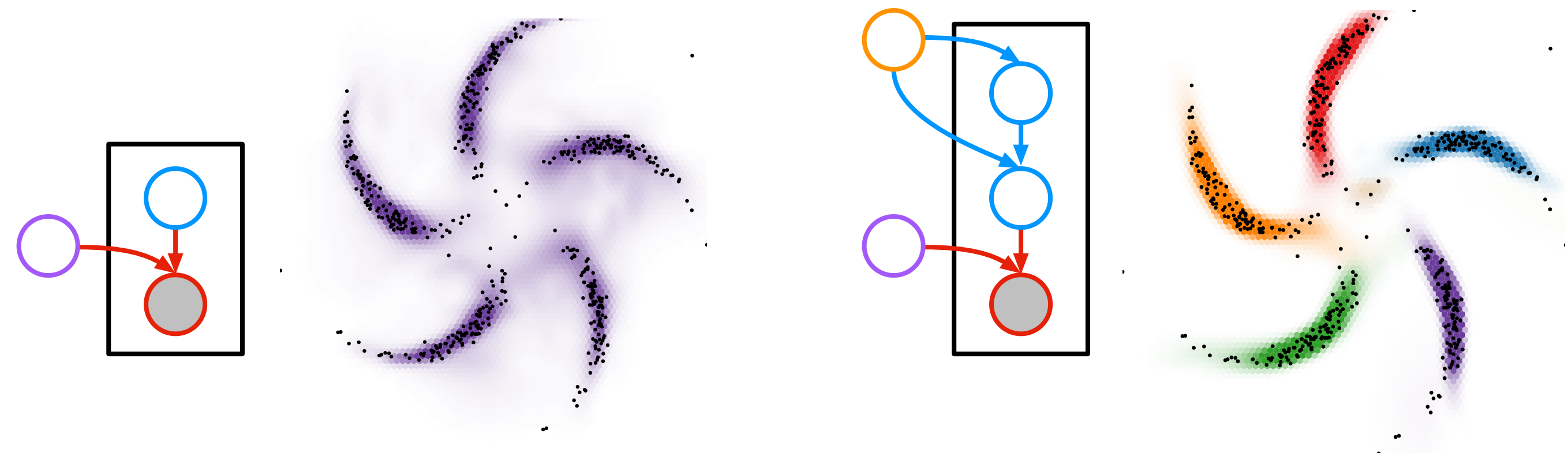


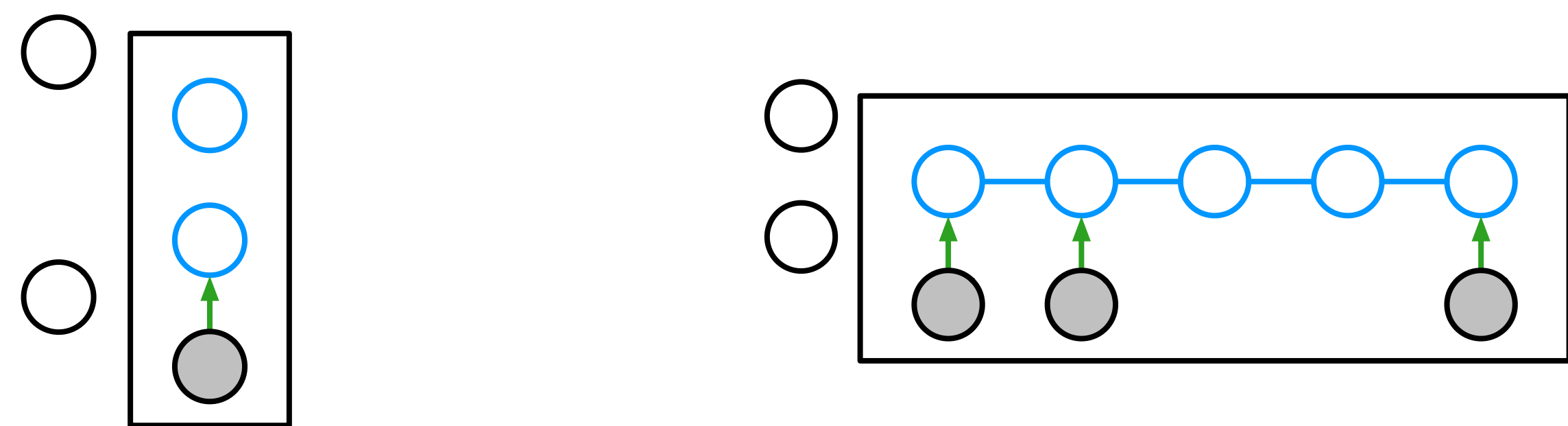
motivation

TL;DR variational autoencoders + latent graphical models

modeling idea use PGM priors to **organize** the latent space, along with neural net observation models for **flexible representations**



inference idea use PGMs to **synthesize information** from **recognition nets** instead of making a single inference net do everything



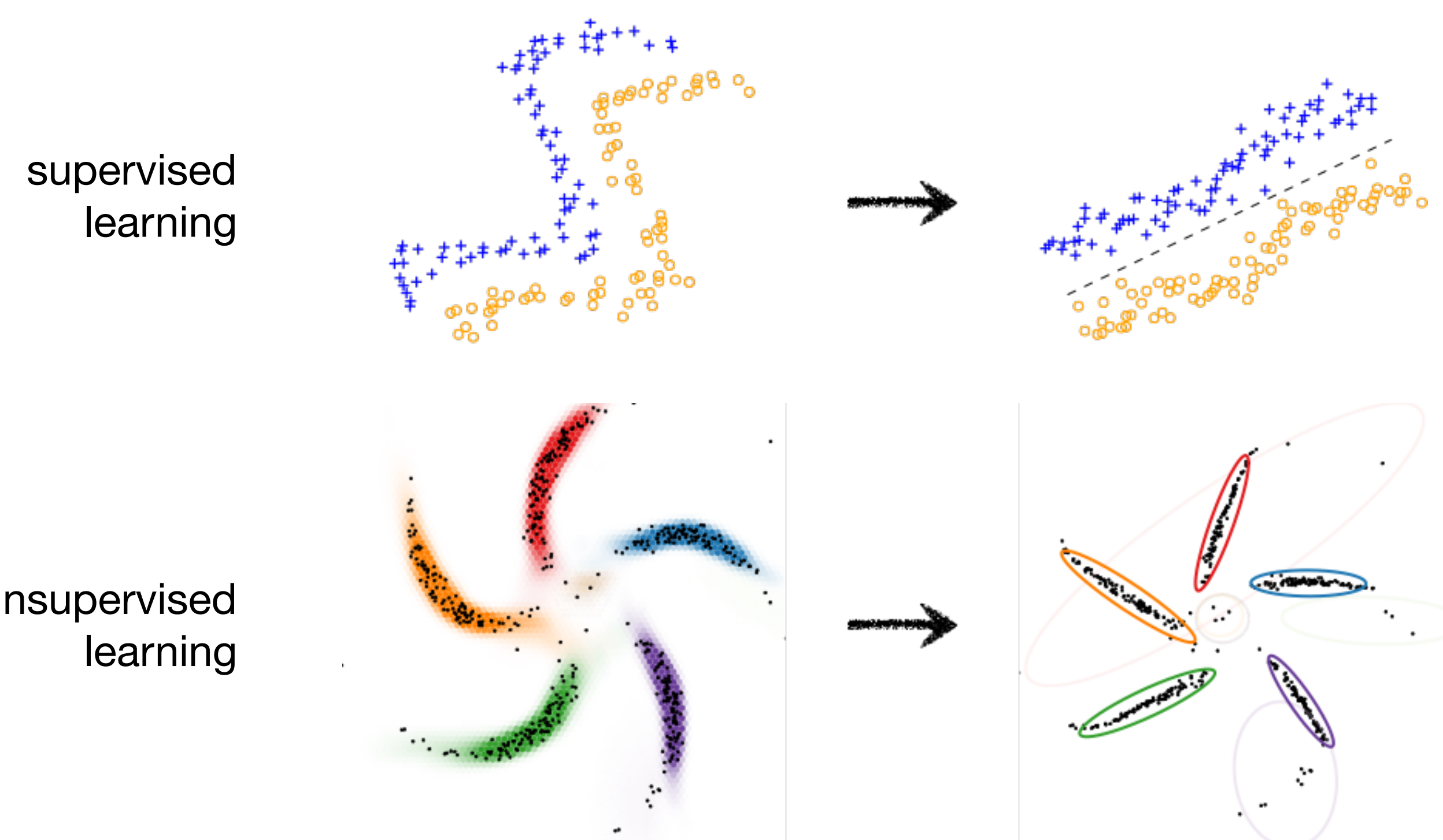
probabilistic graphical models

- + structured representations
- + priors and uncertainty
- + data and computational efficiency within rigid model classes
- rigid assumptions may not fit
- feature engineering
- more flexible models can require slow top-down inference

deep neural networks

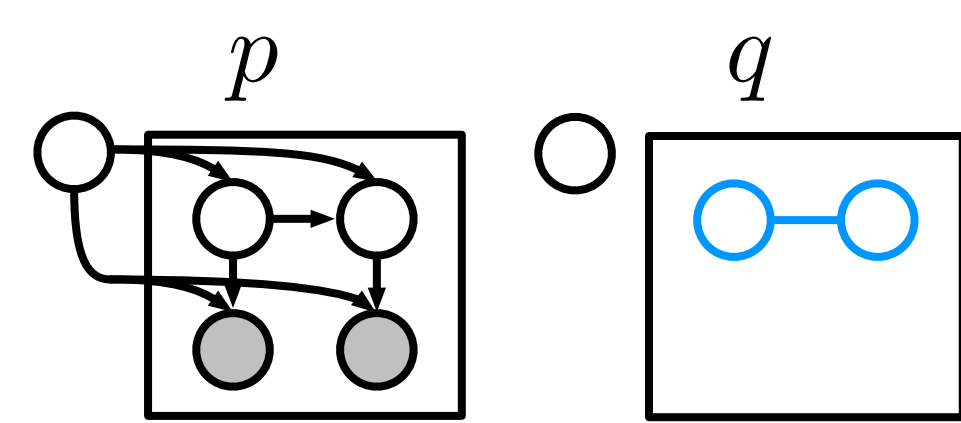
- neural net "goo"
- difficult parameterization
- can require lots of data
- + flexible, high capacity
- + feature learning
- + recognition networks for fast bottom-up inference

automatically learn representations in which structured PGMs fit well



inference

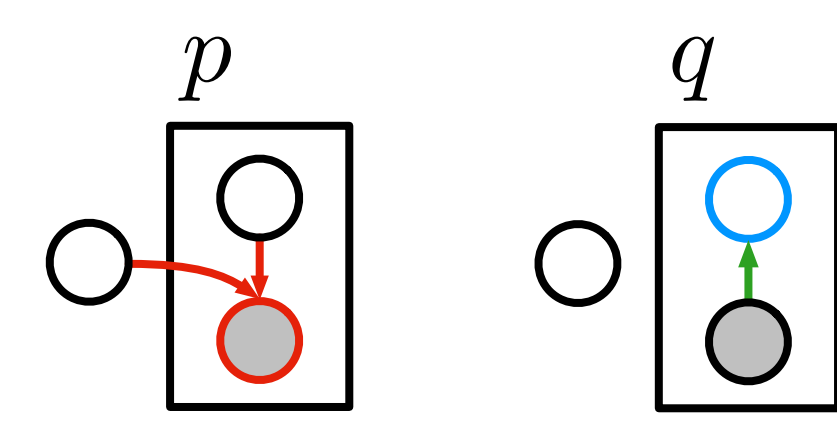
natural gradient SVI



$$q^*(x) \triangleq \arg \max_{q(x)} \mathcal{L}[q(\theta)q(x)]$$

- + optimal local factor
- expensive for general obs.
- + exploit conj. graph structure
- + arbitrary inference queries
- + natural gradients

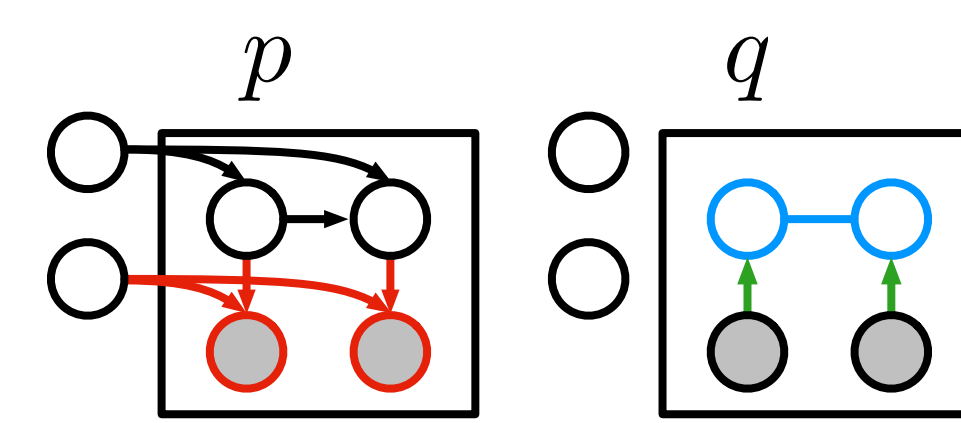
variational autoencoders



$$q^*(x) \triangleq \mathcal{N}(x|\mu(y; \phi), \Sigma(y; \phi))$$

- suboptimal local factor
- + fast for general obs.
- ϕ does all local inference
- limited inference queries
- no natural gradients

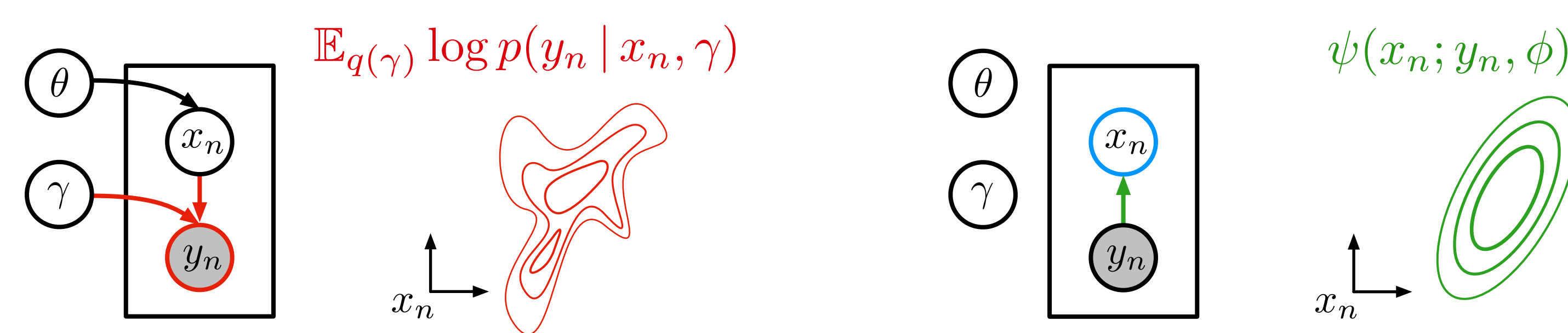
structured VAEs (this work)



$$q^*(x) \triangleq ?$$

- \pm optimal given conj. evidence
- + fast for general obs.
- + exploit conj. graph structure
- + arbitrary inference queries
- + some natural gradients

main idea learn to summarize complicated evidence with simple conjugate potentials (as in CRFs)



step 1

- compute local evidence with recognition networks
- 1(a) sample minibatch
 - 1(b) apply recognition networks...
 - 1(c) ...to get PGM potentials

step 2

- run fast PGM inference
- 2(a) run local mean field
 - 2(b) and message passing
 - 2(c) to fixed point

step 3

- compute unbiased ELBO gradients with respect to all parameters
- 3(a) sample and compute flat gradients w.r.t ϕ and γ
 - 3(b) use already-computed values to get natural gradient w.r.t. η_θ

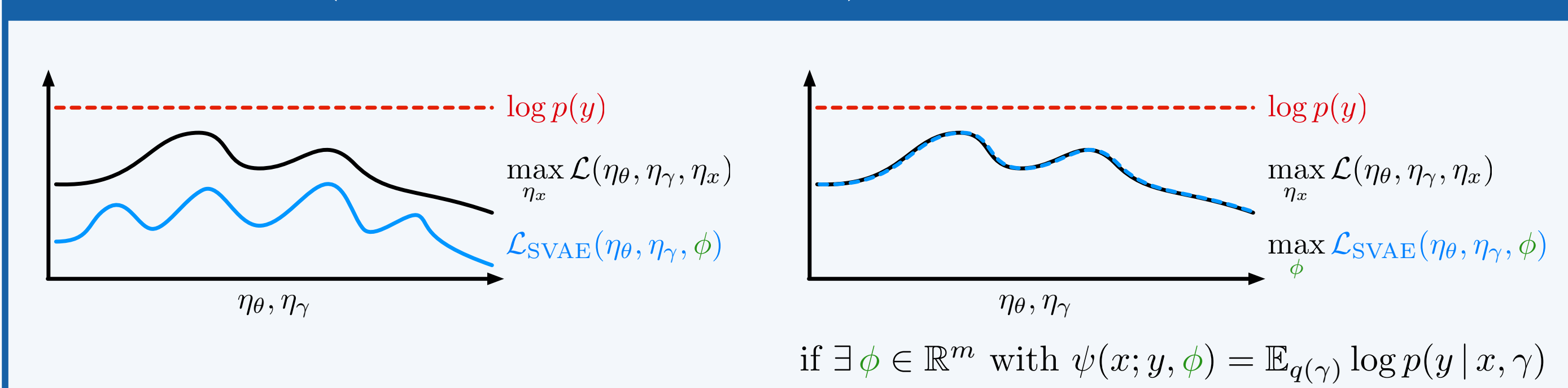
$$\mathcal{L}(\eta_\theta, \eta_\gamma, \eta_x) \triangleq \mathbb{E}_{q(\theta)q(\gamma)q(x)} \left[\log \frac{p(\theta, \gamma, x)p(y|x, \gamma)}{q(\theta)q(\gamma)q(x)} \right]$$

$$\hat{\mathcal{L}}(\eta_\theta, \eta_x, \phi) \triangleq \mathbb{E}_{q(\theta)q(\gamma)q(x)} \left[\log \frac{p(\theta, \gamma, x) \exp\{\psi(x; y, \phi)\}}{q(\theta)q(\gamma)q(x)} \right]$$

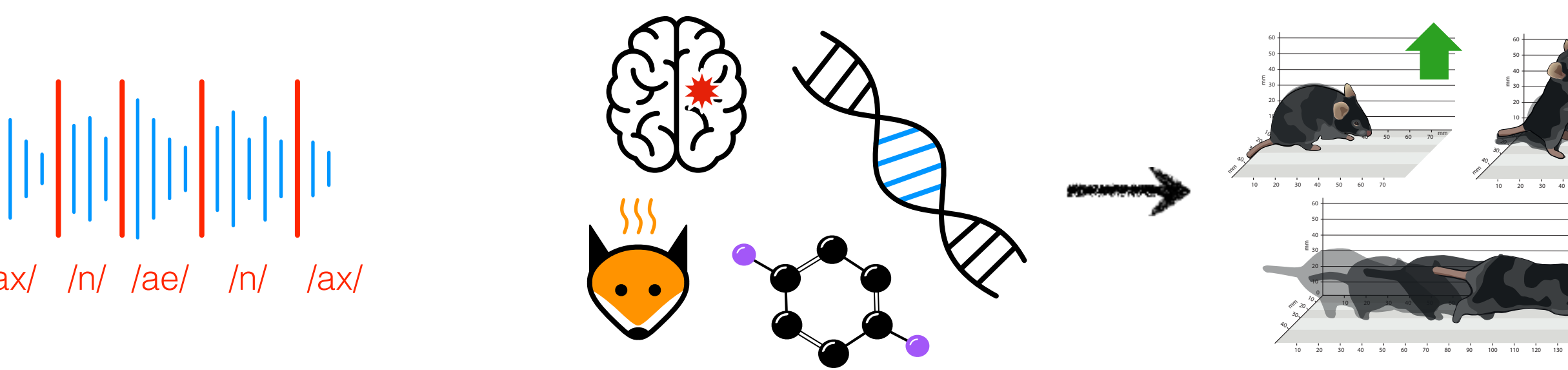
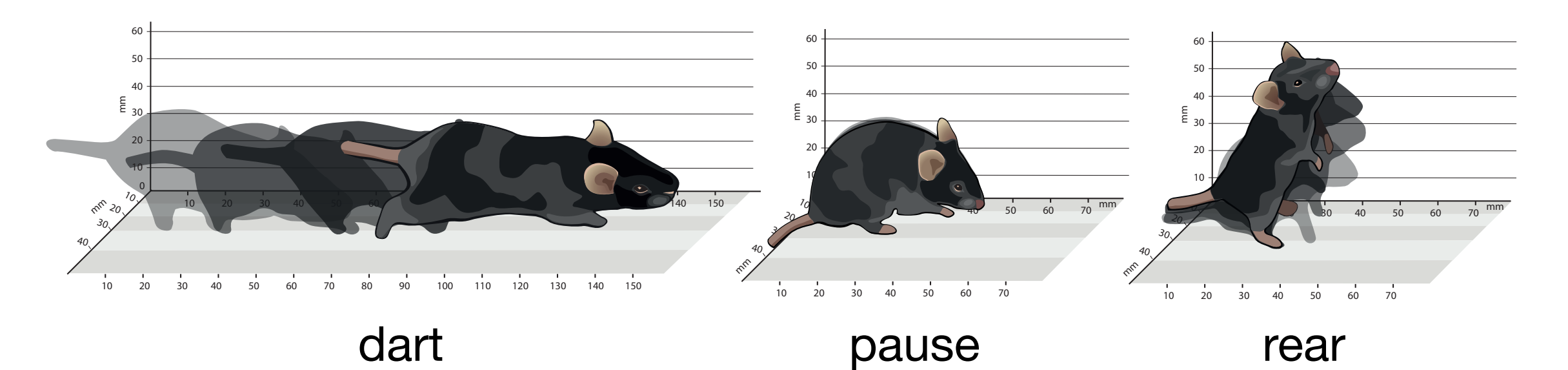
where $\phi(x; y, \phi)$ is a conjugate potential for $p(x|\theta)$

$$\eta_x^*(\eta_\theta, \phi) \triangleq \arg \max_{\eta_x} \hat{\mathcal{L}}(\eta_\theta, \eta_x, \phi) \quad \mathcal{L}_{\text{SVAE}}(\eta_\theta, \eta_\gamma, \phi) \triangleq \mathcal{L}(\eta_\theta, \eta_\gamma, \eta_x^*(\eta_\theta, \phi))$$

Proposition (log evidence lower bound)



learning to parse mouse behavior from depth video



SVAE approach fit a latent switching linear dynamical system (SLDS) and a neural network image model for observations

