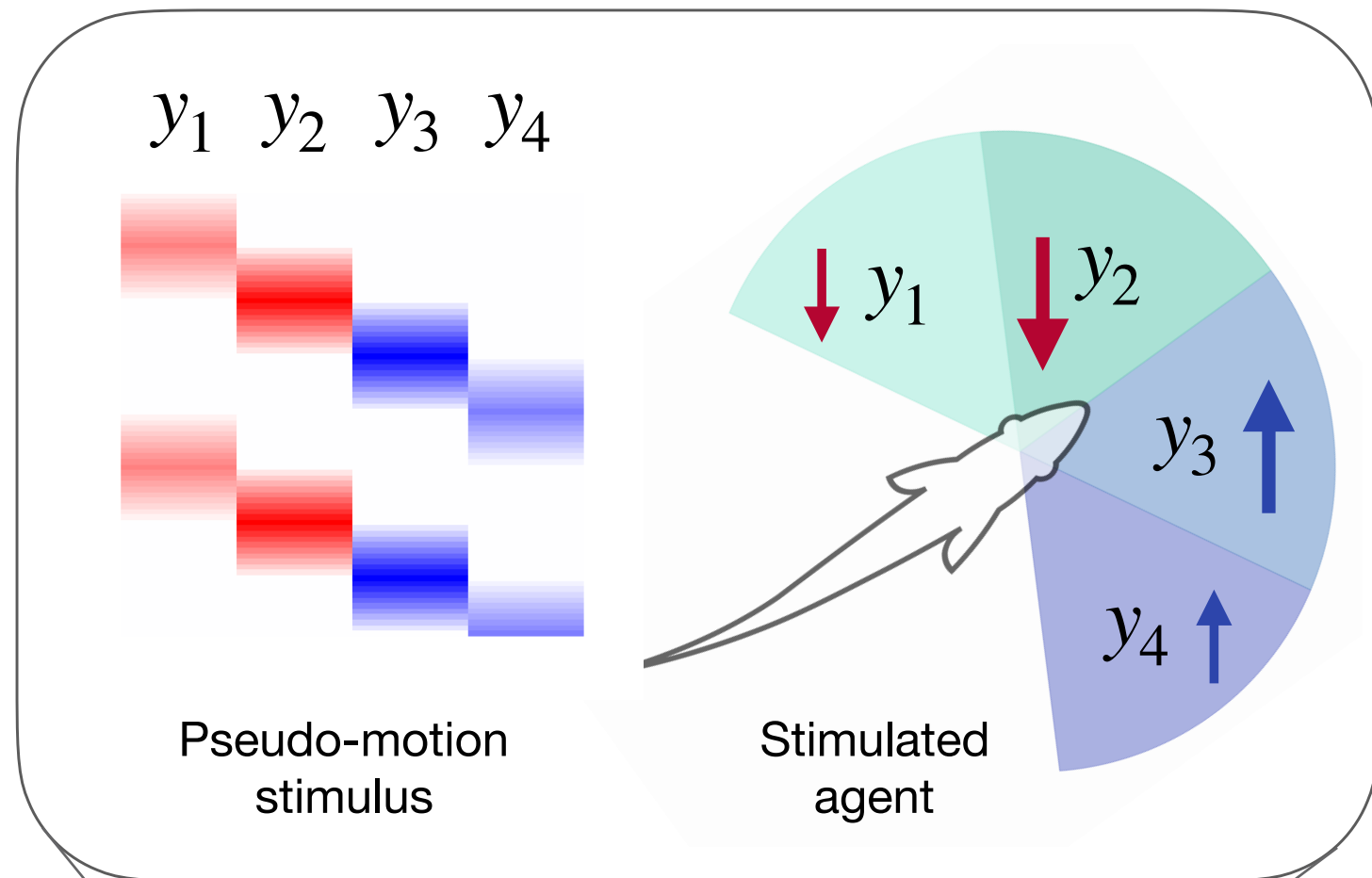


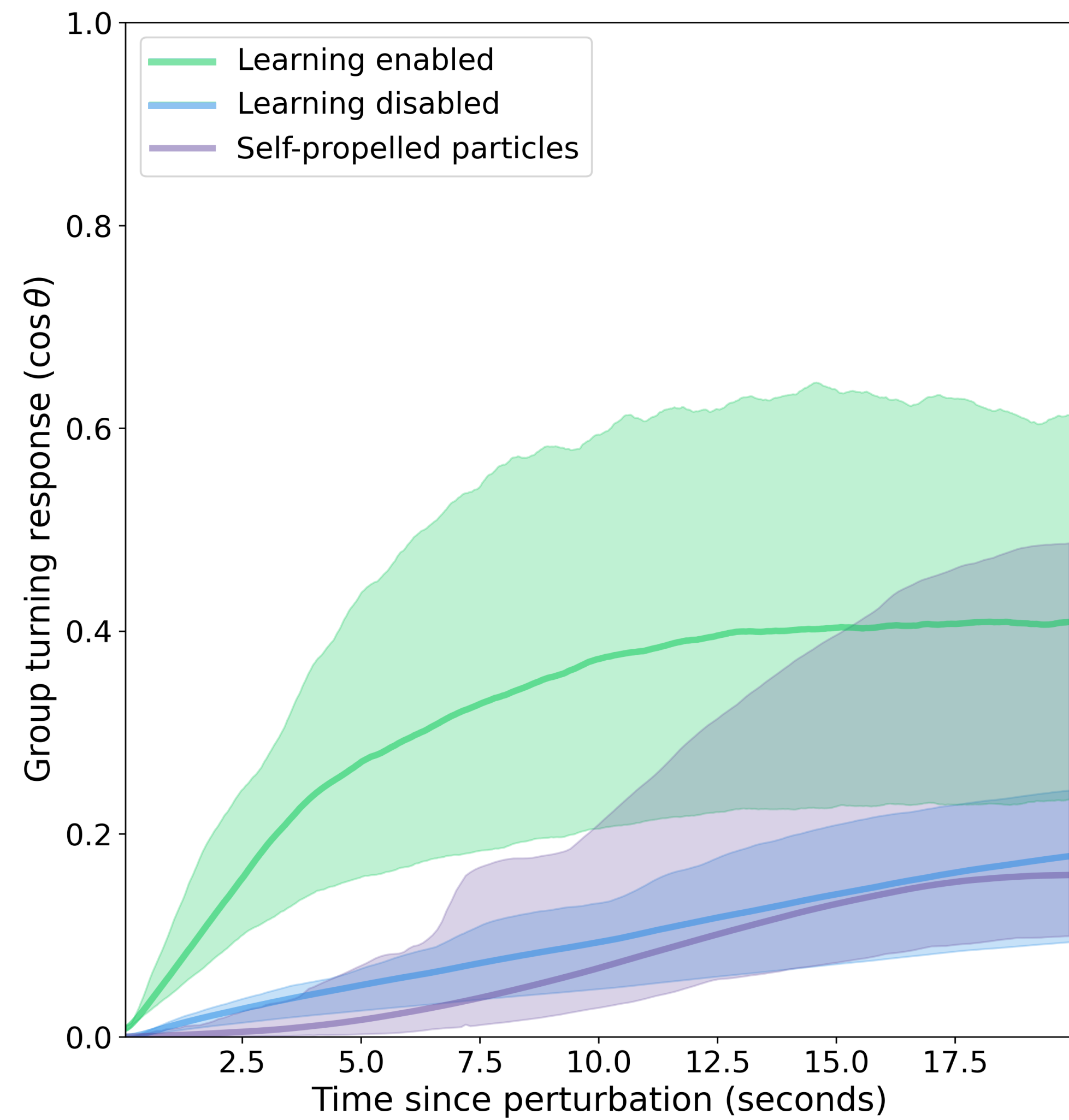
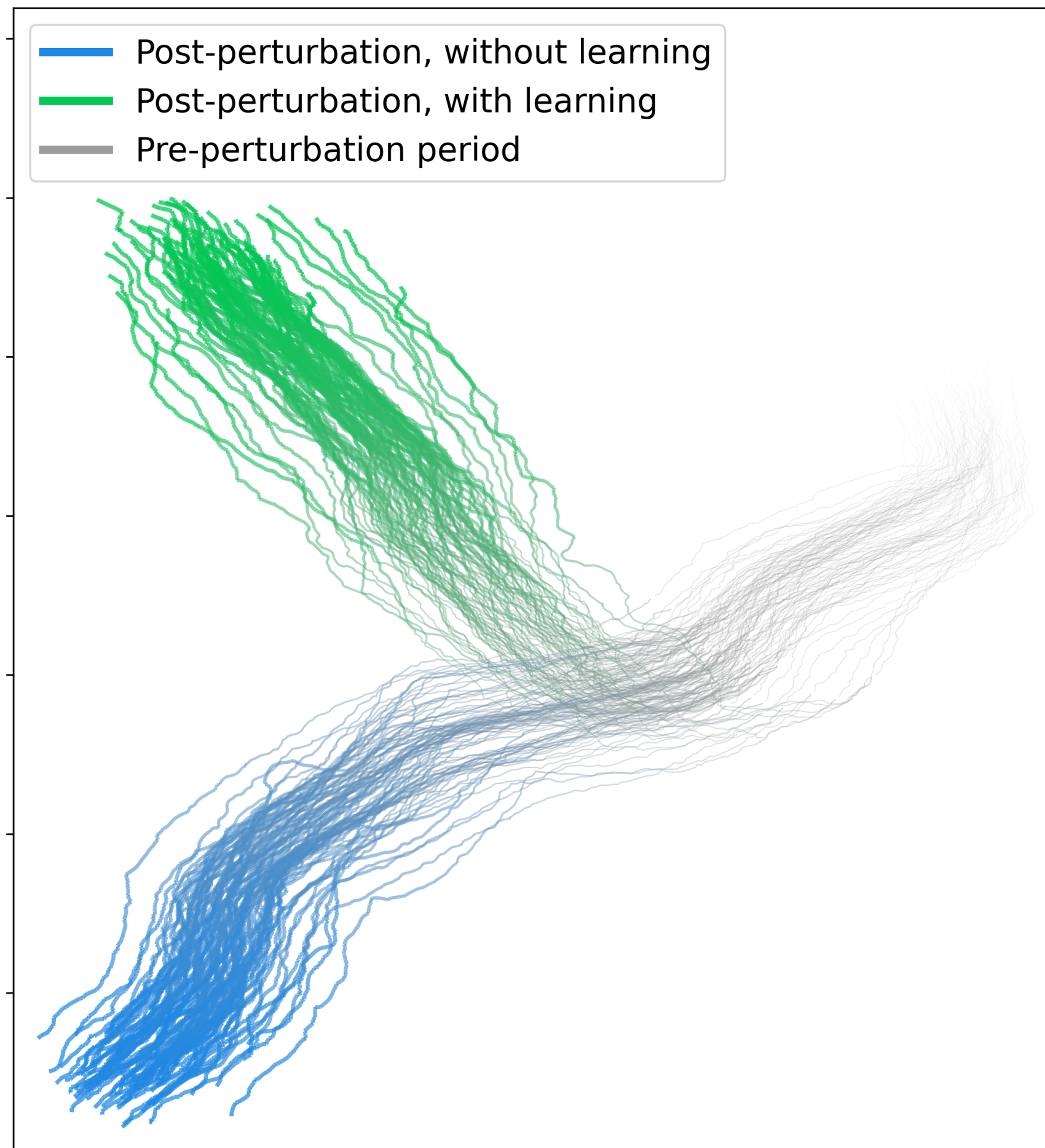
$T = 0$

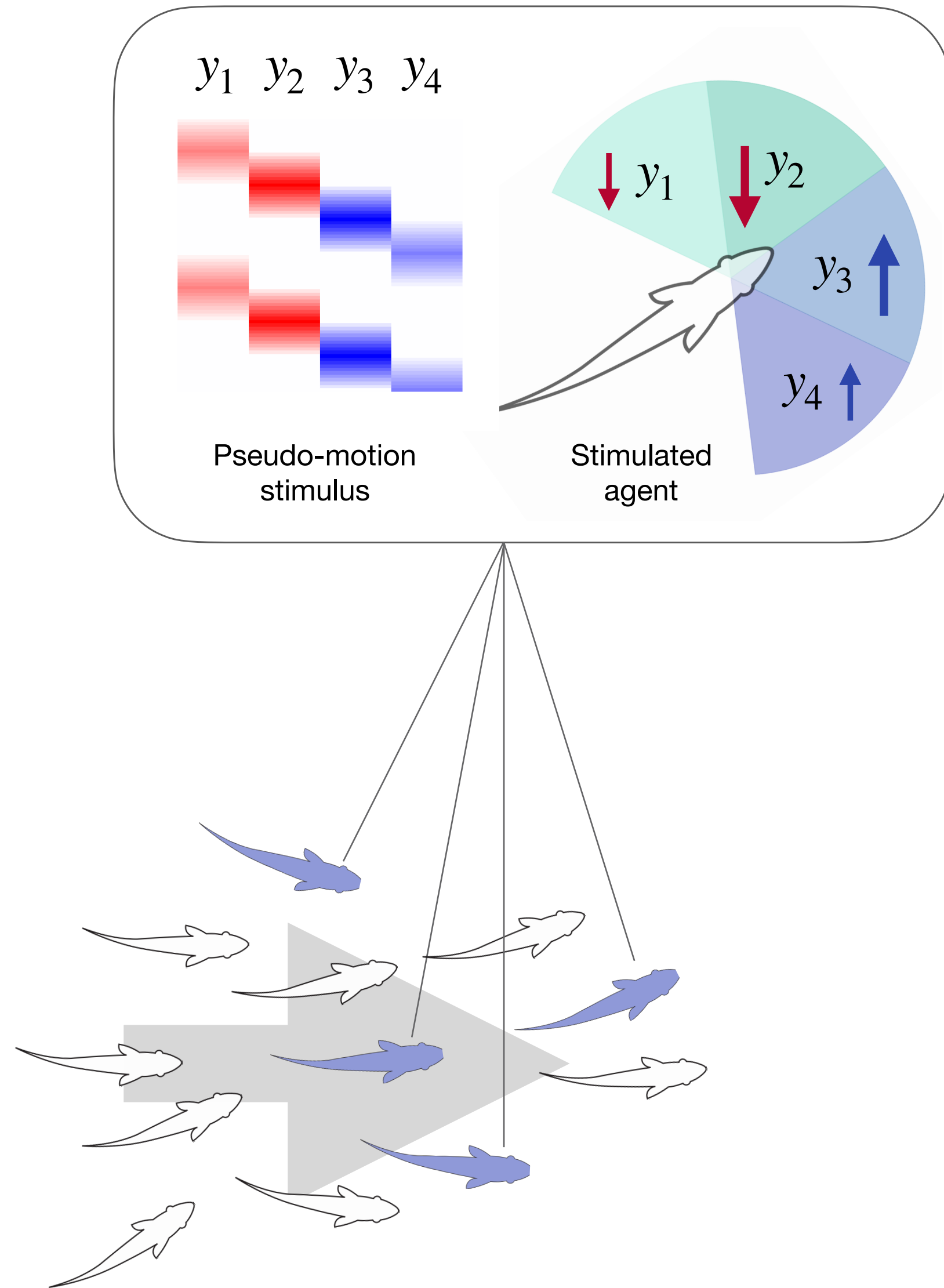


$T = 1000$



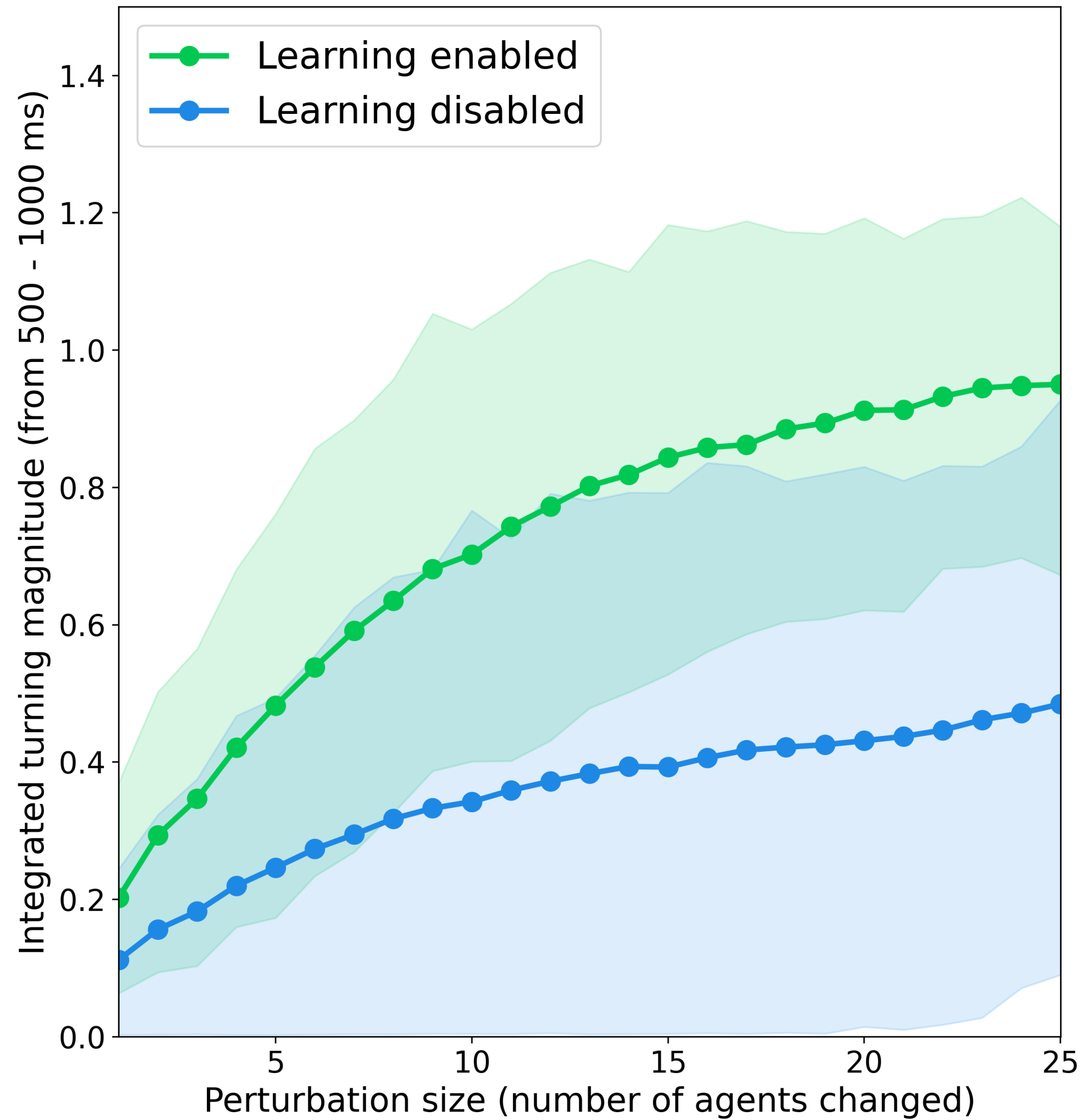




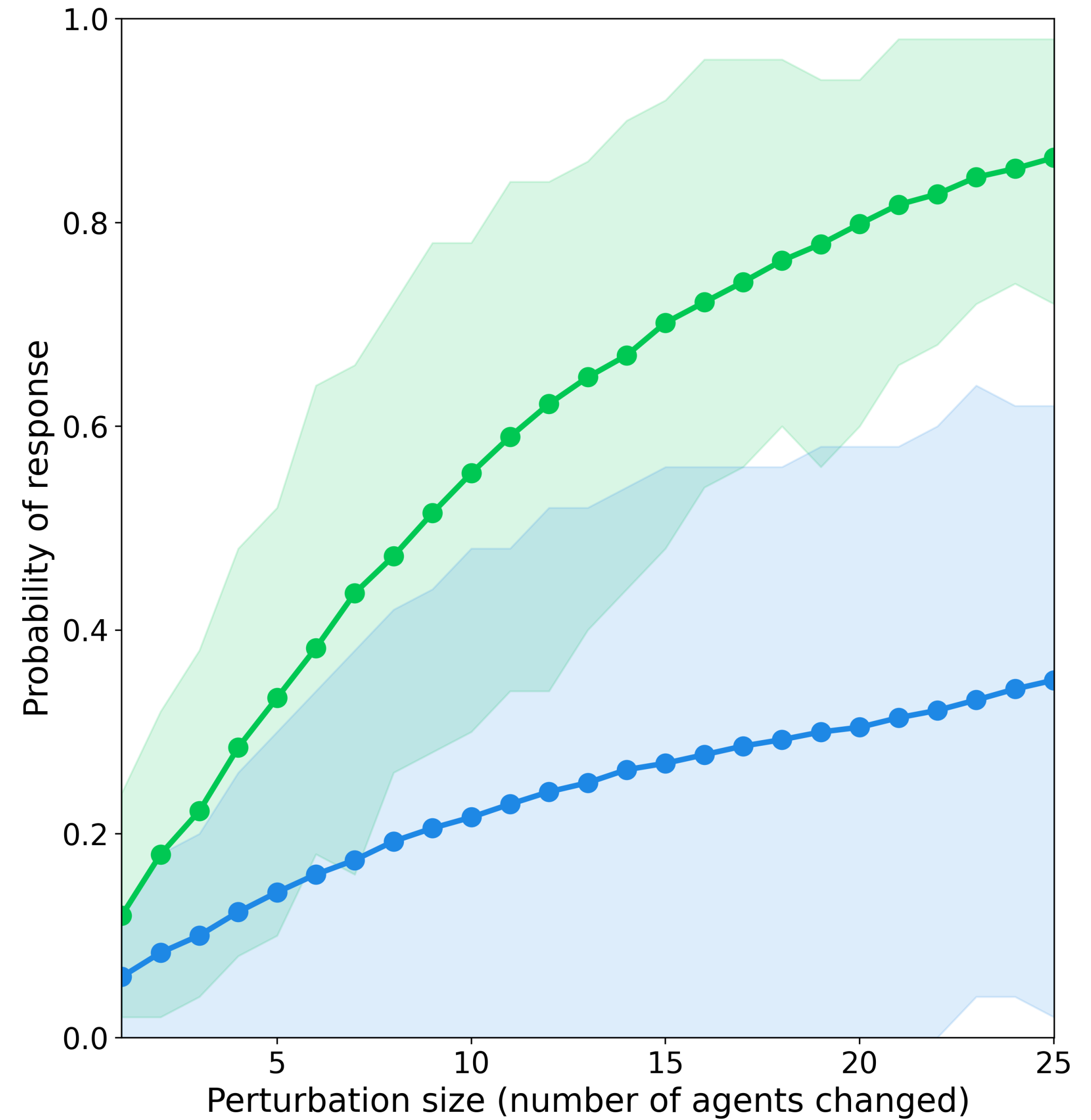


Vary perturbation size  
and measure response

# Total turning magnitude



# Probability of response (large turn)



# Conclusion

## Collective behavior from surprise minimization

Conor Heins<sup>a,b,c,e,1</sup>, Beren Millidge<sup>d</sup>, Lancelot Da Costa<sup>e,f,g</sup>, Richard P. Mann<sup>h</sup>, Karl J. Friston<sup>e,g</sup>, and Iain D. Couzin<sup>a,b,c</sup>

# Thanks to the team behind this work

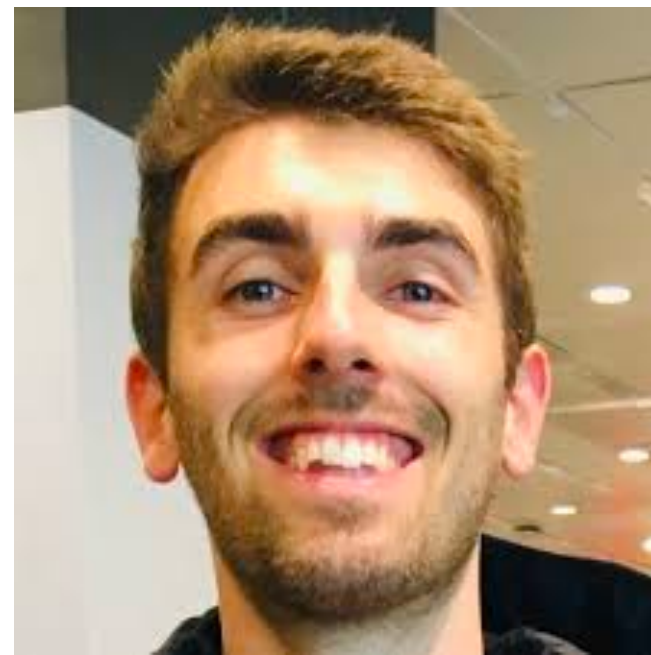
Iain Couzin



Karl Friston



Lance Da Costa



Beren Millidge



Richard Mann



# Bayesian Machine Learning @ VERSES

- Scaling active inference and Bayesian neural networks to modern machine learning contexts
- Variational Bayes for Mixture Models
- Can we train deep nets faster, with less data, while also being Bayesian (i.e., quantifying uncertainty)?
- (Transformer) Attention as Inference



Come chat to me if you want to learn more!

# The goal of model-building

$y$	Data
$\theta$	Parameters for model $m$
$m$	Model

At the end of the day (in my opinion),  
science is about maximizing **model evidence**

$$p(y | m) = \sum_{\theta} p(y | \theta, m) p(\theta | m)$$

Also known as “marginal likelihood”

(log) Model evidence = Accuracy — Complexity

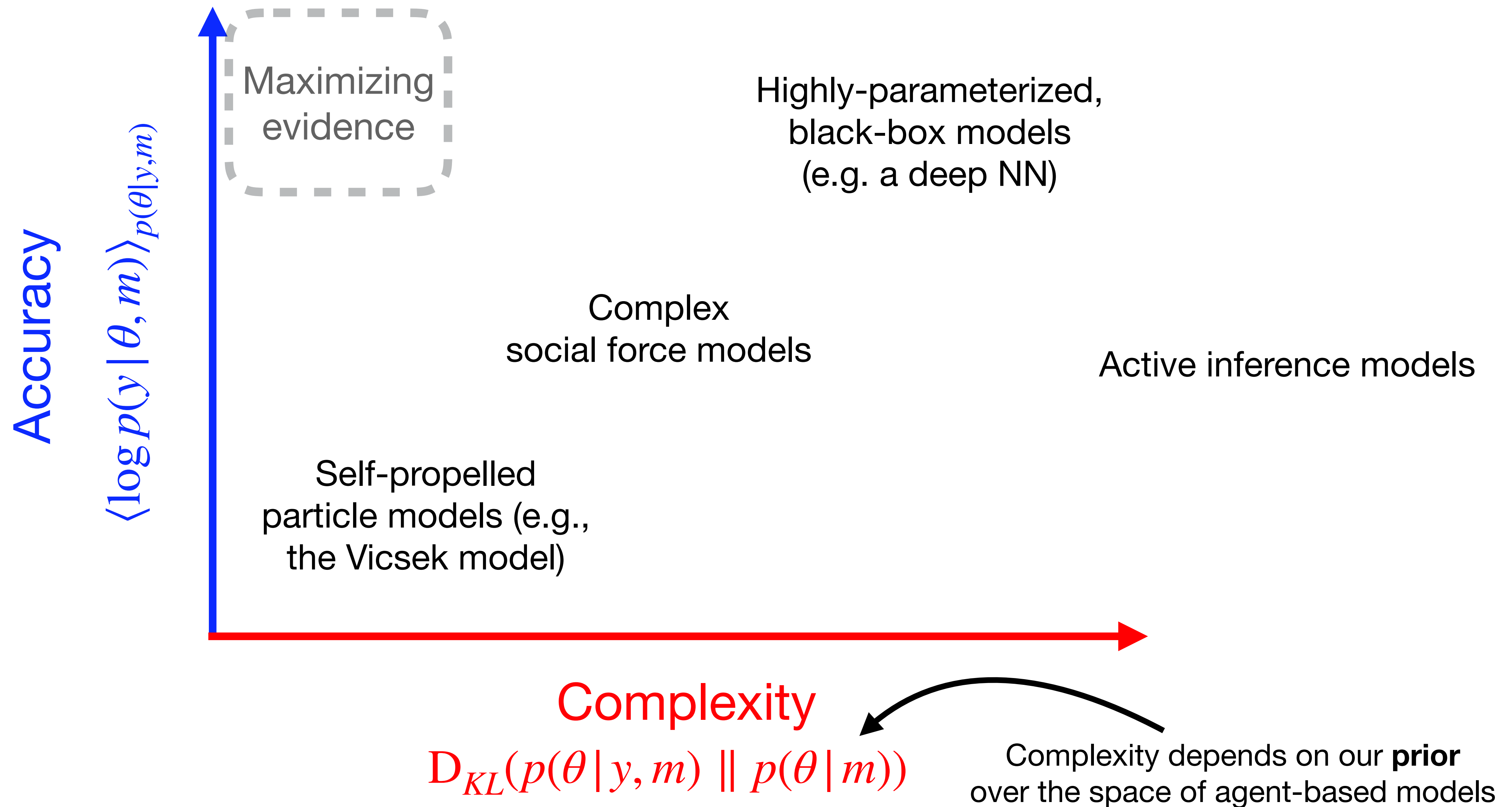
How well I fit the data

How many extra bits my explanation need to encode,  
relative to my “baseline” expectations about explanations

$$\log p(y | m) = \langle \log p(y | \theta, m) \rangle_{p(\theta|y,m)} - D_{KL}(p(\theta | y, m) \parallel p(\theta | m))$$



# The space of agent-based models of collective phenomena



# Variational free energy

$$\mathcal{F} = - \int q(\theta) \ln \frac{p(y, \theta)}{q(\theta)}$$

$$\mathcal{F} = \underbrace{D_{KL} (q(\theta) \parallel p(\theta))}_{\text{Complexity}} - \underbrace{\mathbb{E}_{q(\theta)} [\log p(y | \theta)]}_{\text{Accuracy}}$$