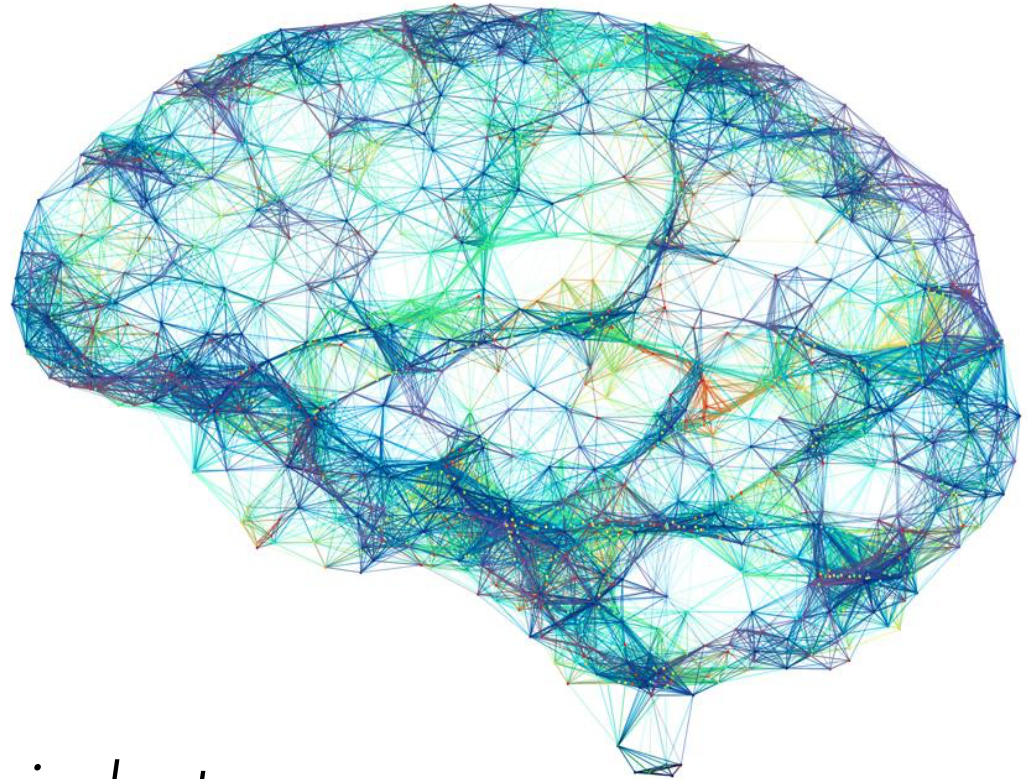


Connecting Variational Autoencoders *Back to the Brain*

arXiv:2011.07464

Joseph Marino

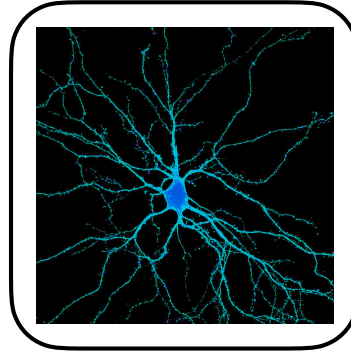
California Institute of Technology



*what is the relationship between
deep learning & neuroscience?*

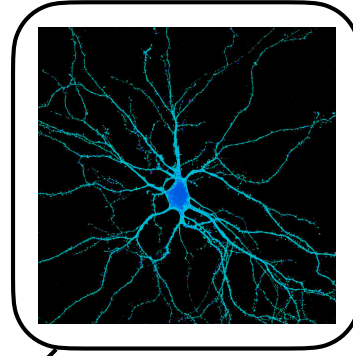
TALK OUTLINE

I. background

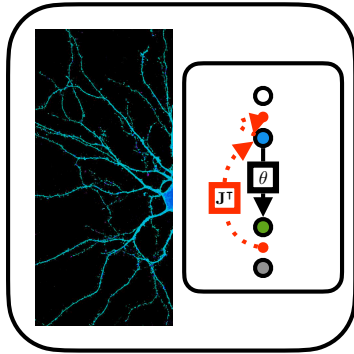


TALK OUTLINE

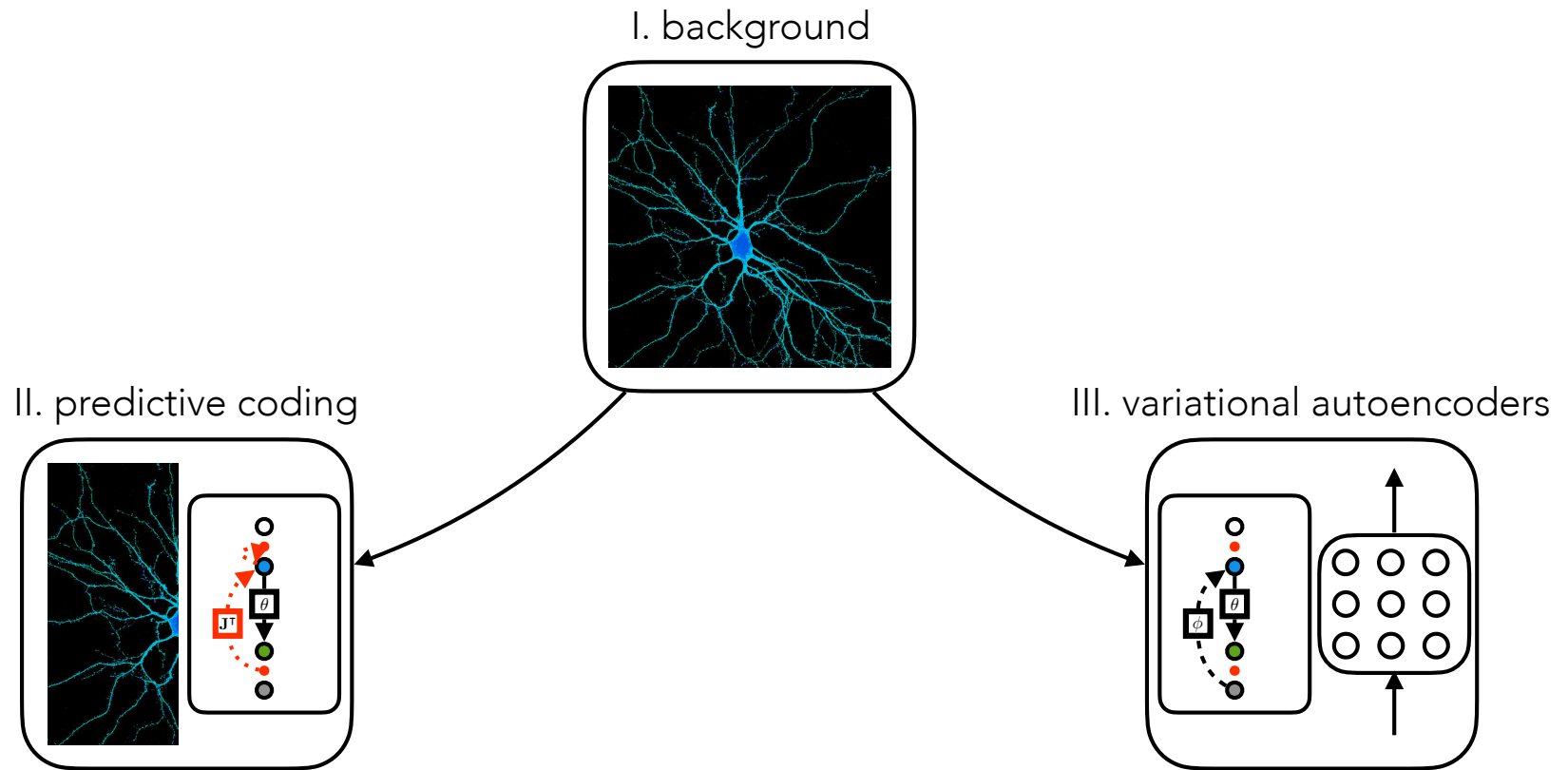
I. background



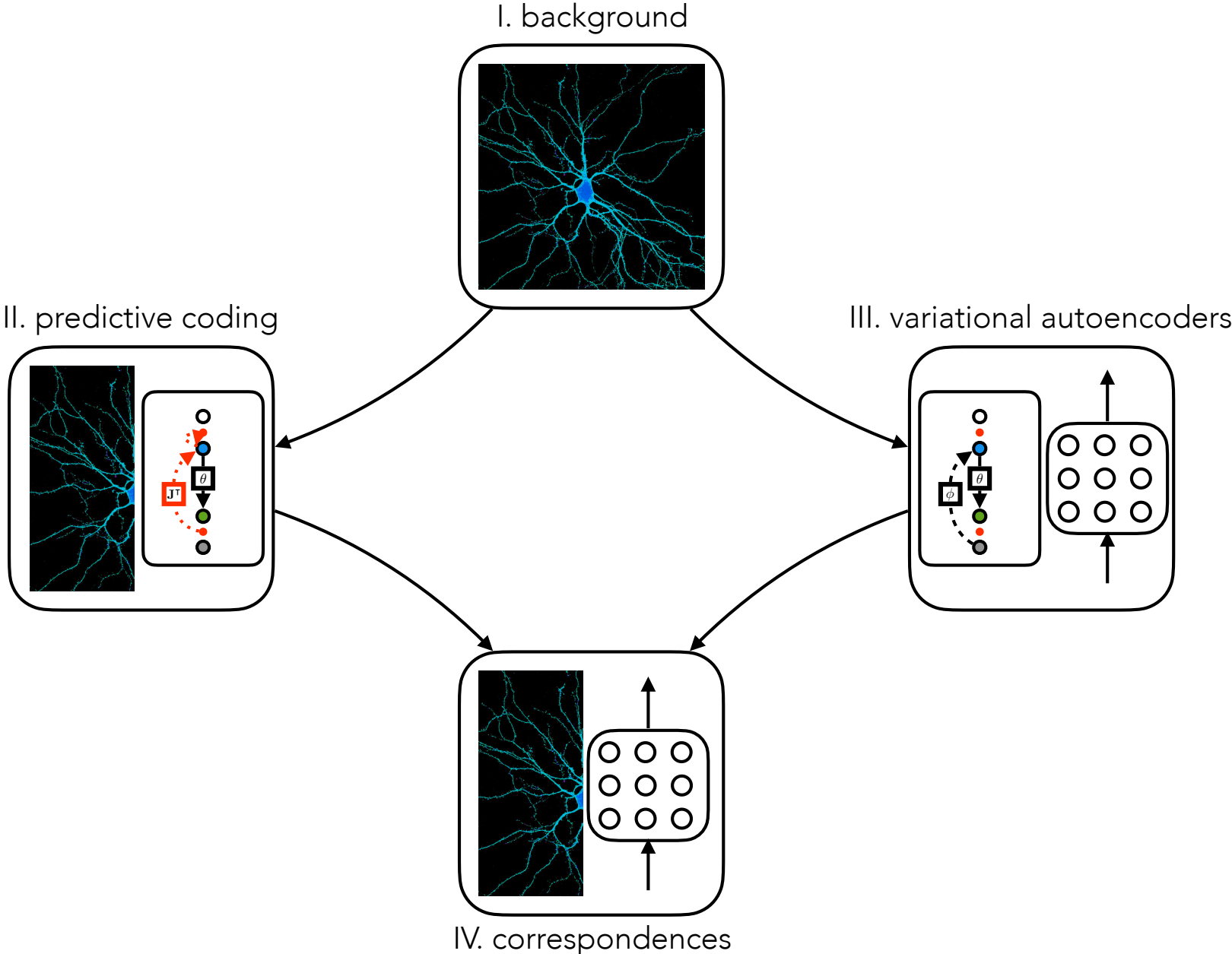
II. predictive coding



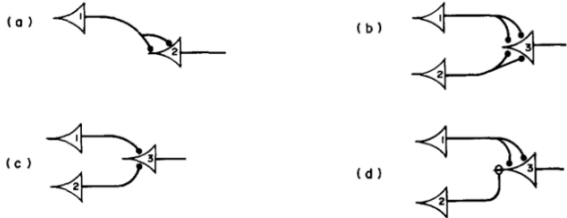
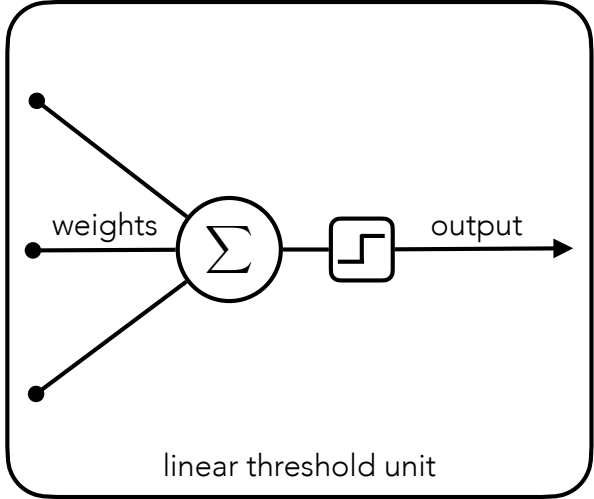
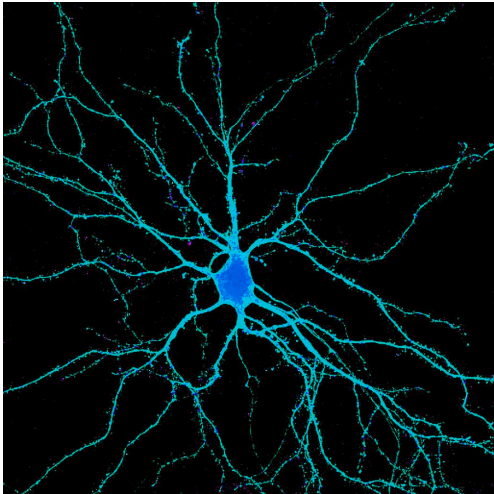
TALK OUTLINE



TALK OUTLINE



ARTIFICIAL NEURAL NETWORKS



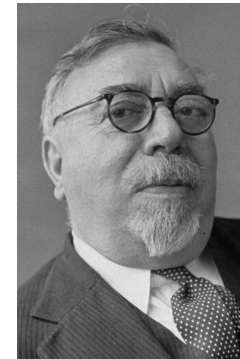
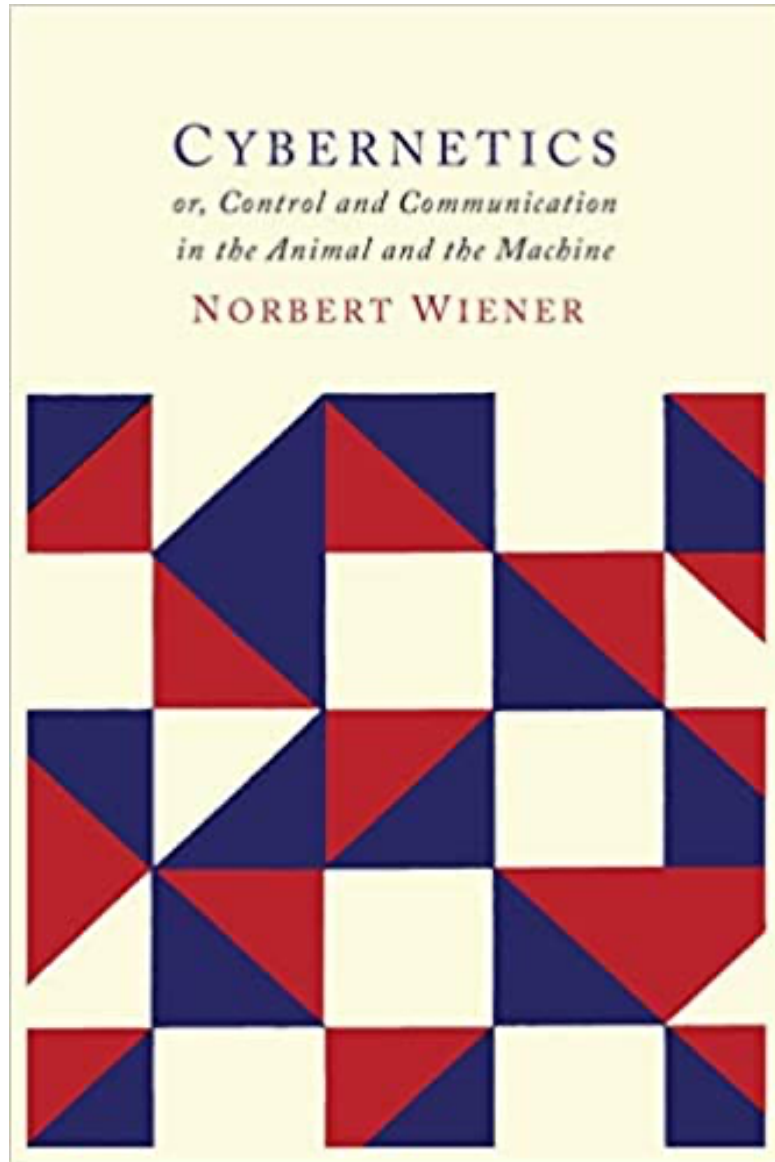
Logical Calculus for Nervous Activity
McCulloch & Pitts, 1943



Walter Pitts Warren McCulloch

cybernetics

CYBERNETICS

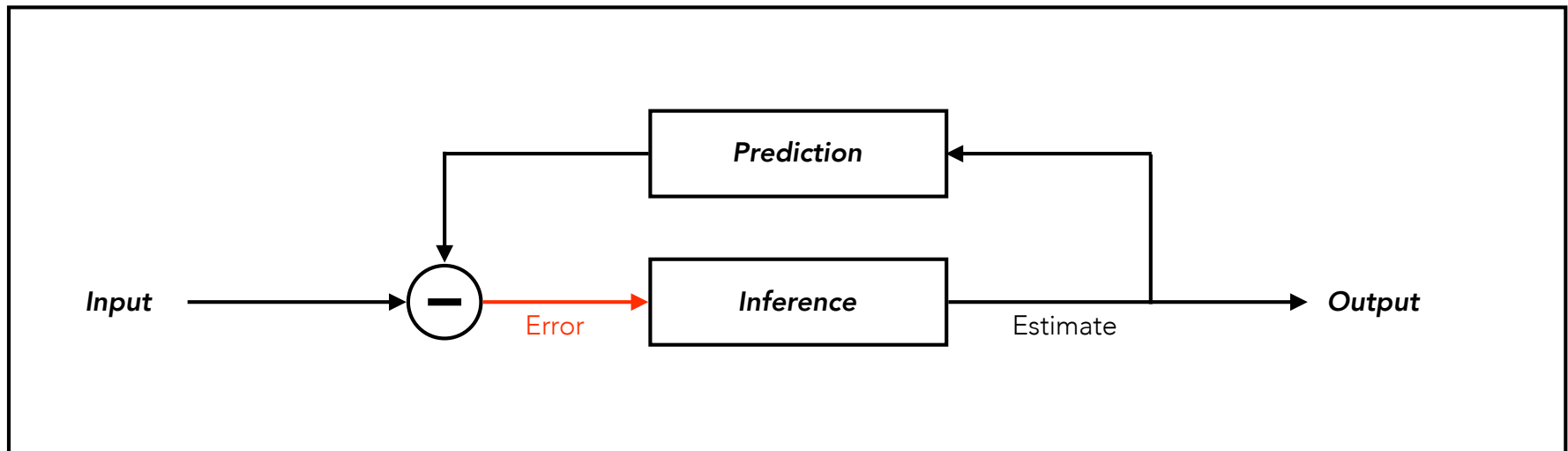


Norbert Wiener

- defined information
- probabilistic models
- variational optimization
- control & perception
- connections to neuroscience

CYBERNETICS

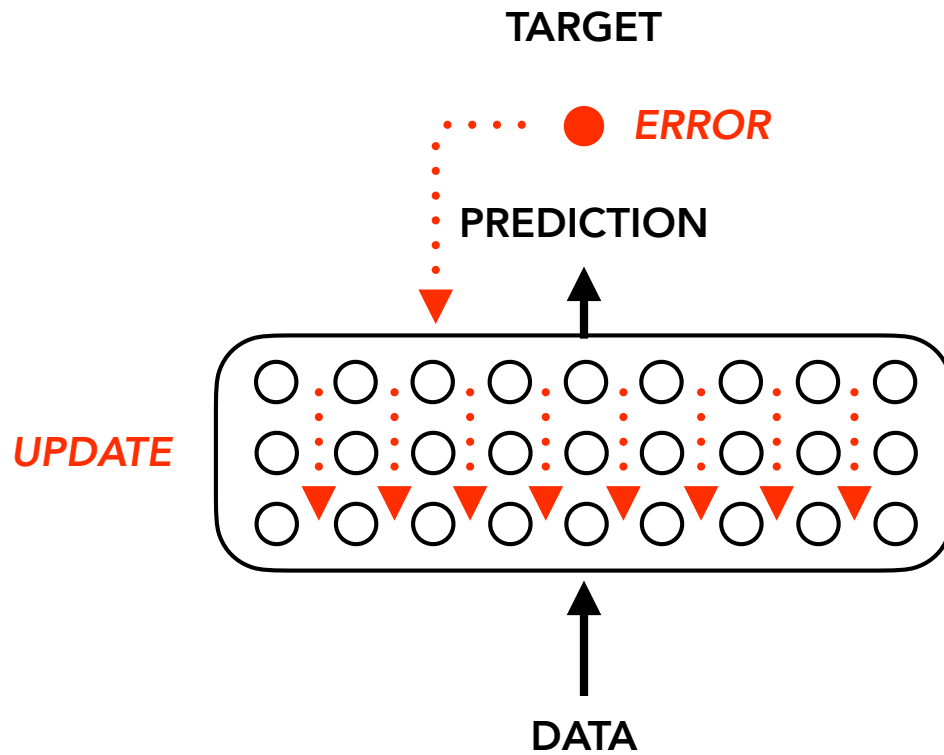
comparator circuit



negative feedback: use errors to correct estimates,
e.g., Kalman filter, linear quadratic regulator, PID

DEEP LEARNING

artificial neural networks + negative feedback



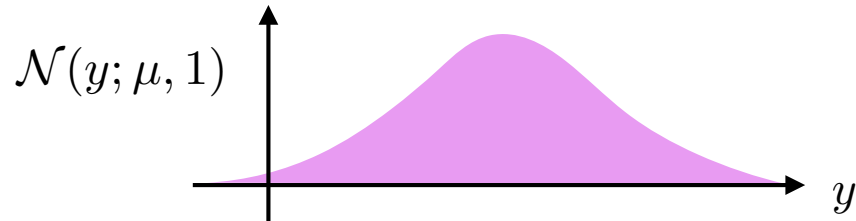
○ artificial neuron

PROBABILISTIC MODELS

why *error*?

gradient of log-probability for exponential family distributions

example: Gaussian density



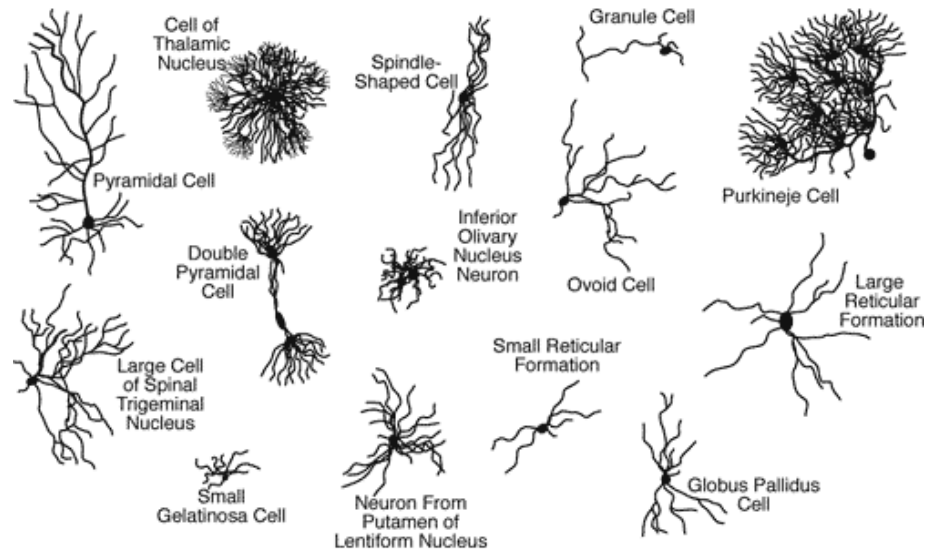
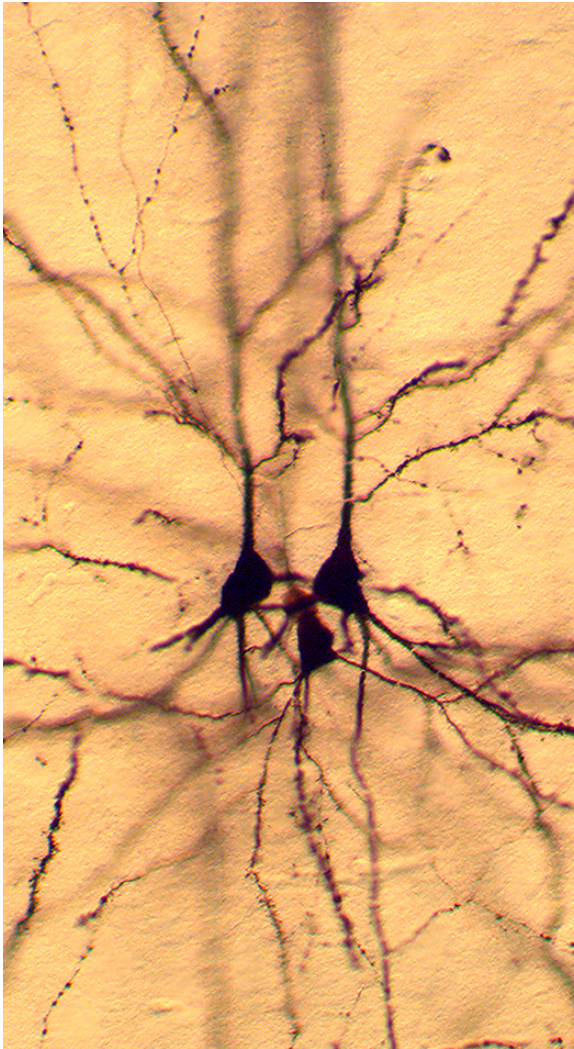
$$\log \mathcal{N}(y; \mu, 1) = \frac{1}{2}(\mu - y)^2 + \text{const.}$$

$$\frac{d}{d\mu} \log \mathcal{N}(y; \mu, 1) = \mu - y$$

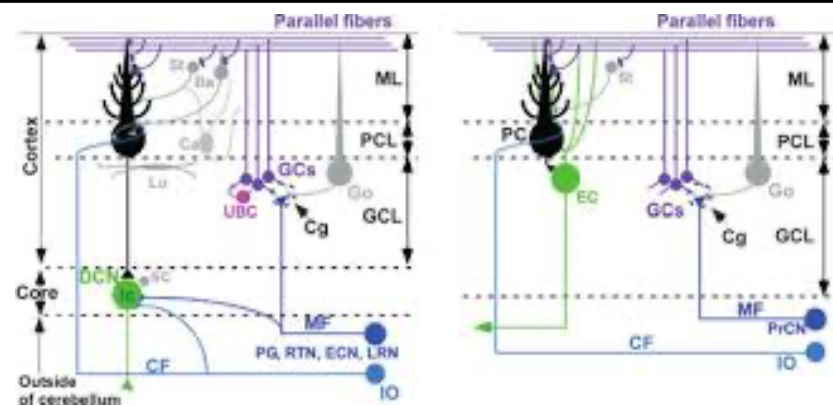
improving predictions requires evaluating and reducing errors

NEUROSCIENCE

neurons are **complex**

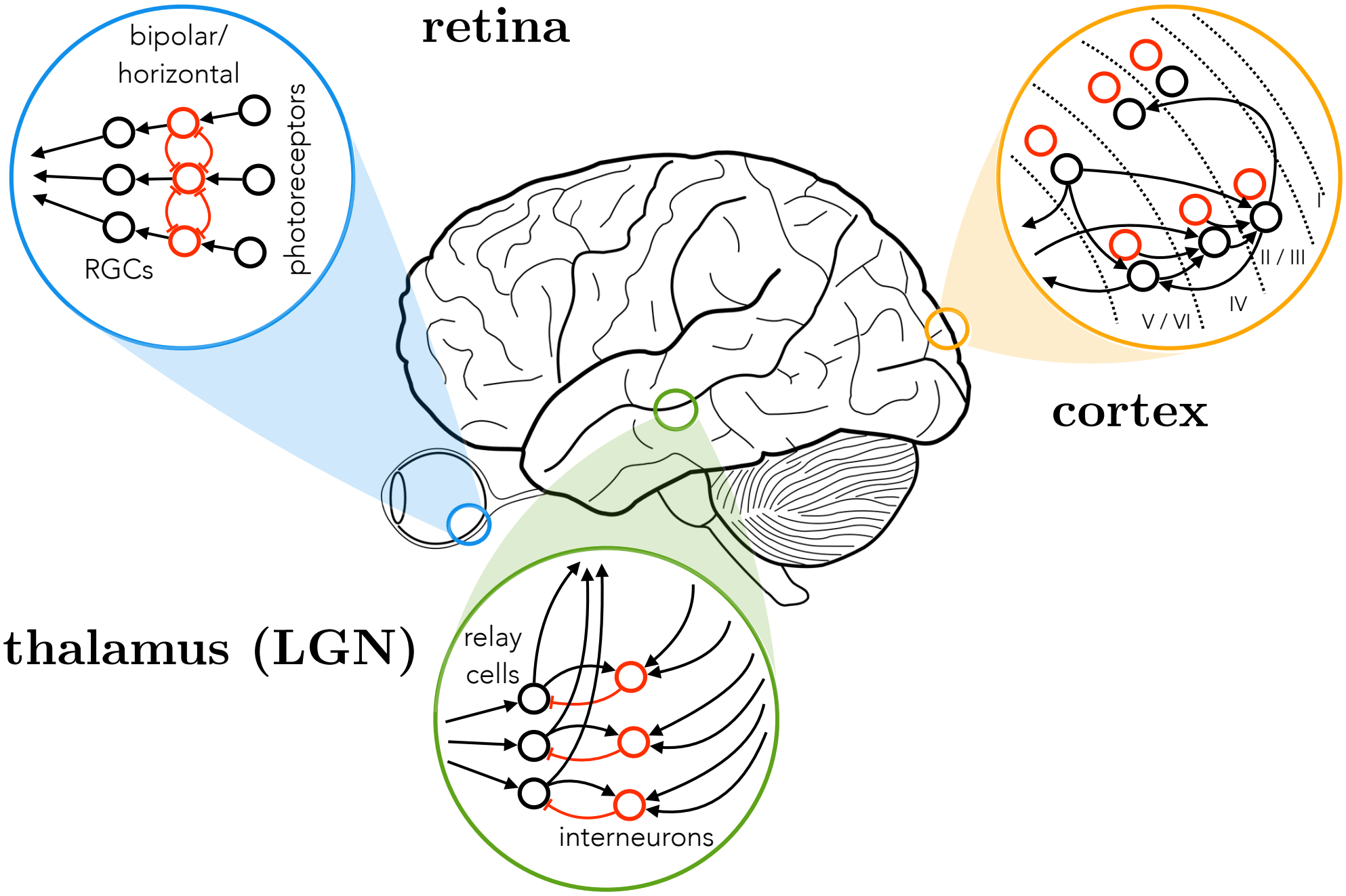


extremely **diverse**



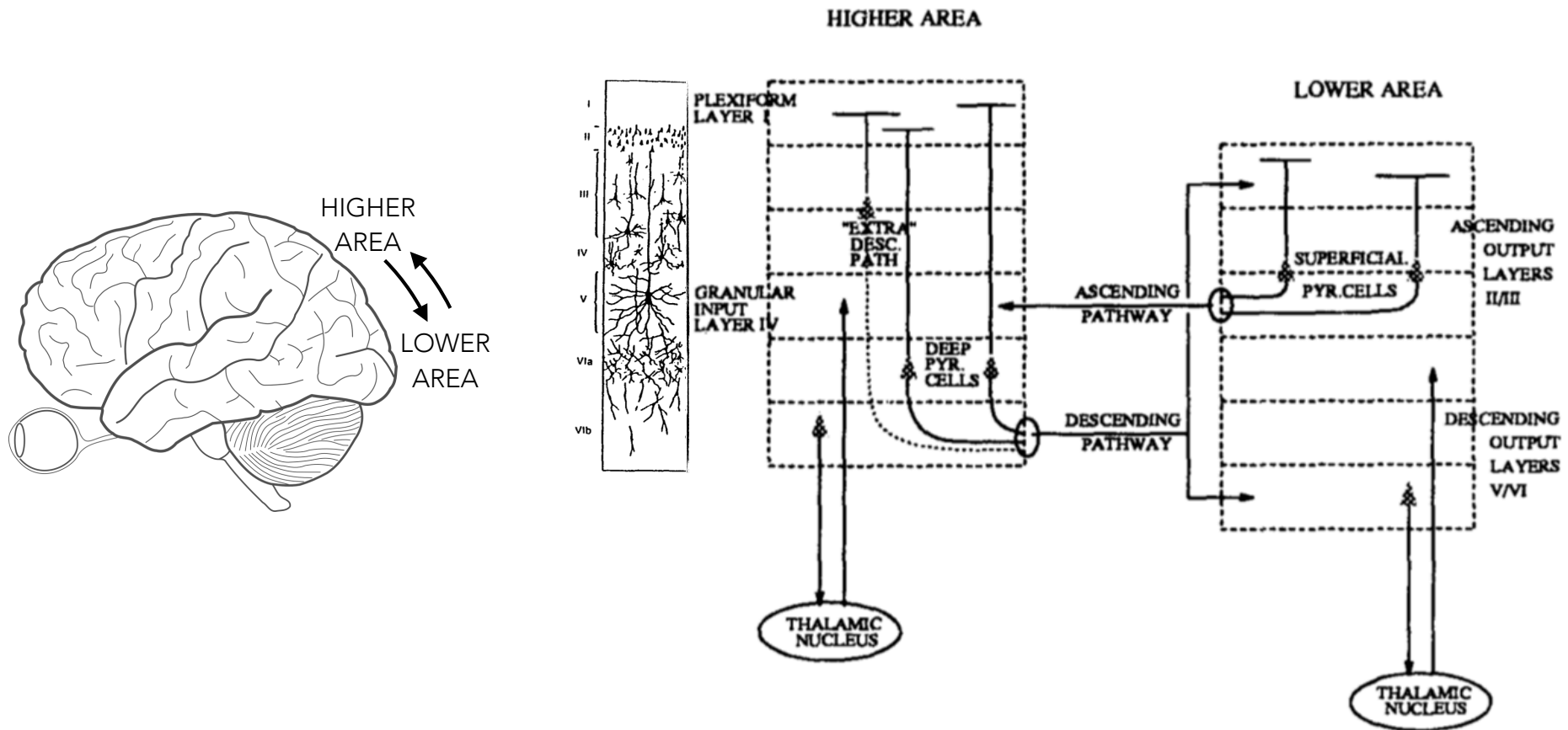
and arranged in intricate **circuits**

VISUAL PATHWAY



II. PREDICTIVE CODING

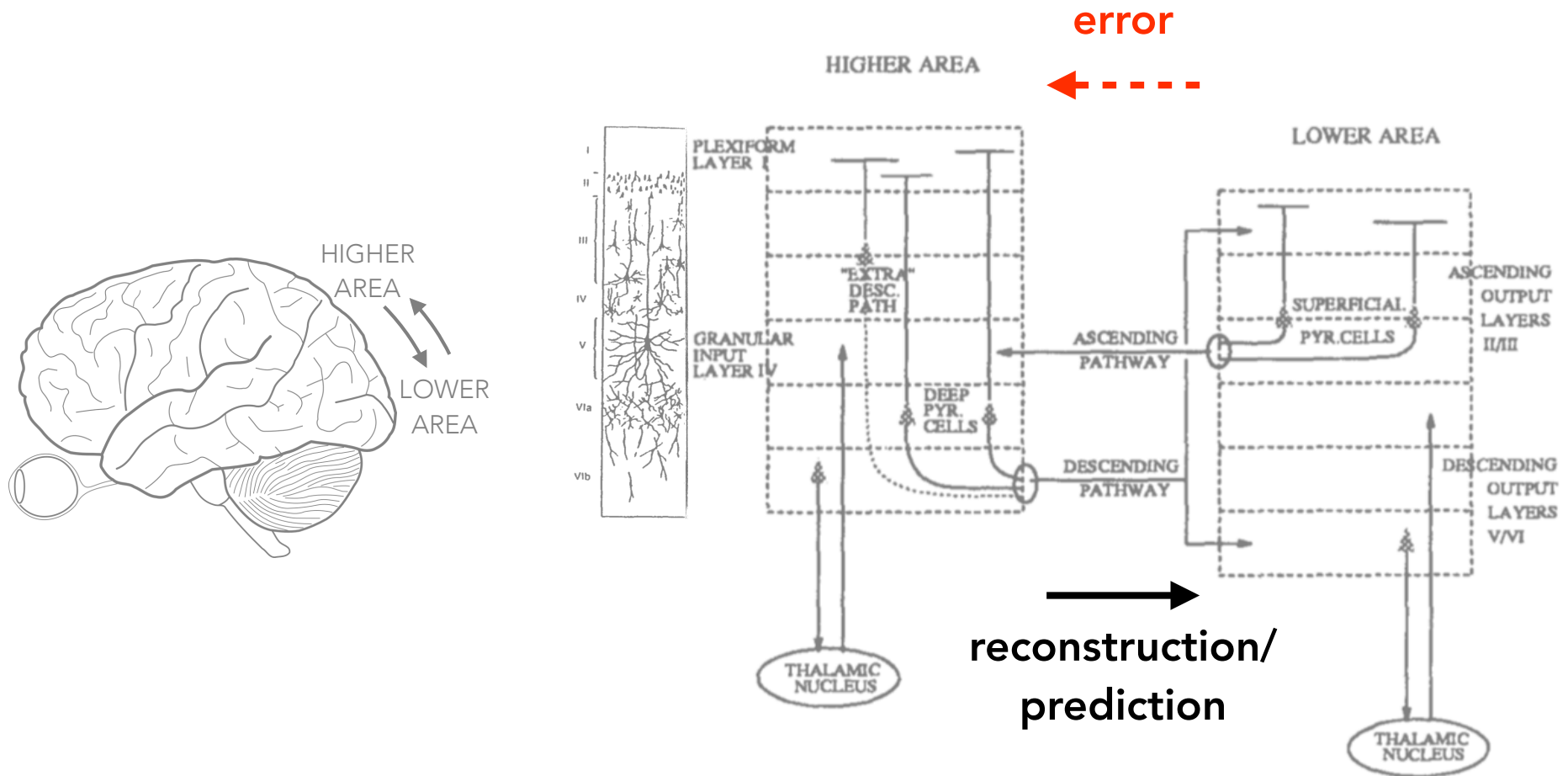
HIERARCHICAL PREDICTIVE CODING



Thalamus "plays the role of an 'active blackboard' on which the current best **reconstruction** of some aspect of the world is always displayed."

Mumford, *Biological Cybernetics* (1991)

HIERARCHICAL PREDICTIVE CODING

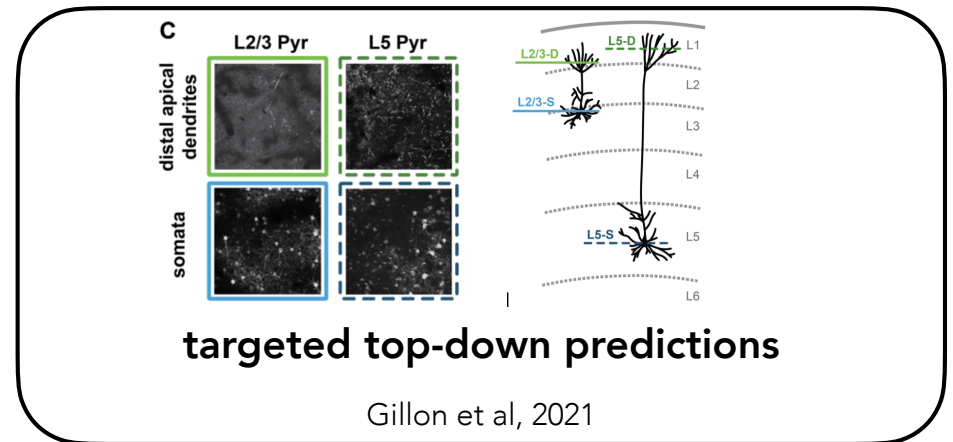
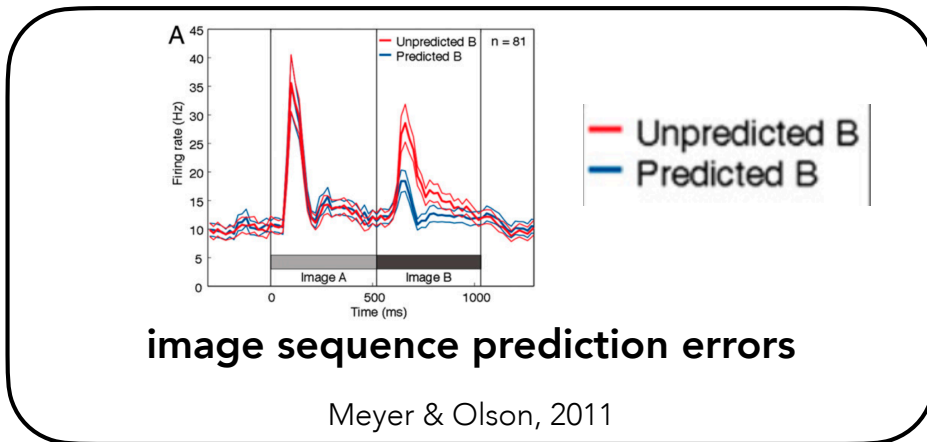
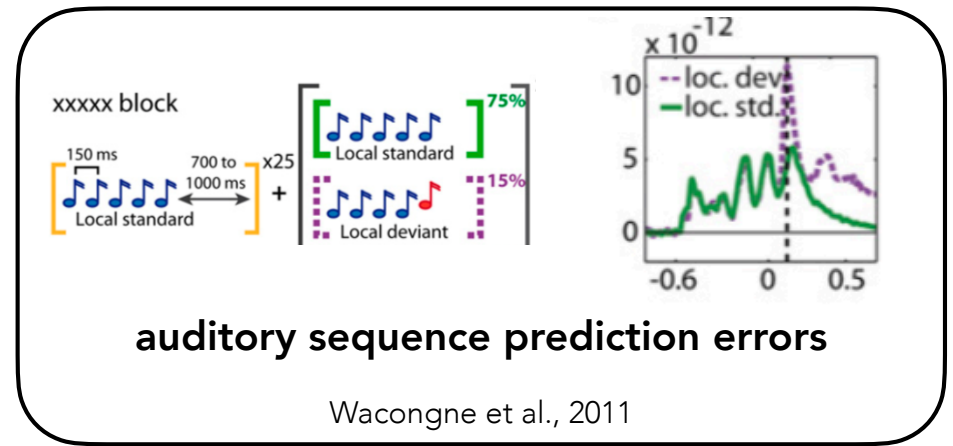
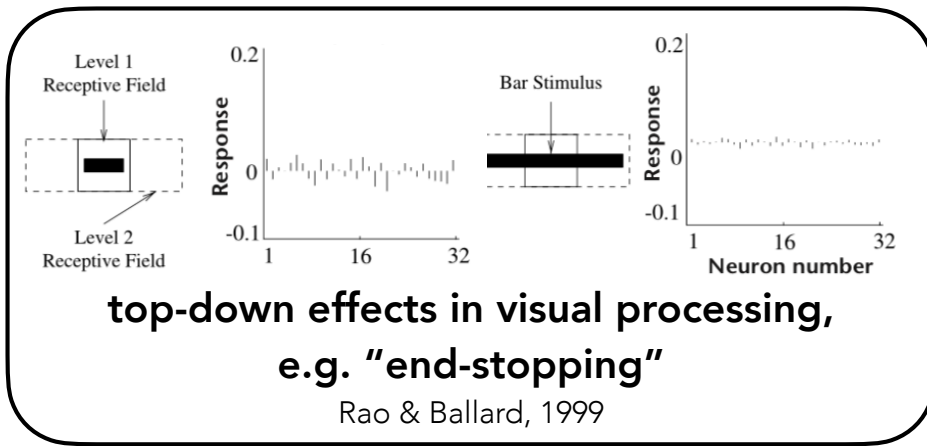


Thalamus "plays the role of an 'active blackboard' on which the current best **reconstruction** of some aspect of the world is always displayed."

Mumford, *Biological Cybernetics* (1991)

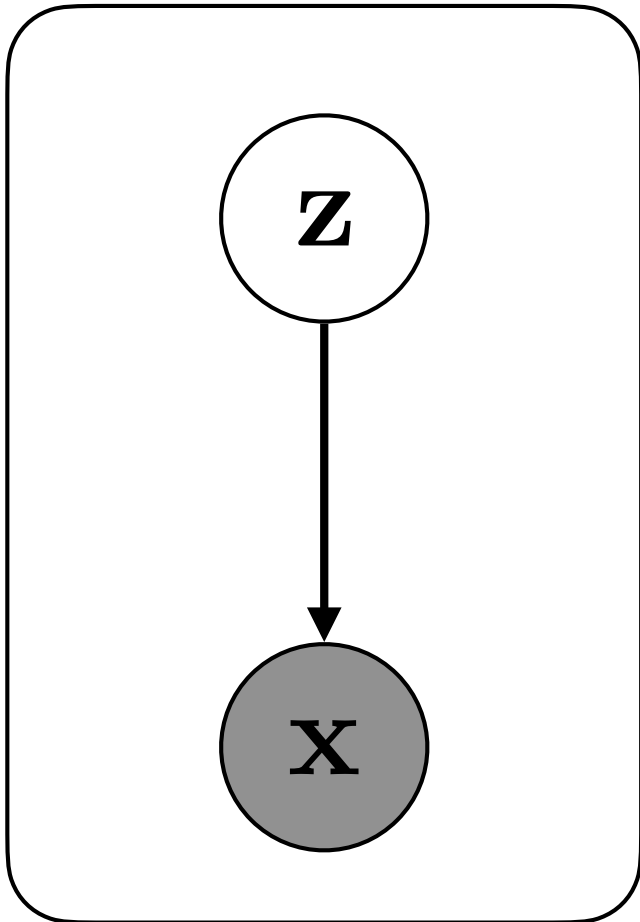
HIERARCHICAL PREDICTIVE CODING

some evidence in support of predictive coding



HIERARCHICAL PREDICTIVE CODING

we can formalize this process as *probabilistic modeling & inference*

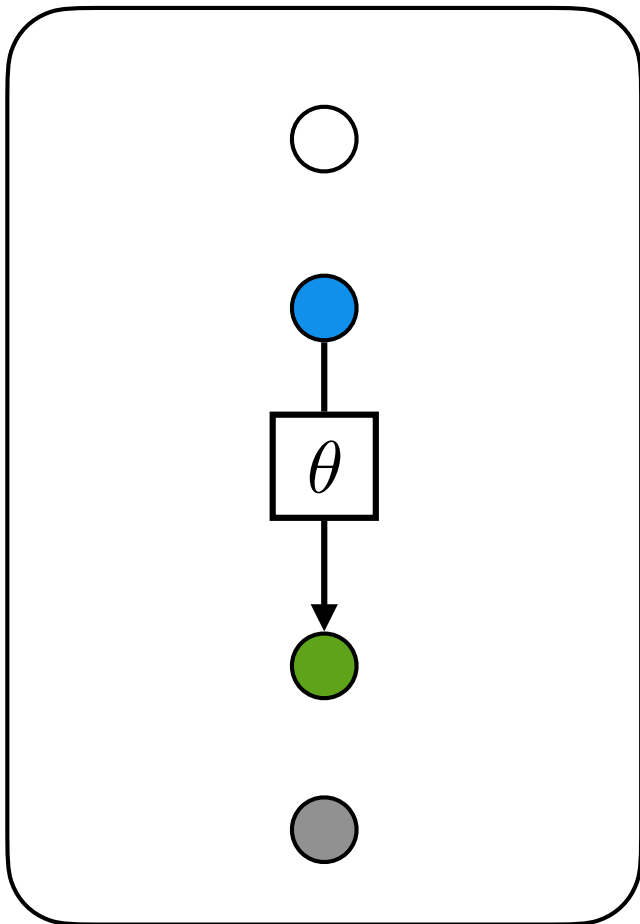


$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})$$

latent variable model

HIERARCHICAL PREDICTIVE CODING

we can formalize this process as *probabilistic modeling & inference*



○ $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z}}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{z}}^2))$

● $q(\mathbf{z}|\mathbf{x}) = \delta(\hat{\mathbf{z}})$



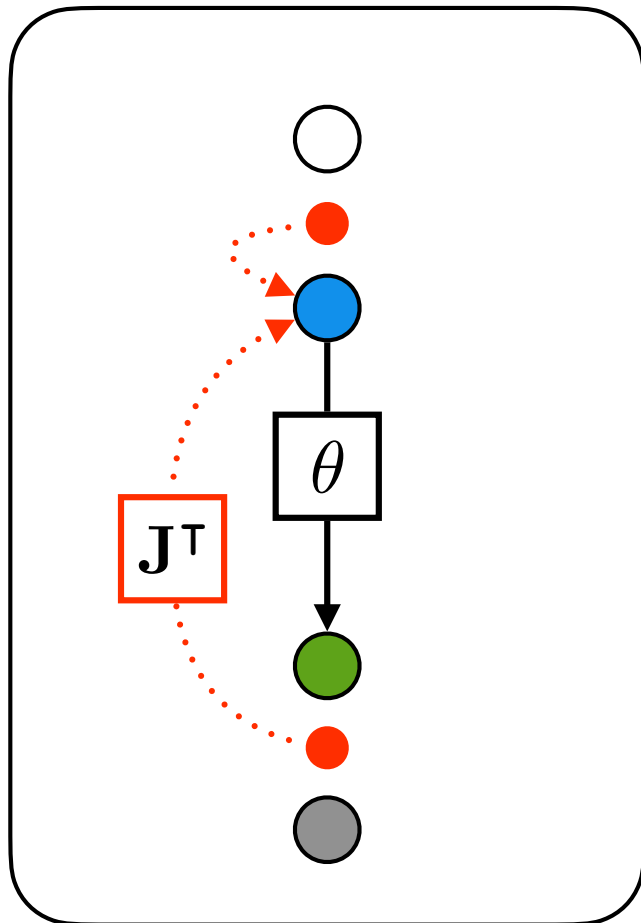
prediction/reconstruction

● $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\mathbf{x}}(\mathbf{z}), \text{diag}(\boldsymbol{\sigma}_{\mathbf{x}}^2))$

● $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$

HIERARCHICAL PREDICTIVE CODING

we can formalize this process as *probabilistic modeling & inference*



\bigcirc $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mu_{\mathbf{z}}, \text{diag}(\sigma_{\mathbf{z}}^2))$

\bullet $q(\mathbf{z}|\mathbf{x}) = \delta(\hat{\mathbf{z}})$



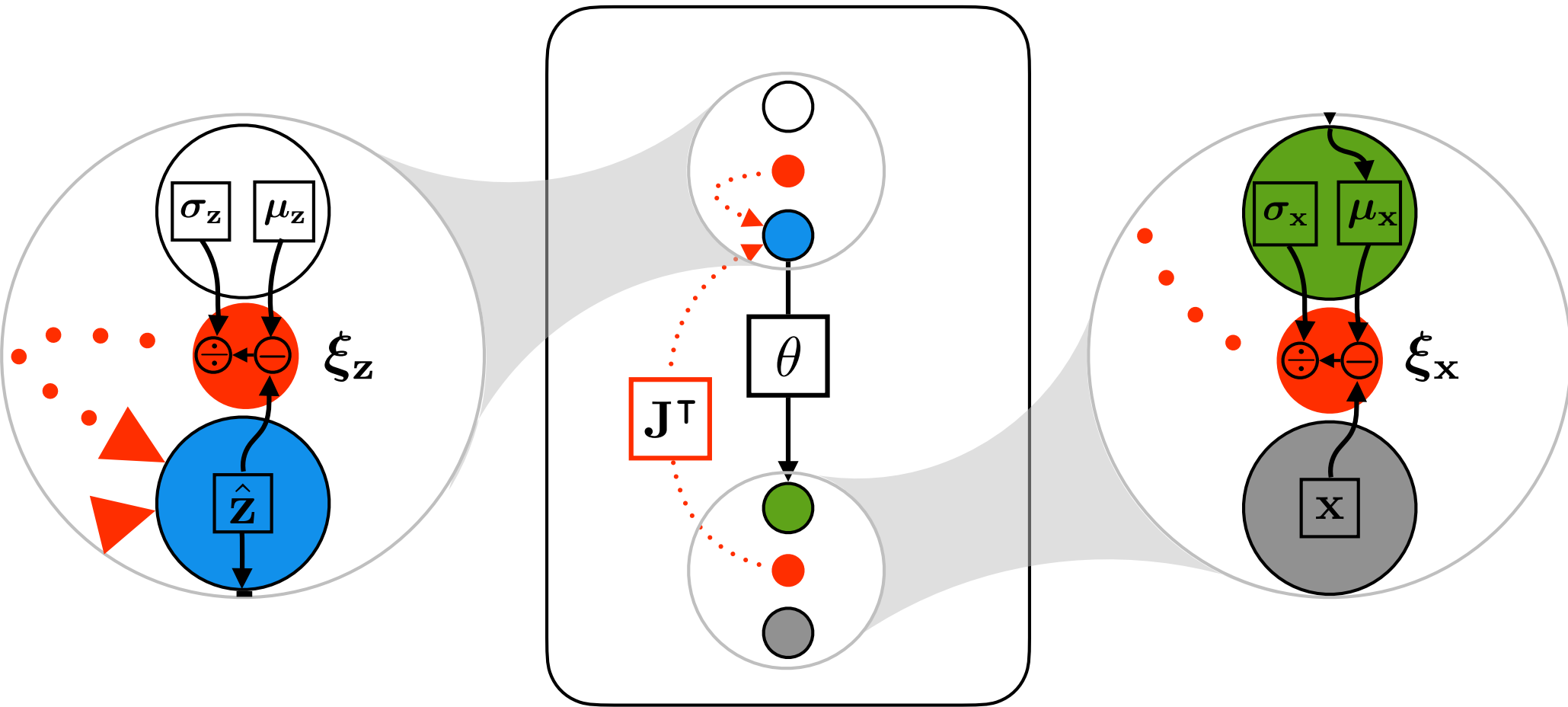
prediction/reconstruction

\bigcirc $p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mu_{\mathbf{x}}(\mathbf{z}), \text{diag}(\sigma_{\mathbf{x}}^2))$

\bigcirc $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$

inference: maximize a (variational) objective \mathcal{L} w.r.t. \bullet

HIERARCHICAL PREDICTIVE CODING



inference involves a sum of (prediction) errors

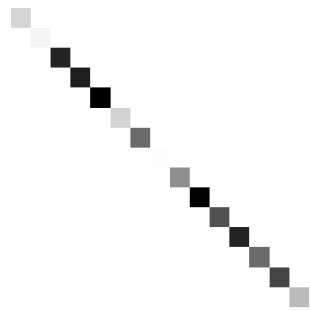
$$\nabla_{\hat{z}} \mathcal{L} = \mathbf{J}^T \xi_x - \xi_z$$

SPATIOTEMPORAL PREDICTIVE CODING

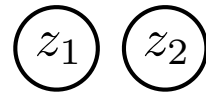
lateral inhibition implements a covariance matrix/normalization

independent dimensions: $\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$

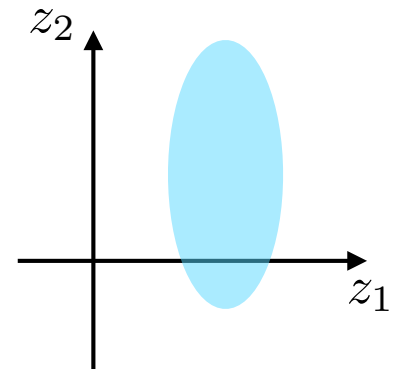
$\text{diag}(\boldsymbol{\sigma}^2) =$



in 2D:

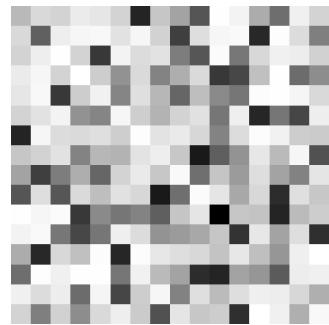


*independent
sampling*

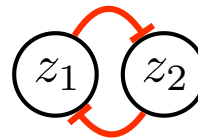


linearly-dependent dimensions: $\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$

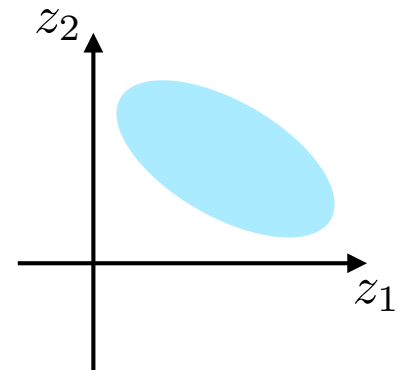
$\boldsymbol{\Sigma} =$



in 2D:

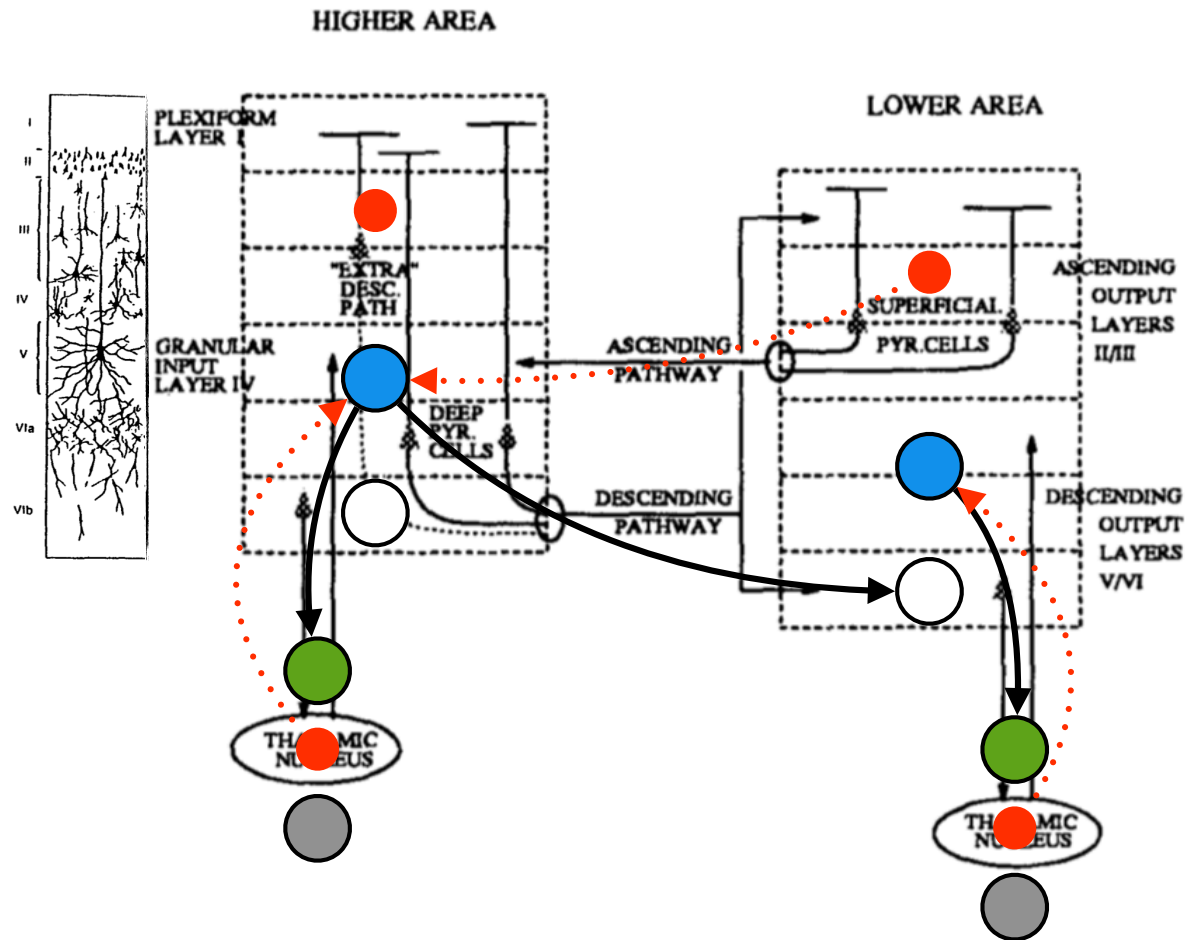
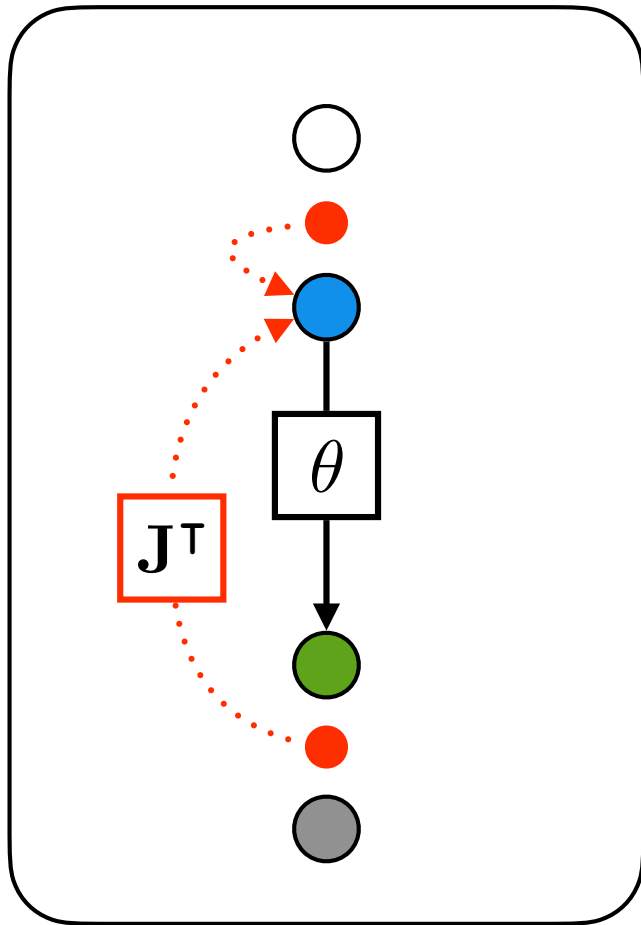


*linearly-dependent
sampling*



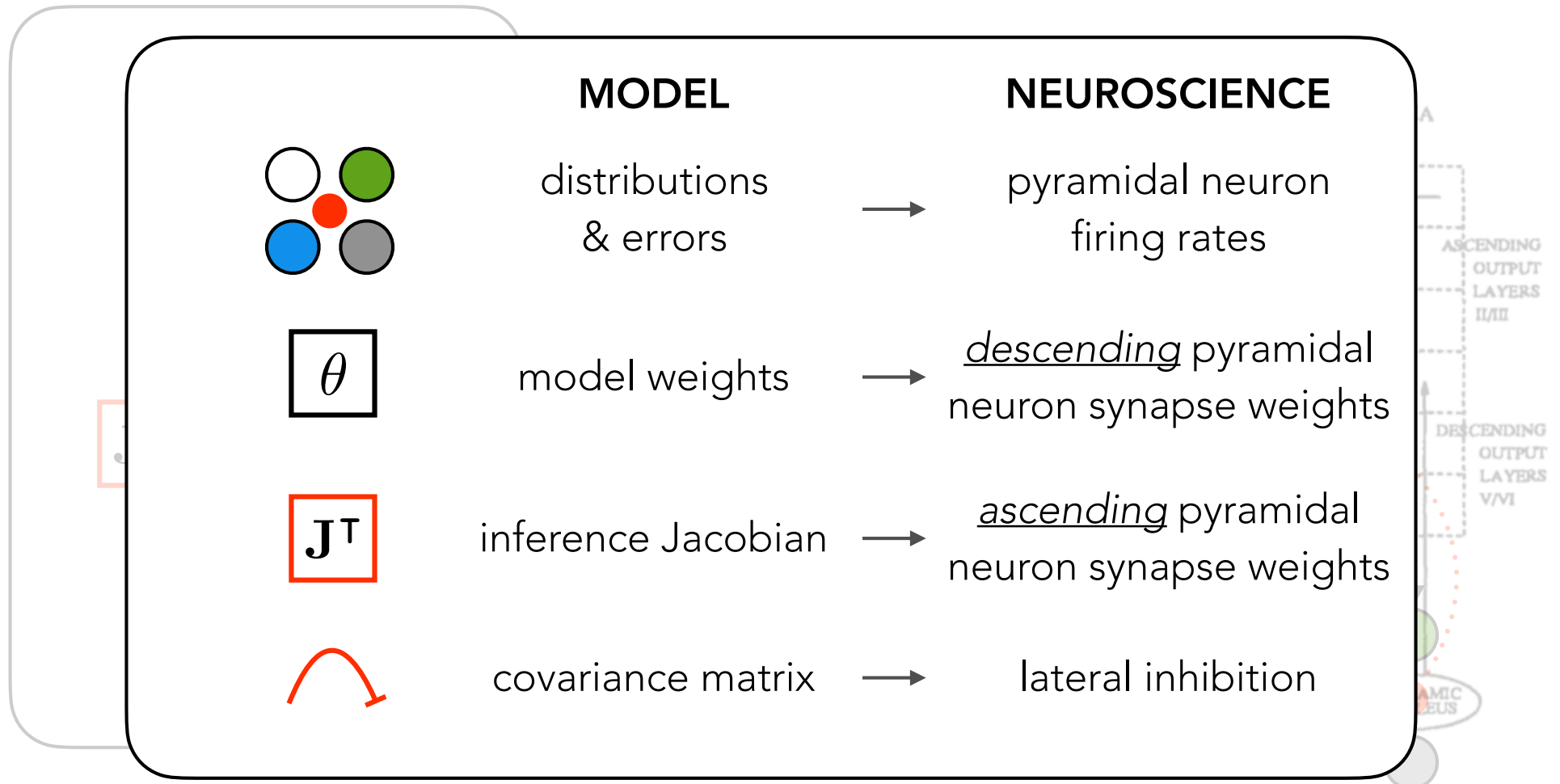
PREDICTIVE CODING

(proposed) biological correspondences



PREDICTIVE CODING

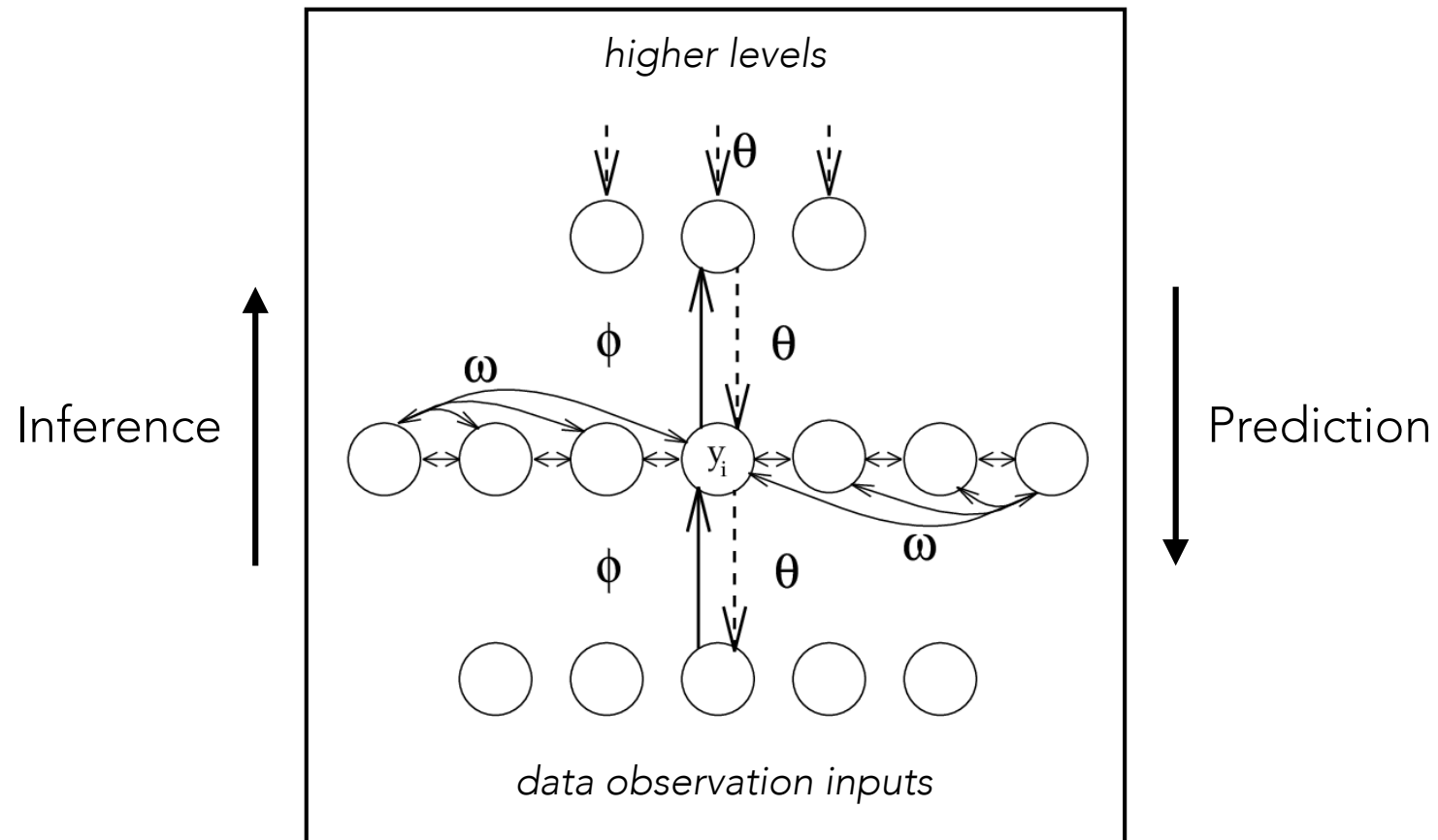
(proposed) biological correspondences



III. VARIATIONAL AUTOENCODERS

DEEP LATENT VARIABLE MODELS

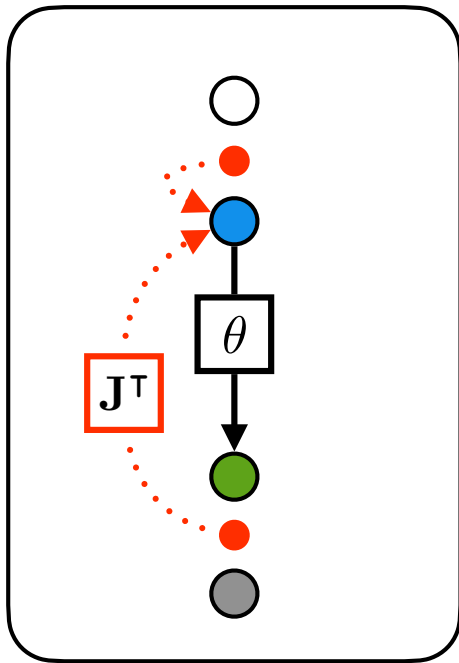
The Helmholtz Machine: *learn a separate inference model*



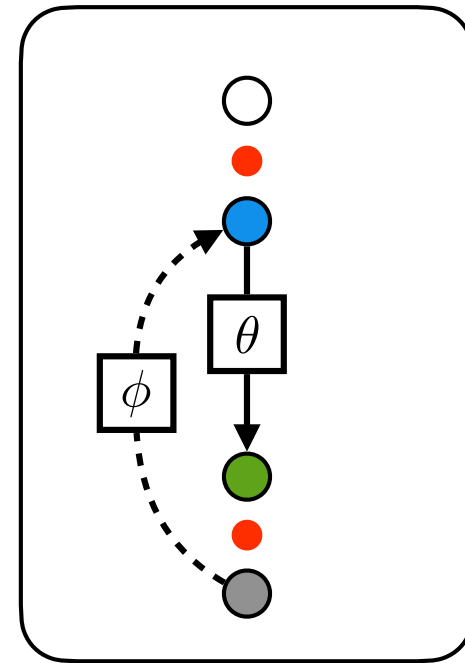
AMORTIZED INFERENCE

amortized inference:

spread out inference costs by learning a separate model (learning to infer)



gradient-based inference



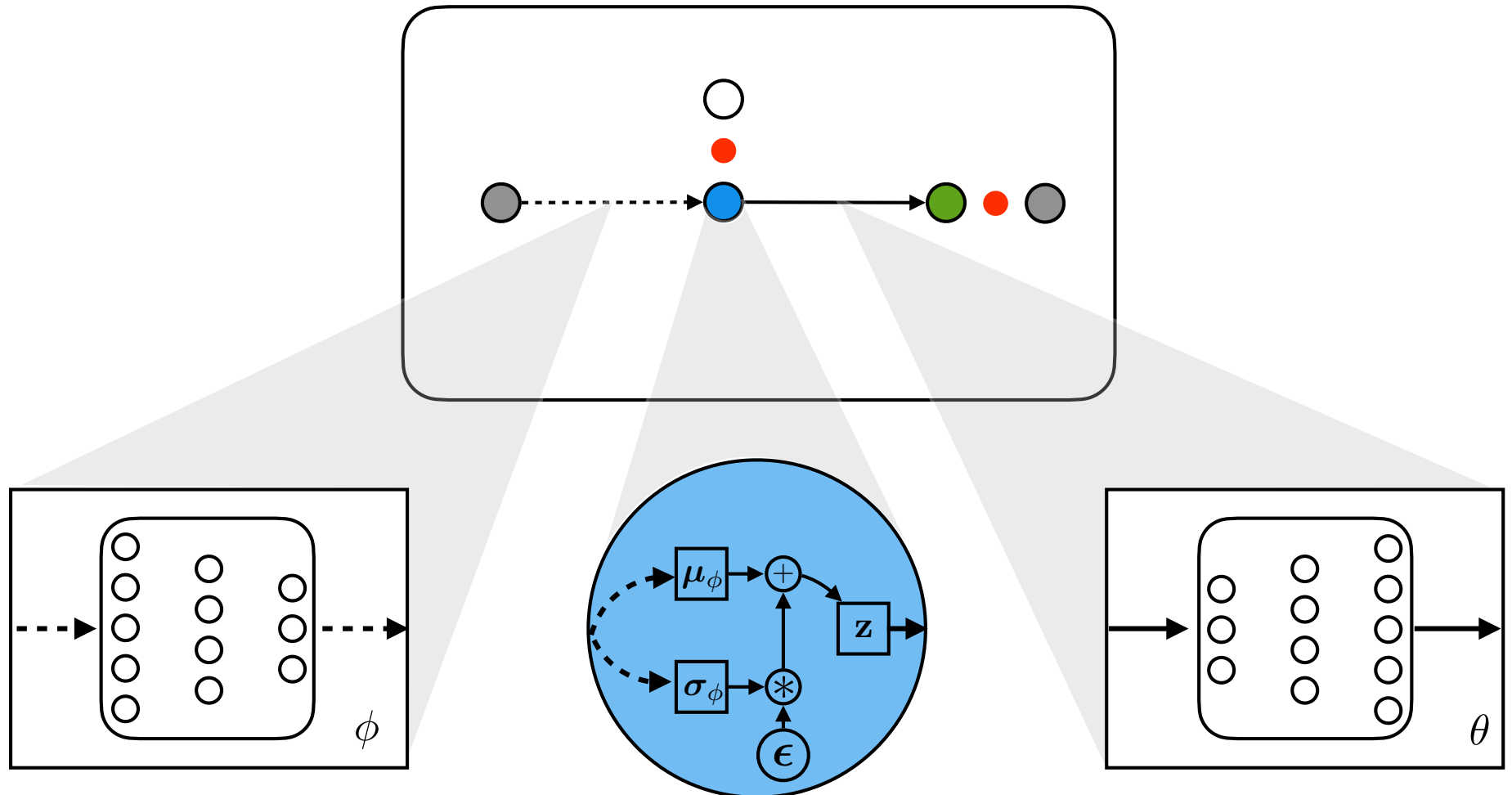
amortized inference

substantially more efficient!

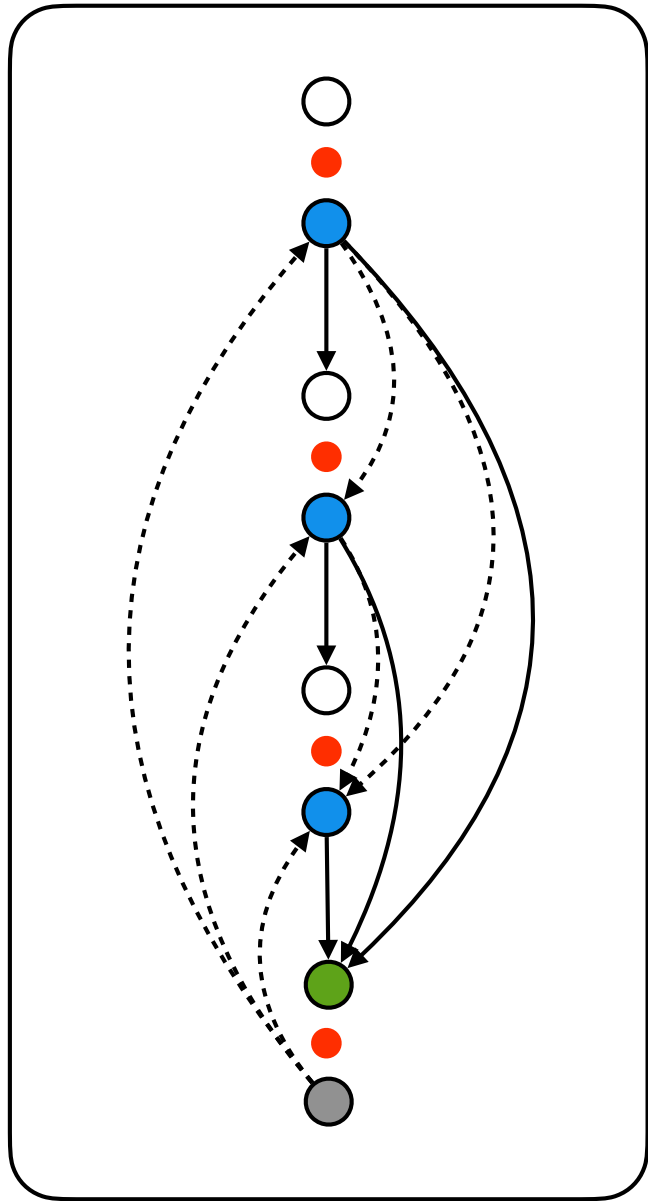
VARIATIONAL AUTOENCODERS

Variational Autoencoder (VAE):

deep latent variable model + variational inference + direct encoder + *reparameterized* Gaussian



HIERARCHICAL VAES



256 x 256 samples

Child, 2020

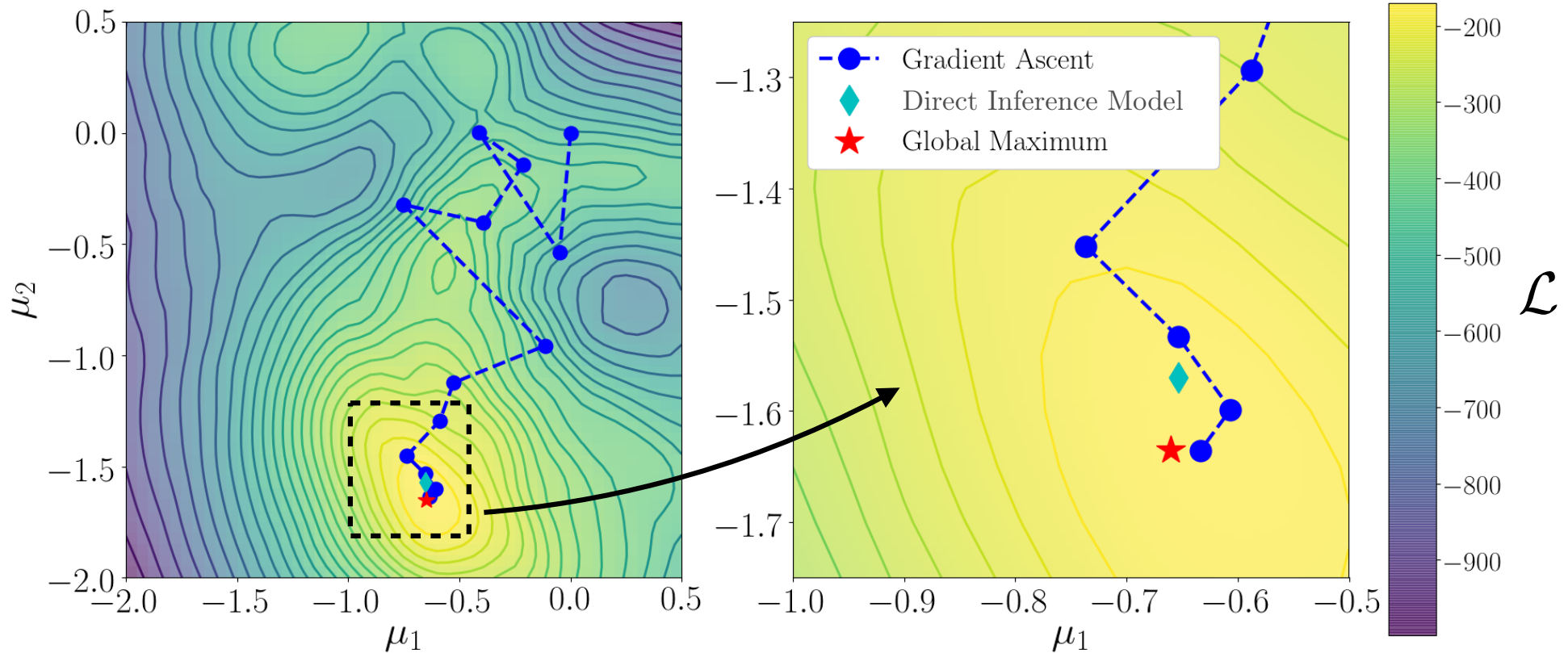


up the hierarchy
increasingly abstract

INFERENCE SUBOPTIMALITY

direct inference models provide suboptimal estimates

"amortization gap"



see also Cremer et al., 2018

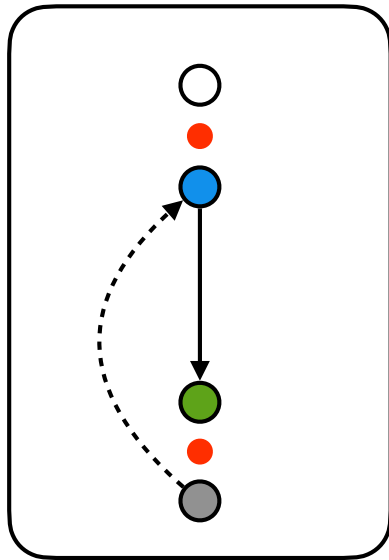
ITERATIVE AMORTIZED INFERENCE

perform inference via a learned *iterative* procedure

encode inference gradient or errors

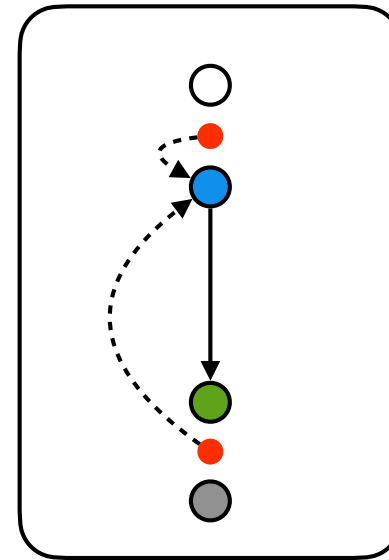
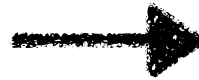
$$\nabla_{\hat{z}} \mathcal{L}$$

$$\xi_z \quad \xi_x$$



Direct

Amortization



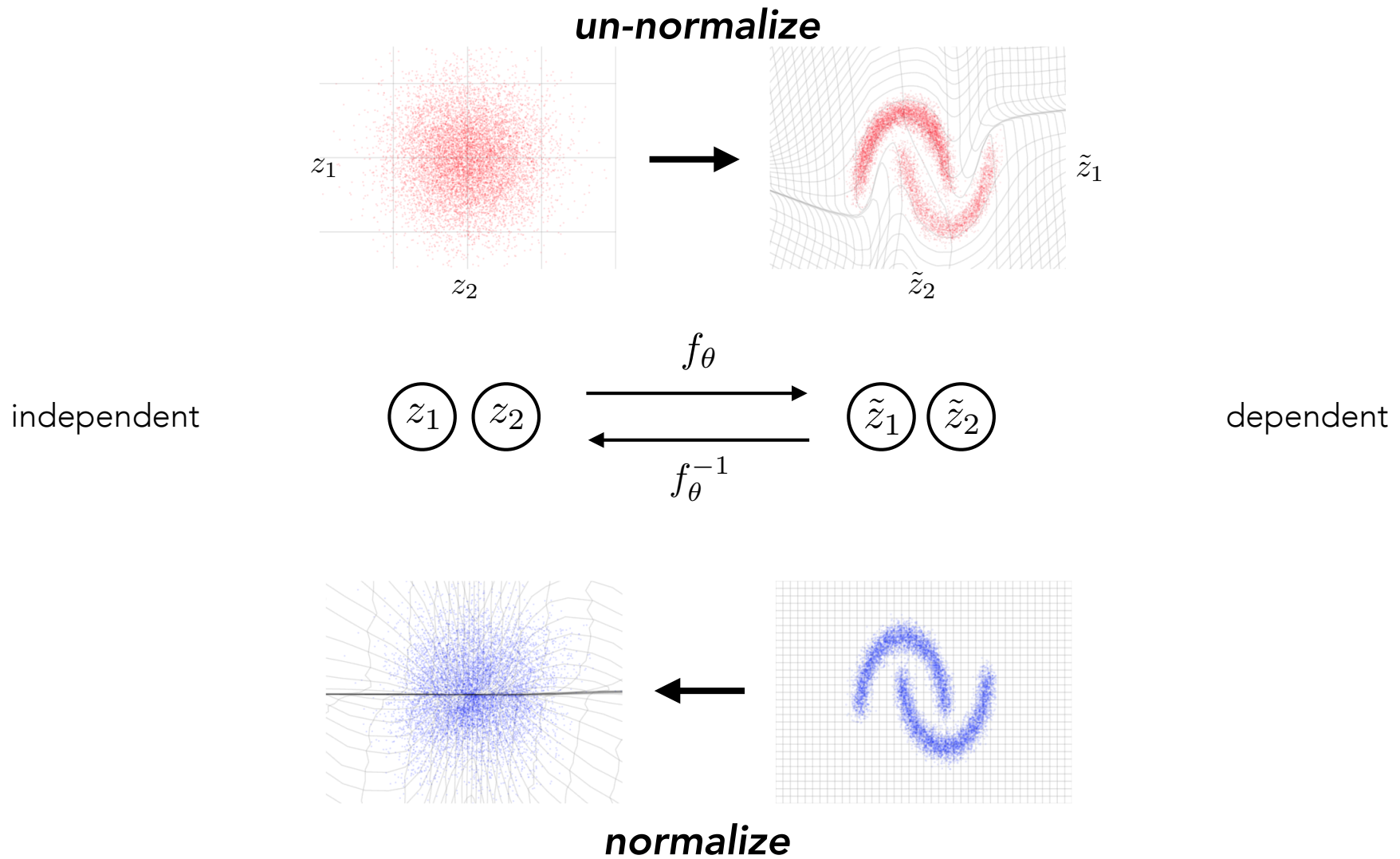
Iterative

Amortization

- + more accurate
- + more general

NORMALIZING FLOWS

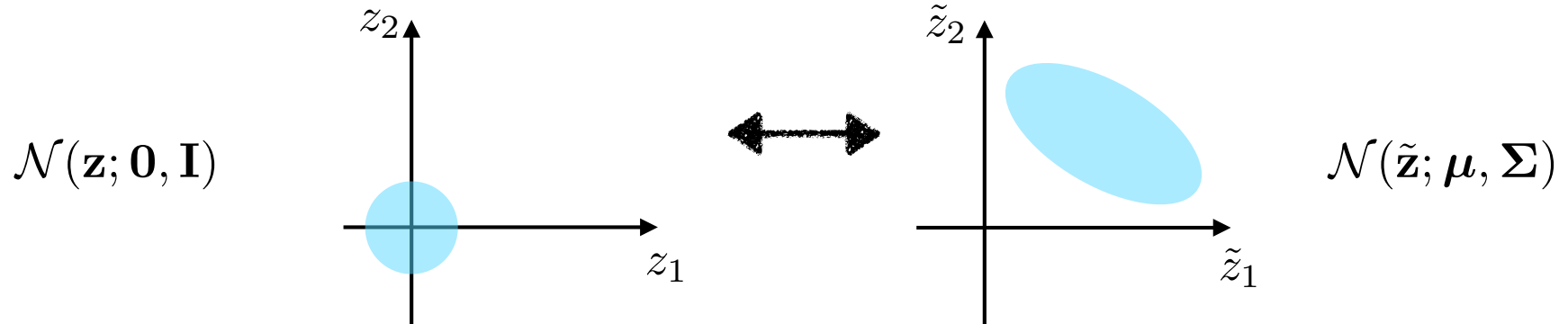
normalizing flows add/remove dependencies



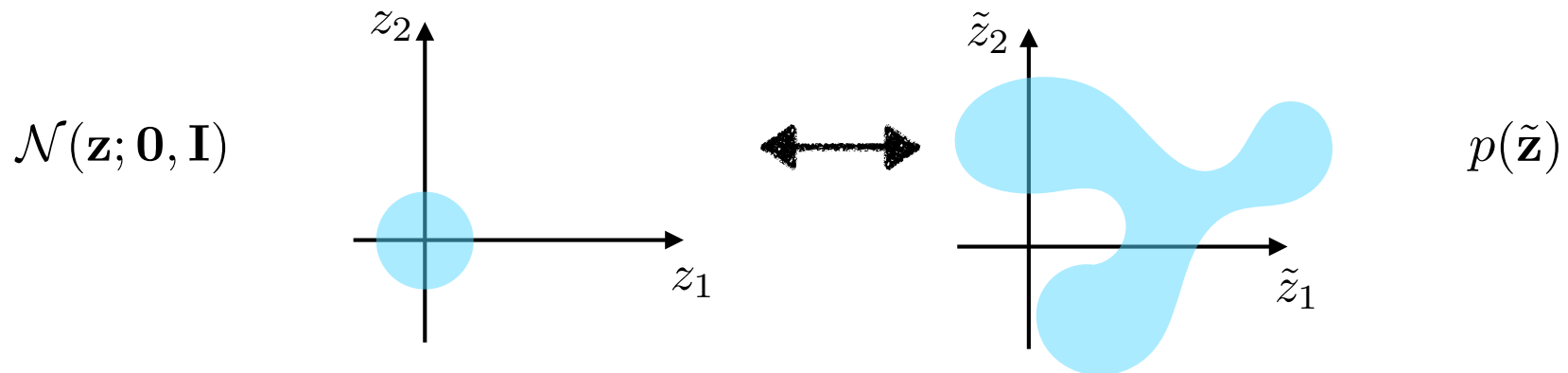
NORMALIZING FLOWS

a covariance matrix is an affine normalizing flow with *linear* dependencies

LINEAR (covariance)



NON-LINEAR (non-linear normalizing flow)

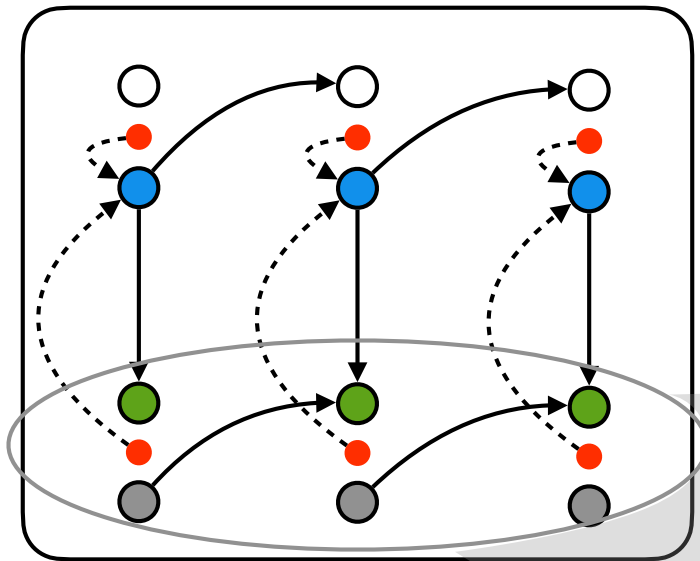


NORMALIZING FLOWS

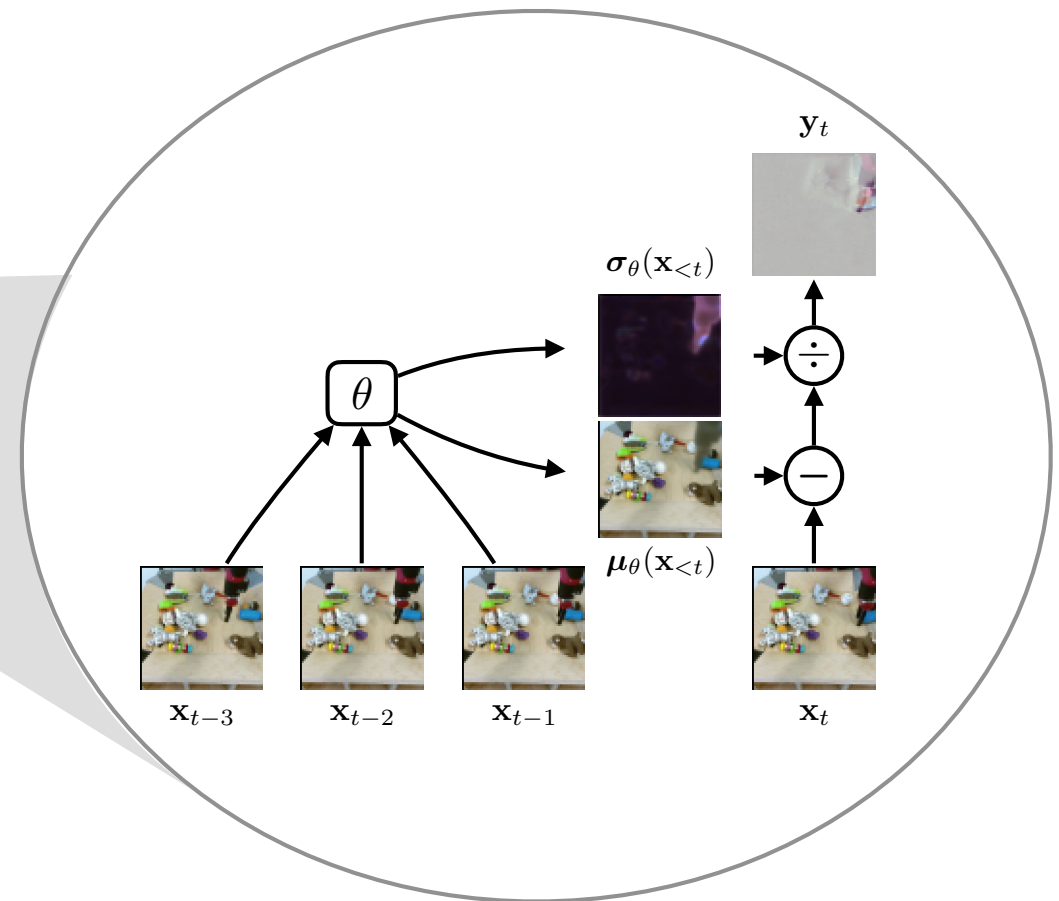
normalizing flows can be applied to any parametric distribution

sequential autoregressive flows

remove dependencies across time

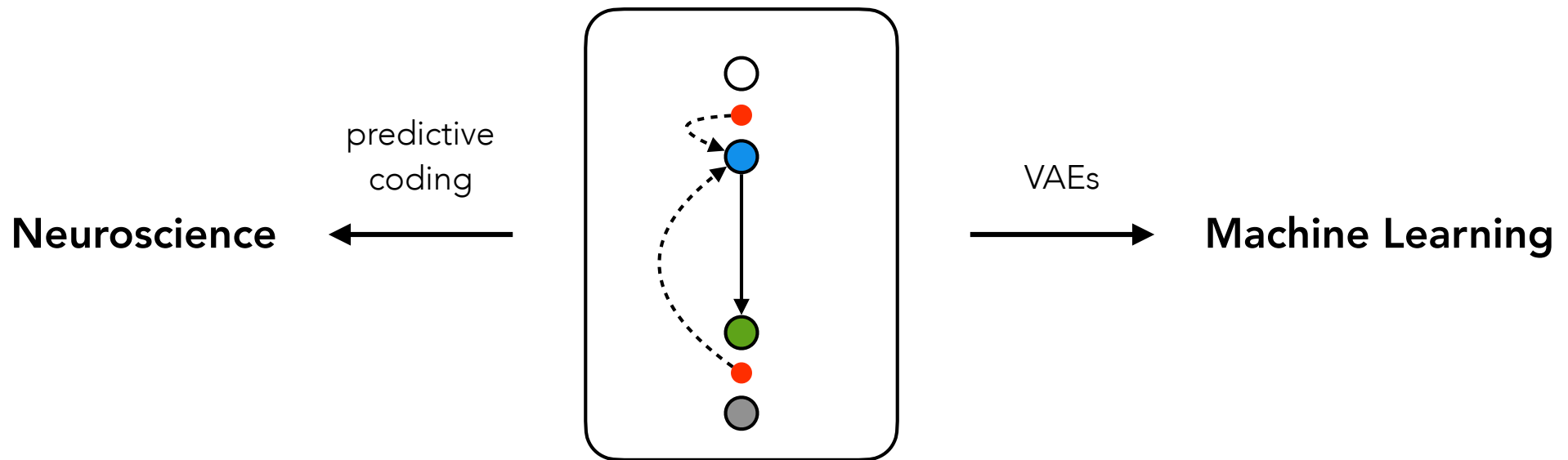
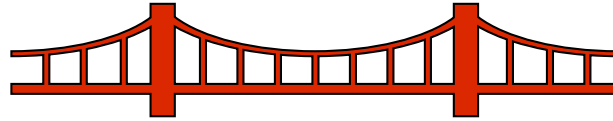


improved performance
& generalization

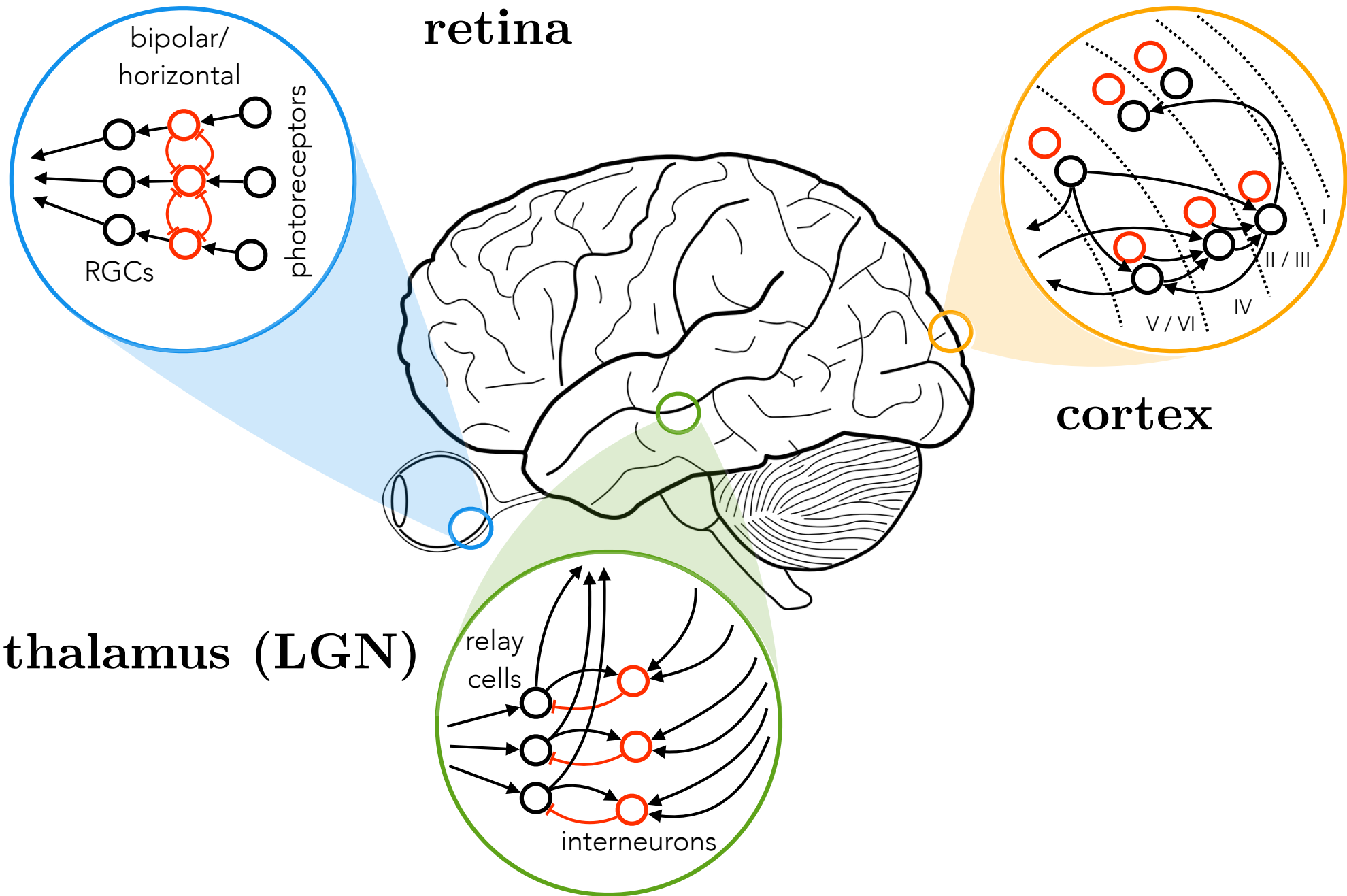


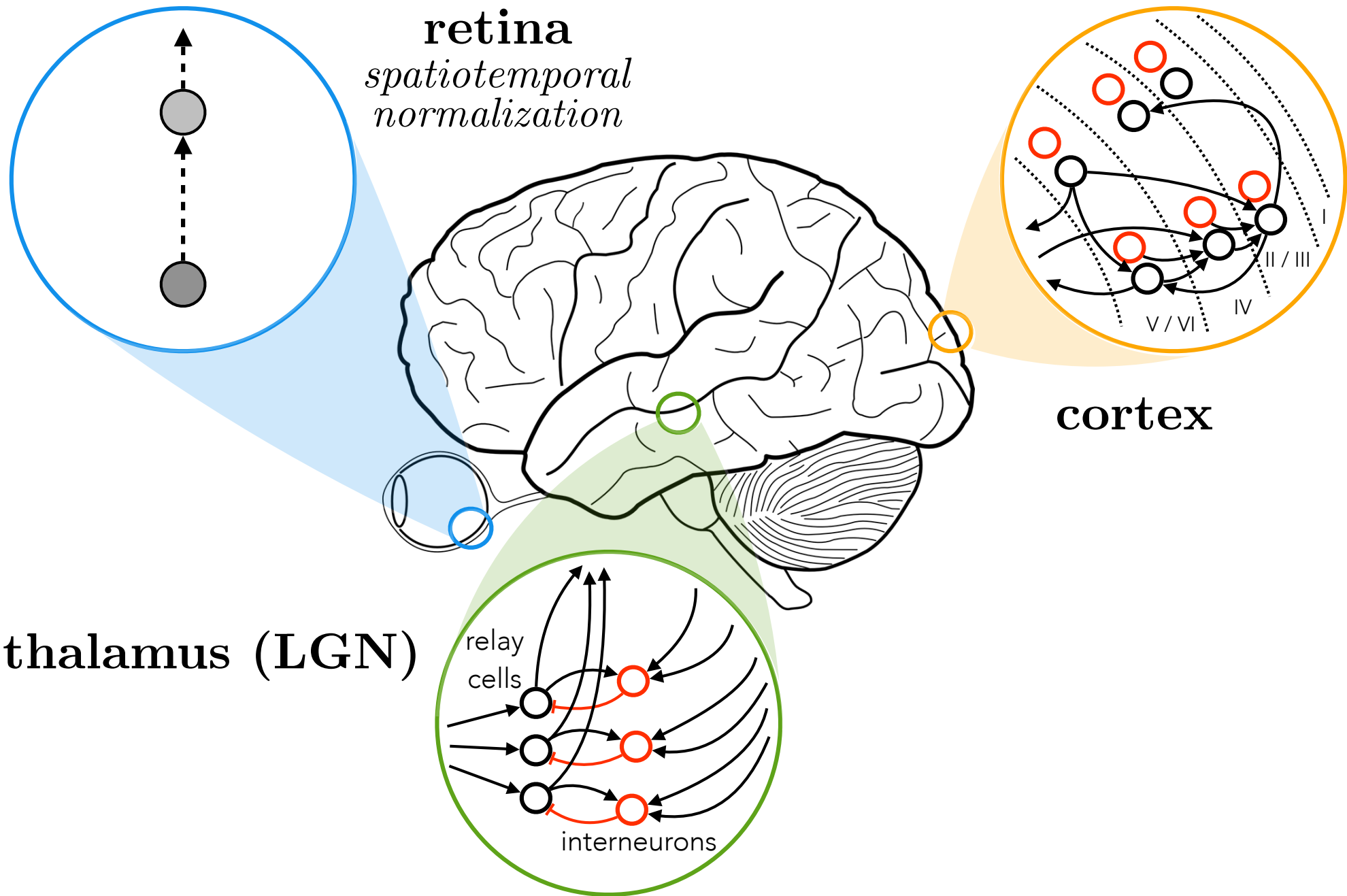
IV. CORRESPONDENCES

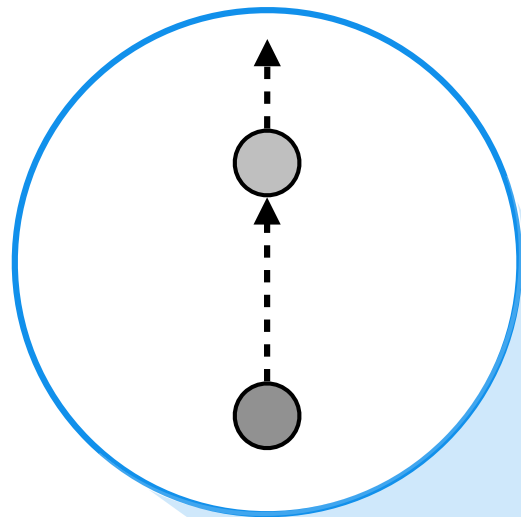
a bridge between neuroscience and machine learning



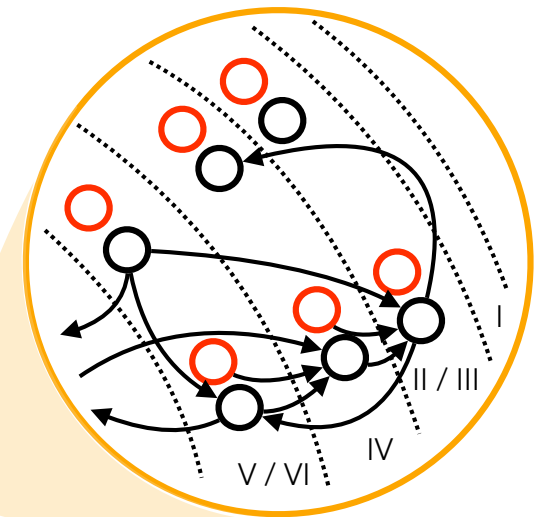
providing new correspondences...



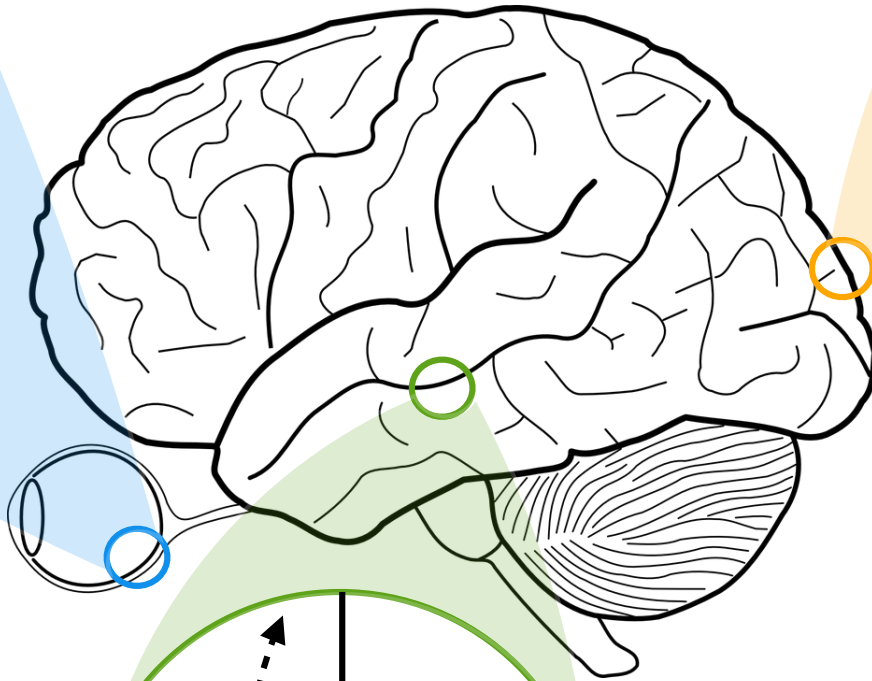




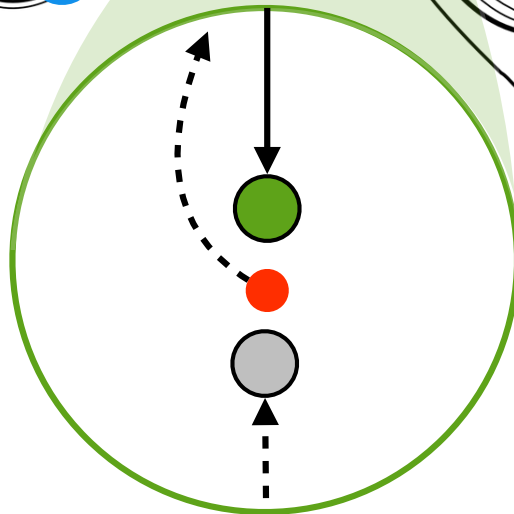
retina
*spatiotemporal
normalization*

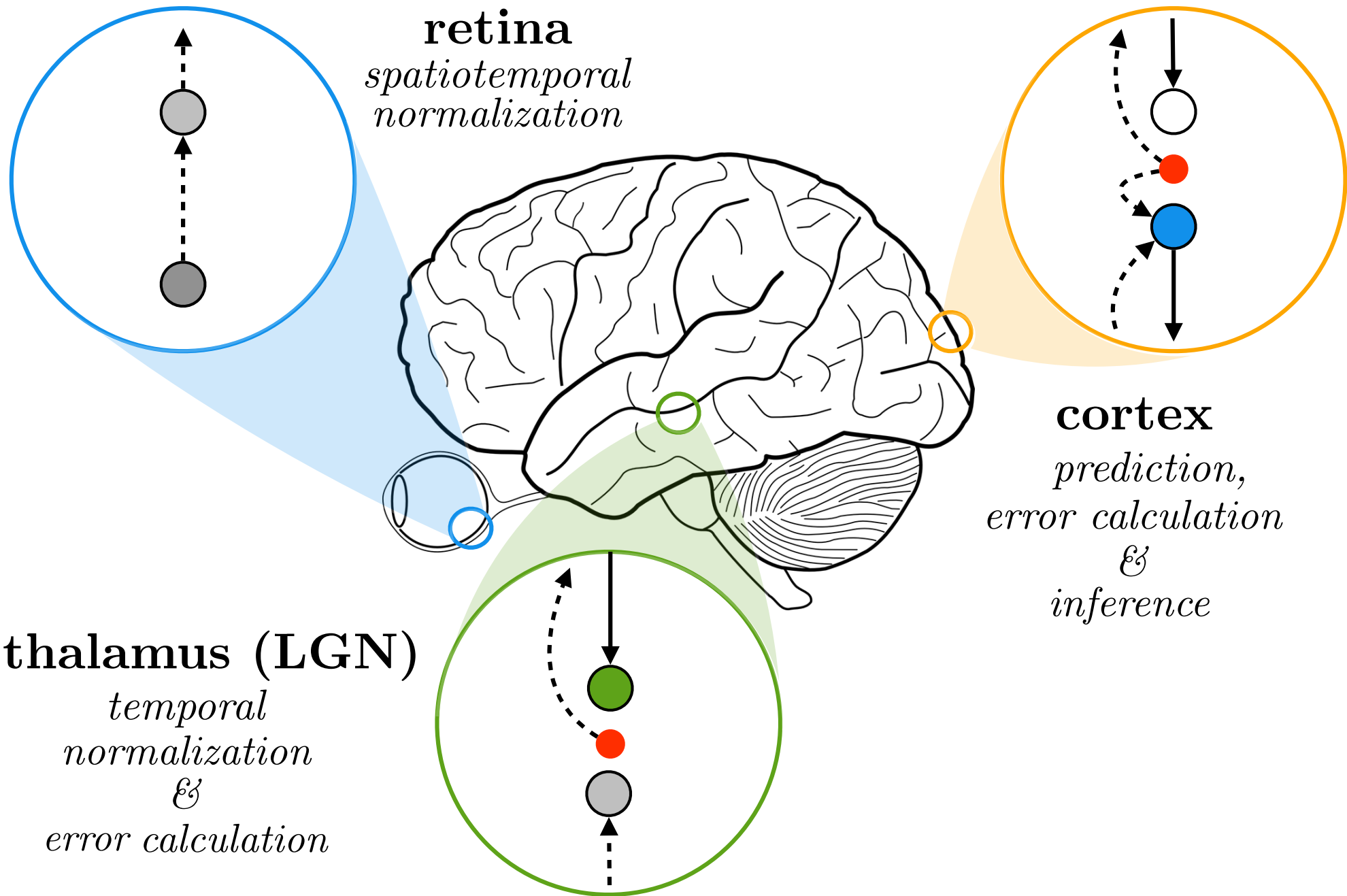


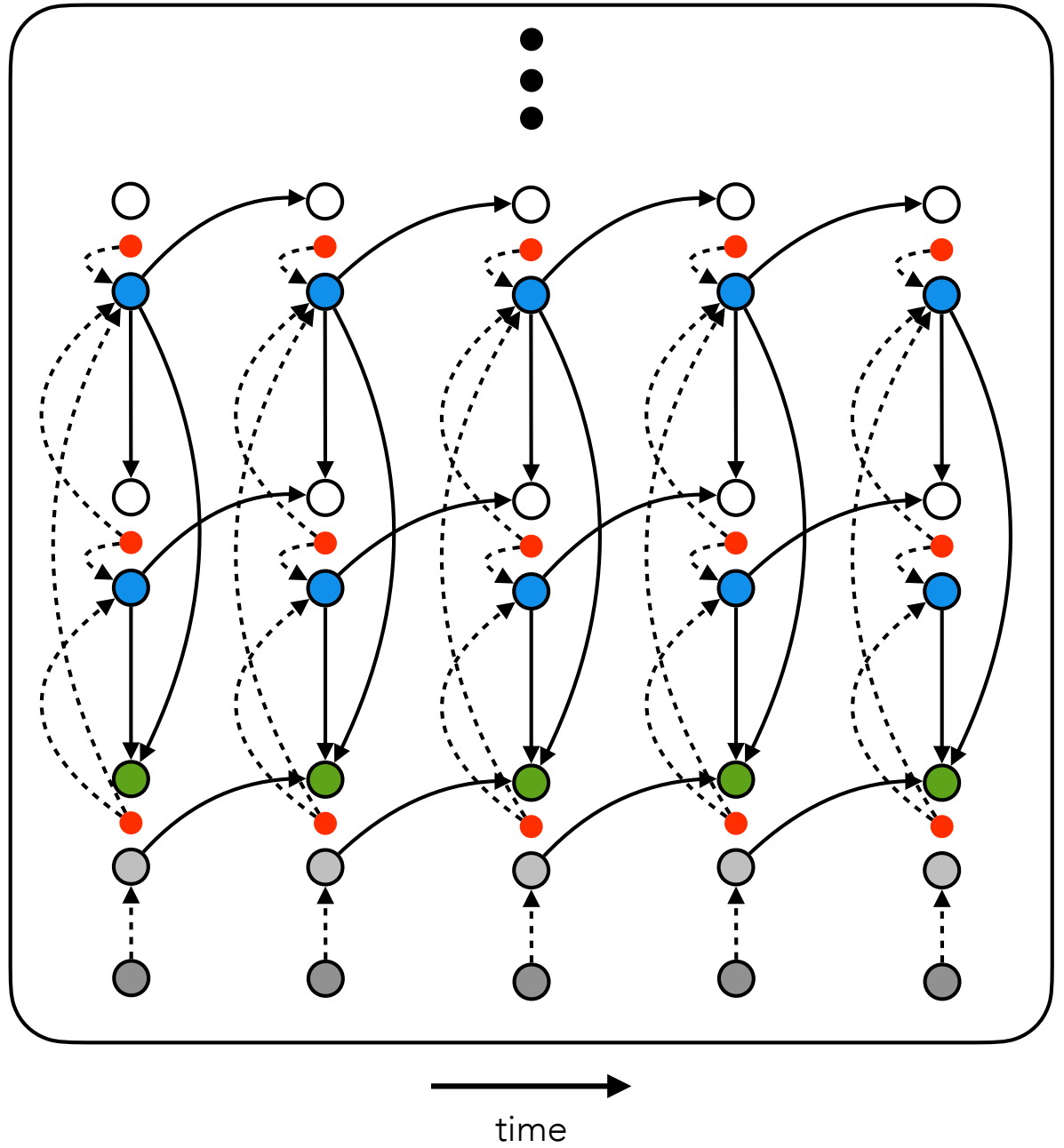
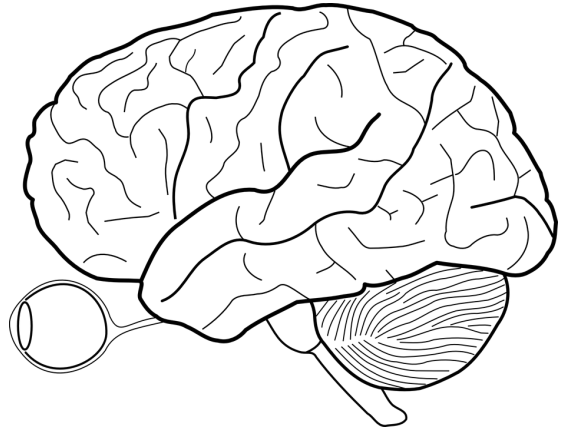
cortex

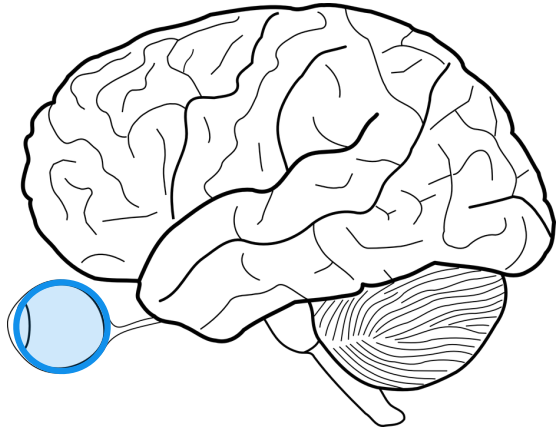


thalamus (LGN)
*temporal
normalization
&
error calculation*

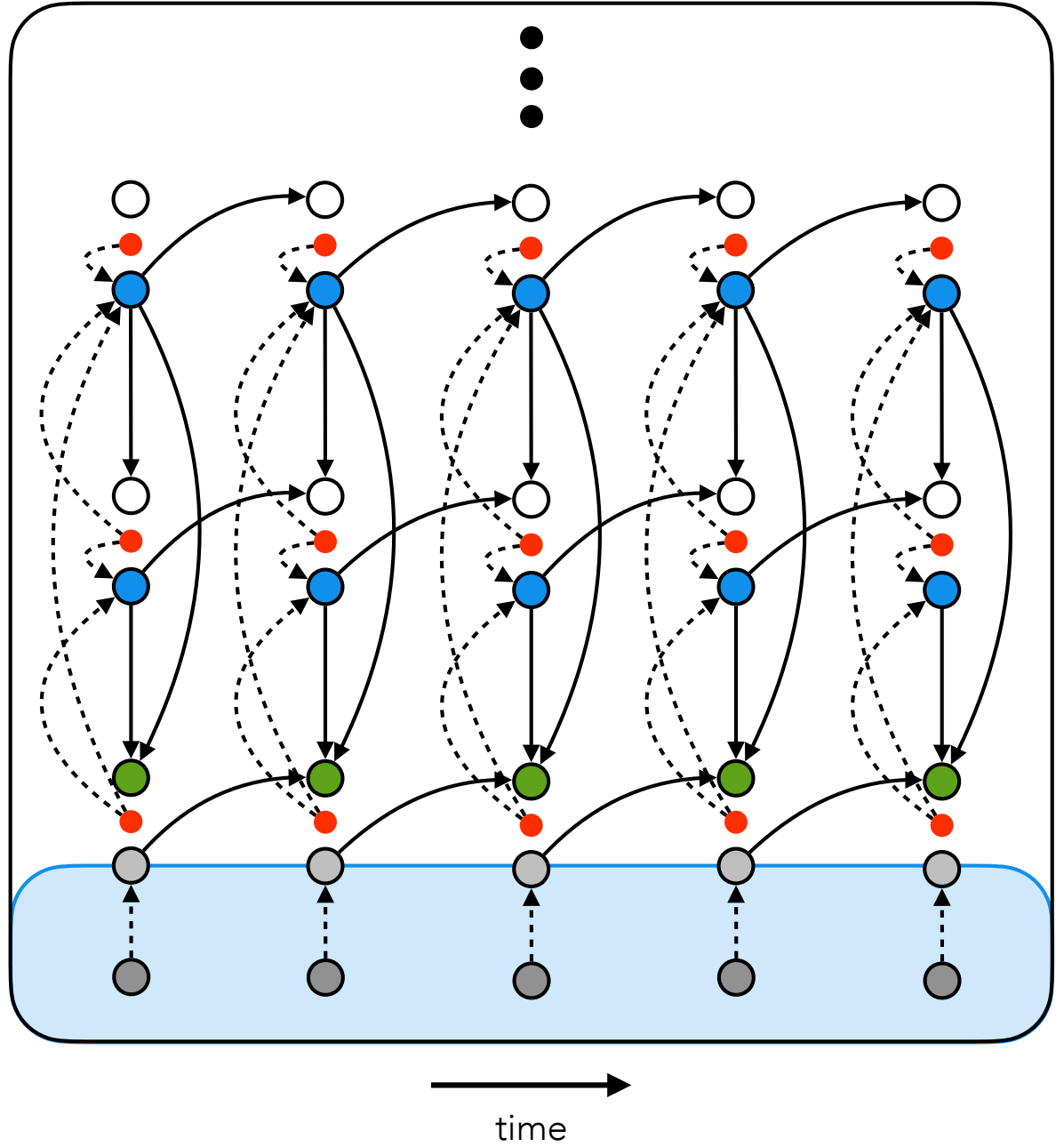


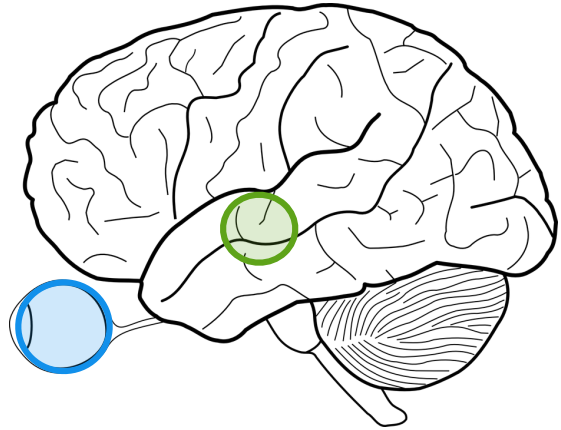






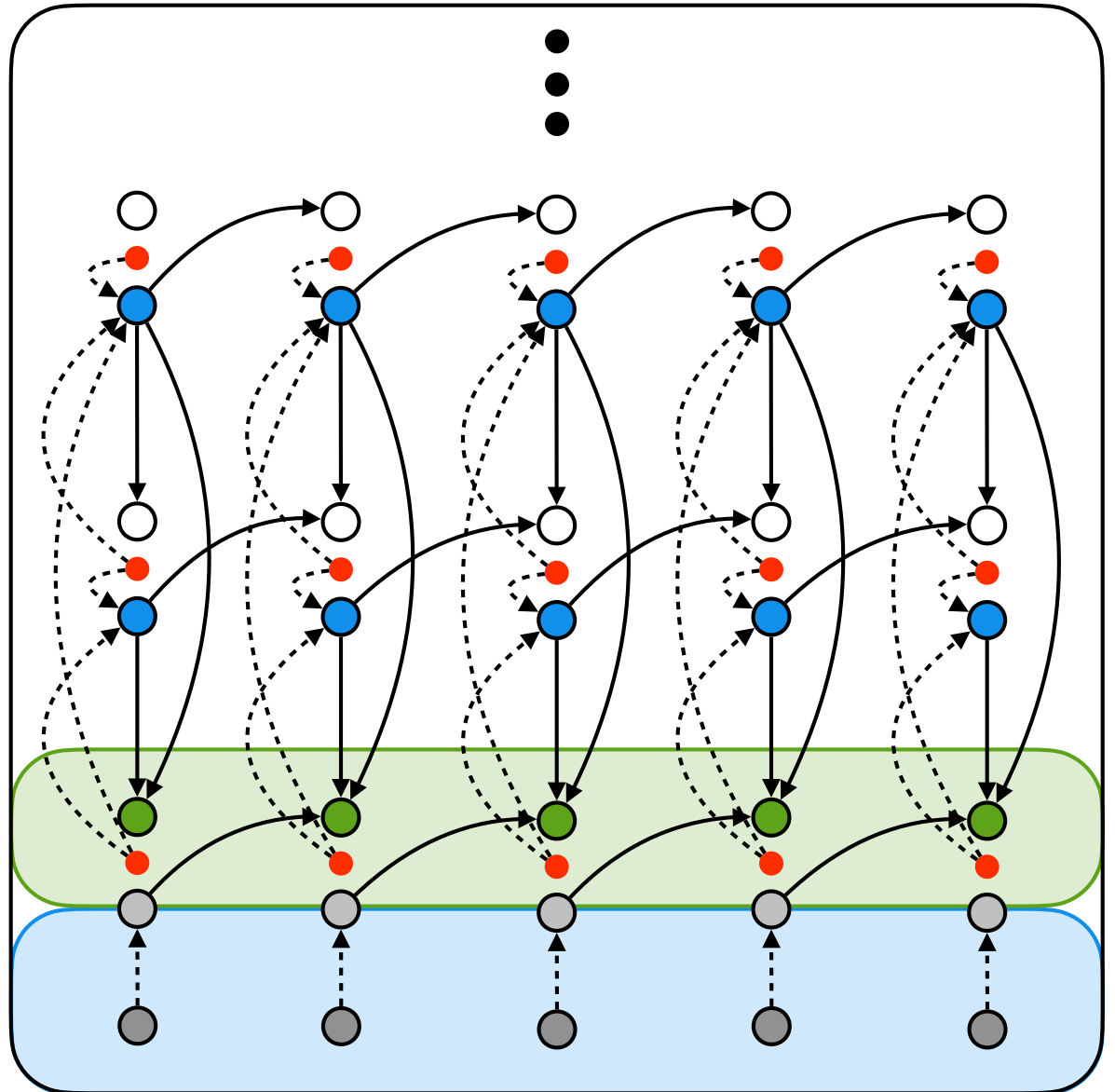
retina

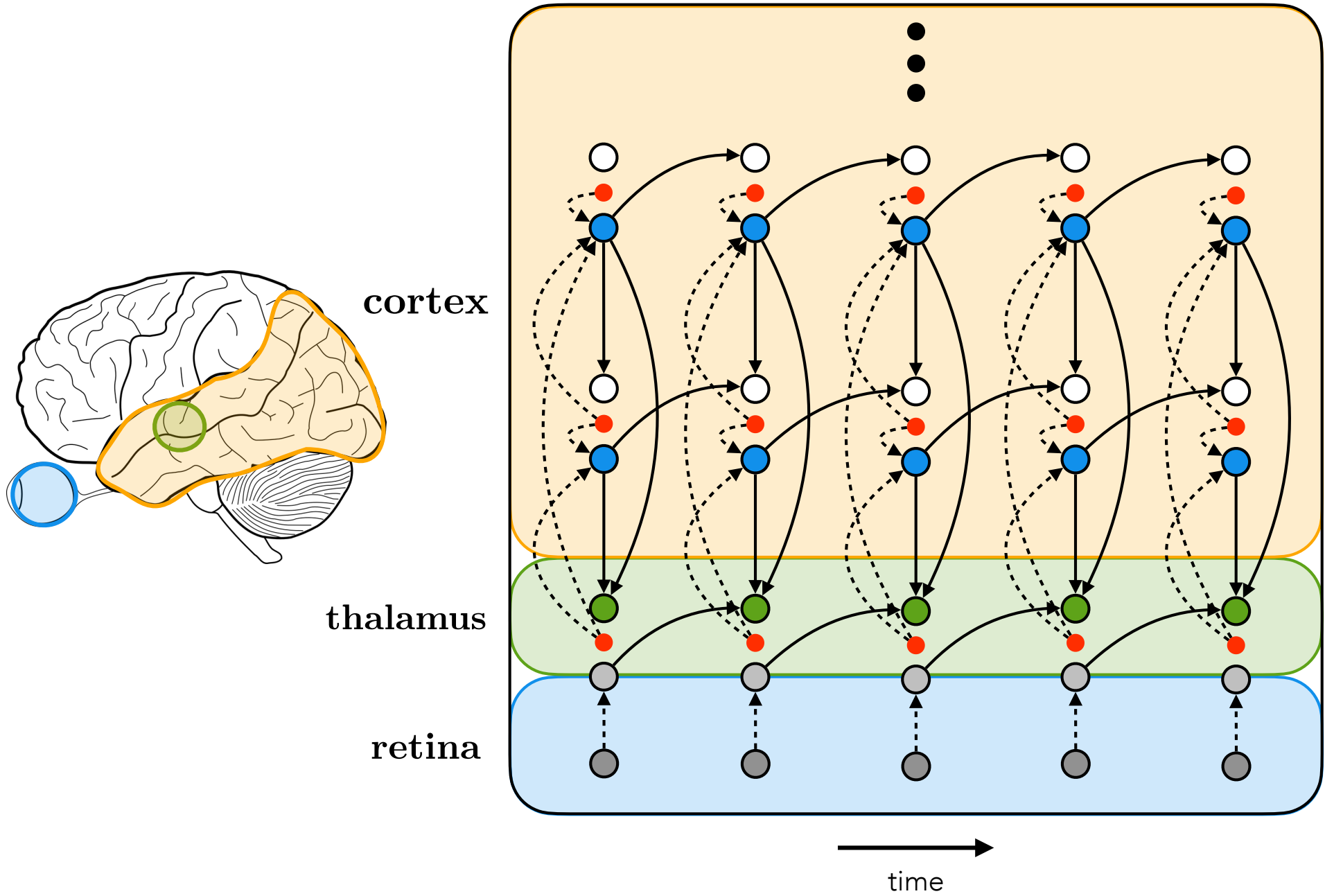




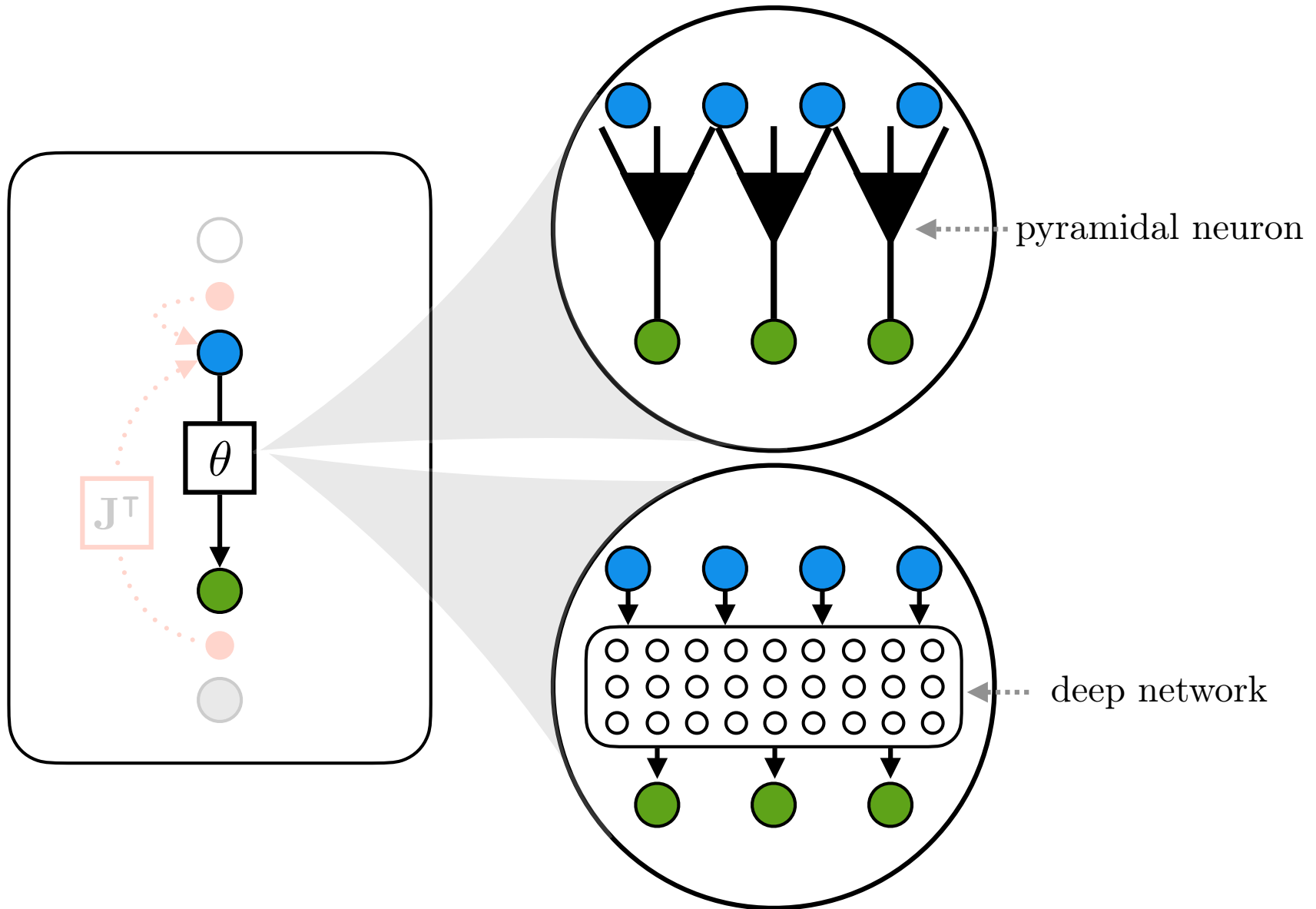
thalamus

retina

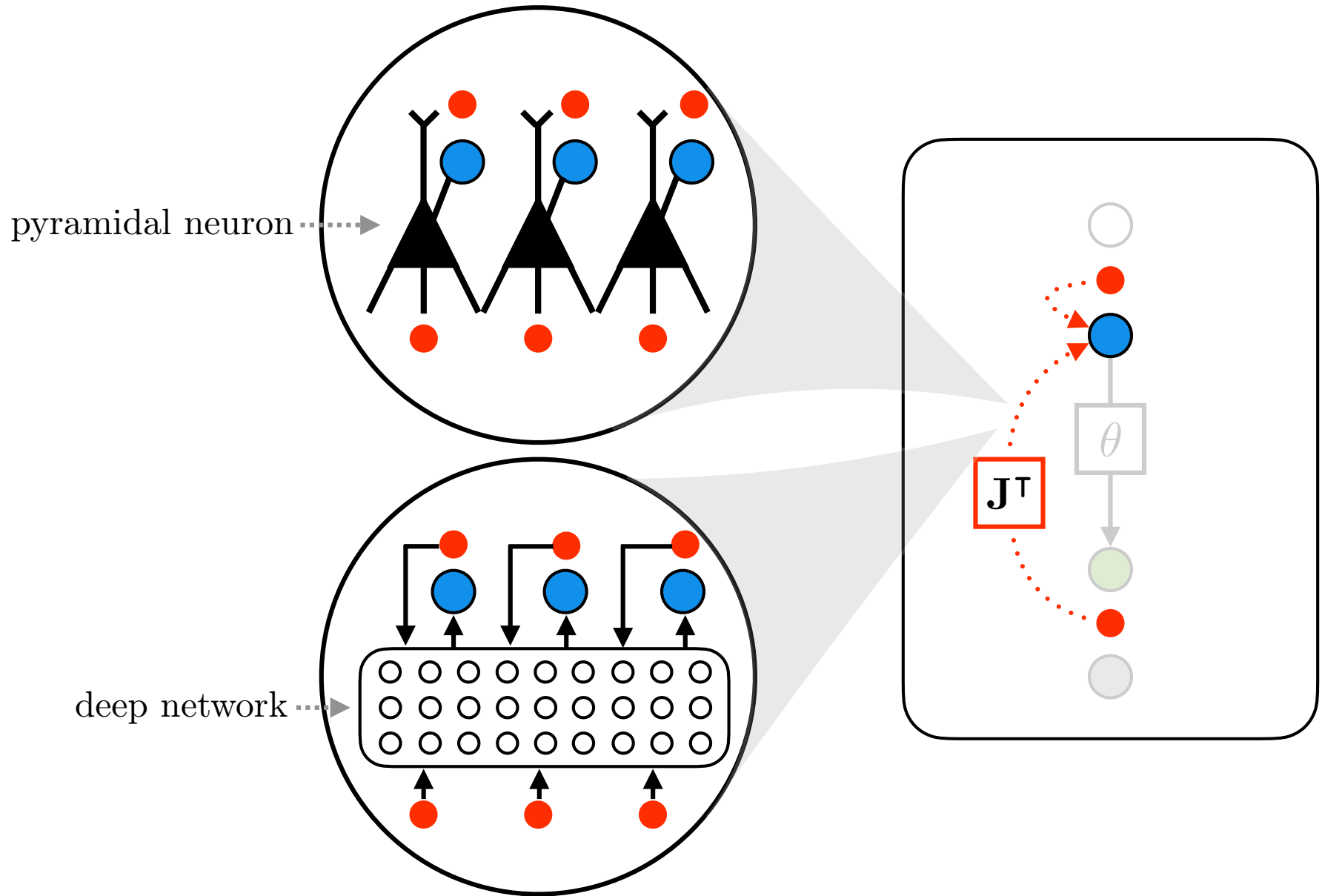




PYRAMIDAL NEURONS AS DEEP NETWORKS

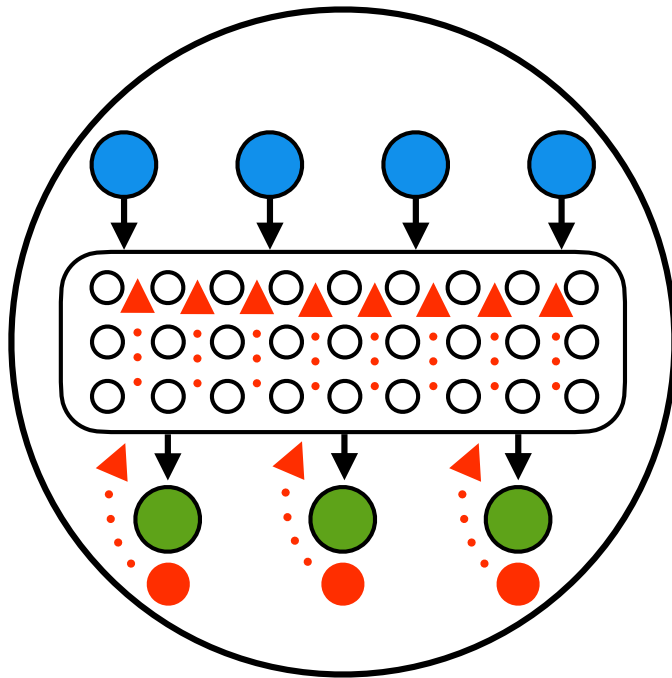


PYRAMIDAL NEURONS AS DEEP NETWORKS

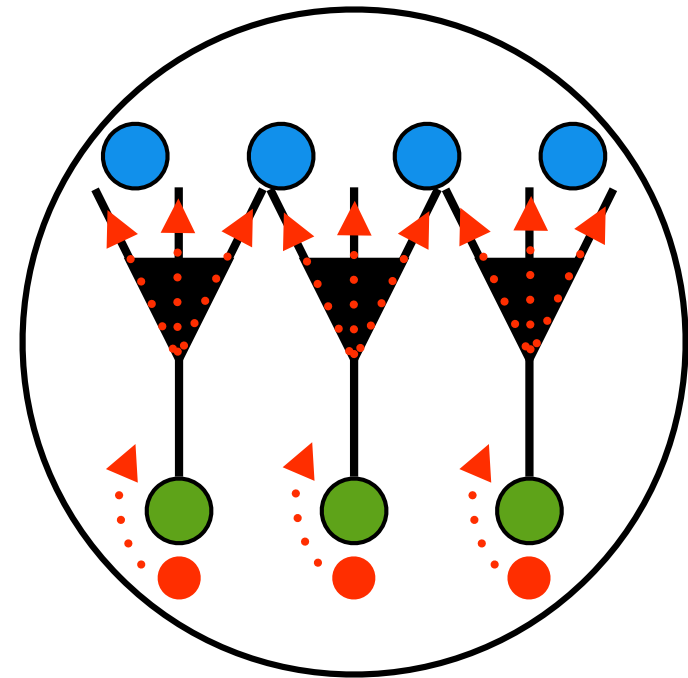


BACKPROPAGATION WITHIN NEURONS

in hierarchical latent variable models, errors provide a **local** training signal
see, e.g., target propagation (Bengio, 2014)

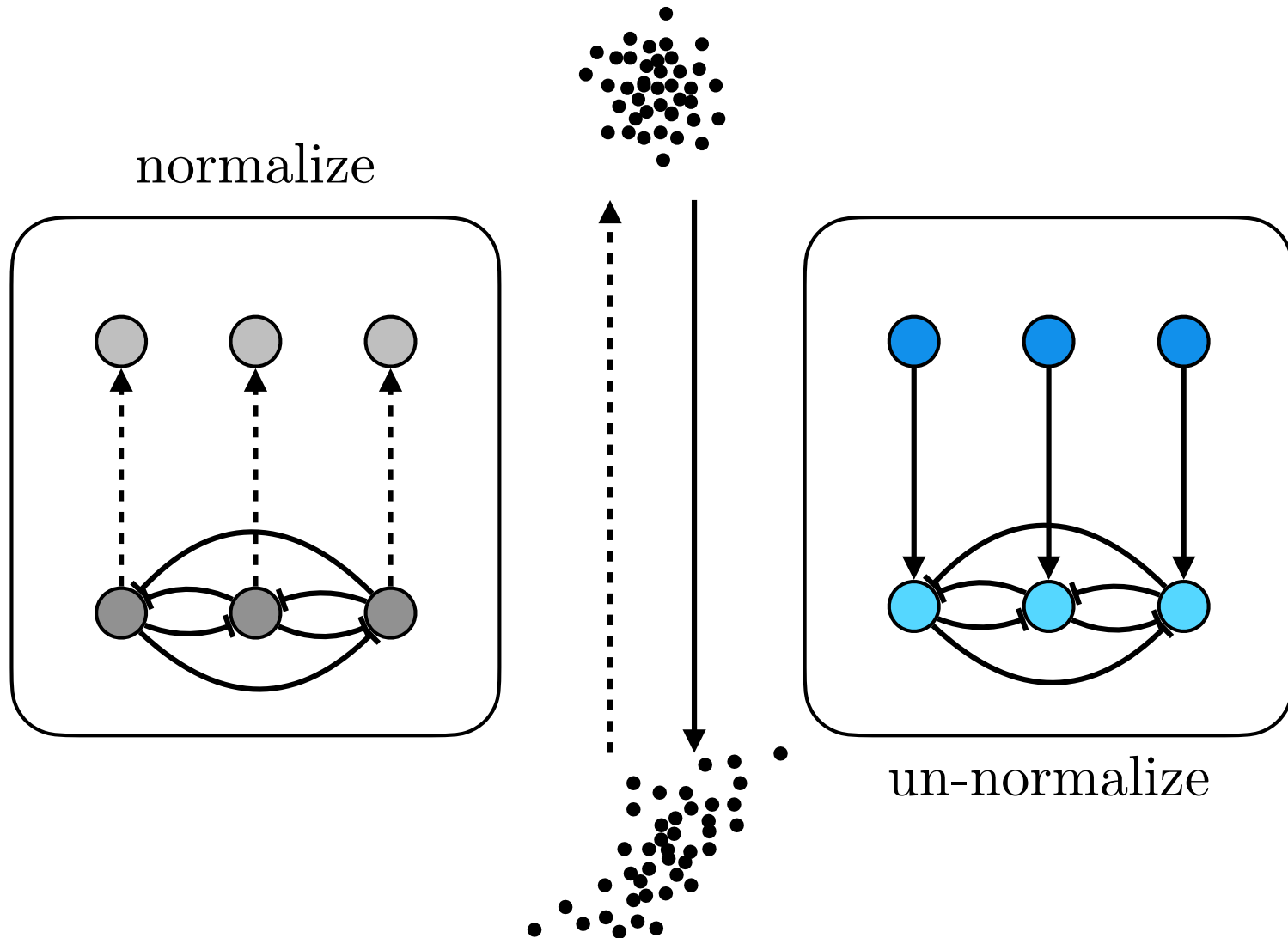


only need to
backpropagate gradients
between variables

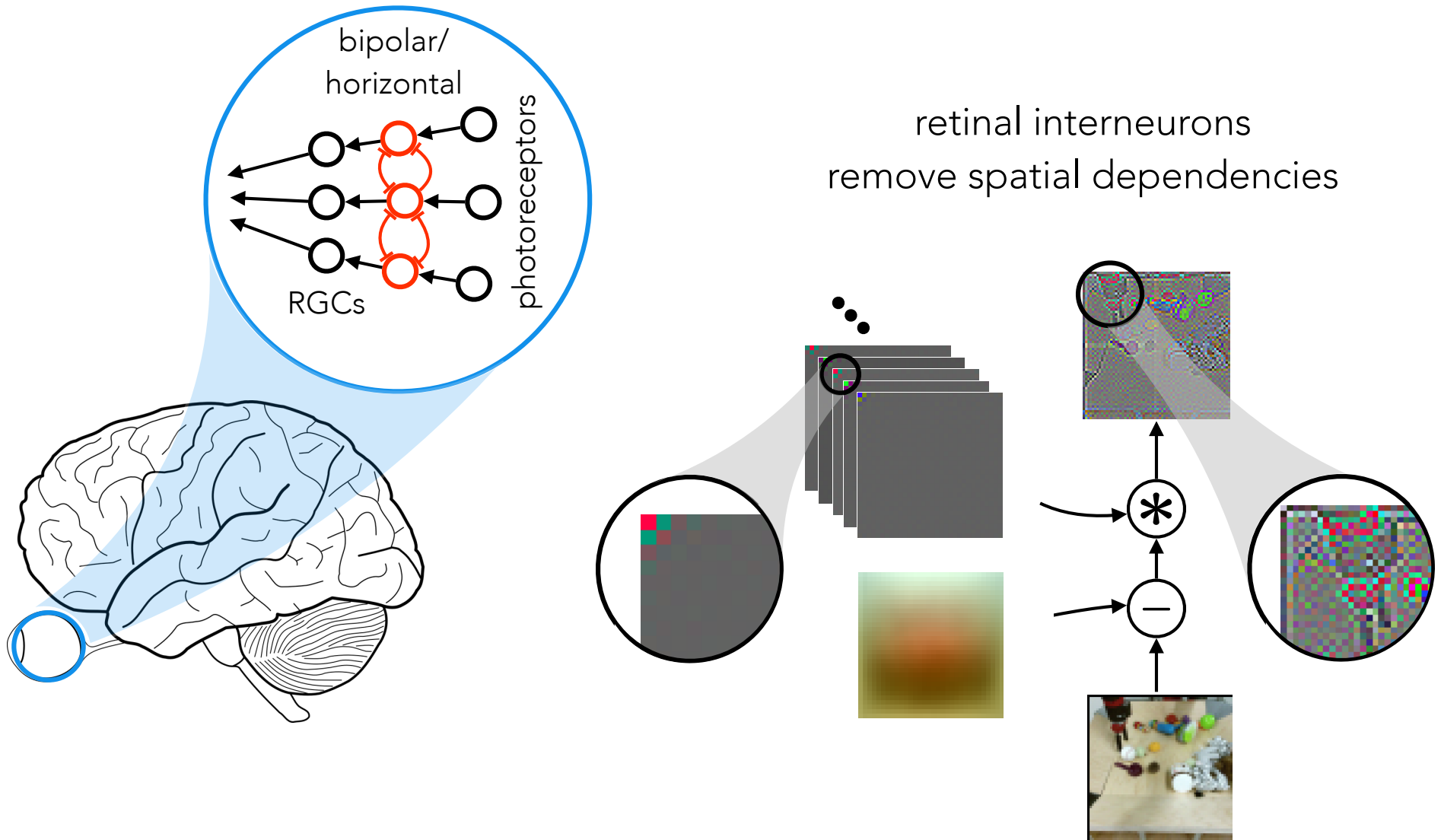


corresponding process
would occur within neurons

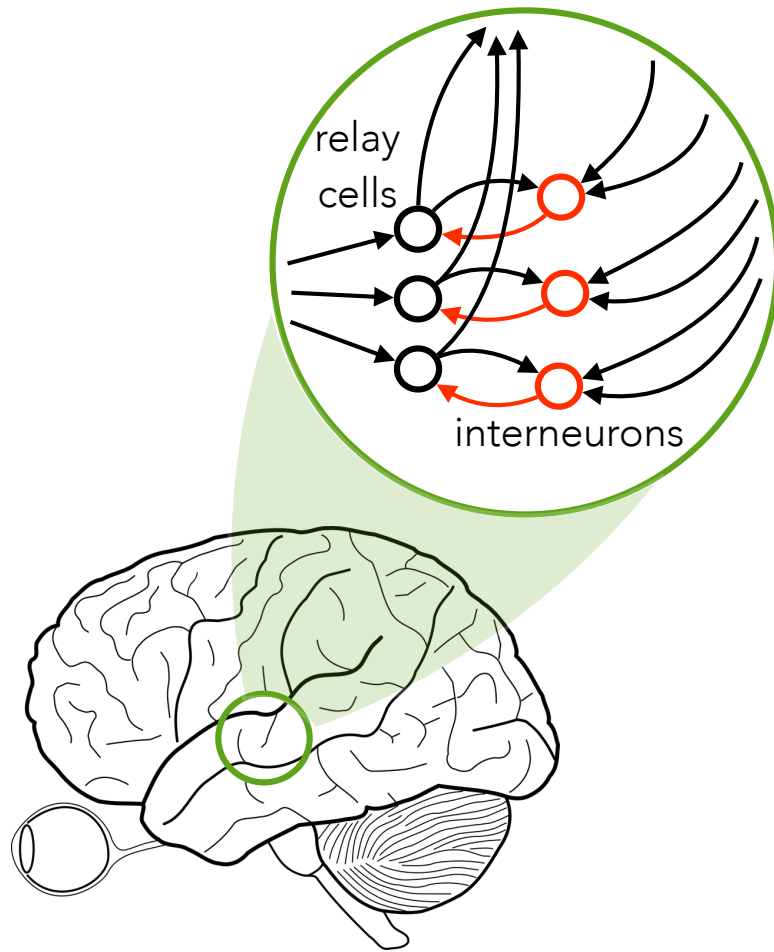
interneurons can add or remove spatiotemporal dependencies between neurons



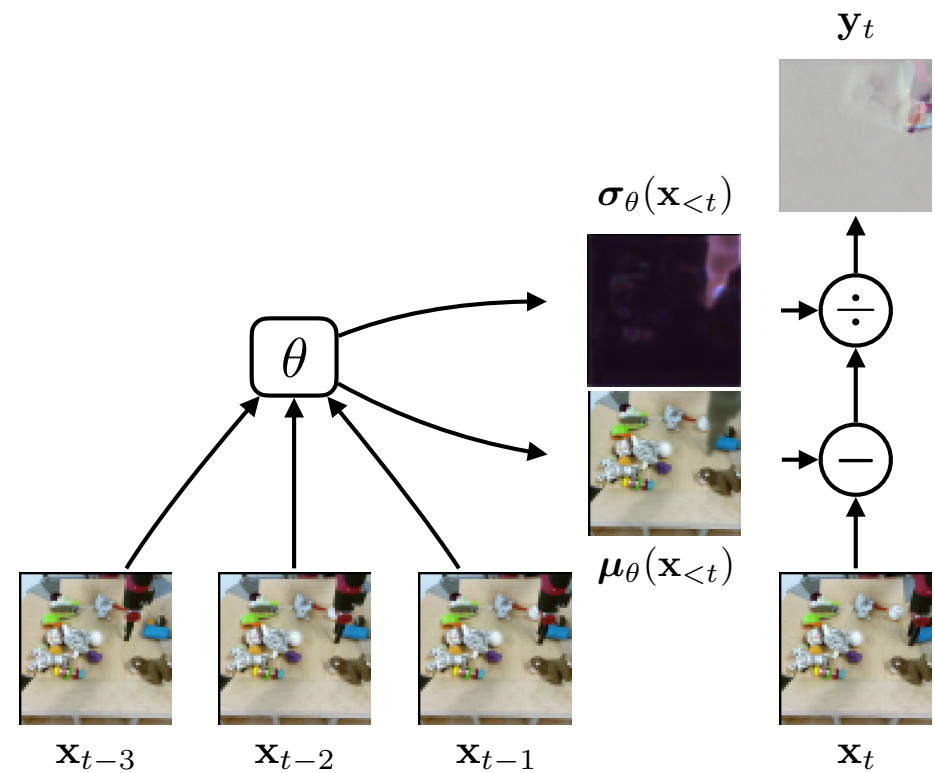
LATERAL INHIBITION & NORMALIZING FLOWS



LATERAL INHIBITION & NORMALIZING FLOWS

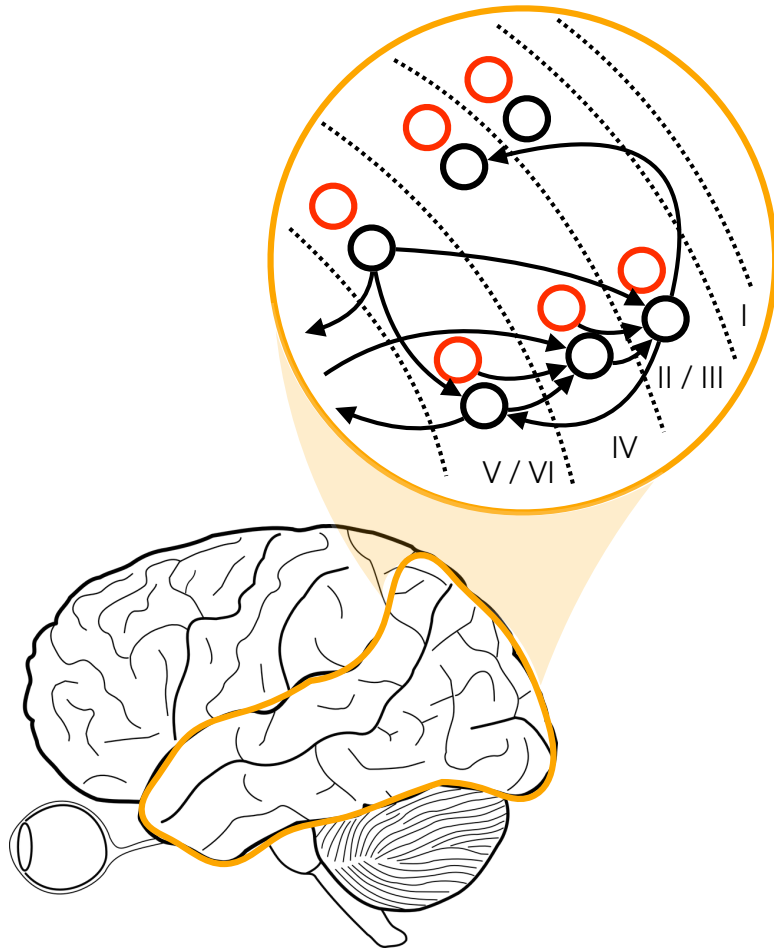


thalamic interneurons
remove low-level temporal dependencies

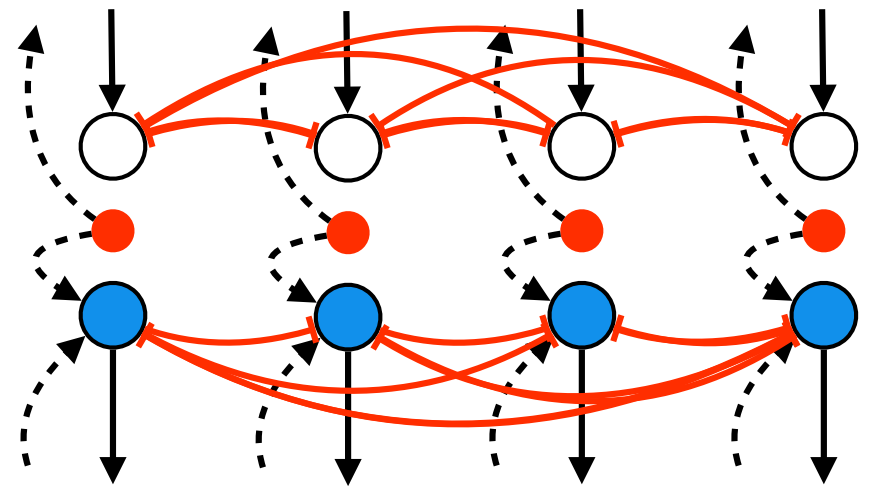


temporal whitening:
inverting an affine normalizing flow

LATERAL INHIBITION & NORMALIZING FLOWS



cortical interneurons
add/remove spatiotemporal dependencies
between columns

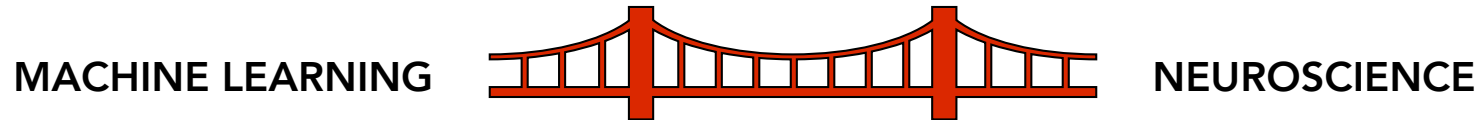


cortical columns

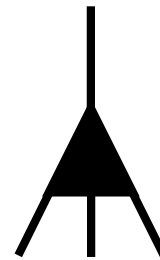
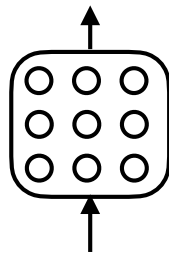
forward/inverse normalizing flows

SUMMARY OF CORRESPONDENCES

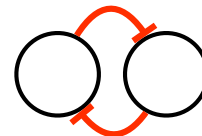
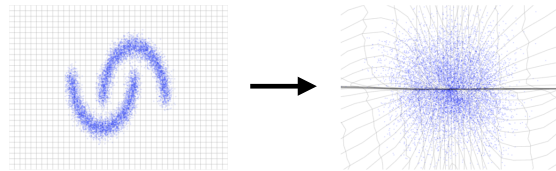
by traversing the bridge between predictive coding & VAES...



deep networks & pyramidal neurons

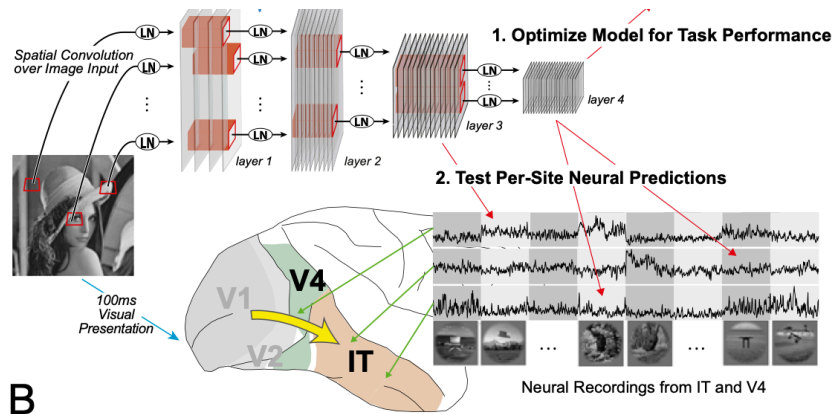


normalizing flows & lateral inhibition



A NEW COMPUTATIONAL MODEL FOR NEUROSCIENCE

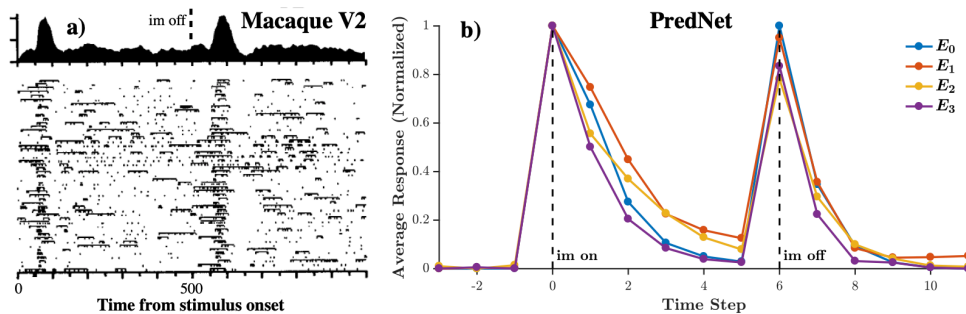
current efforts to compare deep networks and the brain are limited



*compare with supervised learning
on brief image presentation
ignore feedback + dynamics*

Yamins et al., 2014

large-scale VAEs provide a route toward testing predictive coding

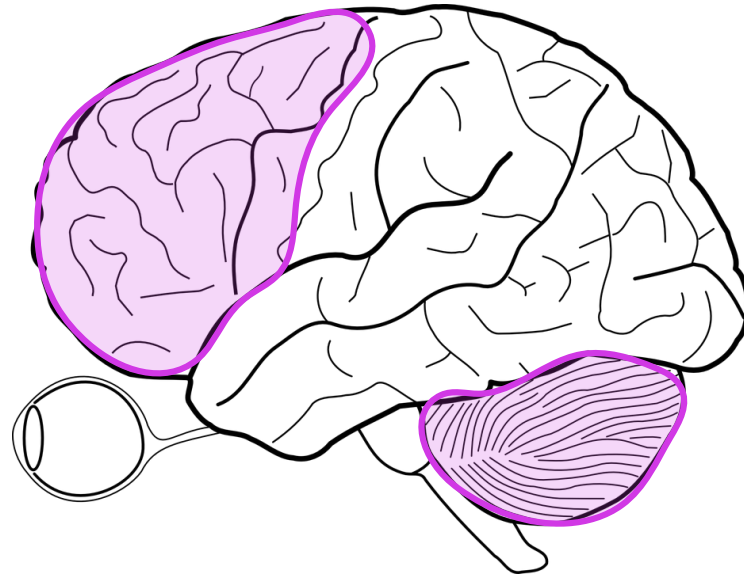


Lotter et al., 2018



Child, 2020

MOTOR CONTROL / EXECUTIVE FUNCTION



can apply the same computational techniques to control/RL

→ prediction errors drive **feedback control**

Iterative Amortized Policy Optimization, Marino et al., 2020

arXiv:2010.10670

Thank You