# Composing graphical models with neural networks for structured representations and fast inference
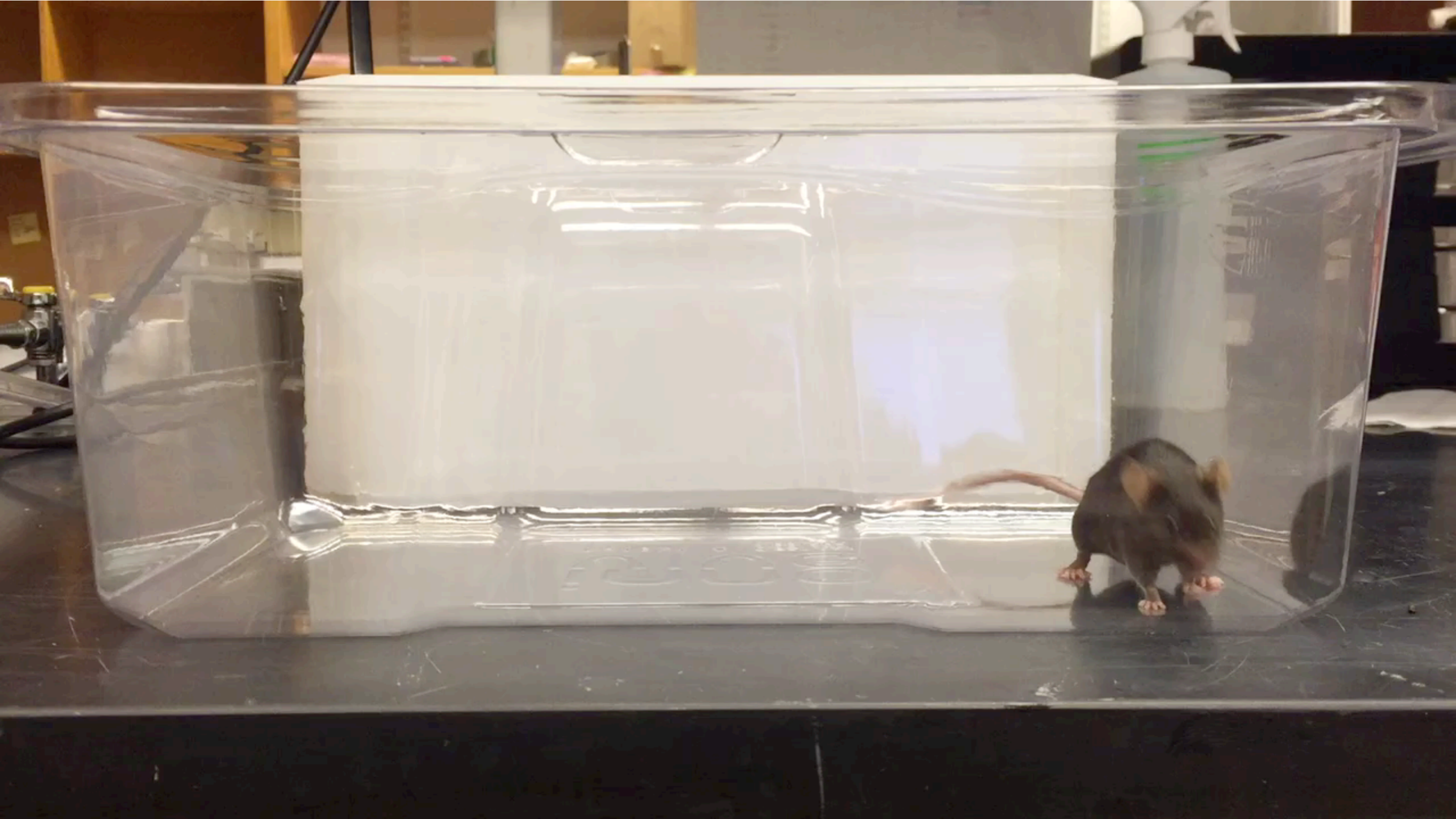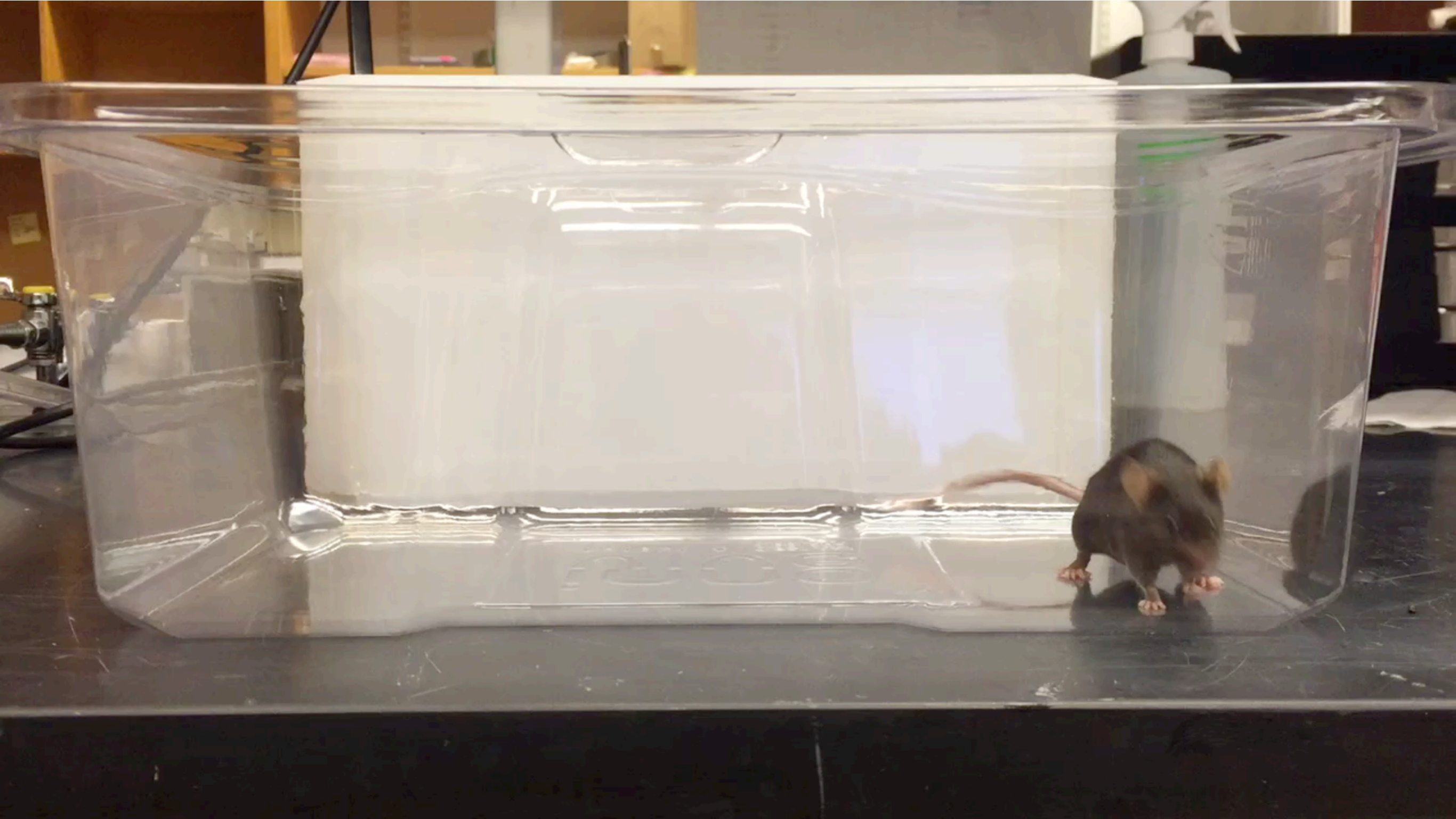
Matthew J Johnson ([mattjj@google.com](mailto:mattjj@google.com))
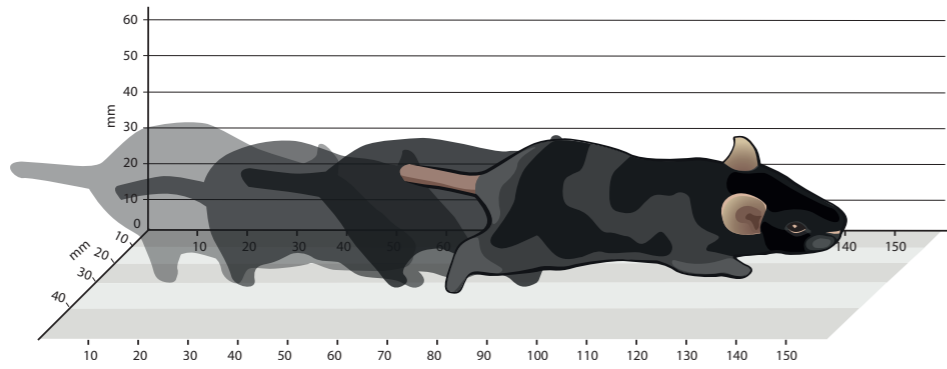Deep Learning Summer School
Montreal 2017

dart          pause          rear

dart          pause          rear

dart

pause

rear

/b/ /ax/ /n/ /ae/ /n/ /ax/
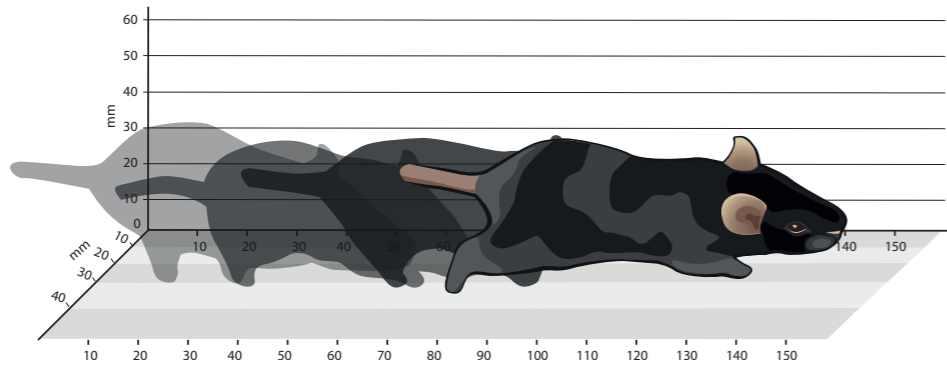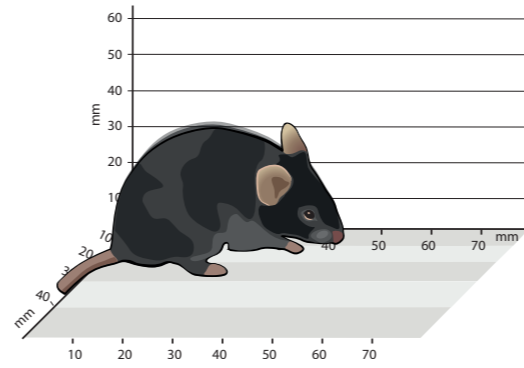
[1,2]

[1] Lee and Glass. A Nonparametric Bayesian Approach to Acoustic Model Discovery. ACL 2012.
[2] Lee. Discovering Linguistic Structures in Speech: Models and Applications. MIT Ph.D. Thesis 2014.

dart          pause          rear

/b/  /ax/  /n/  /ae/  /n/  /ax/     [1,2]
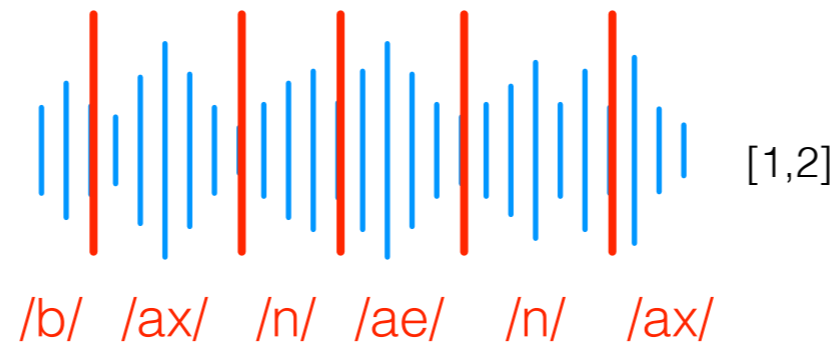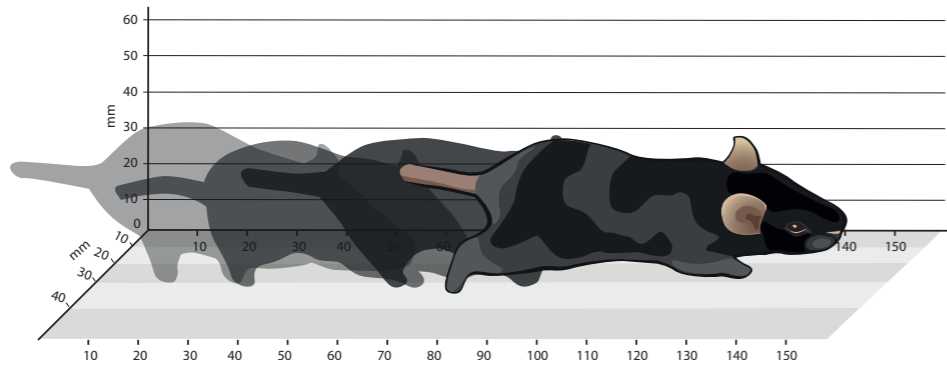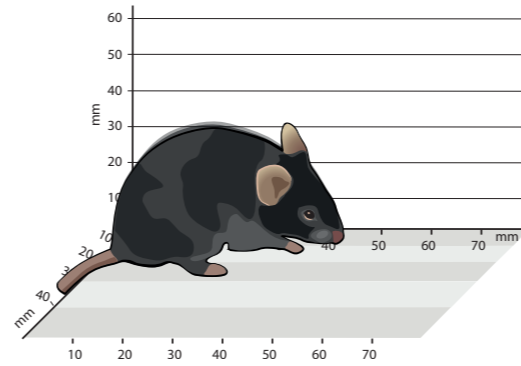
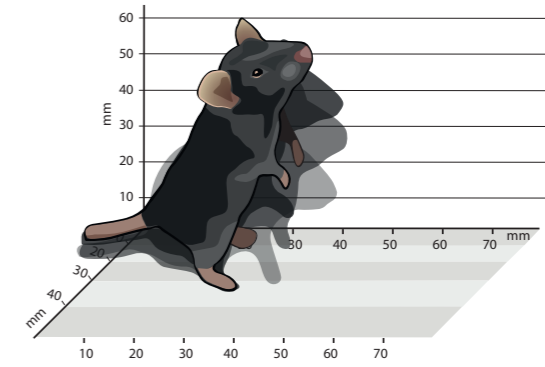[1] Lee and Glass. A Nonparametric Bayesian Approach to Acoustic Model Discovery. ACL 2012.
[2] Lee. Discovering Linguistic Structures in Speech: Models and Applications. MIT Ph.D. Thesis 2014.

Frame 0

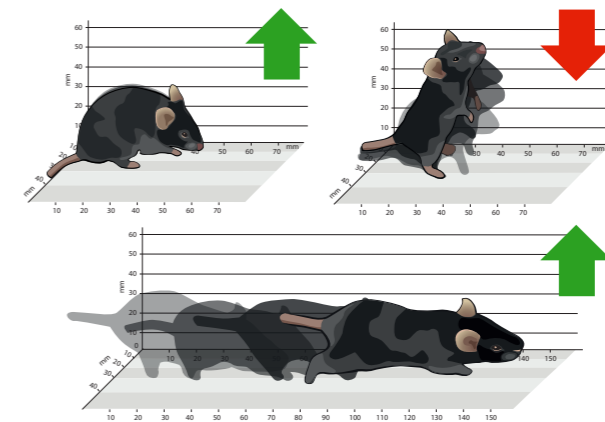Alexander Wiltschko, **Matthew Johnson**, et al., Neuron 2015.

# Frame 0

image
manifold

depth
video

image
manifold

depth
video

image
manifold

depth
video

image
manifold

depth
video

image
manifold

depth
video

image
manifold

depth
video

image
manifold

depth video

image manifold

manifold coordinates

dart    rear

# Recurrent neural networks? [1,2,3]



Figure 1. LSTM unit



Figure 2. LSTM Autoencoder Model

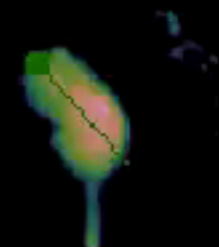[1] Srivastava, Mansimov, Salakhutdinov. Unsupervised learning of video representations using LSTMs. ICML 2015.
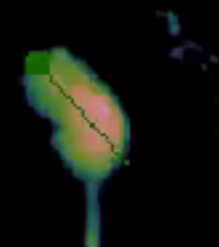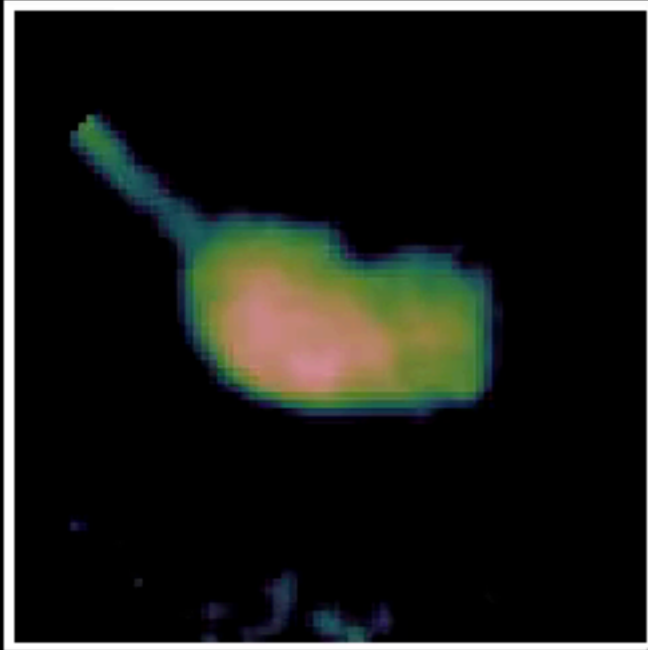[2] Ranzato, MarcAurelio, et al. Video (language) modeling: a baseline for generative models of natural videos. Preprint 2015.
[3] Sutskever, Hinton, and Taylor. The Recurrent Temporal Restricted Boltzmann Machine. NIPS 2008.

# Recurrent neural networks? [1,2,3]



Figure 1. LSTM unit



Figure 2. LSTM Autoencoder Model

# Probabilistic graphical models? [4,5,6]

[1] Srivastava, Mansimov, Salakhutdinov. Unsupervised learning of video representations using LSTMs. ICML 2015.
[2] Ranzato, MarcAurelio, et al. Video (language) modeling: a baseline for generative models of natural videos. Preprint 2015.
[3] Sutskever, Hinton, and Taylor. The Recurrent Temporal Restricted Boltzmann Machine. NIPS 2008.
[4] Fox, Sudderth, Jordan, Willsky. Bayesian nonparametric inference of switching dynamic linear models. IEEE TSP 2011.
[5] **Johnson** and Willsky. Bayesian nonparametric hidden semi-Markov models. JMLR 2013.
[6] Murphy. Machine learning: a probabilistic perspective. MIT Press 2012.

supervised
learning

supervised learning

unsupervised learning

Probabilistic graphical models                    Deep learning

# Probabilistic graphical models          Deep learning

**+** structured representations

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

## Deep learning

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

## Deep learning

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

**–** feature engineering

## Deep learning

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

**–** feature engineering

**+** arbitrary inference queries

## Deep learning

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

**–** feature engineering

**+** arbitrary inference queries

**+** data and computational efficiency
within rigid model classes

## Deep learning

## Probabilistic graphical models                    ## Deep learning

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

**–** feature engineering

**+** arbitrary inference queries

**+** data and computational efficiency
within rigid model classes

**–** more flexible models can require
slow top-down inference

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

**–** rigid assumptions may not fit

**–** feature engineering

**+** arbitrary inference queries

**+** data and computational efficiency
within rigid model classes

**–** more flexible models can require
slow top-down inference

## Deep learning

**–** neural net "goo"

**–** difficult parameterization

**+** flexible, high capacity

**+** feature learning

## Probabilistic graphical models

**+** structured representations

**+** priors and uncertainty

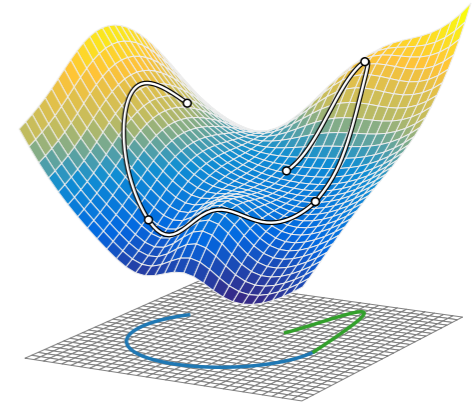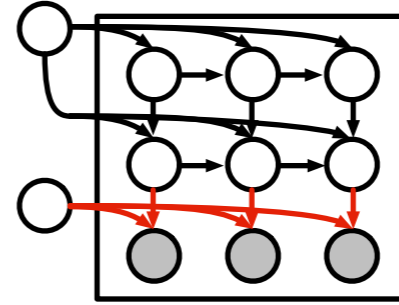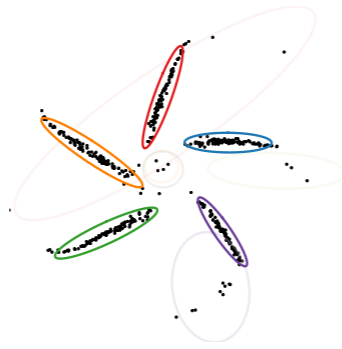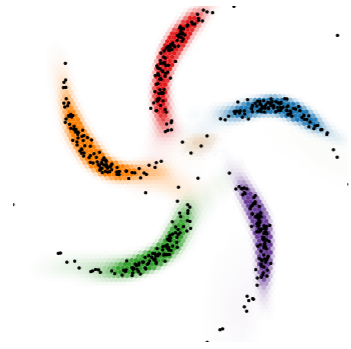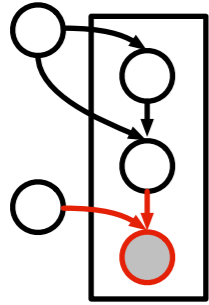**–** rigid assumptions may not fit

**–** feature engineering

**+** arbitrary inference queries

**+** data and computational efficiency
within rigid model classes

**–** more flexible models can require
slow top-down inference

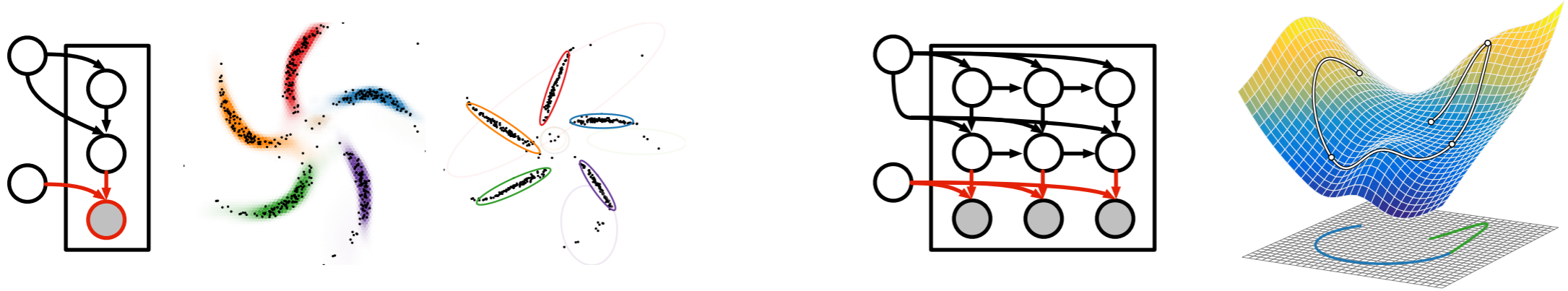## Deep learning

**–** neural net "goo"

**–** difficult parameterization

**+** flexible, high capacity

**+** feature learning

**–** limited inference queries

**–** data- and compute-hungry

**+** recognition networks learn
how to do inference

**Modeling idea:** graphical models on latent variables,
neural network models for observations

**Modeling idea:** graphical models on latent variables, neural network models for observations



**Inference:** recognition networks output conjugate potentials, then apply fast graphical model inference
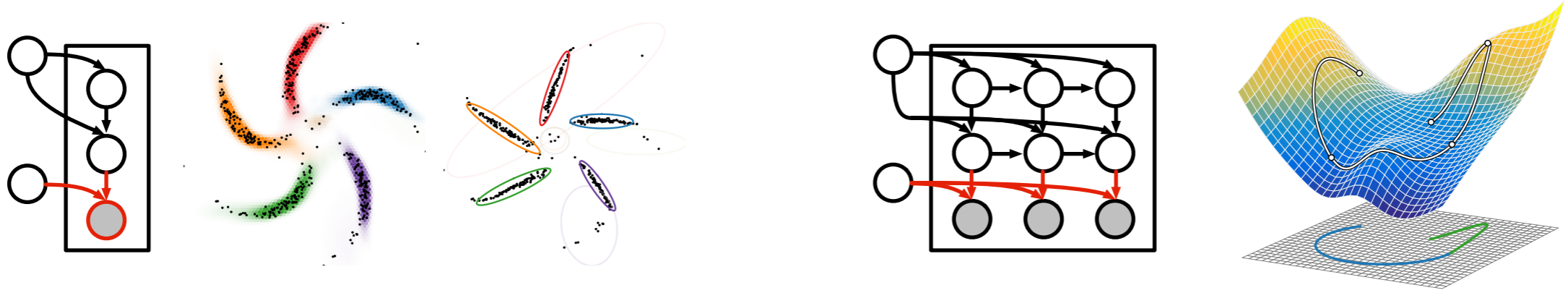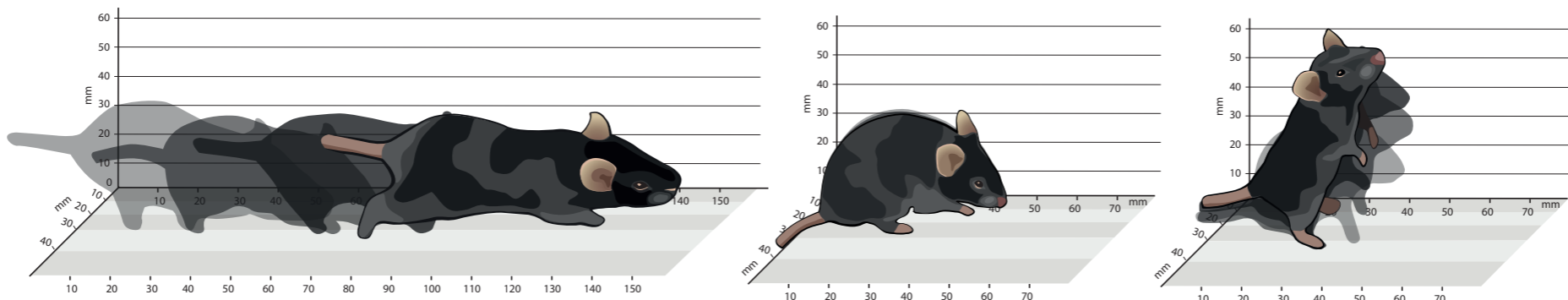
**Modeling idea:** graphical models on latent variables,
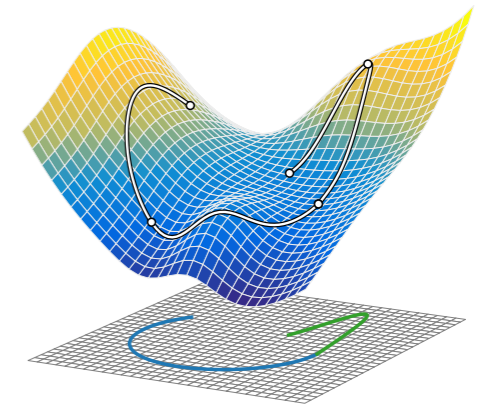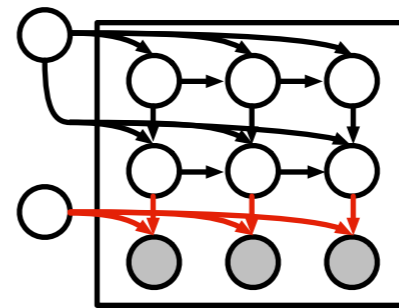neural network models for observations
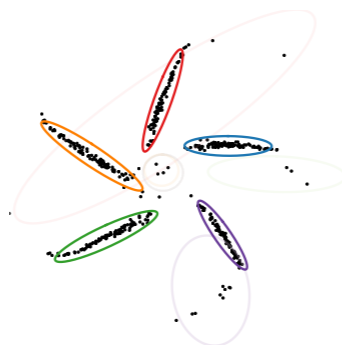


**Inference:** recognition networks output conjugate potentials,
then apply fast graphical model inference



**Application:** learn syllable representation of behavior from video

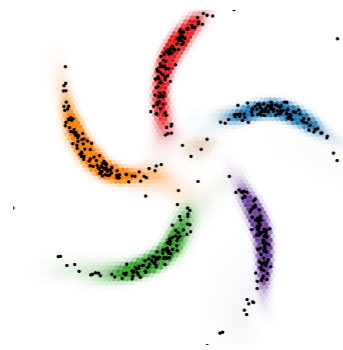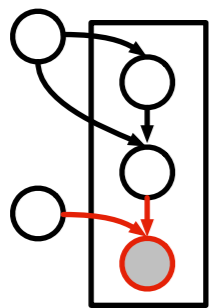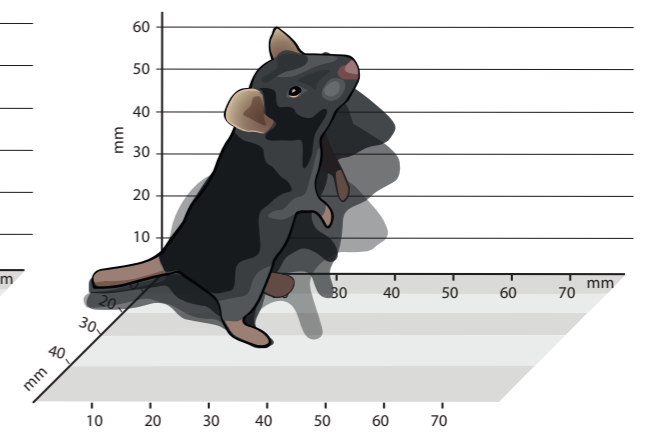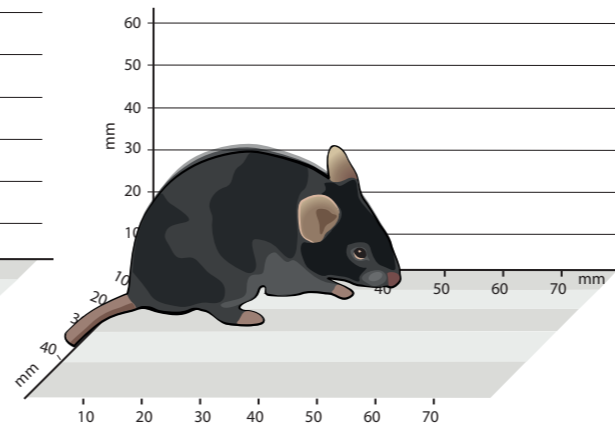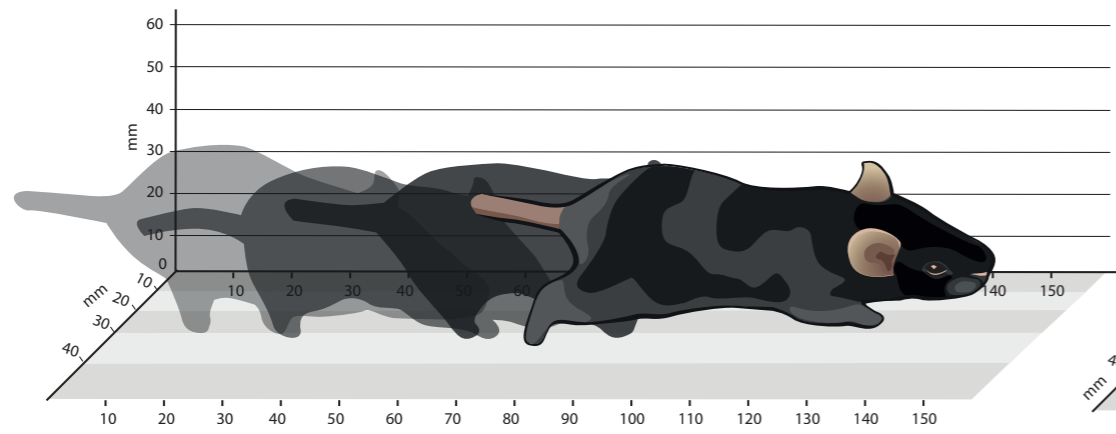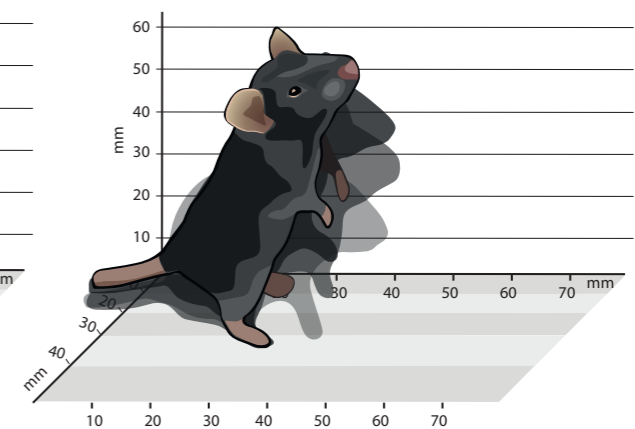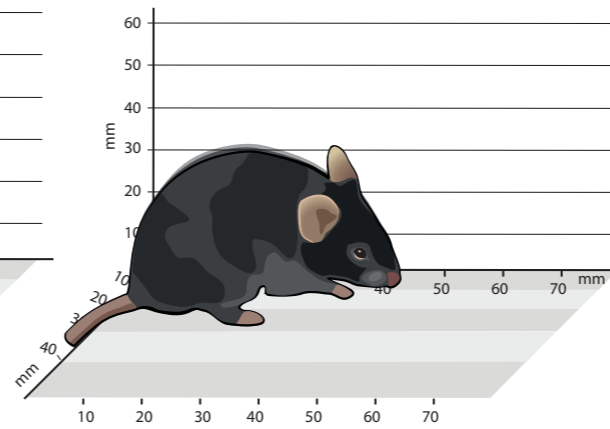**Modeling idea:** graphical models on latent variables, neural network models for observations

$$\pi = \begin{array}{c} \color{blue}\blacksquare \\ \color{red}\blacksquare \\ \color{green}\blacksquare \end{array} \left[ \begin{array}{ccc} \color{blue}\blacksquare & \color{red}\blacksquare & \color{green}\blacksquare \\ \rule{1em}{0.4pt} & \pi^{(1)} & \rule{1em}{0.4pt} \\ \rule{1em}{0.4pt} & \pi^{(2)} & \rule{1em}{0.4pt} \\ \rule{1em}{0.4pt} & \pi^{(3)} & \rule{1em}{0.4pt} \end{array} \right]$$

$$z_{t+1} \sim \pi^{(z_t)}$$

$$\pi = \begin{bmatrix} \underline{\quad} & \pi^{(1)} & \underline{\quad} \\ \underline{\quad} & \pi^{(2)} & \underline{\quad} \\ \underline{\quad} & \pi^{(3)} & \underline{\quad} \end{bmatrix}$$

$$z_{t+1} \sim \pi^{(z_t)}$$

$$A^{(1)} \quad A^{(2)} \quad A^{(3)}$$

$$B^{(1)} \quad B^{(2)} \quad B^{(3)}$$

$$x_{t+1} = A^{(z_t)} x_t + B^{(z_t)} u_t \qquad u_t \overset{\text{iid}}{\sim} \mathcal{N}(0, I)$$

$$\pi = \begin{bmatrix} \text{——} & \pi^{(1)} & \text{——} \\ \text{——} & \pi^{(2)} & \text{——} \\ \text{——} & \pi^{(3)} & \text{——} \end{bmatrix}$$

$A^{(1)}$  $A^{(2)}$  $A^{(3)}$
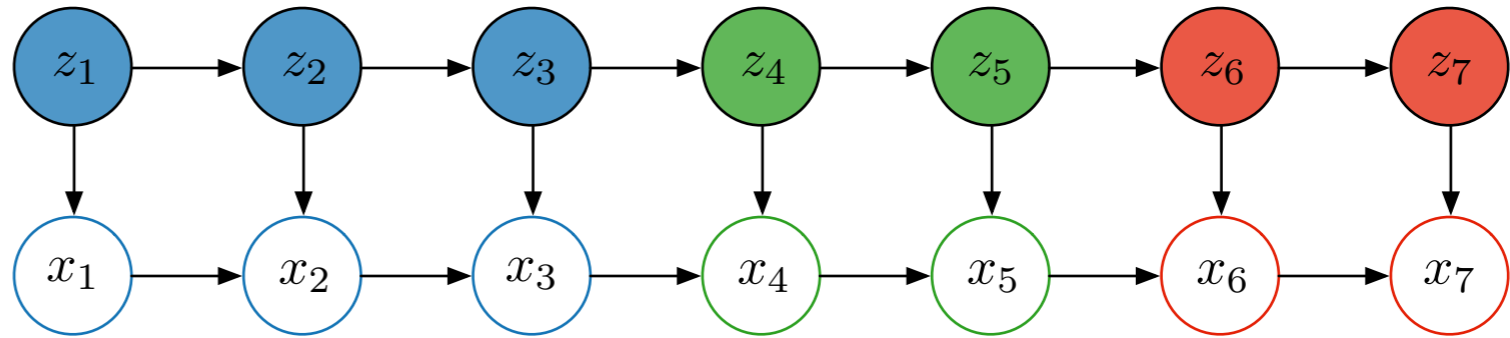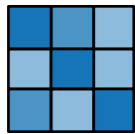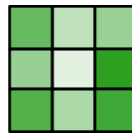
$B^{(1)}$  $B^{(2)}$  $B^{(3)}$

$z_1 \rightarrow z_2 \rightarrow z_3 \rightarrow z_4 \rightarrow z_5 \rightarrow z_6 \rightarrow z_7$

$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow x_4 \rightarrow x_5 \rightarrow x_6 \rightarrow x_7$

$$y_t \mid x_t, \gamma \; \sim \; \mathcal{N}(\mu(x_t; \gamma), \; \Sigma(x_t; \gamma))$$

$$y_t \mid x_t, \gamma \ \sim \ \mathcal{N}(\mu(x_t; \gamma), \ \Sigma(x_t; \gamma))$$

$p(\theta)$
$p(x \mid \theta)$

conjugate prior on global variables
exponential family on local variables

$p(\theta)$    conjugate prior on global variables

$p(x \mid \theta)$    exponential family on local variables

$p(\gamma)$    any prior on observation parameters

$p(y \mid x, \gamma)$    neural network observation model

Gaussian mixture model [1]  Linear dynamical system [2]  Hidden Markov model [3]  Switching LDS [4]

[1] Palmer, Wipf, Kreutz-Delgado, and Rao. Variational EM algorithms for non-Gaussian latent variable models. NIPS 2005.
[2] Ghahramani and Beal. Propagation algorithms for variational Bayesian learning. NIPS 2001.
[3] Beal. Variational algorithms for approximate Bayesian inference, Ch. 3. U of London Ph.D. Thesis 2003.
[4] Ghahramani and Hinton. Variational learning for switching state-space models. Neural Computation 2000.

Gaussian mixture model [1]

Linear dynamical system [2]

Hidden Markov model [3]

Switching LDS [4]

Mixture of Experts [5]

Driven LDS [2]

IO-HMM [6]

Factorial HMM [7]

[1] Palmer, Wipf, Kreutz-Delgado, and Rao. Variational EM algorithms for non-Gaussian latent variable models. NIPS 2005.

[2] Ghahramani and Beal. Propagation algorithms for variational Bayesian learning. NIPS 2001.

[3] Beal. Variational algorithms for approximate Bayesian inference, Ch. 3. U of London Ph.D. Thesis 2003.

[4] Ghahramani and Hinton. Variational learning for switching state-space models. Neural Computation 2000.

[5] Jordan and Jacobs. Hierarchical Mixtures of Experts and the EM algorithm. Neural Computation 1994.

[6] Bengio and Frasconi. An Input Output HMM Architecture. NIPS 1995.

[7] Ghahramani and Jordan. Factorial Hidden Markov Models. Machine Learning 1997.

Gaussian mixture model

Linear dynamical system

Hidden Markov model

Switching LDS

Mixture of Experts

Driven LDS

IO-HMM

Factorial HMM

Canonical correlations analysis

admixture / LDA / NMF

[1] Palmer, Wipf, Kreutz-Delgado, and Rao. Variational EM algorithms for non-Gaussian latent variable models. NIPS 2005.

[2] Ghahramani and Beal. Propagation algorithms for variational Bayesian learning. NIPS 2001.

[3] Beal. Variational algorithms for approximate Bayesian inference, Ch. 3. U of London Ph.D. Thesis 2003.

[4] Ghahramani and Hinton. Variational learning for switching state-space models. Neural Computation 2000.

[5] Jordan and Jacobs. Hierarchical Mixtures of Experts and the EM algorithm. Neural Computation 1994.

[6] Bengio and Frasconi. An Input Output HMM Architecture. NIPS 1995.

[7] Ghahramani and Jordan. Factorial Hidden Markov Models. Machine Learning 1997.

[8] Bach and Jordan. A probabilistic interpretation of Canonical Correlation Analysis. Tech. Report 2005.

[9] Archambeau and Bach. Sparse probabilistic projections. NIPS 2008.

[10] Hoffman, Bach, Blei. Online learning for Latent Dirichlet Allocation. NIPS 2010.

**Inference?**

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

Natural gradient SVI
for nice exp. fam. PGMs [1,2]

[1] Hoffman, Bach, Blei. Online learning for Latent Dirichlet Allocation. NIPS 2010.
[2] Hoffman, Blei, Wang, and Paisley. Stochastic variational inference. JMLR 2013.

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$p(x \,|\, \theta)$ is a linear dynamical system
$p(y \,|\, x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$$q(\theta)q(x) \approx p(\theta, x \,|\, y)$$

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$$q(\theta)q(x) \approx p(\theta, x \mid y)$$

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

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$$q(\theta)q(x) \approx p(\theta, x \mid y)$$

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

$$\eta_x^*(\eta_\theta) \triangleq \arg\max_{\eta_x} \mathcal{L}(\eta_\theta, \eta_x) \qquad \mathcal{L}_{\mathrm{SVI}}(\eta_\theta) \triangleq \mathcal{L}(\eta_\theta, \eta_x^*(\eta_\theta))$$

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$$q(\theta)q(x) \approx p(\theta, x \mid y)$$

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

$$\eta_x^*(\eta_\theta) \triangleq \arg\max_{\eta_x} \mathcal{L}(\eta_\theta, \eta_x) \qquad \mathcal{L}_{\text{SVI}}(\eta_\theta) \triangleq \mathcal{L}(\eta_\theta, \eta_x^*(\eta_\theta))$$

**Proposition (natural gradient SVI of Hoffman et al. 2013)**

$$\widetilde{\nabla} \mathcal{L}_{\text{SVI}}(\eta_\theta) = \eta_\theta^0 + \mathbb{E}_{q^*(x)}(t_{xy}(x, y), 1) - \eta_\theta$$

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \theta)$ is a linear-Gaussian observation
$p(\theta)$ is a conjugate prior

$$q(\theta)q(x) \approx p(\theta, x \mid y)$$

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

$$\eta_x^*(\eta_\theta) \triangleq \arg\max_{\eta_x} \mathcal{L}(\eta_\theta, \eta_x) \qquad \mathcal{L}_{\mathrm{SVI}}(\eta_\theta) \triangleq \mathcal{L}(\eta_\theta, \eta_x^*(\eta_\theta))$$

**Proposition (natural gradient SVI of Hoffman et al. 2013)**

$$\widetilde{\nabla}\mathcal{L}_{\mathrm{SVI}}(\eta_\theta) = \eta_\theta^0 + \sum_{n=1}^{N} \mathbb{E}_{q^*(x_n)}(t_{xy}(x_n, y_n), 1) - \eta_\theta$$

# Step 1: compute evidence potentials

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

# Step 1: compute evidence potentials

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

Step 1: compute evidence potentials

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

Step 1: compute evidence potentials

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

Step 1: compute evidence potentials

Step 2: run fast message passing

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

Step 1: compute evidence potentials

Step 2: run fast message passing

[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

## Step 1: compute evidence potentials

## Step 2: run fast message passing
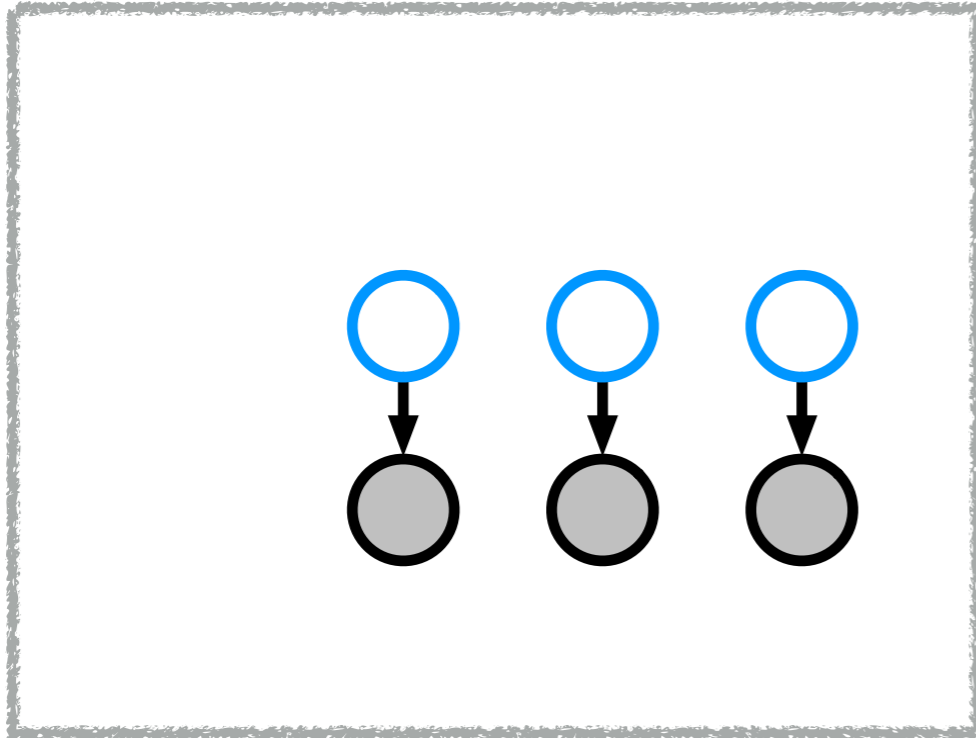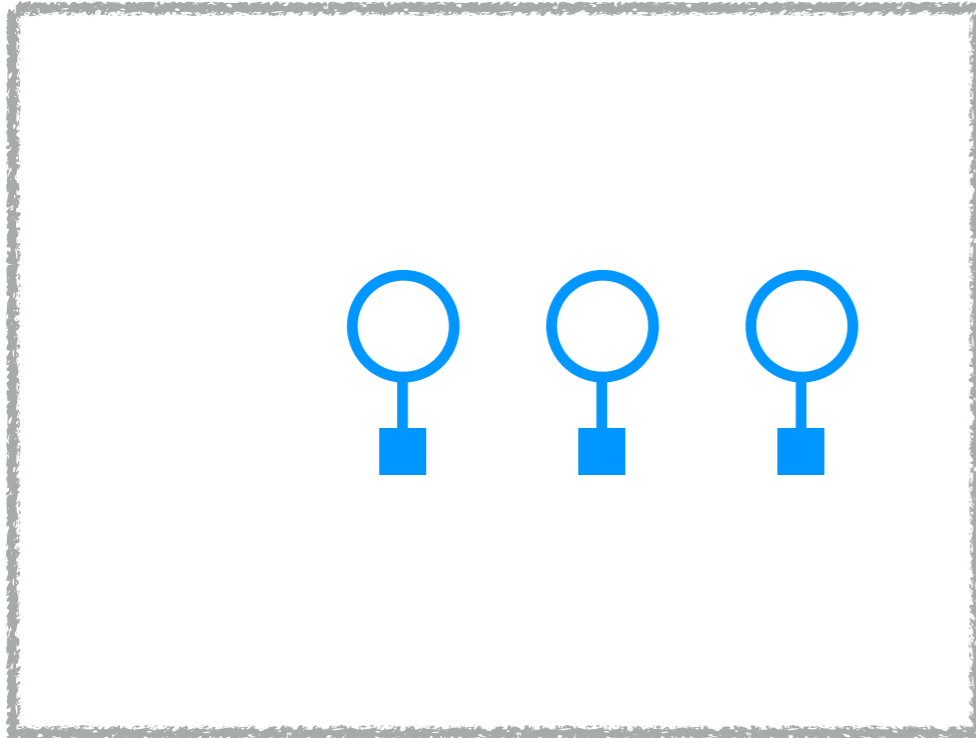
## Step 3: compute natural gradient
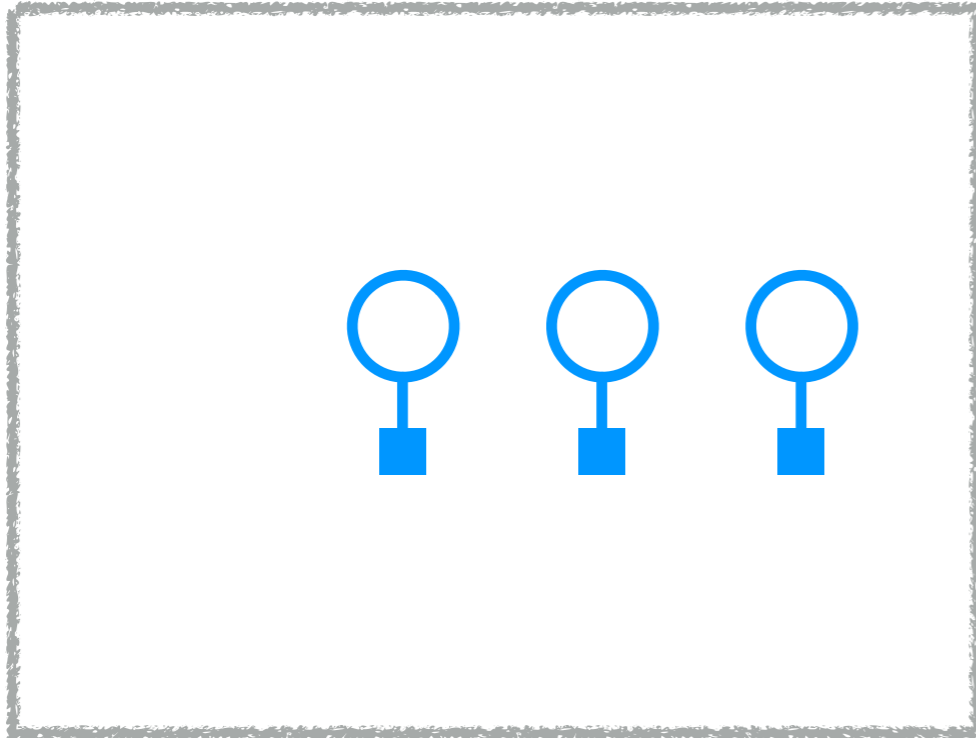
[1] **Johnson** and Willsky. Stochastic variational inference for Bayesian time series models. ICML 2014.
[2] Foti, Xu, Laird, and Fox. Stochastic variational inference for hidden Markov models. NIPS 2014.

arbitrary inference queries

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \gamma)$ is a neural network decoder
$p(\theta)$ is a conjugate prior, $p(\gamma)$ is generic

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \gamma)$ is a neural network decoder
$p(\theta)$ is a conjugate prior, $p(\gamma)$ is generic

$$q(\theta)q(\gamma)q(x) \approx p(\theta, \gamma, x \mid y)$$

$p(x \mid \theta)$ is a linear dynamical system
$p(y \mid x, \gamma)$ is a neural network decoder
$p(\theta)$ is a conjugate prior, $p(\gamma)$ is generic

$$q(\theta)q(\gamma)q(x) \approx p(\theta, \gamma, x \mid y)$$

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

$p(x \mid \theta)$ is a linear dynamical system
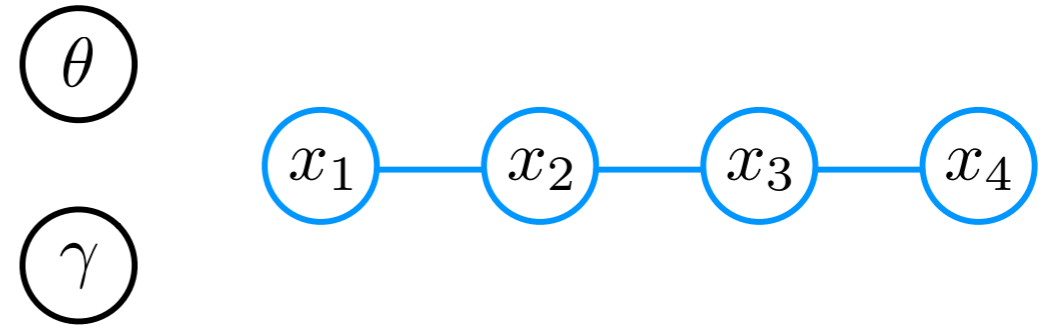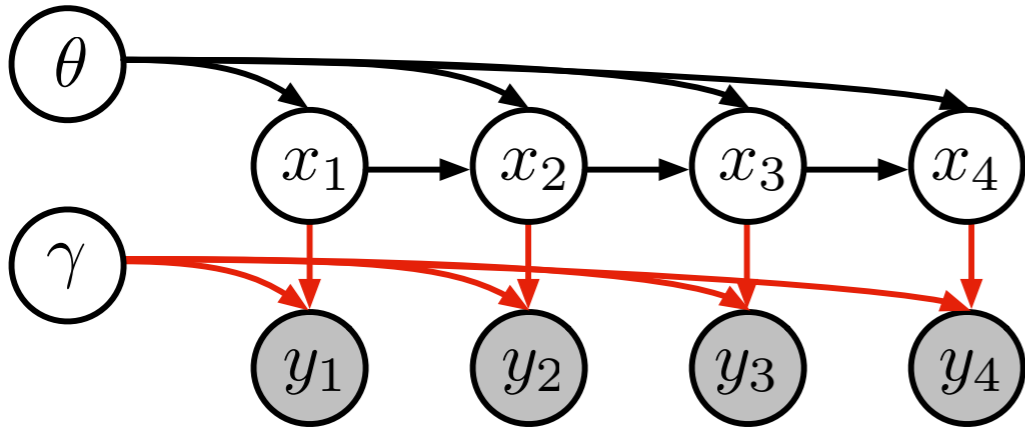$p(y \mid x, \gamma)$ is a neural network decoder
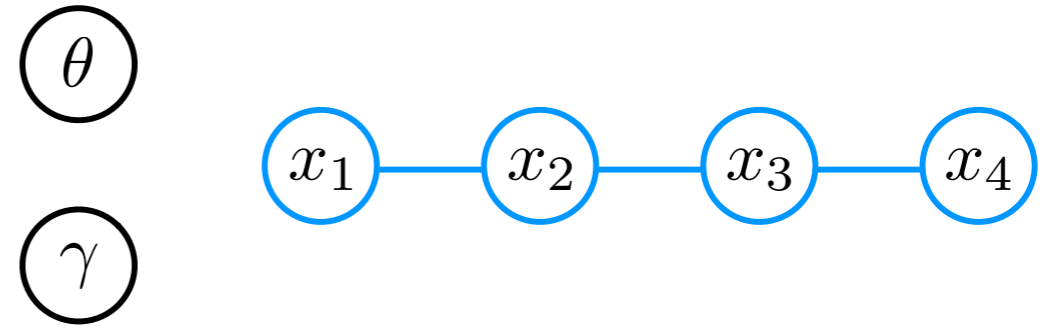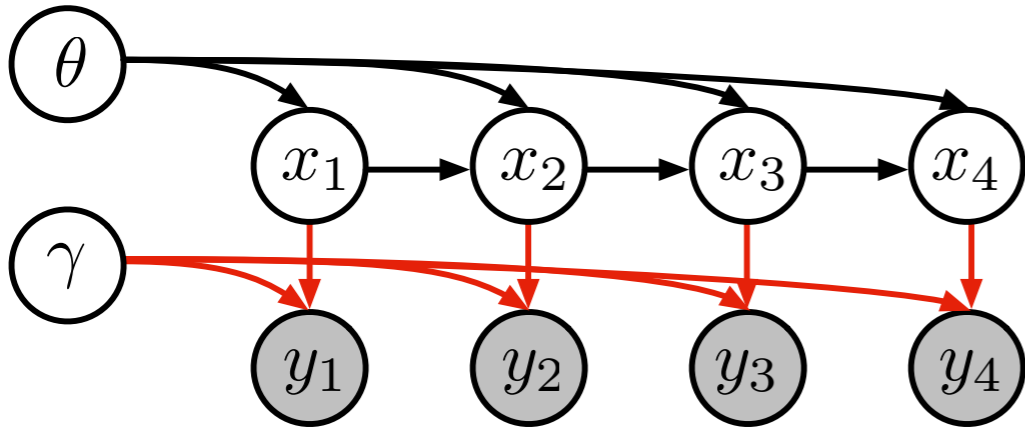$p(\theta)$ is a conjugate prior, $p(\gamma)$ is generic

$$q(\theta)q(\gamma)q(x) \approx p(\theta, \gamma, x \mid y)$$

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

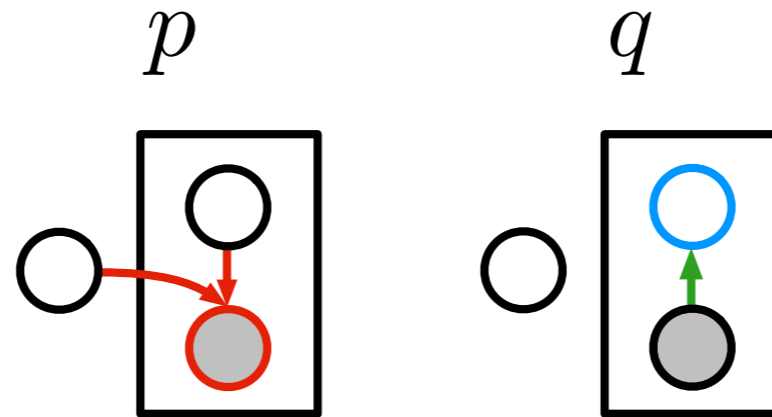$$\eta_x^\star(\eta_\theta, \eta_\gamma) \triangleq \arg\max_{\eta_x} \mathcal{L}(\eta_\theta, \eta_\gamma, \eta_x)$$

$$\mathcal{L}_{\mathrm{SVI}}(\eta_\theta, \eta_\gamma) \triangleq \mathcal{L}(\eta_\theta, \eta_\gamma, \eta_x^\star(\eta_\theta, \eta_\gamma))$$

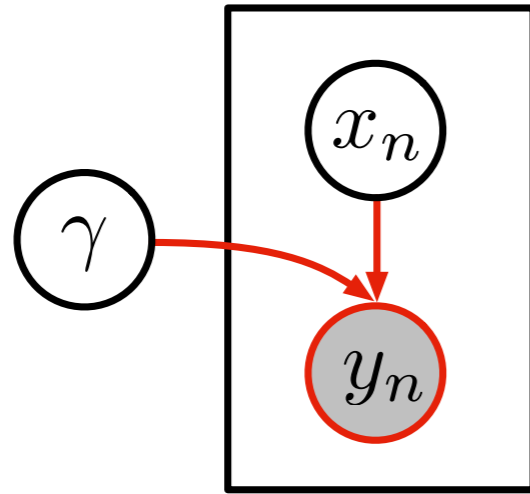$$q^*(x) \triangleq \mathcal{N}(x \mid \mu(y; \phi), \Sigma(y; \phi))$$

Variational autoencoders
and amortized inference [1,2]

[1] Kingma and Welling. Auto-encoding variational Bayes. ICLR 2014.
[2] Rezende, Mohamed, and Wierstra. Stochastic backpropagation and approximate inference in deep generative models. ICML 2014

$$q^{\star}(x_n) \triangleq \mathcal{N}(x_n \mid \mu(y_n; \phi), \Sigma(y_n; \phi))$$

$$q^{\star}(x_n) \triangleq \mathcal{N}(x_n \mid \mu(y_n; \phi), \Sigma(y_n; \phi))$$

$$q^{\star}(x_n) \triangleq \mathcal{N}(x_n \mid \mu(y_n; \phi), \Sigma(y_n; \phi))$$

$$\mathcal{L}_{\text{VAE}}(\eta_\gamma, \phi) \triangleq \mathcal{L}(\eta_\gamma, \eta_x^{\star}(\phi))$$

$$\mu_t(y_t; \phi_\mu)$$

$$J_{t,t}(y_t; \phi_D)$$

$$J_{t,t+1}(y_t, y_{t+1}; \phi_B)$$

[1,2]

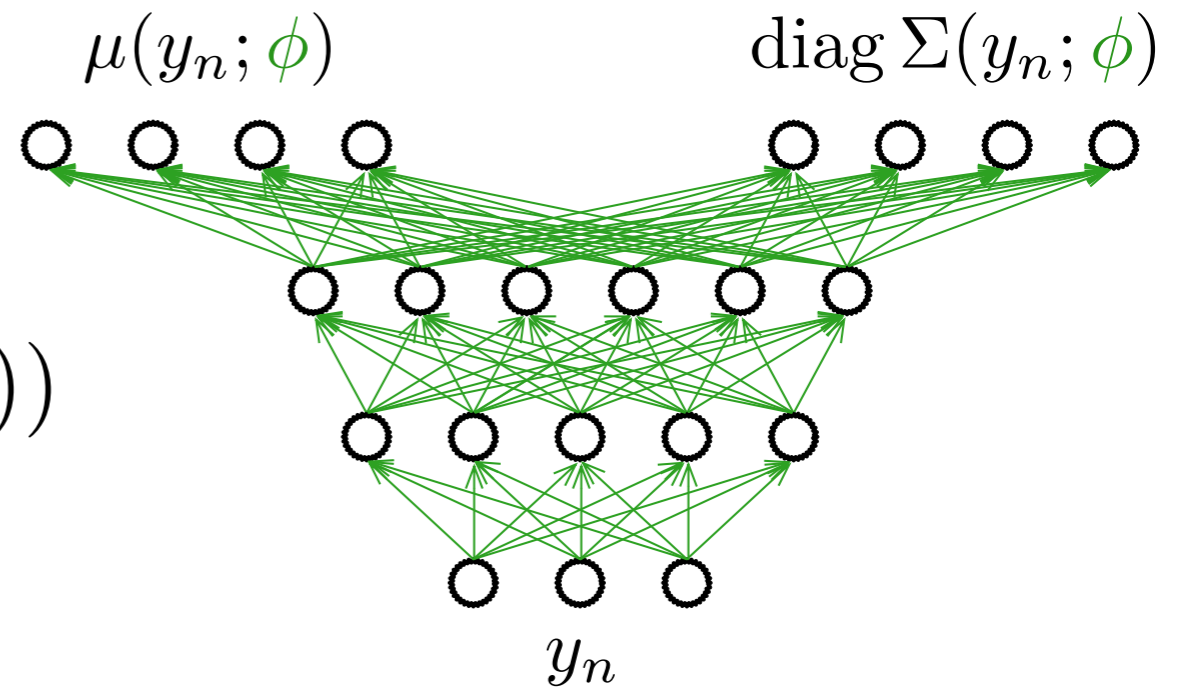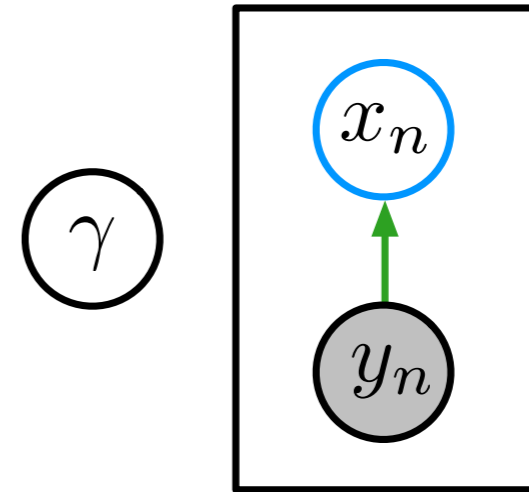[1] Archer, Park, Buesing, Cunningham, Paninski. Black box variational inference for state space models. ICLR 2016 Workshops.
[2] Gao*, Archer*, Paninski, Cunningham. Linear dynamical neural population models through nonlinear embeddings. NIPS 2016.

$$\mu_t(y_t; \phi_\mu)$$
$$J_{t,t}(y_t; \phi_D)$$
$$J_{t,t+1}(y_t, y_{t+1}; \phi_B)$$

[1,2]

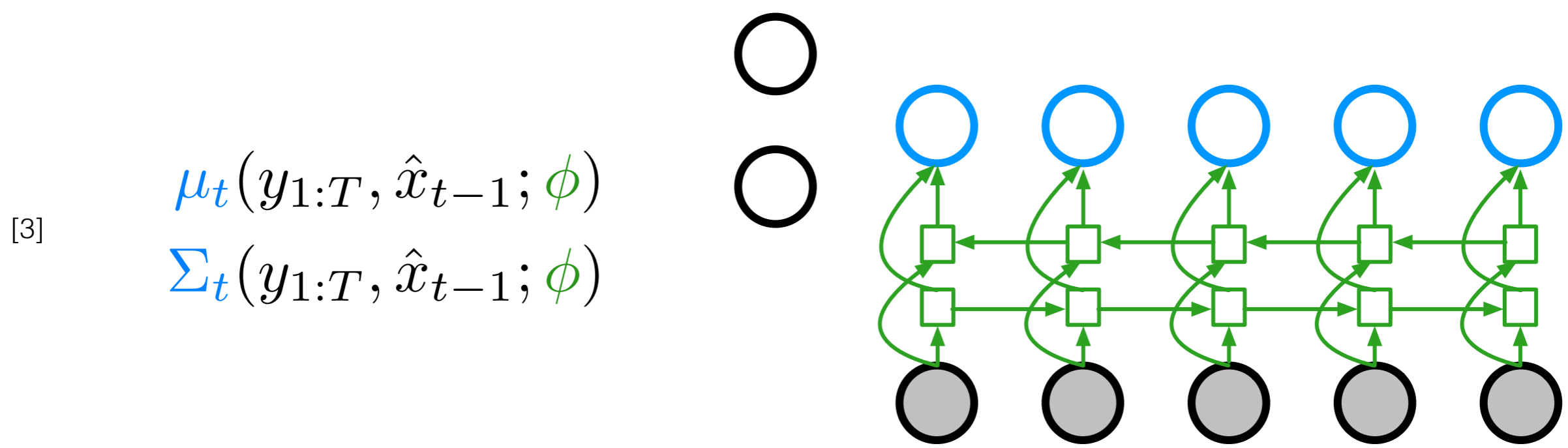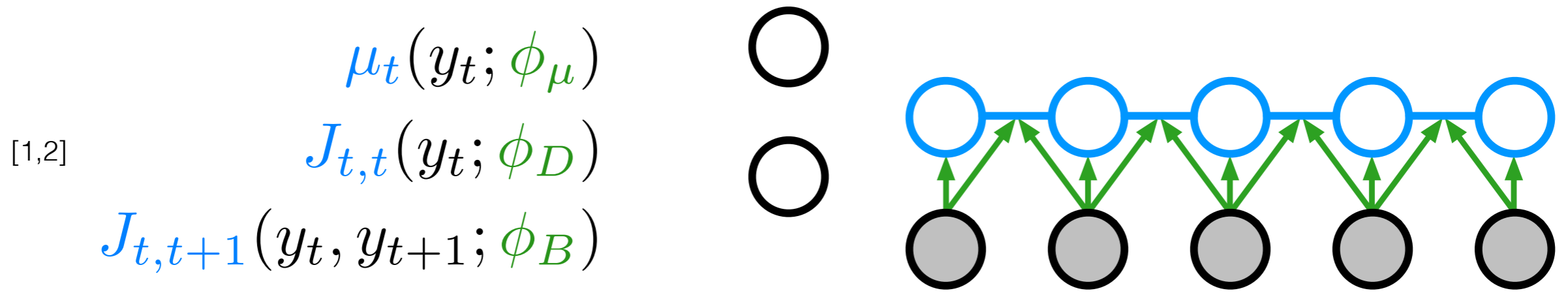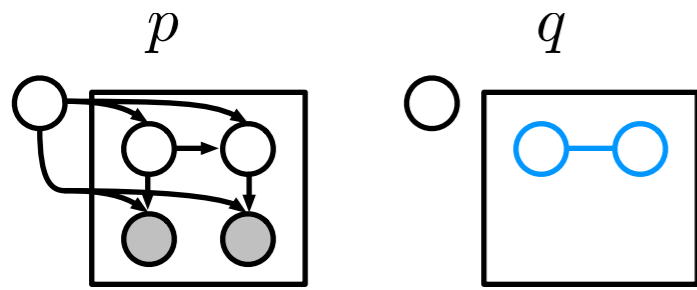$$\mu_t(y_{1:T}, \hat{x}_{t-1}; \phi)$$
$$\Sigma_t(y_{1:T}, \hat{x}_{t-1}; \phi)$$

[3]

[1] Archer, Park, Buesing, Cunningham, Paninski. Black box variational inference for state space models. ICLR 2016 Workshops.
[2] Gao*, Archer*, Paninski, Cunningham. Linear dynamical neural population models through nonlinear embeddings. NIPS 2016.
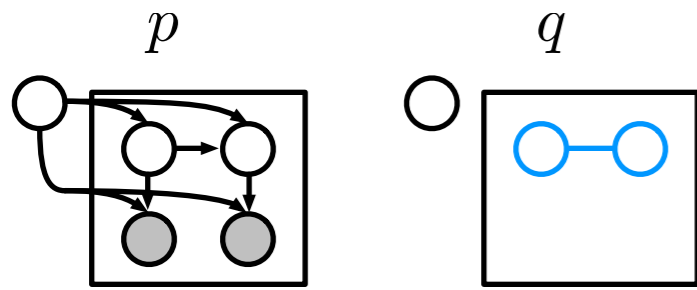[3] Krishnan, Shalit, Sontag. Structured inference networks for nonlinear state space models. AISTATS 2017.

$$q^*(x) \triangleq \underset{q(x)}{\arg\max}\, \mathcal{L}[\, q(\theta)q(x)\,]$$
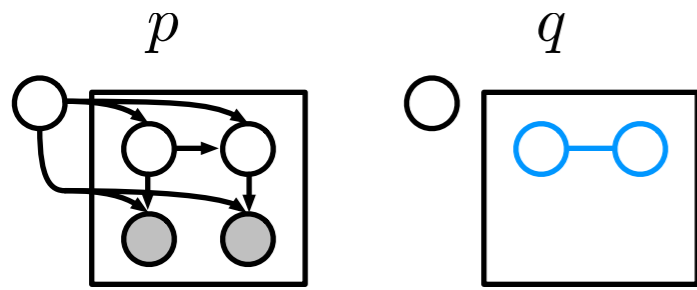
Natural gradient SVI

$p$ $q$

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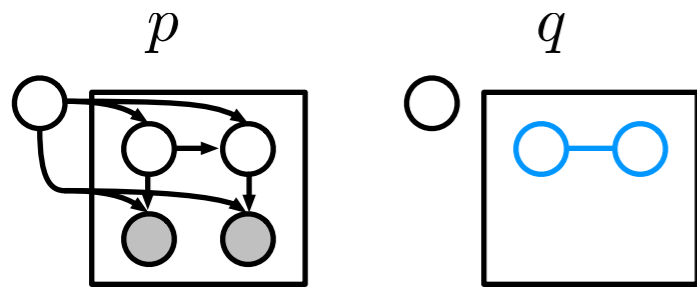
Natural gradient SVI

— expensive for general obs.

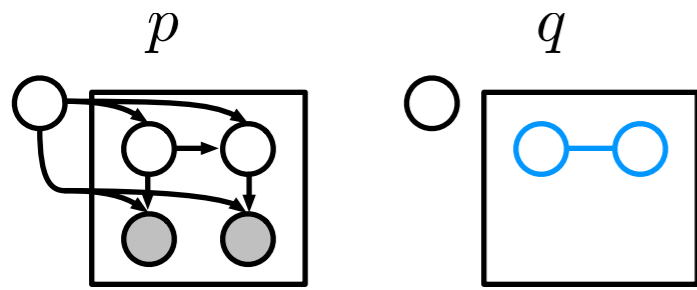$$q^*(x) \triangleq \arg\max_{q(x)} \mathcal{L}[\, q(\theta)q(x)\,]$$

Natural gradient SVI

**−** expensive for general obs.

**+** optimal local factor

$p$ $q$

$$q^*(x) \triangleq \underset{q(x)}{\arg\max}\, \mathcal{L}[\,q(\theta)q(x)\,]$$
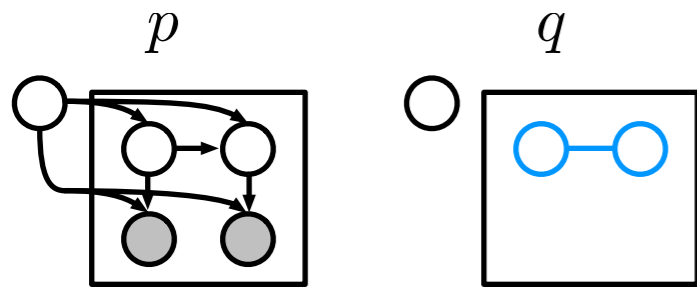
Natural gradient SVI

− expensive for general obs.

+ optimal local factor

+ exploits conj. graph structure

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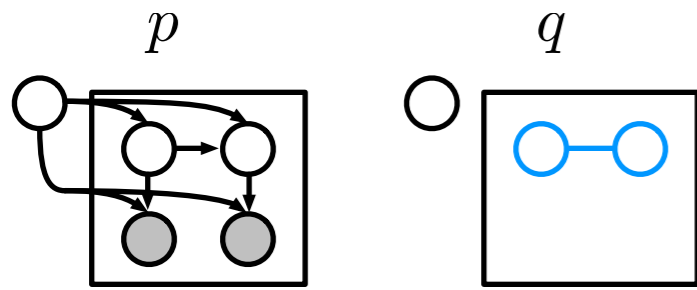
Natural gradient SVI

– expensive for general obs.

+ optimal local factor

+ exploits conj. graph structure

+ arbitrary inference queries

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

Natural gradient SVI
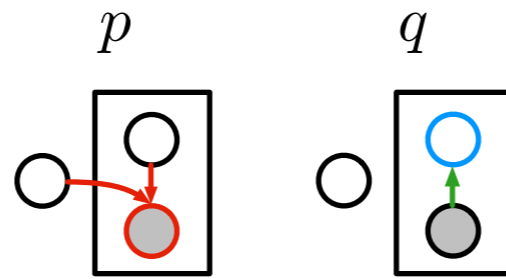
− expensive for general obs.

+ optimal local factor

+ exploits conj. graph structure

+ arbitrary inference queries

+ natural gradients

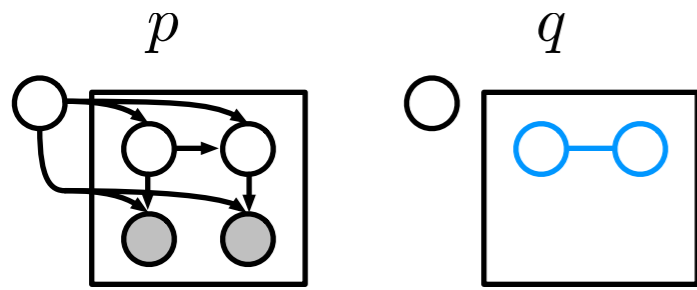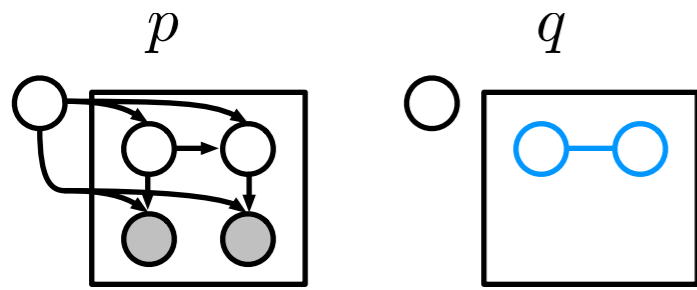$$q^*(x) \triangleq \arg\max_{q(x)} \mathcal{L}[\, q(\theta)q(x)\,]$$

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

Natural gradient SVI      Variational autoencoders

**−** expensive for general obs.

**+** optimal local factor

**+** exploits conj. graph structure

**+** arbitrary inference queries

**+** natural gradients

$$q^*(x) \triangleq \underset{q(x)}{\arg\max} \, \mathcal{L}[\, q(\theta)q(x) \,]$$
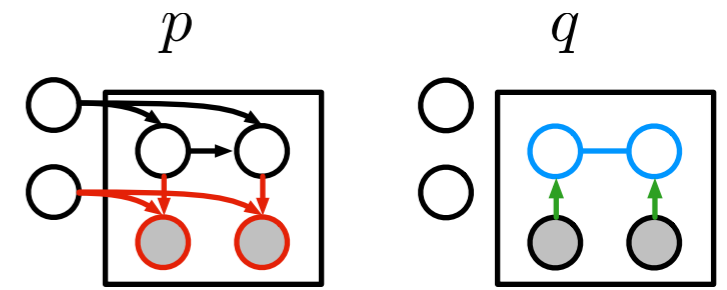
$$q^*(x) \triangleq \mathcal{N}(x \mid \mu(y; \phi), \Sigma(y; \phi))$$

Natural gradient SVI

Variational autoencoders

− expensive for general obs.

+ fast for general obs.

+ optimal local factor

− suboptimal local inference

+ exploits conj. graph structure

− $\phi$ does all local inference

+ arbitrary inference queries

− limited inference queries

+ natural gradients

− no cheap natural gradients

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

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

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

Natural gradient SVI

Variational autoencoders

Structured VAEs [1]
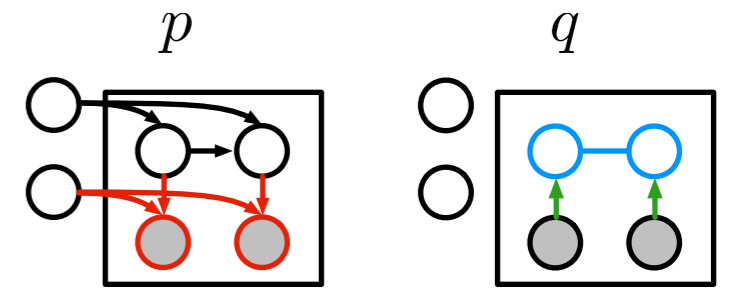
− expensive for general obs.

+ optimal local factor

+ exploits conj. graph structure

+ arbitrary inference queries

+ natural gradients

+ fast for general obs.

− suboptimal local inference

− $\phi$ does all local inference
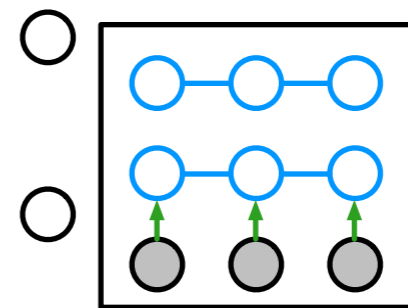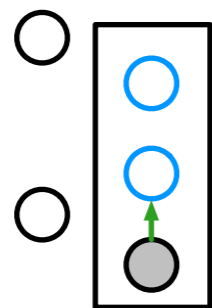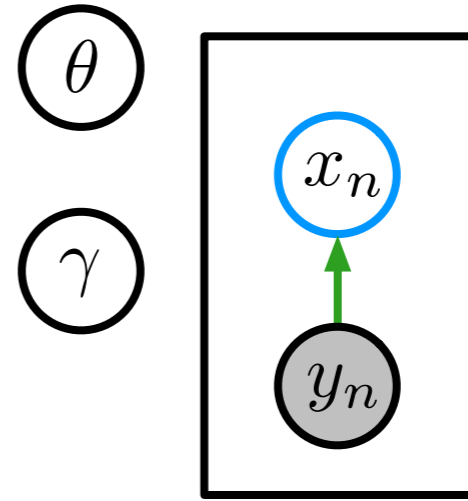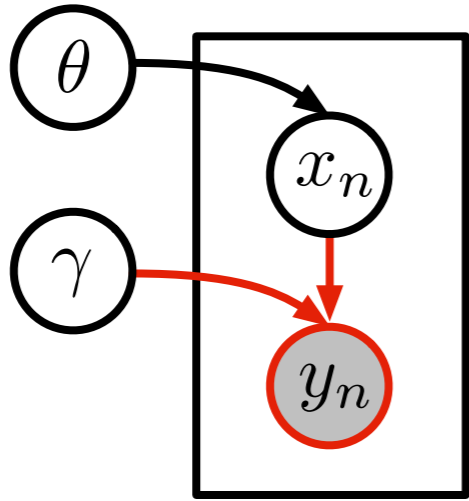
− limited inference queries

− no cheap natural gradients

[1] **Johnson**, Duvenaud, Wiltschko, Datta, and Adams. Composing graphical models and neural networks. NIPS 2016.

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

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

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

| Natural gradient SVI | Variational autoencoders | Structured VAEs [1] |
|---|---|---|
| − expensive for general obs. | + fast for general obs. | + fast for general obs. |
| + optimal local factor | − suboptimal local inference | ± optimal given conj. evidence |
| + exploits conj. graph structure | − $\phi$ does all local inference | + exploits conj. graph structure |
| + arbitrary inference queries | − limited inference queries | + arbitrary inference queries |
| + natural gradients | − no cheap natural gradients | + natural gradients on $\eta_\theta$ |

[1] **Johnson**, Duvenaud, Wiltschko, Datta, and Adams. Composing graphical models and neural networks. NIPS 2016.
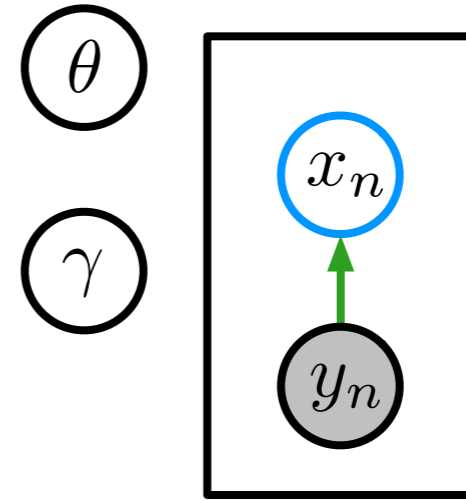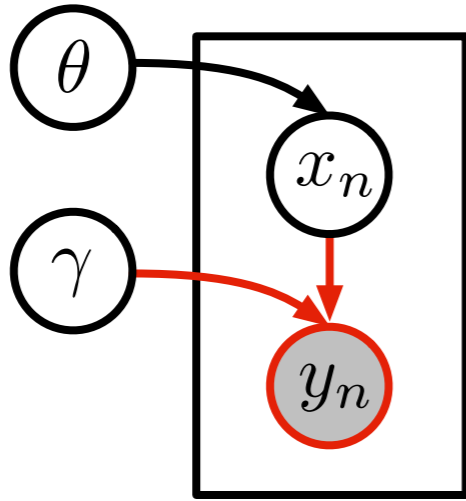
**Inference:** recognition networks output conjugate potentials, then apply fast graphical model inference
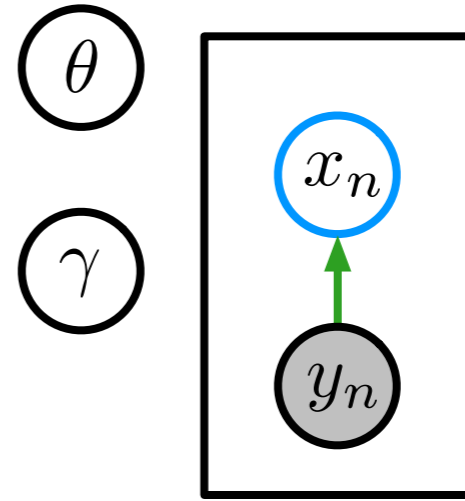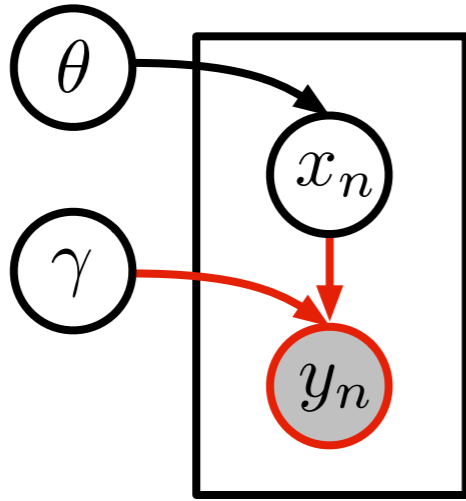
$$\mathcal{L}\big[\, q(\theta)q(\gamma)q(x) \,\big] \triangleq \mathbb{E}_{q(\theta)q(\gamma)q(x)}\left[\log \frac{p(\theta,\gamma,x)\,p(y \mid x,\gamma)}{q(\theta)q(\gamma)q(x)}\right]$$
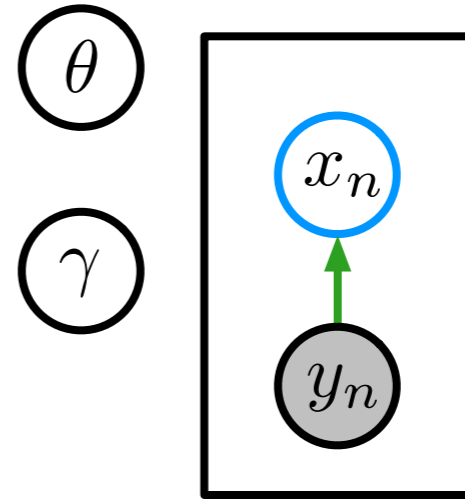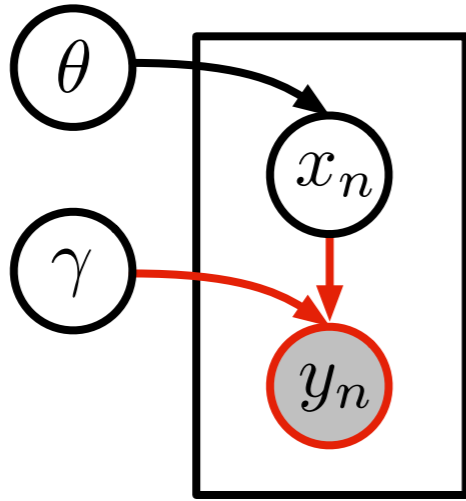
$$\mathcal{L}\big[\,q(\theta)q(\gamma)\textcolor{blue}{q(x)}\,\big] \triangleq \mathbb{E}_{q(\theta)q(\gamma)\textcolor{blue}{q(x)}}\left[\log \frac{p(\theta,\gamma,x)\textcolor{red}{p(y\,|\,x,\gamma)}}{q(\theta)q(\gamma)\textcolor{blue}{q(x)}}\right]$$

$$q(\theta) \leftrightarrow \eta_\theta \qquad q(\gamma) \leftrightarrow \eta_\gamma \qquad \textcolor{blue}{q(x) \leftrightarrow \eta_x}$$
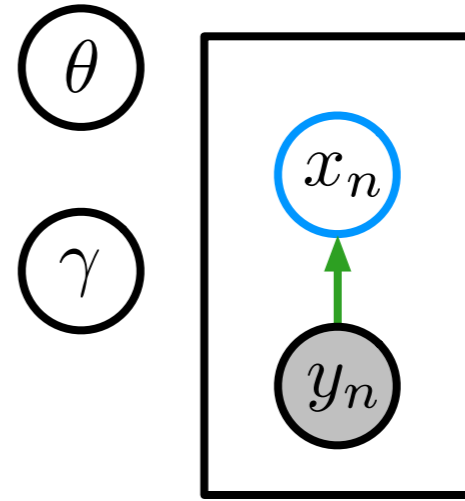
$$\mathcal{L}(\eta_\theta, \eta_\gamma, {\color{blue}\eta_x}) \triangleq \mathbb{E}_{q(\theta)q(\gamma){\color{blue}q(x)}} \left[ \log \frac{p(\theta,\gamma,x){\color{red}p(y \mid x,\gamma)}}{q(\theta)q(\gamma){\color{blue}q(x)}} \right]$$

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

$$\textcolor{red}{\mathbb{E}_{q(\gamma)} \log p(y_t \mid x_t, \gamma)}$$

$x_t$

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

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

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

$$\mathbb{E}_{q(\gamma)} \log p(y_t \mid x_t, \gamma)$$

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

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

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

$$\mathbb{E}_{q(\gamma)} \log p(y_t \mid x_t, \gamma)$$

$$\psi(x_t; y_t, \phi)$$

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

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

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

$$\eta_x^*(\eta_\theta, \phi) \triangleq \arg\max_{\eta_x} \widehat{\mathcal{L}}(\eta_\theta, \eta_x, \phi) \qquad \mathcal{L}_{\mathrm{SVAE}}(\eta_\theta, \eta_\gamma, \phi) \triangleq \mathcal{L}(\eta_\theta, \eta_\gamma, \eta_x^*(\eta_\theta, \phi))$$

# Step 1: apply recognition network

Step 1: apply recognition network

Step 1: apply recognition network

Step 1: apply recognition network

Step 1: apply recognition network

Step 2: run fast PGM algorithms

Step 1: apply recognition network

Step 2: run fast PGM algorithms

Step 1: apply recognition network



Step 2: run fast PGM algorithms



Step 3: sample, compute flat grads

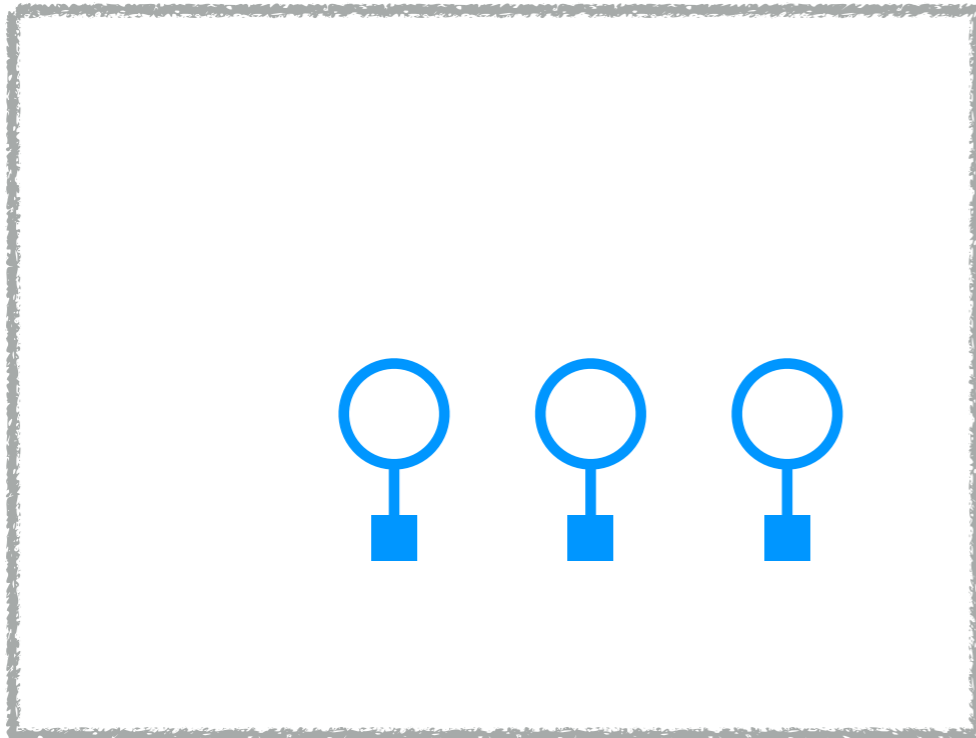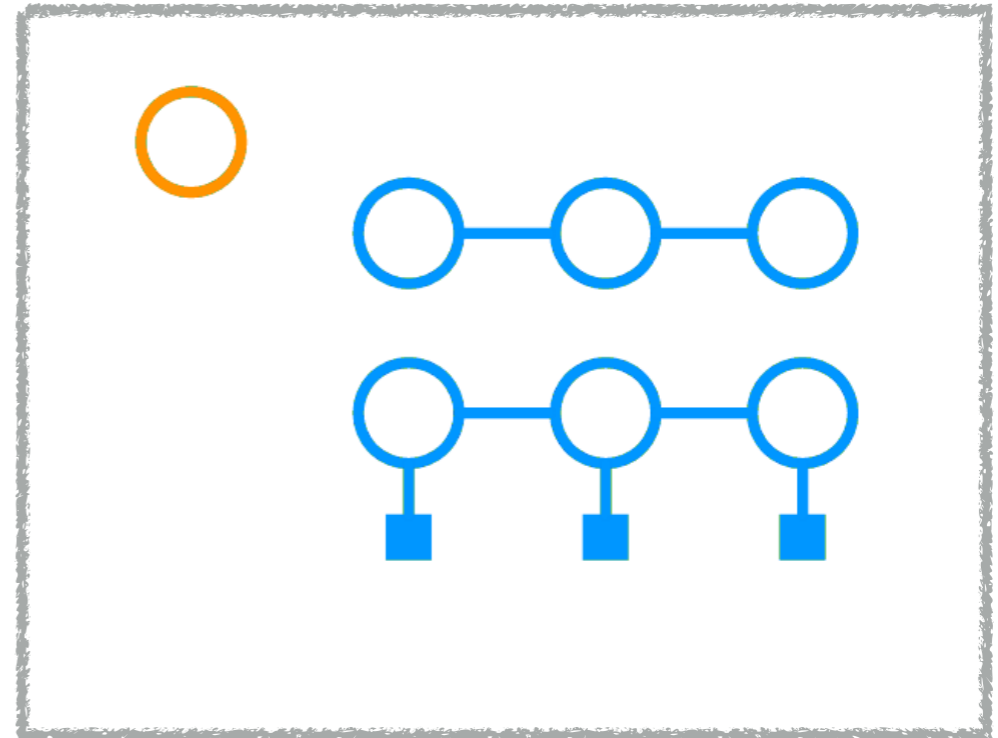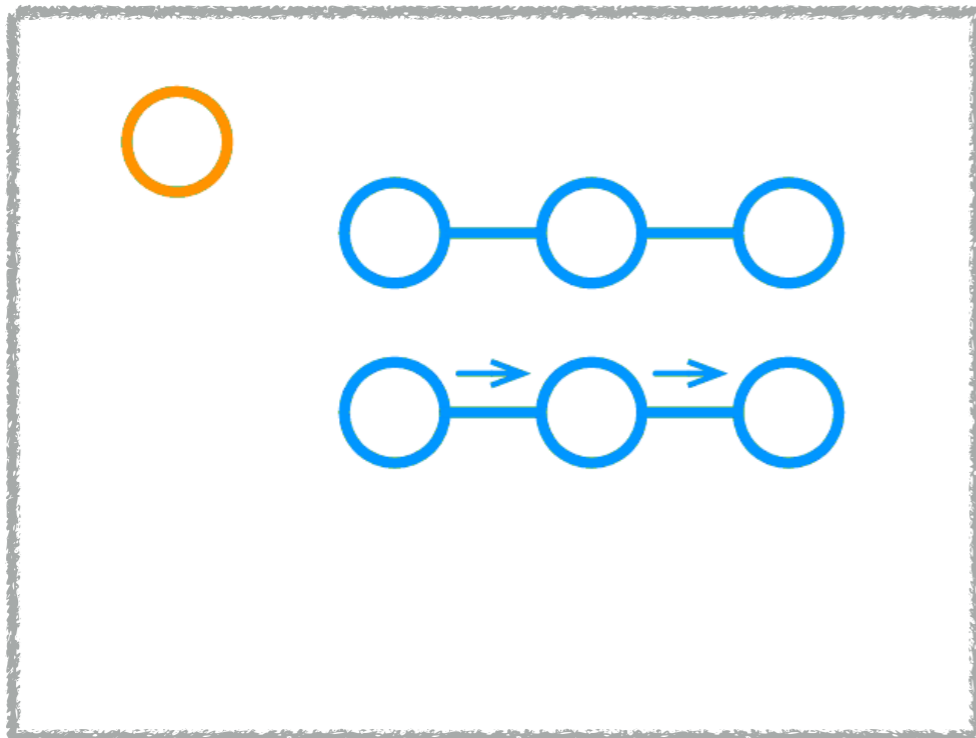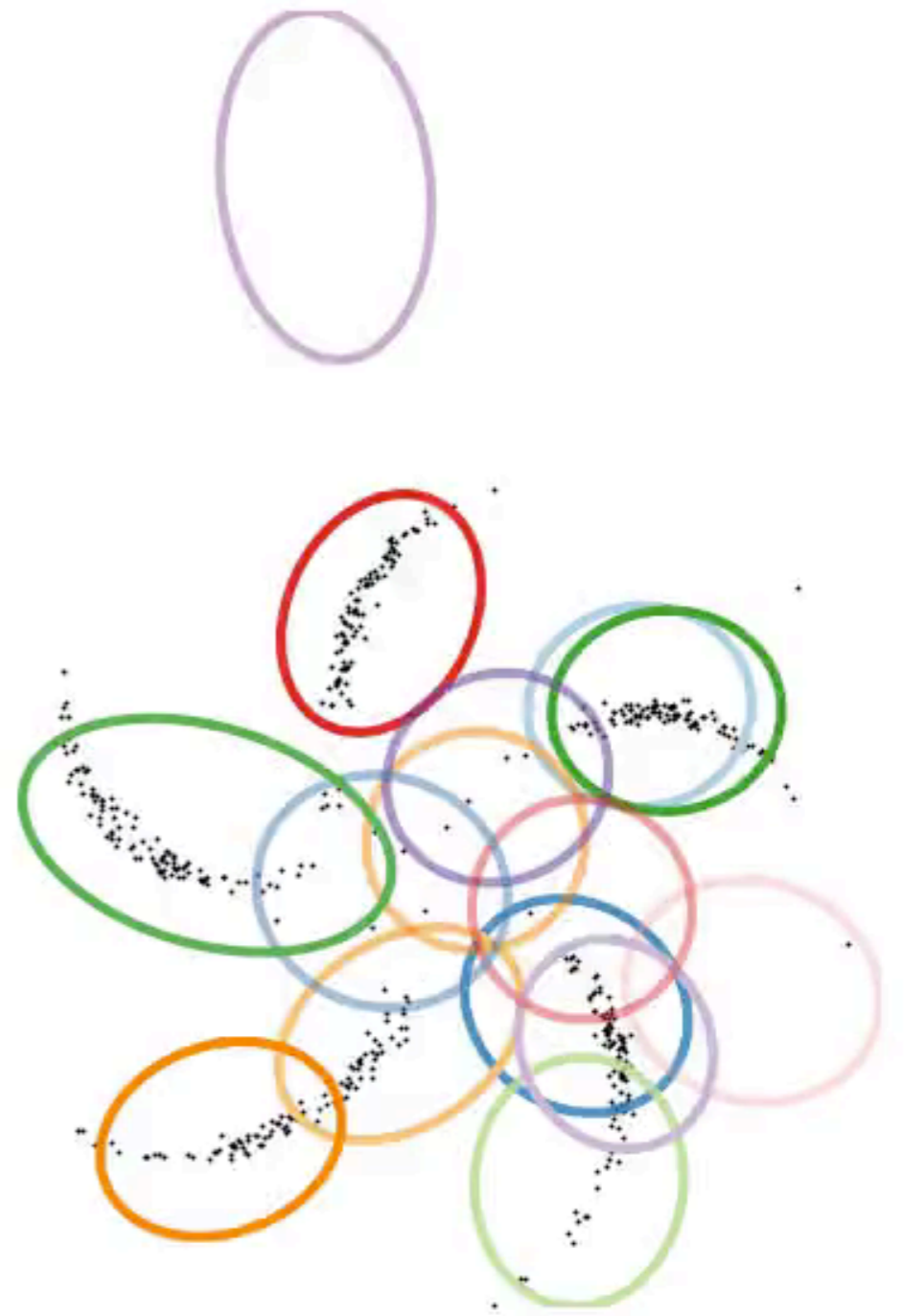Step 1: apply recognition network

Step 2: run fast PGM algorithms

Step 3: sample, compute flat grads

Step 1: apply recognition network

Step 2: run fast PGM algorithms

Step 3: sample, compute flat grads

Step 4: compute natural gradient

Step 1: apply recognition network



Step 2: run fast PGM algorithms



Step 3: sample, compute flat grads



Step 4: compute natural gradient

data space                    latent space

data space

latent space

data

predictions

latent states

frame index

data

predictions

latent states

0   20   40   60   80

frame index

# arbitrary inference queries*



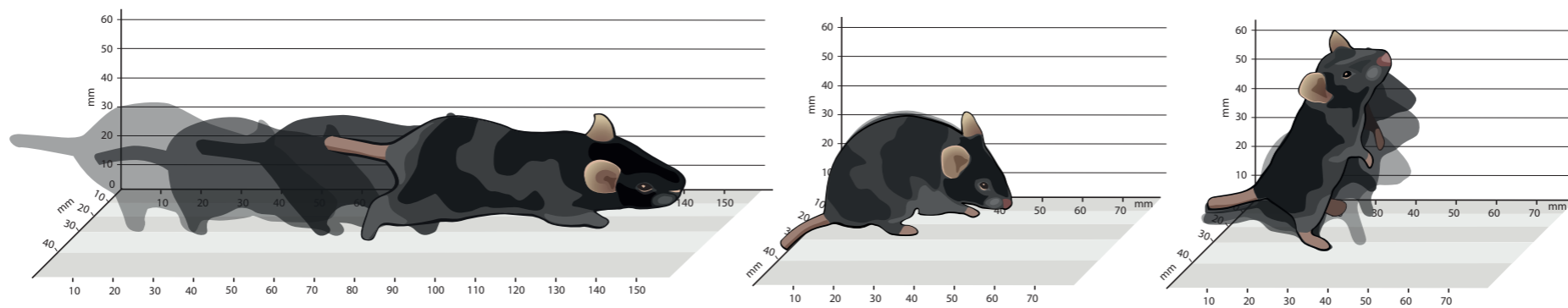*see next slide

# SVAEs can use any inference network architectures

[1] Archer, Park, Buesing, Cunningham, Paninski. Black box variational inference for state space models. ICLR 2016 Workshops.
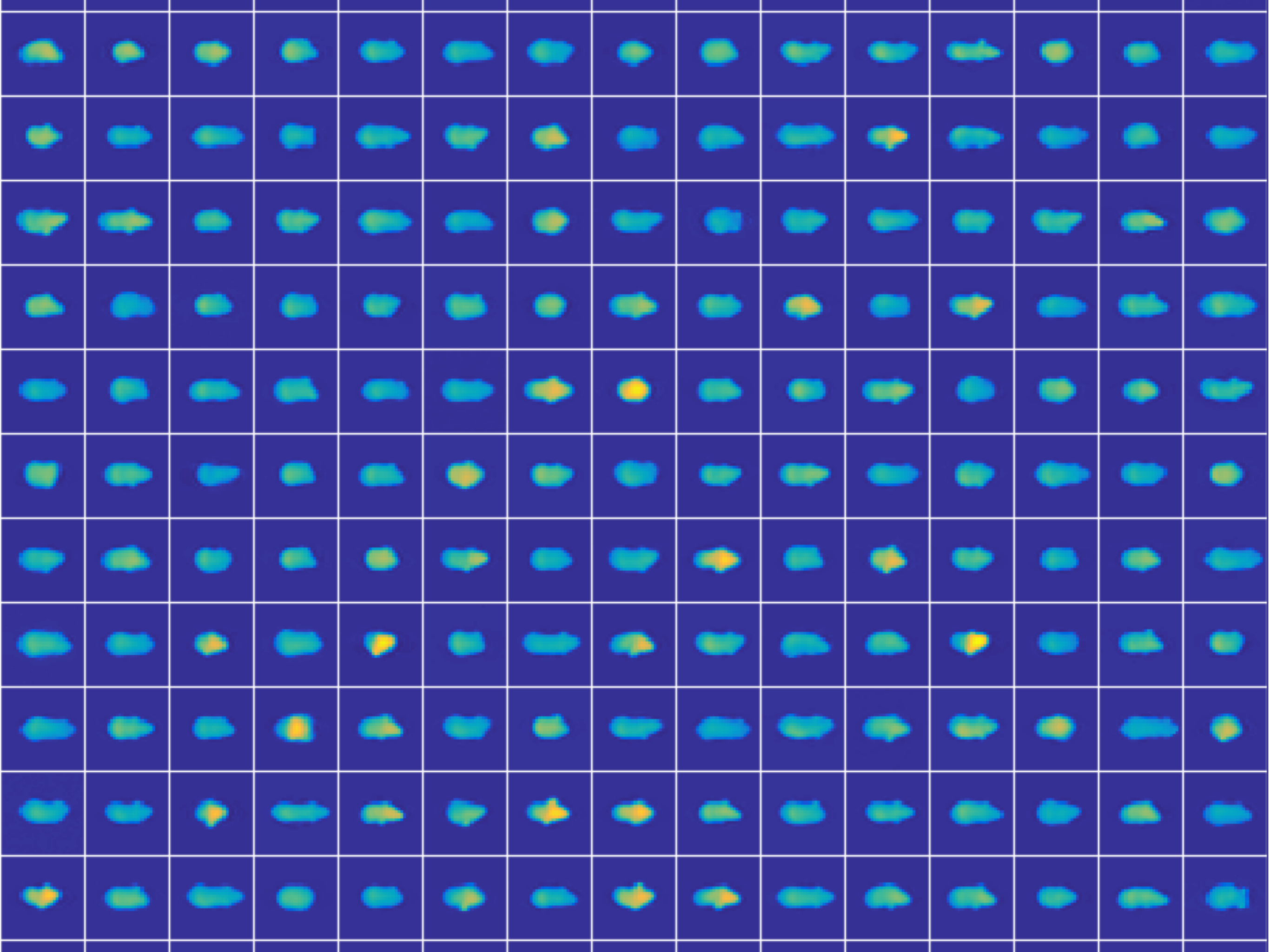[2] Gao*, Archer*, Paninski, Cunningham. Linear dynamical neural population models through nonlinear embeddings. NIPS 2016.
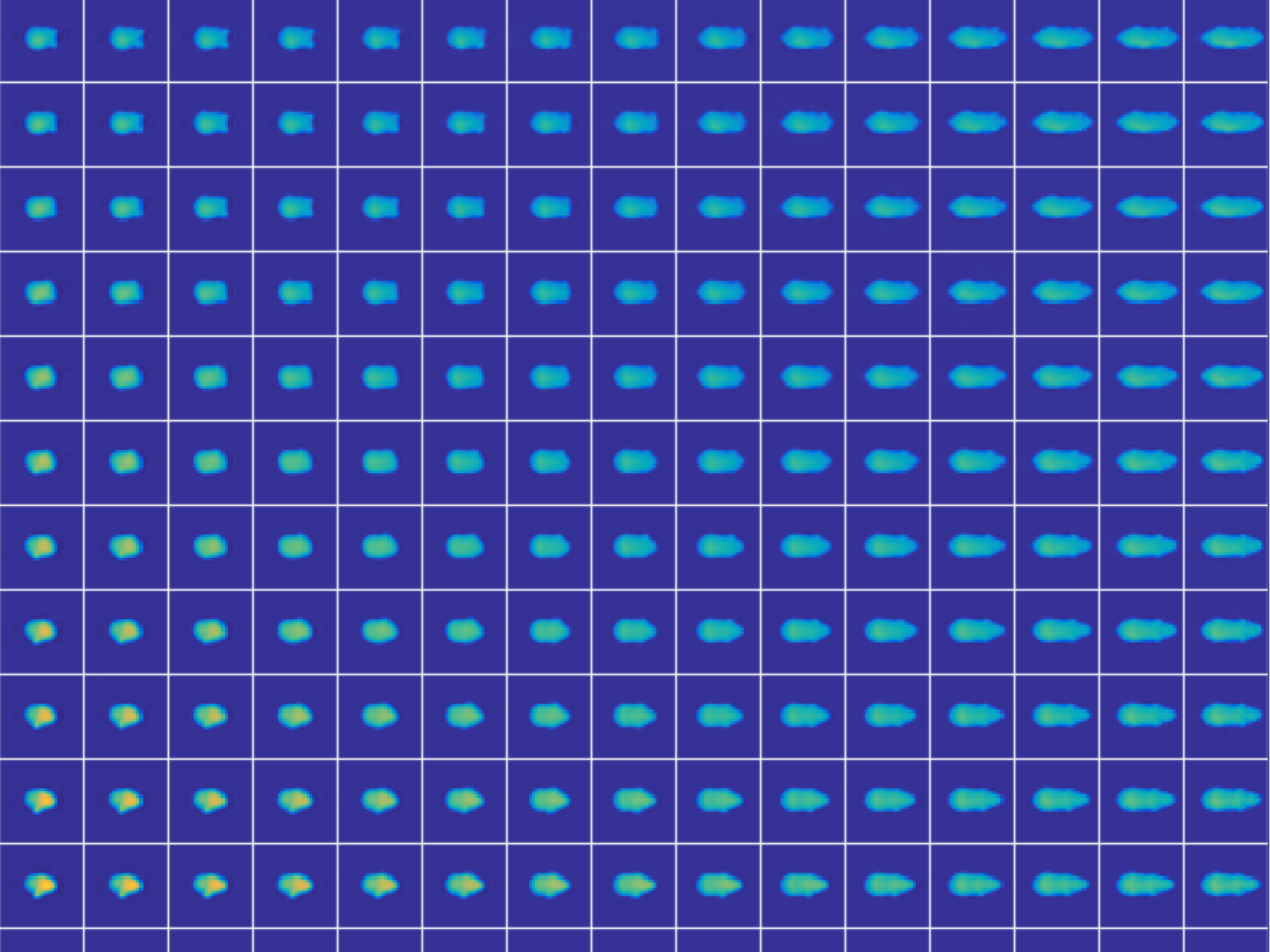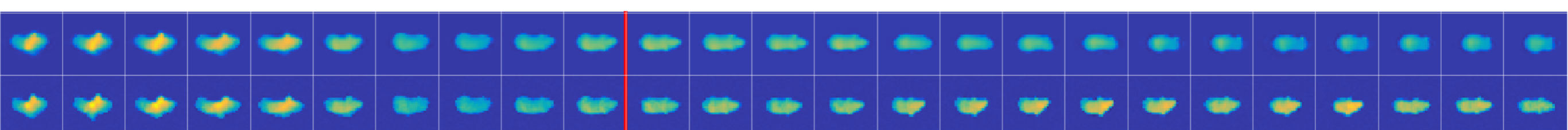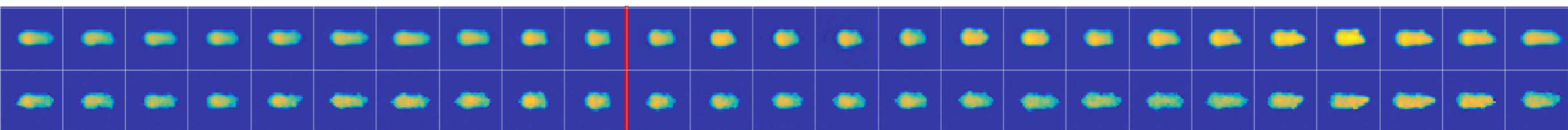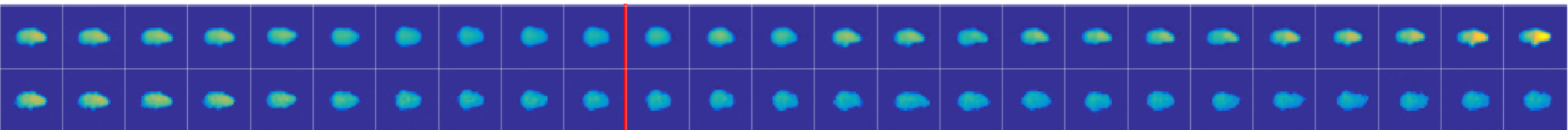
SVAEs
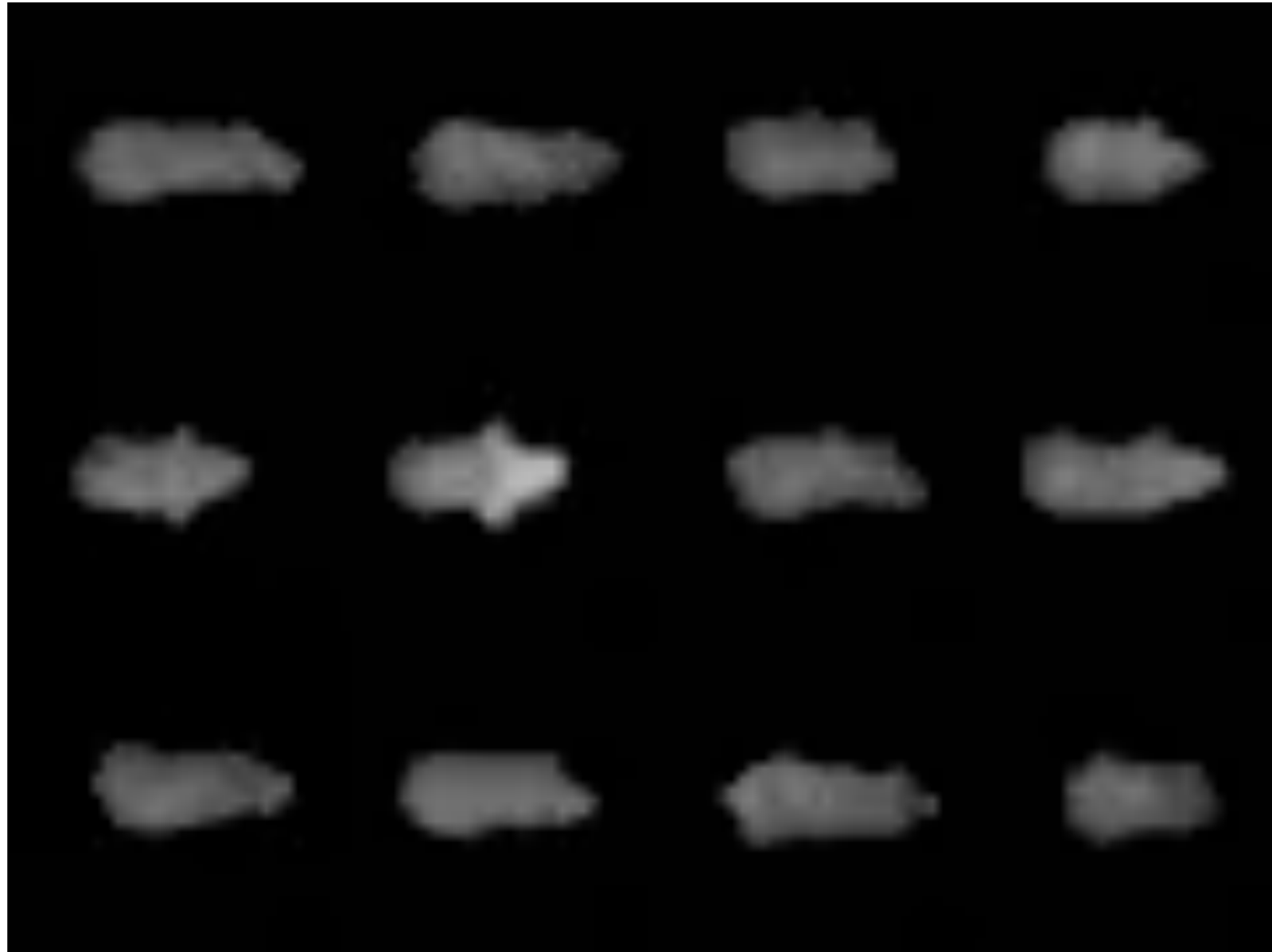
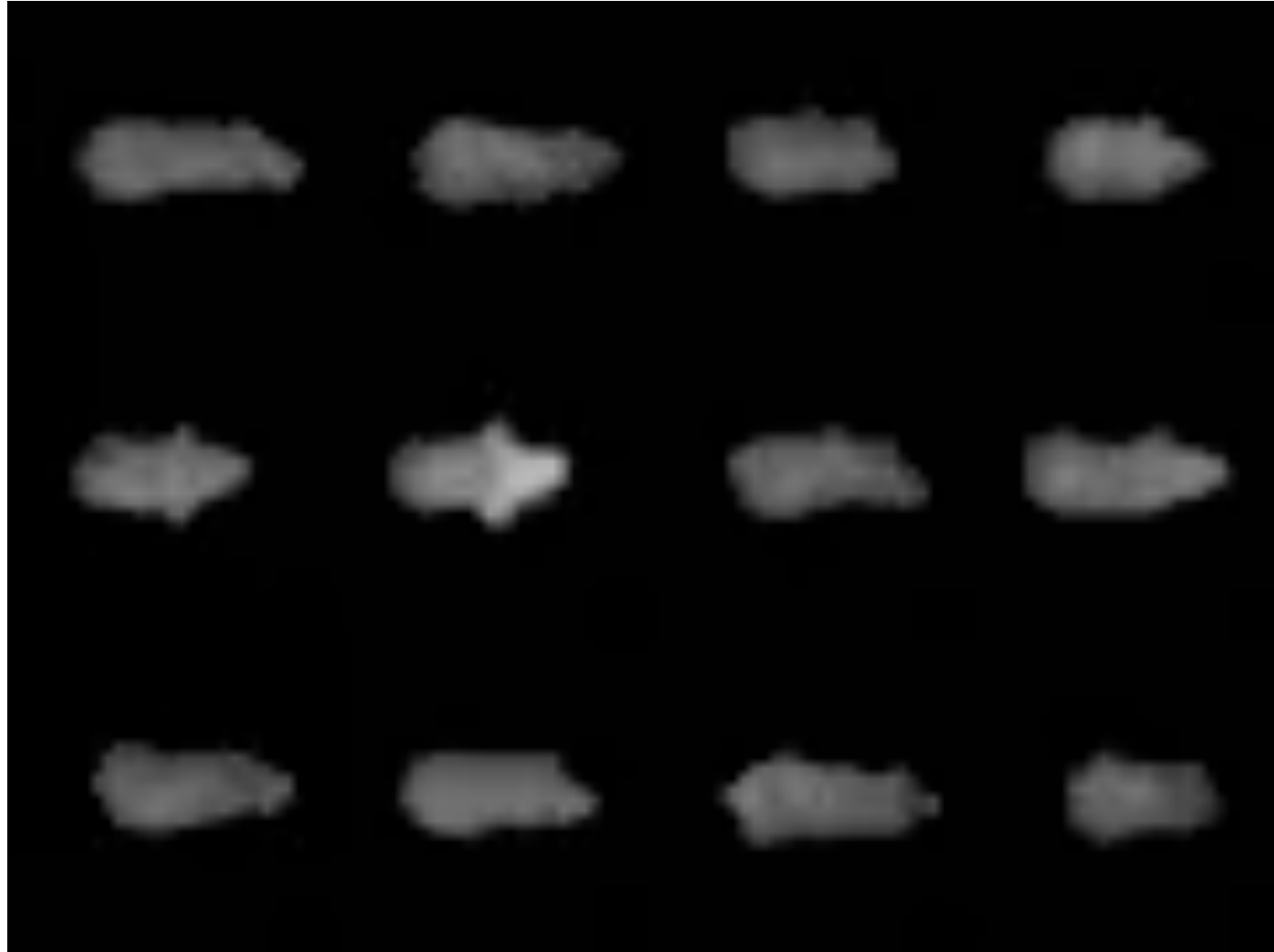**Application:** learn syllable representation of behavior from video
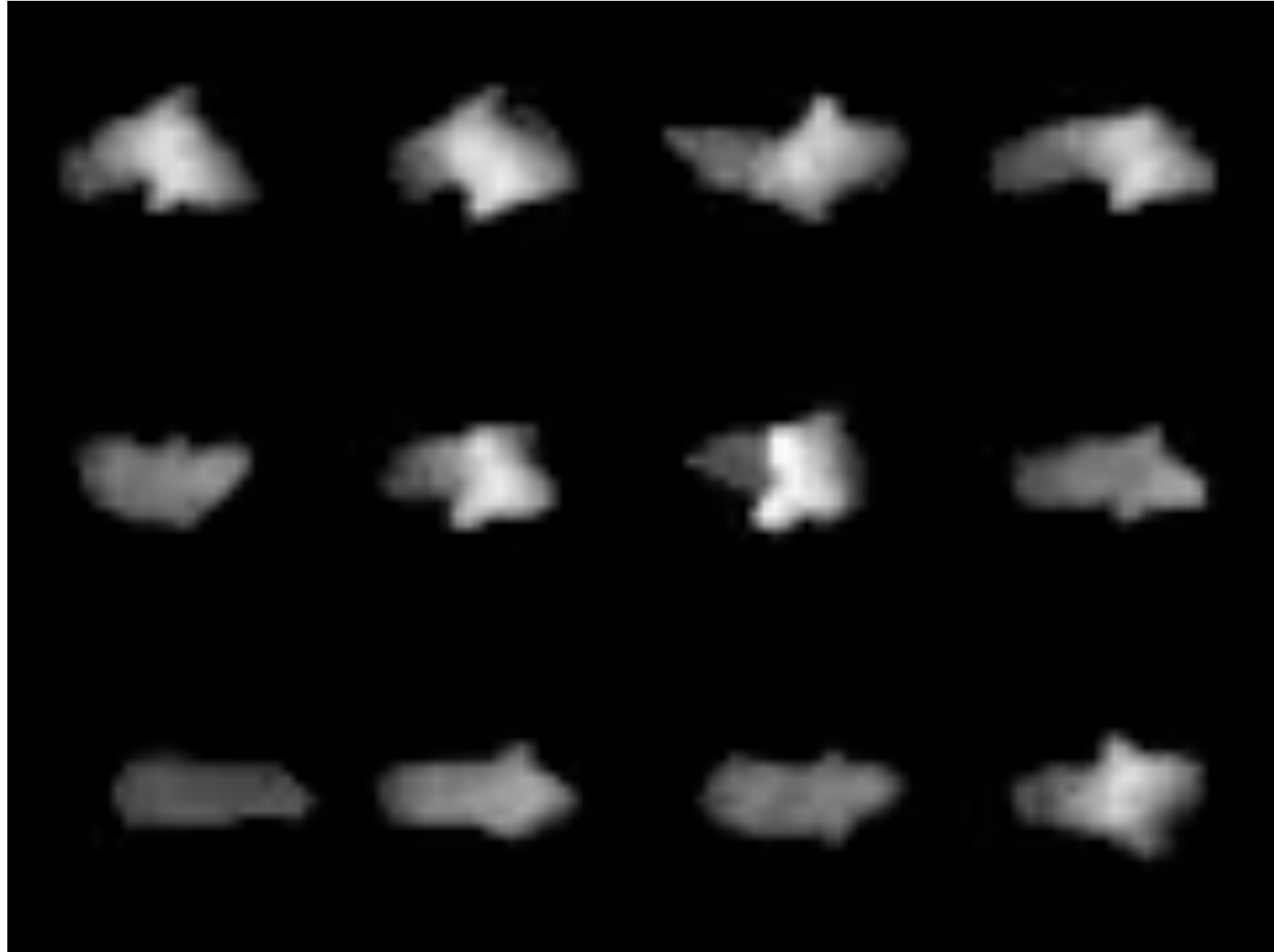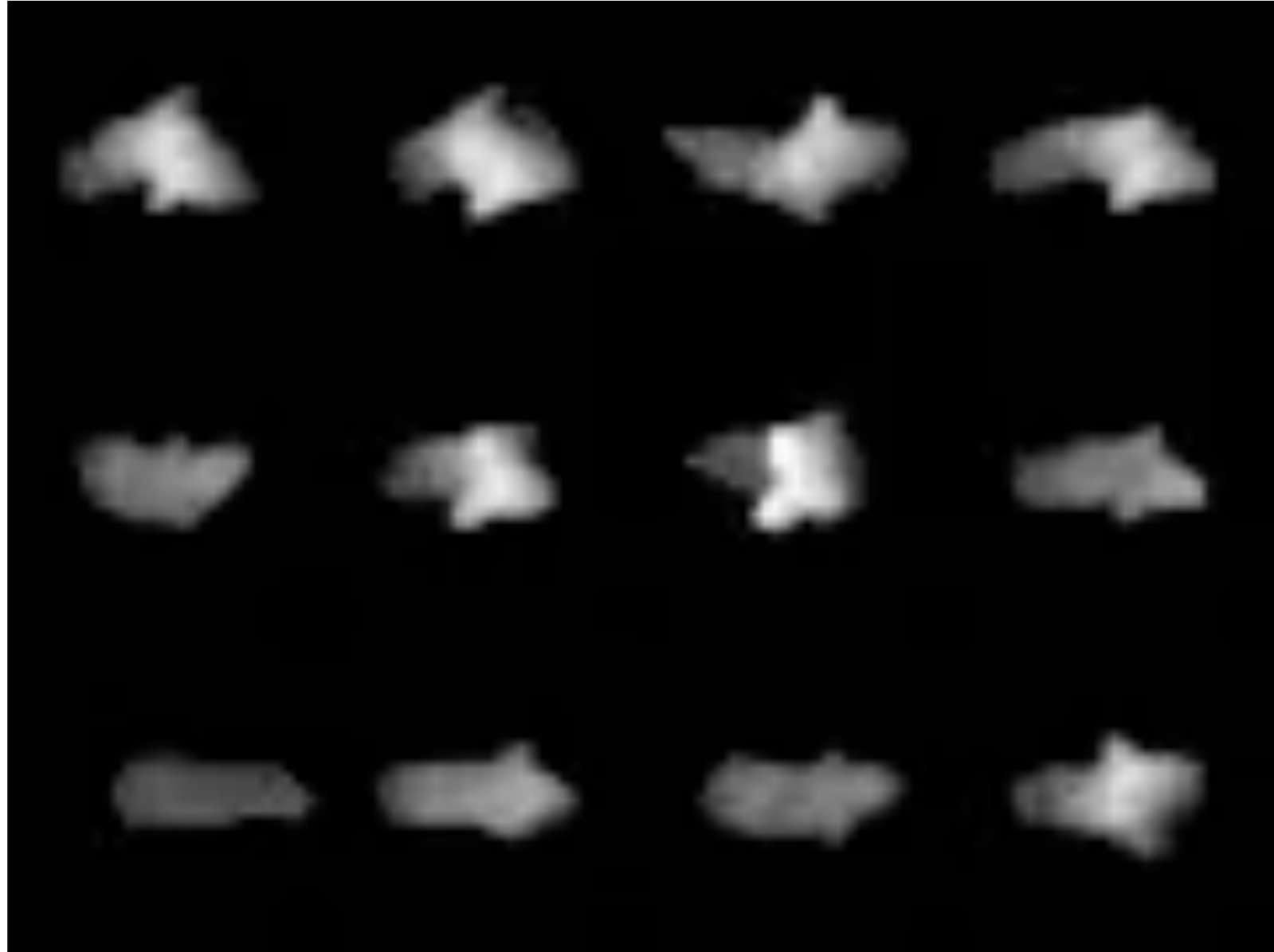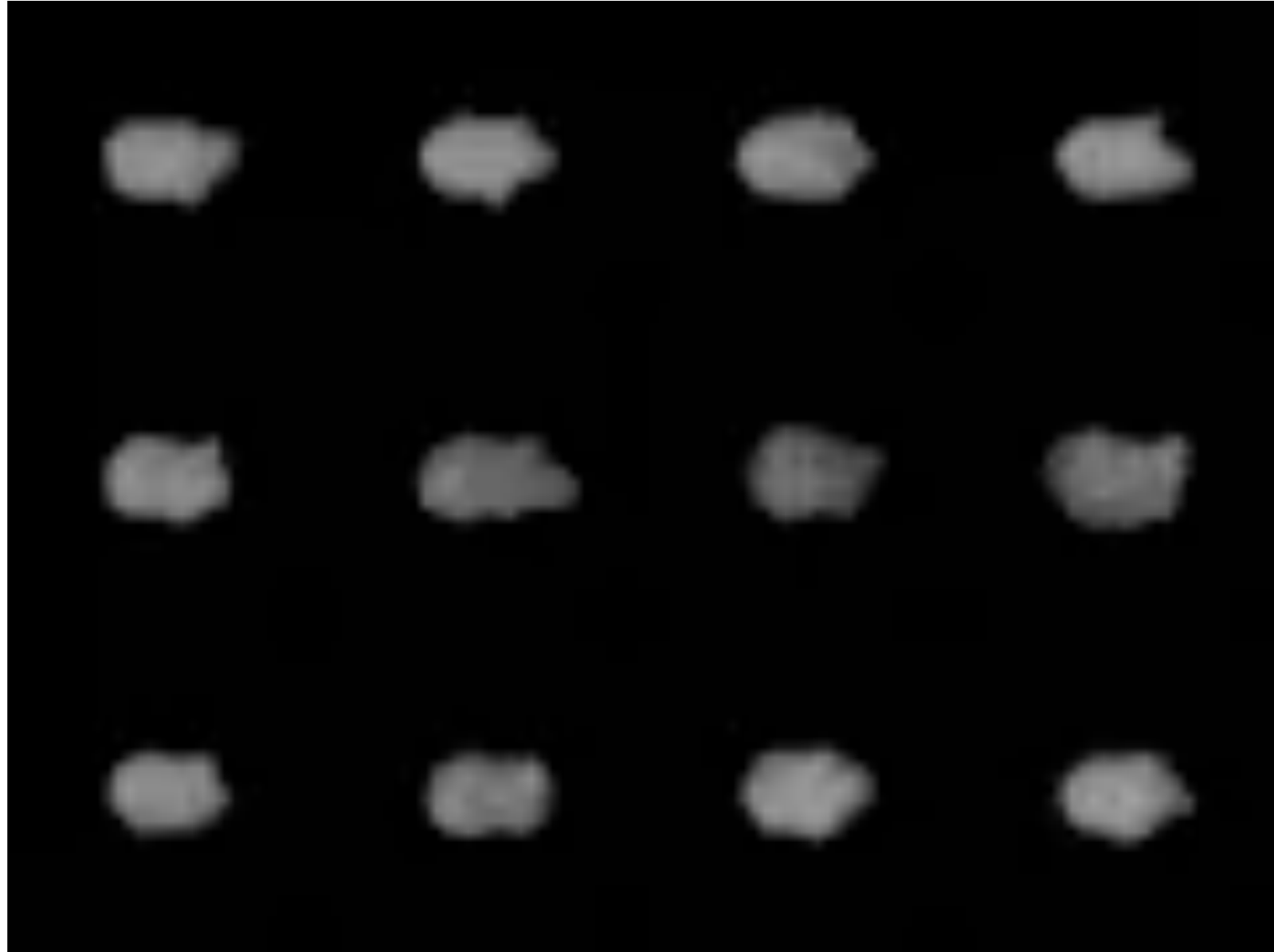
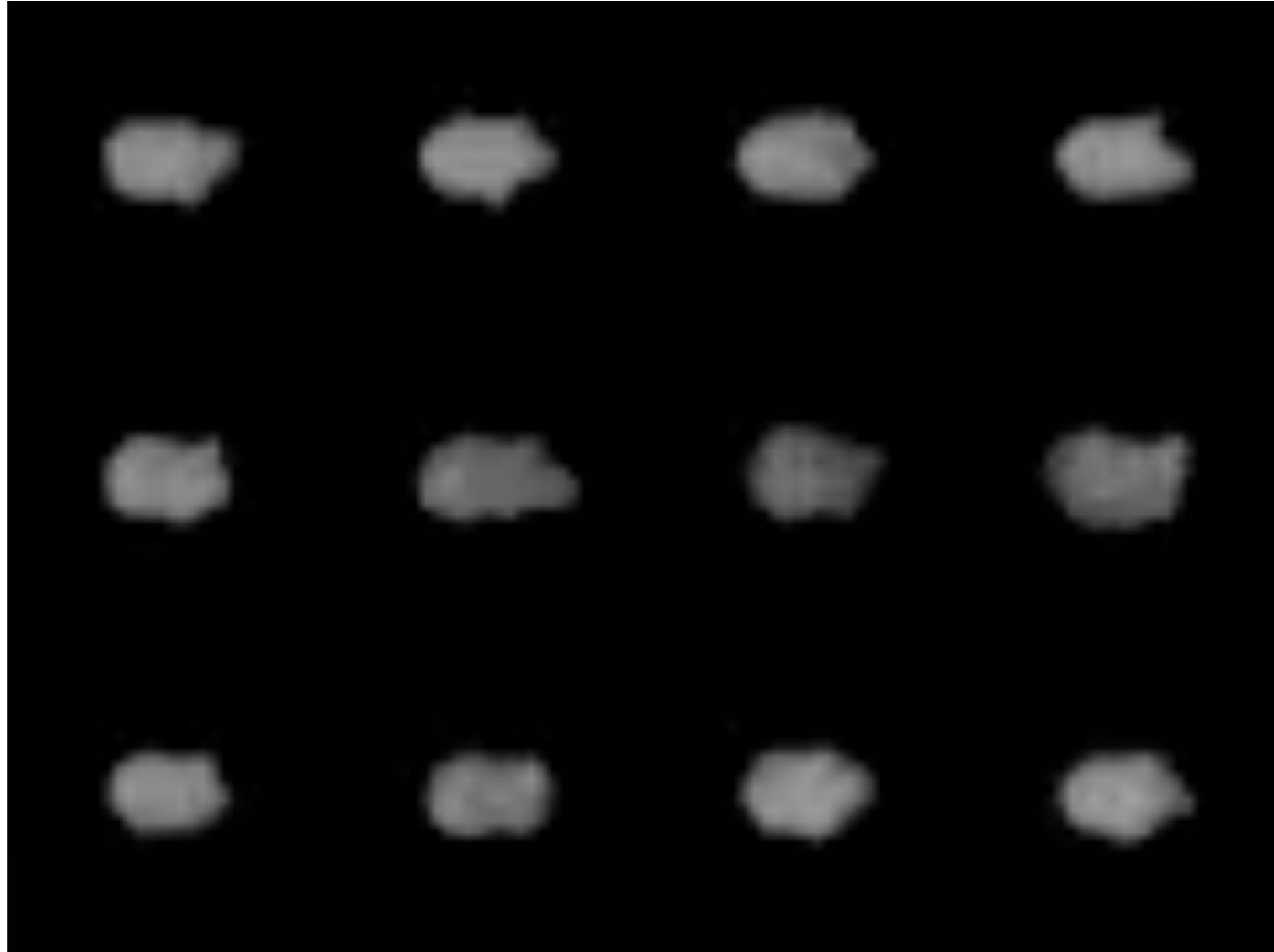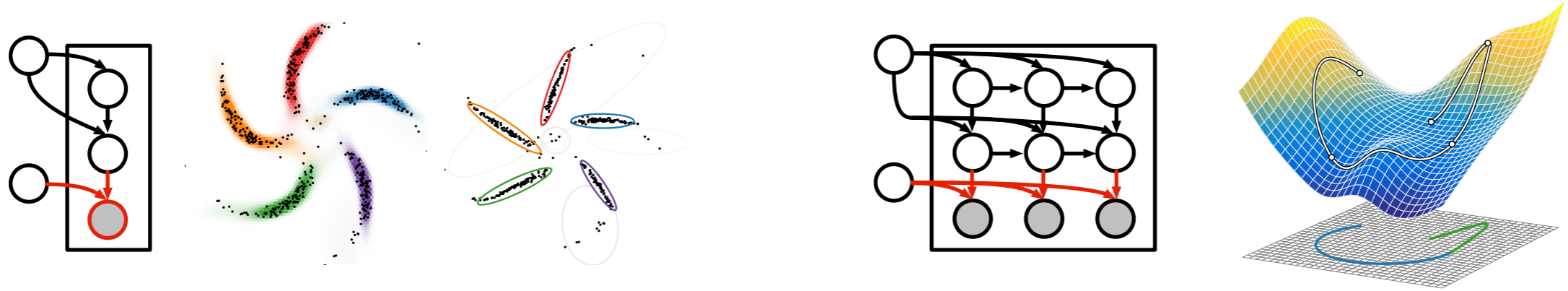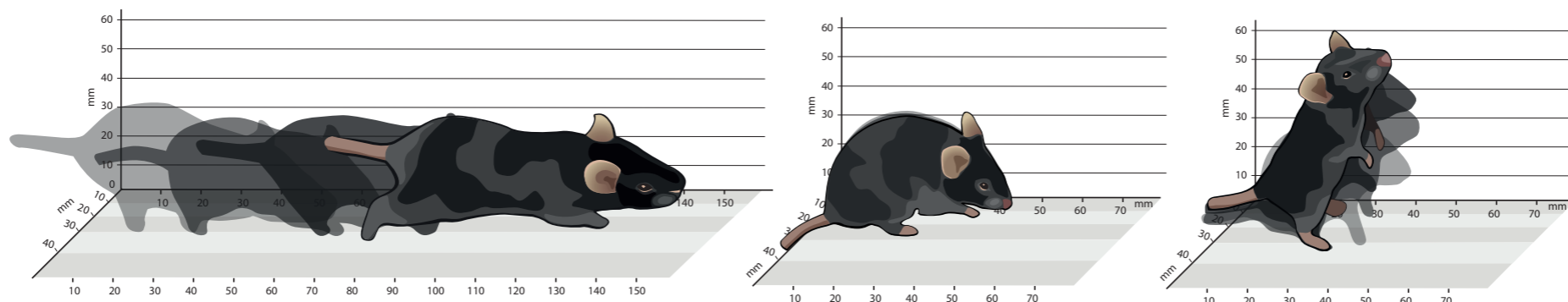start rear

start rear

fall from rear

fall from rear

grooming

grooming

**Modeling idea:** graphical models on latent variables, neural network models for observations



**Inference:** recognition networks output conjugate potentials, then apply fast graphical model inference



**Application:** learn syllable representation of behavior from video

# Thanks!