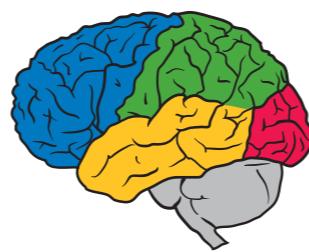
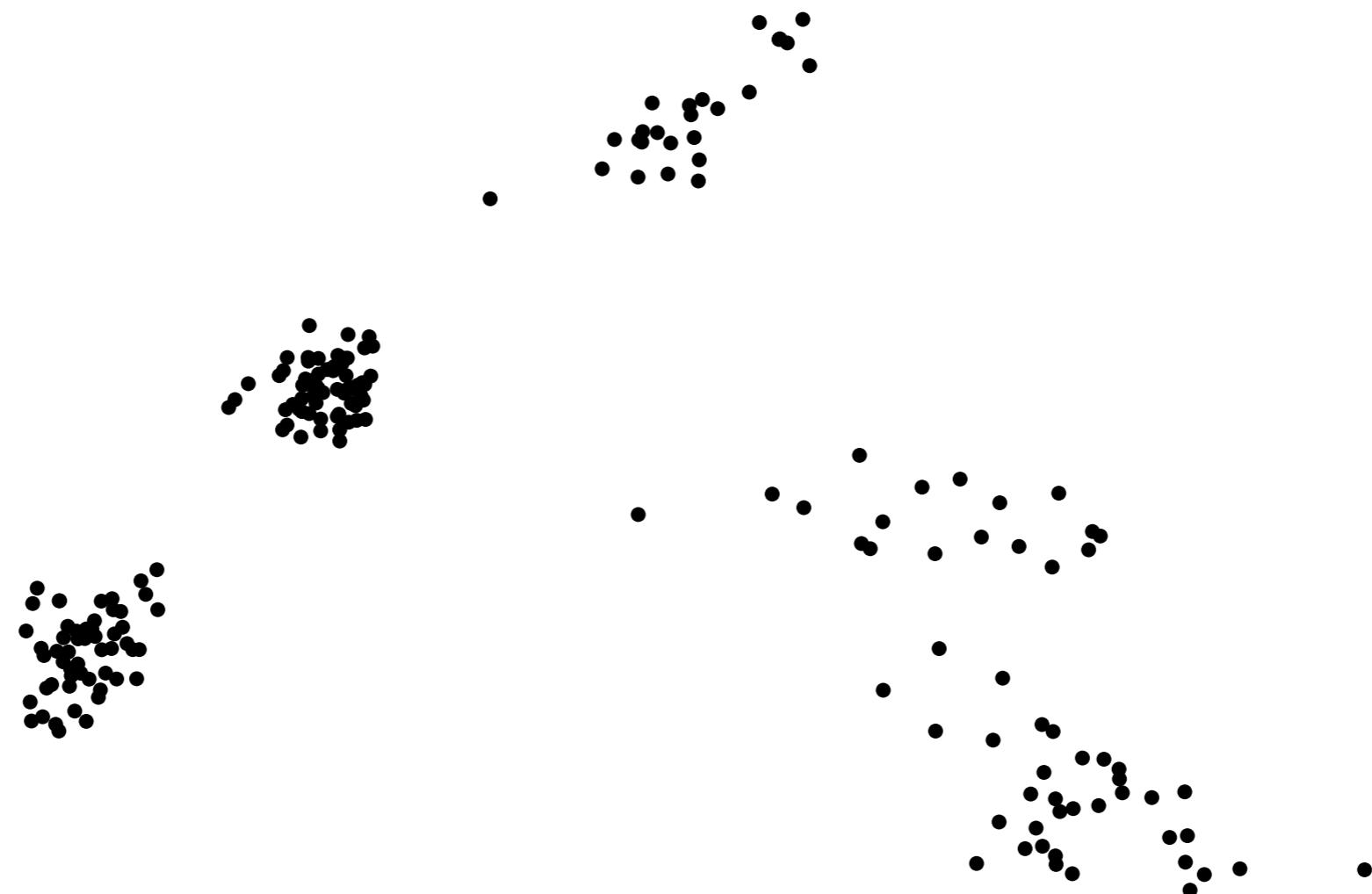
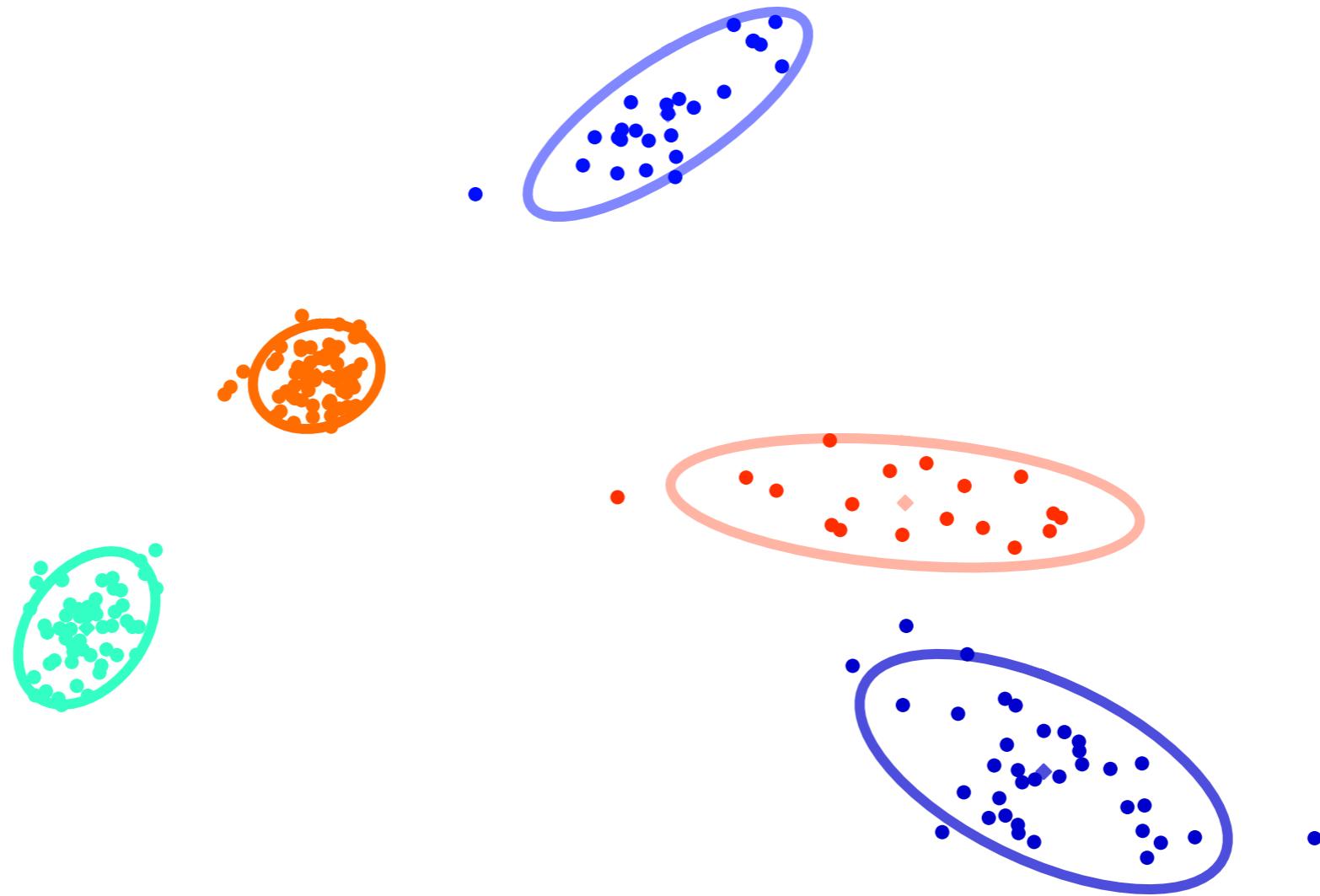


# Composing graphical models with neural networks for structured representations and fast inference

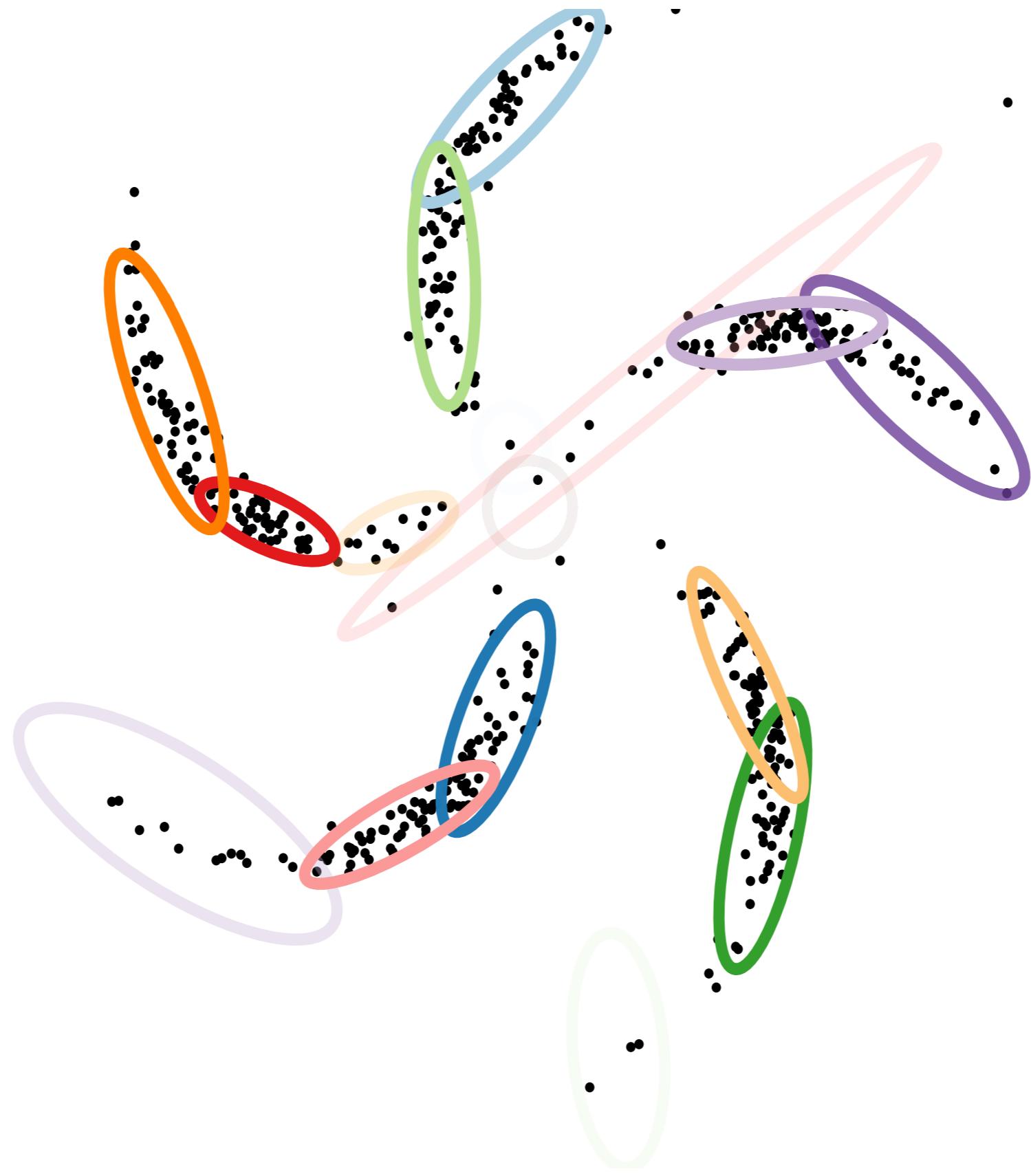
Matt Johnson, David Duvenaud, Alex Wiltschko, Bob Datta, Ryan Adams

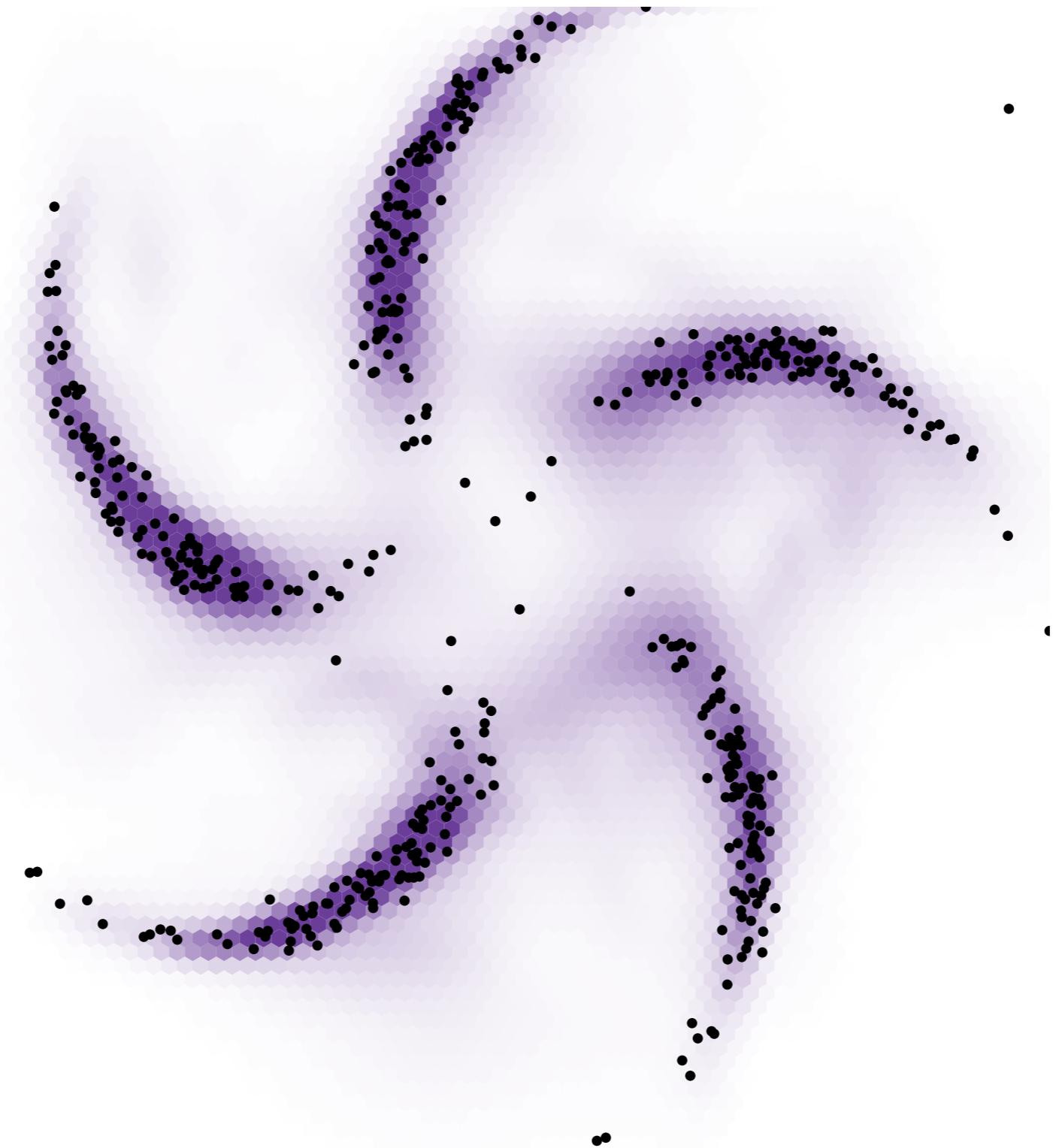


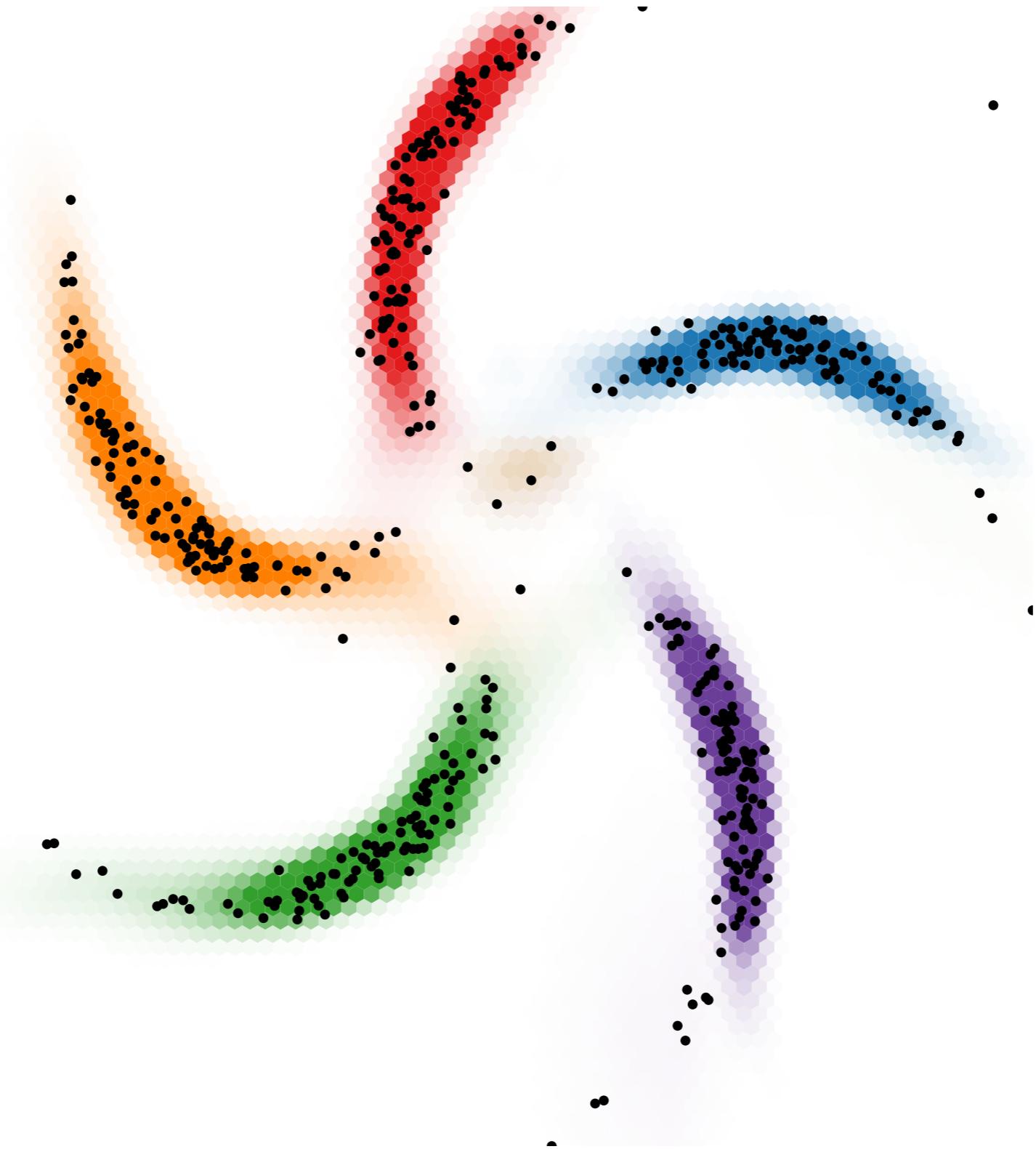


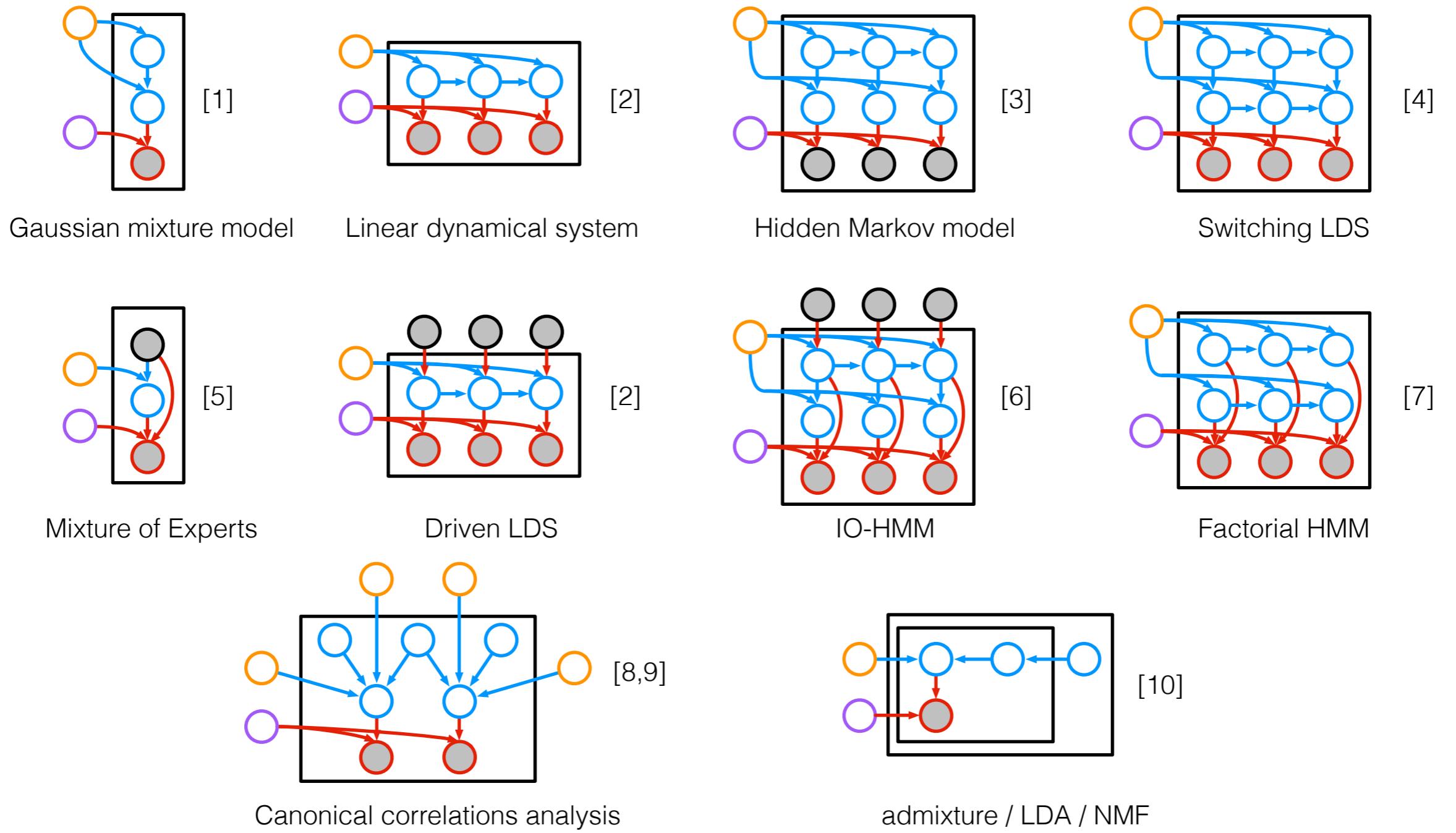












- [1] Palmer, Wipf, Kreutz-Delgado, and Rao. Variational EM algorithms for non-Gaussian latent variable models. NIPS 2005.
- [2] Ghahramani and Beal. Propagation algorithms for variational Bayesian learning. NIPS 2001.
- [3] Beal. Variational algorithms for approximate Bayesian inference, Ch. 3. U of 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.

## Probabilistic graphical models

- + structured representations
- + priors and uncertainty
- + data and computational efficiency
- rigid assumptions may not fit
- feature engineering
- top-down inference

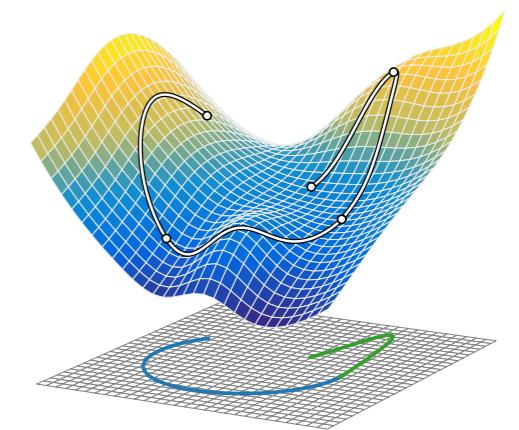
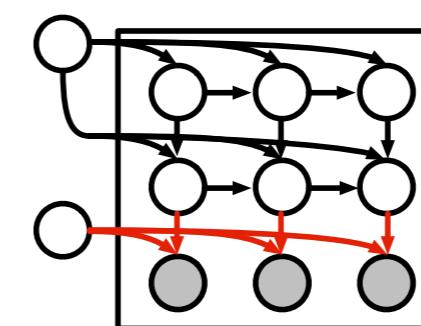
## Deep learning

- neural net “goo”
- difficult parameterization
- can require lots of data
- + flexible
- + feature learning
- + recognition networks

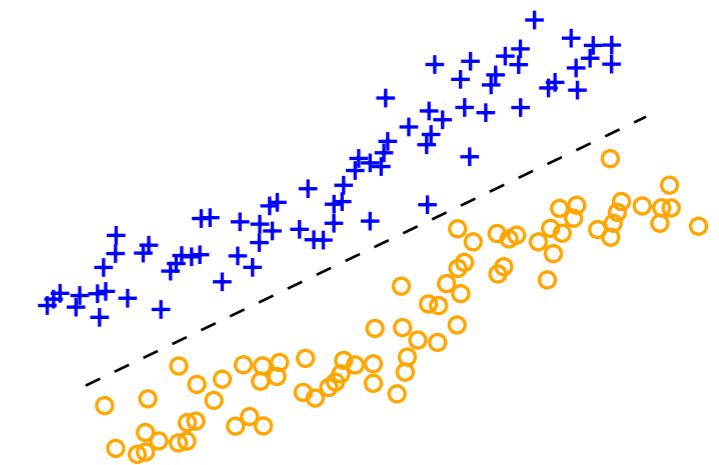
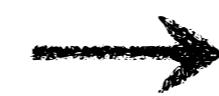
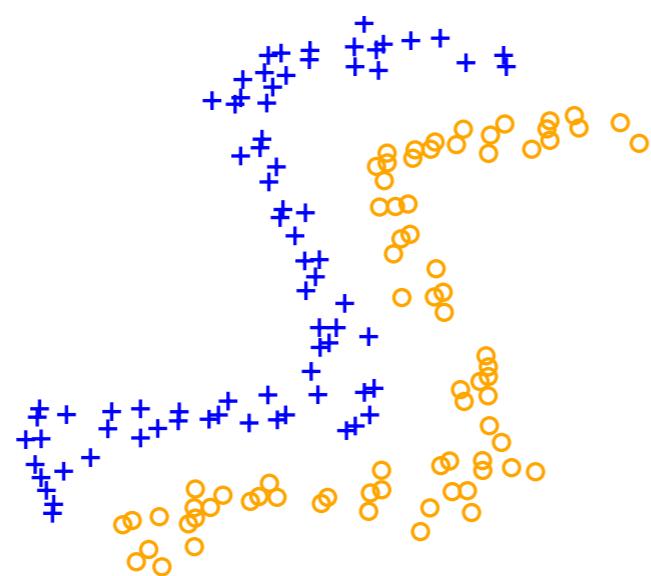


MAKE PGMS  
GREAT AGAIN

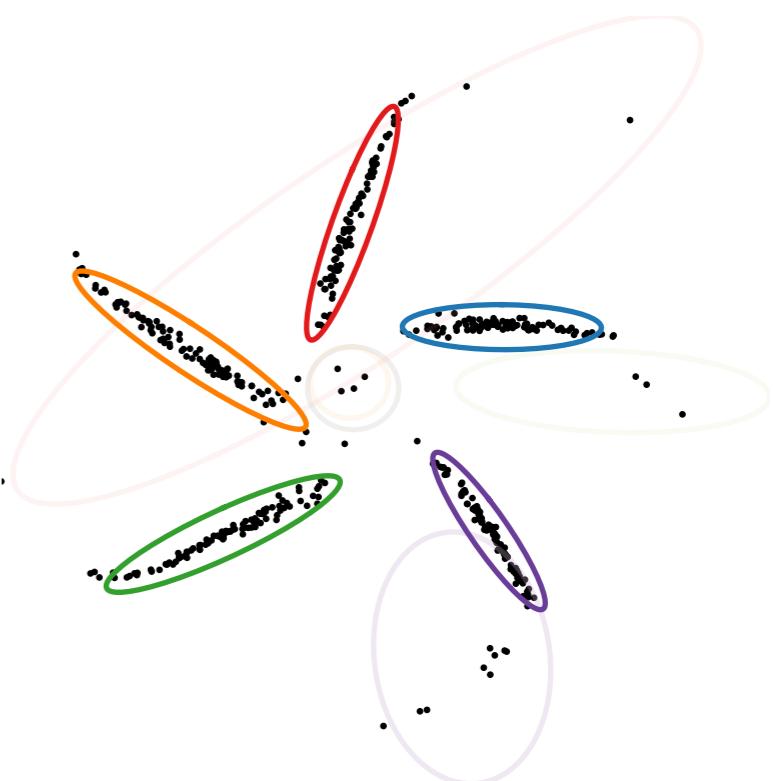
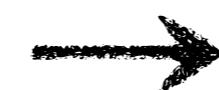
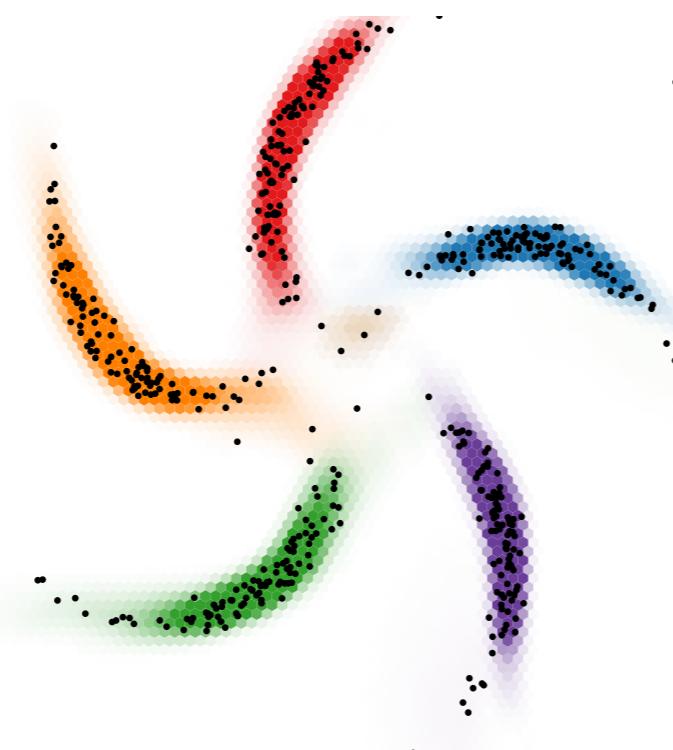
**Modeling idea:** graphical models on latent variables,  
neural network models for observations



supervised  
learning

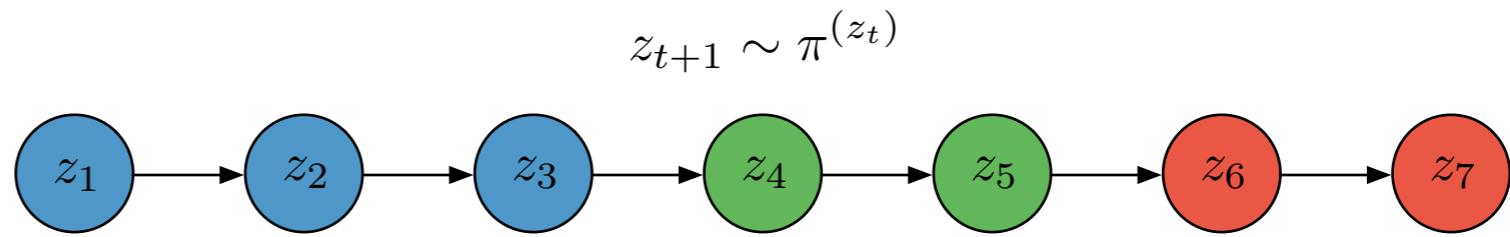


unsupervised  
learning



$$\pi = \begin{bmatrix} \textcolor{blue}{\blacksquare} & \textcolor{red}{\blacksquare} & \textcolor{green}{\blacksquare} \\ \hline \textcolor{blue}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \textcolor{red}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \textcolor{green}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \hline \end{bmatrix} \quad \begin{array}{c} \textcolor{blue}{\blacksquare} \\ \textcolor{red}{\blacksquare} \\ \textcolor{green}{\blacksquare} \end{array} \quad \begin{array}{c} \textcolor{blue}{\blacksquare} \\ \textcolor{red}{\blacksquare} \\ \textcolor{green}{\blacksquare} \end{array}$$

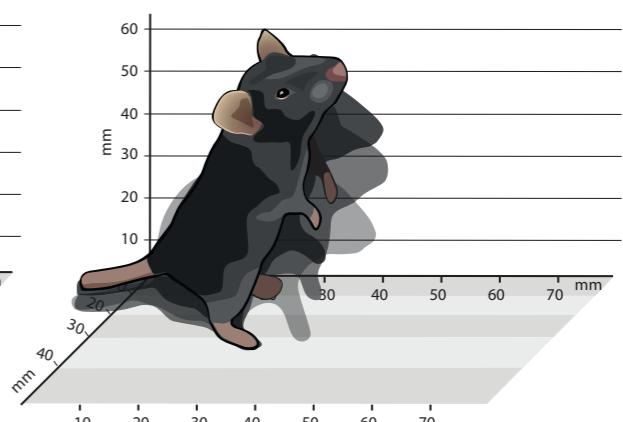
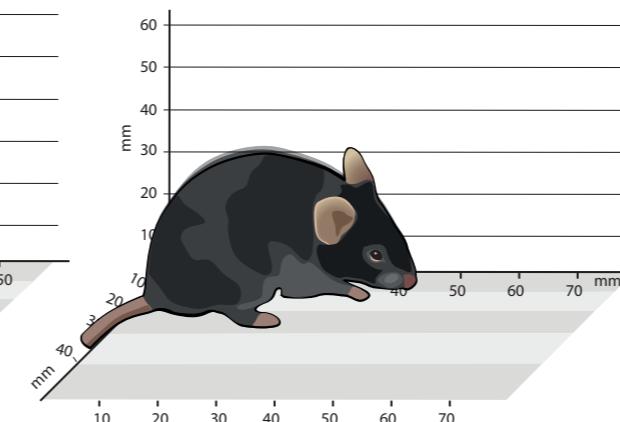
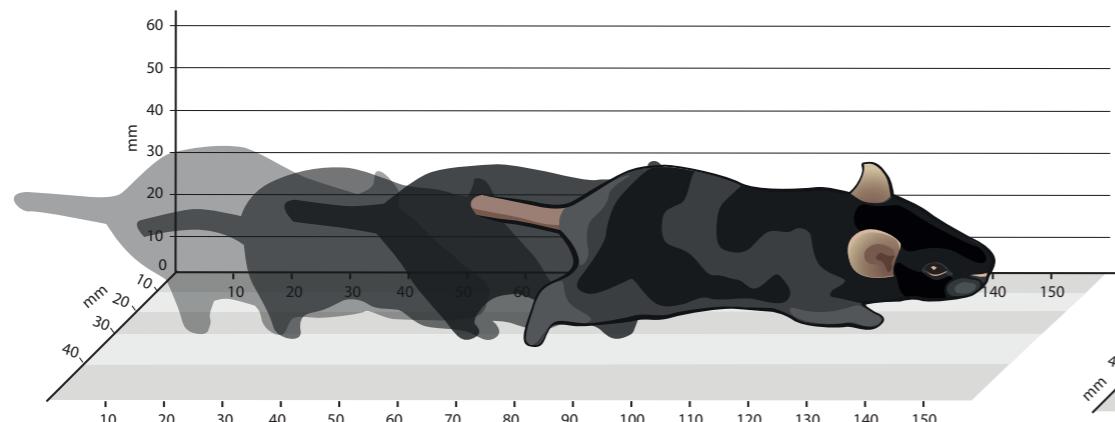
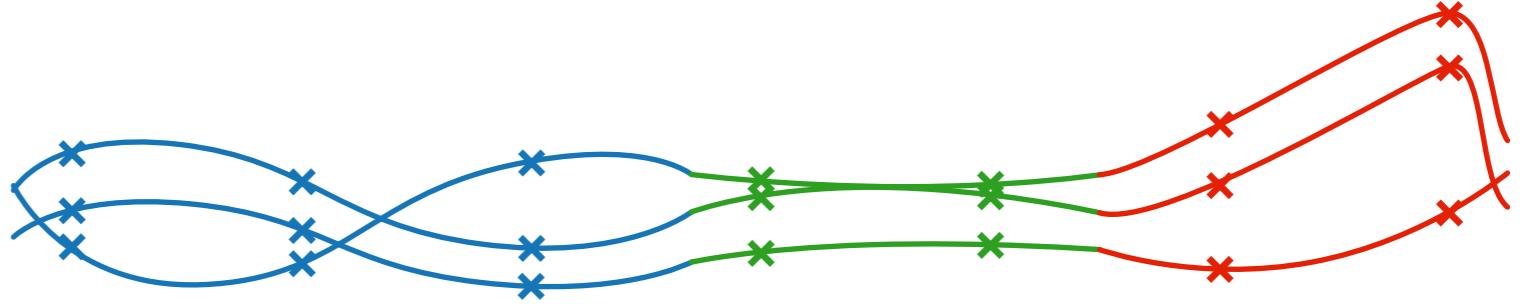
$$\pi^{(1)} \\ \pi^{(2)} \\ \pi^{(3)}$$



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

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

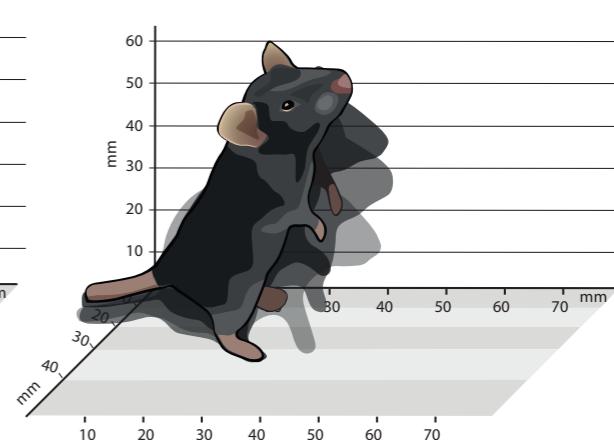
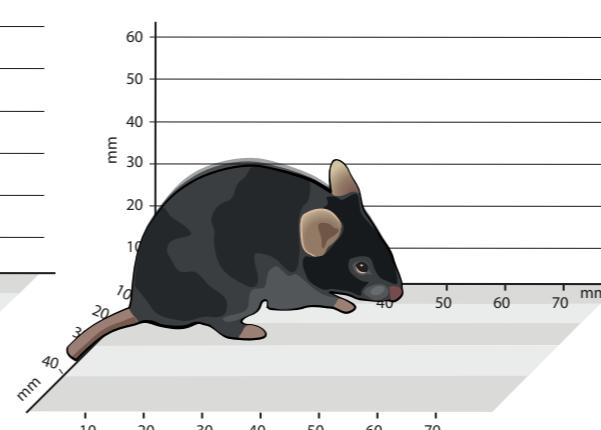
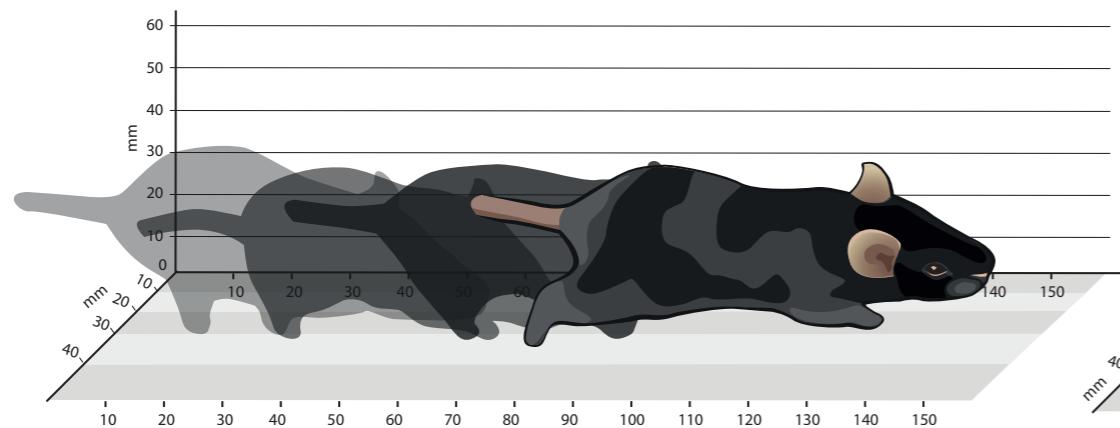
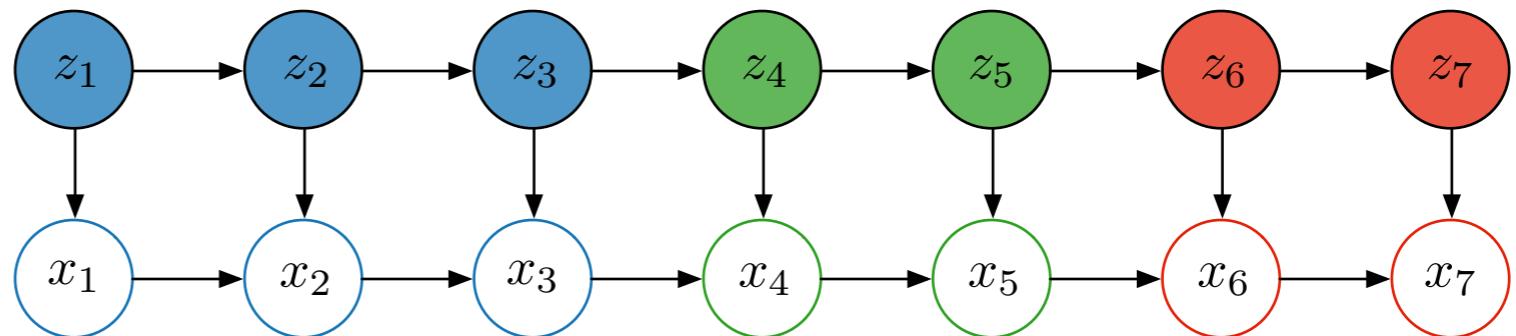
$$x_{t+1} = A^{(z_t)} x_t + B^{(z_t)} u_t \quad u_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I)$$

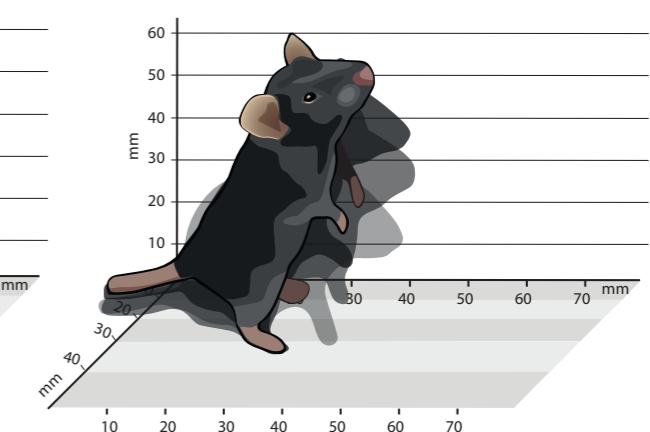
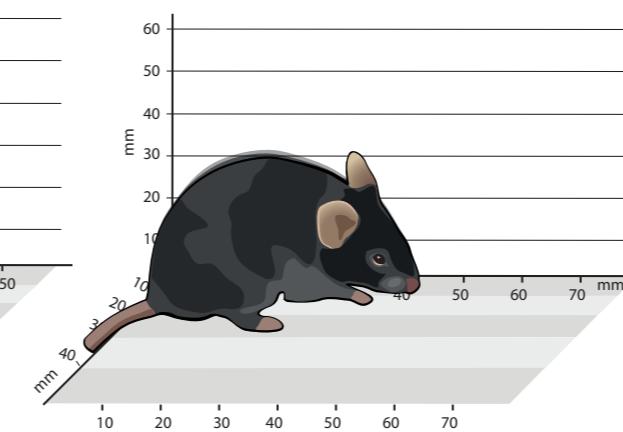
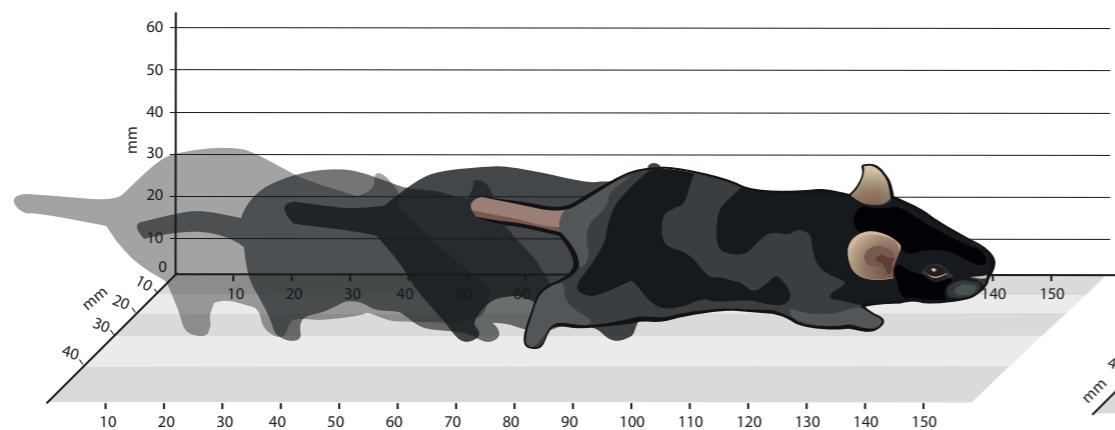
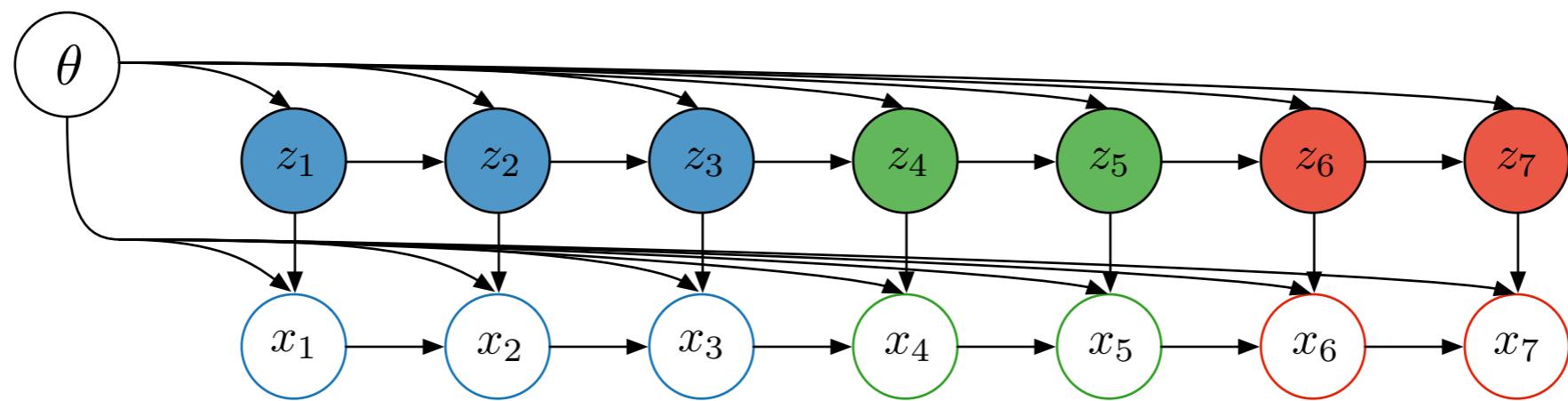


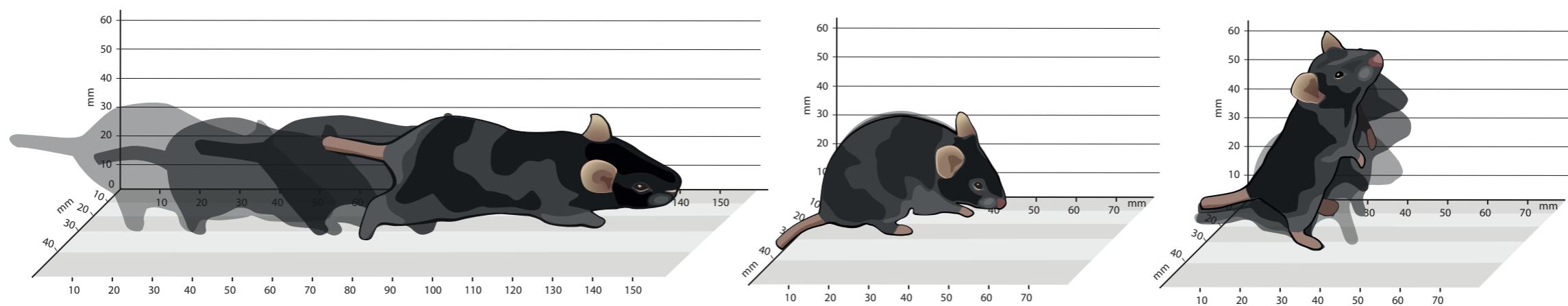
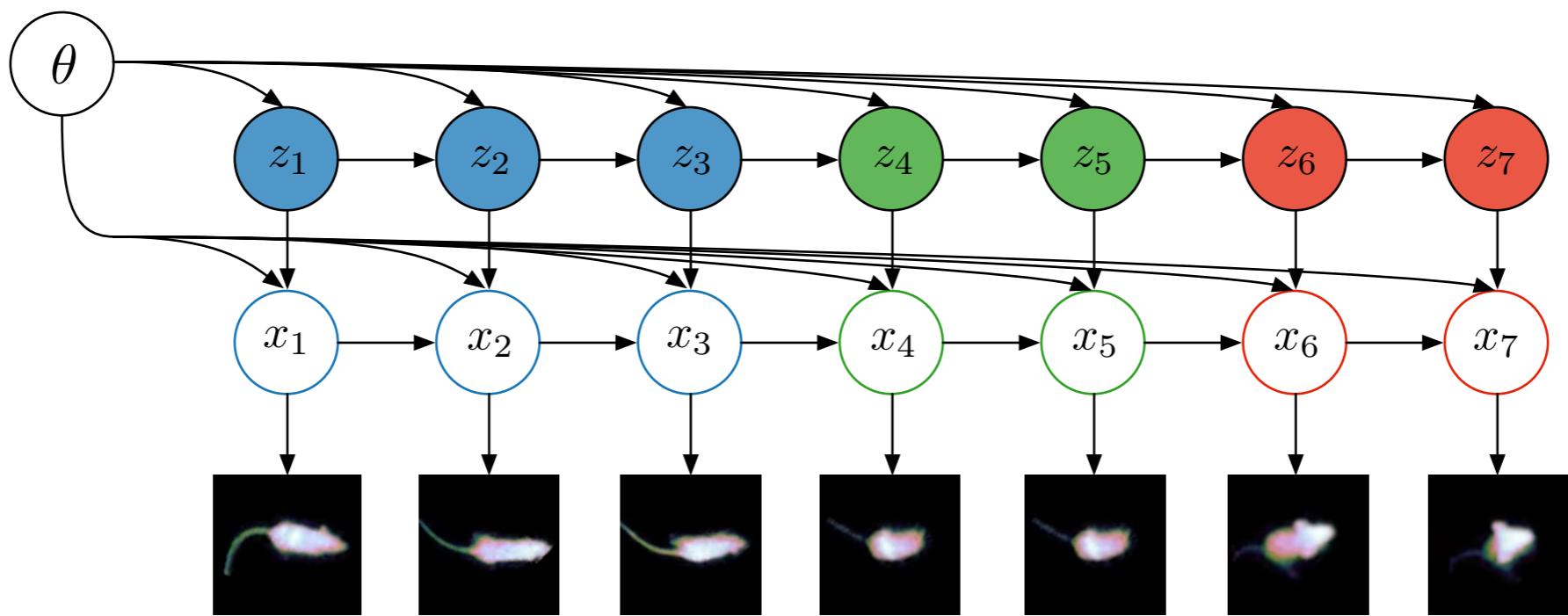
$$\pi = \begin{bmatrix} \textcolor{blue}{\blacksquare} & \textcolor{red}{\blacksquare} & \textcolor{green}{\blacksquare} \\ \hline \textcolor{blue}{\blacksquare} & \textcolor{black}{\rule{0pt}{1em}} & \textcolor{black}{\rule{0pt}{1em}} \\ \textcolor{red}{\blacksquare} & \textcolor{black}{\rule{0pt}{1em}} & \textcolor{black}{\rule{0pt}{1em}} \\ \textcolor{green}{\blacksquare} & \textcolor{black}{\rule{0pt}{1em}} & \textcolor{black}{\rule{0pt}{1em}} \end{bmatrix} \quad \begin{array}{c} \pi^{(1)} \\ \pi^{(2)} \\ \pi^{(3)} \end{array}$$

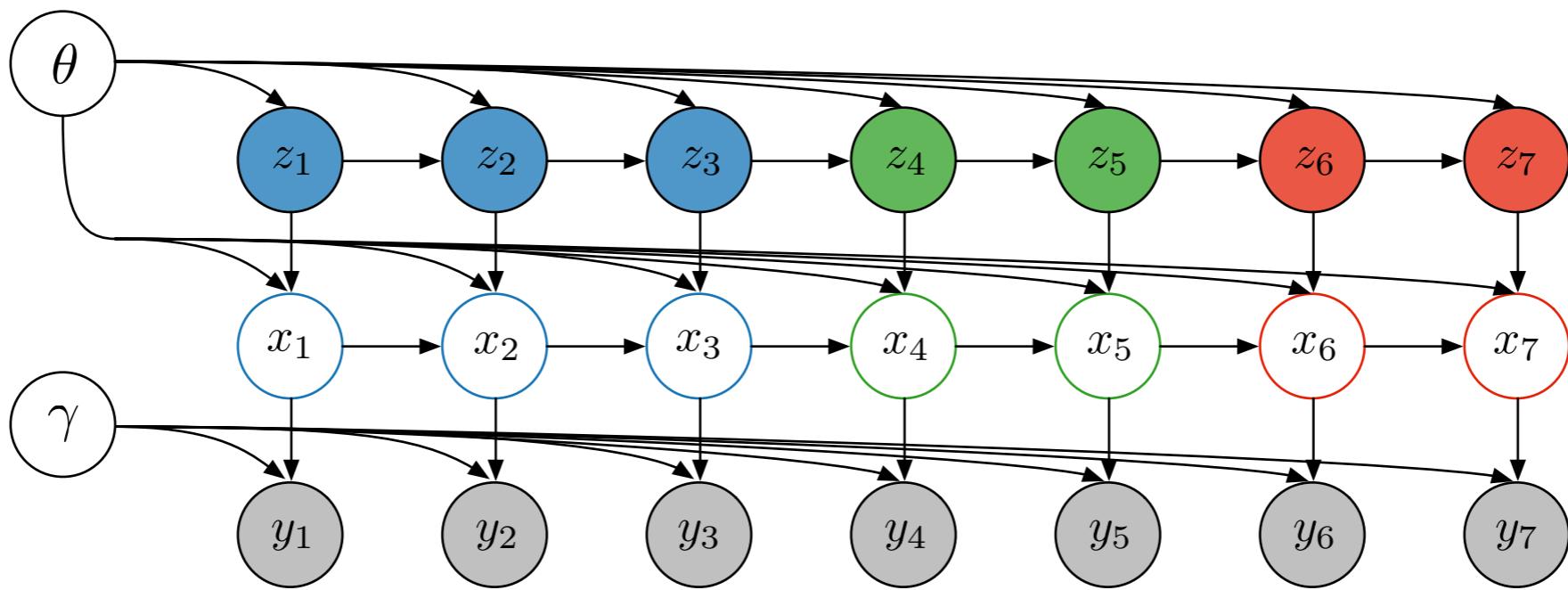
$$A^{(1)} \qquad A^{(2)} \qquad A^{(3)}$$

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

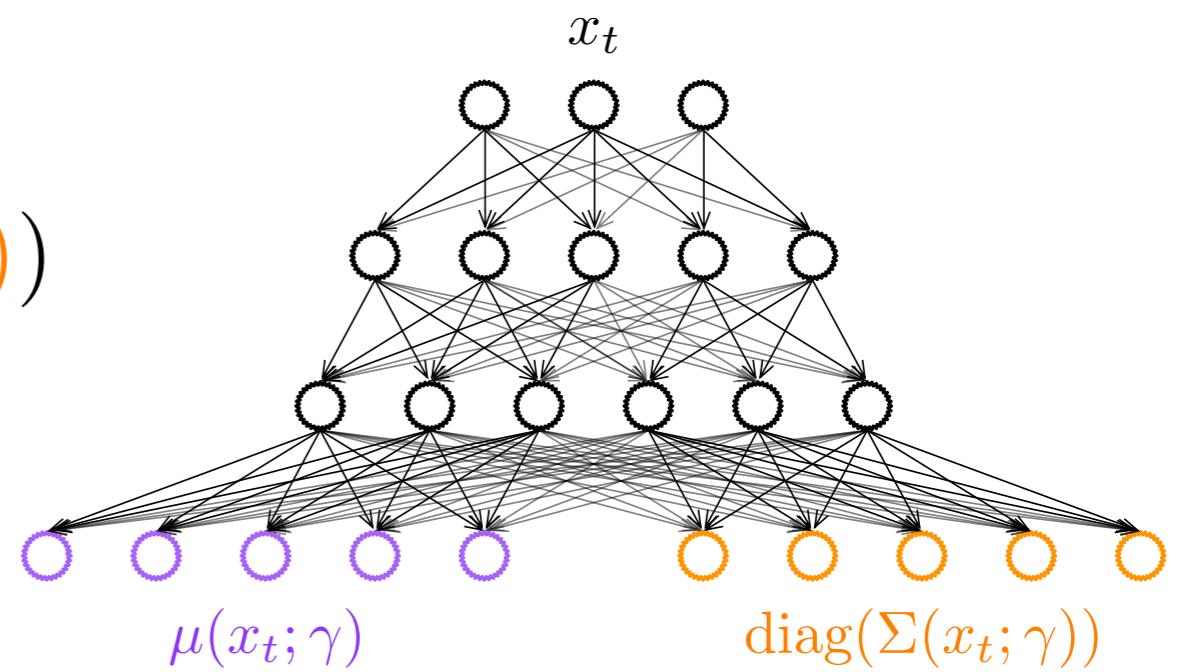


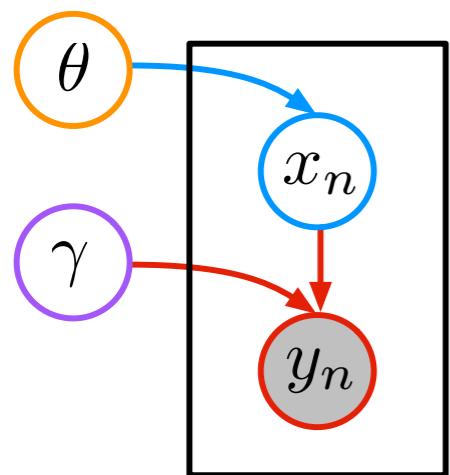
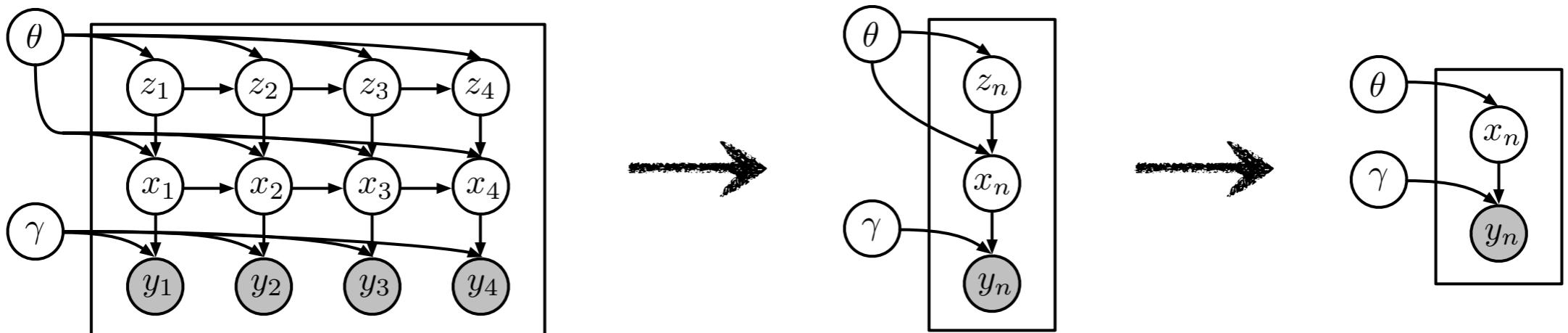






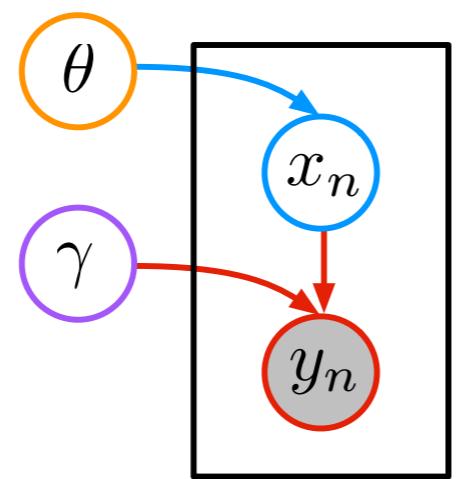
$$y_t \mid x_t, \gamma \sim \mathcal{N}(\mu(x_t; \gamma), \Sigma(x_t; \gamma))$$





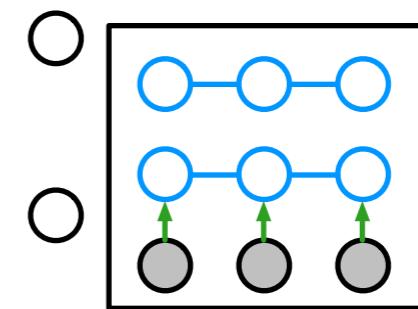
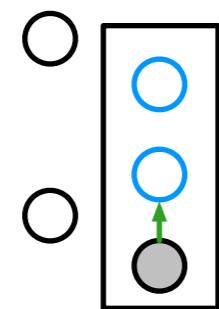
$p(\theta)$   
 $p(x | \theta)$   
 $p(\gamma)$   
 $p(y | x, \gamma)$

conjugate prior on global variables  
 exponential family on local variables  
 any prior on observation parameters  
 neural network observation model

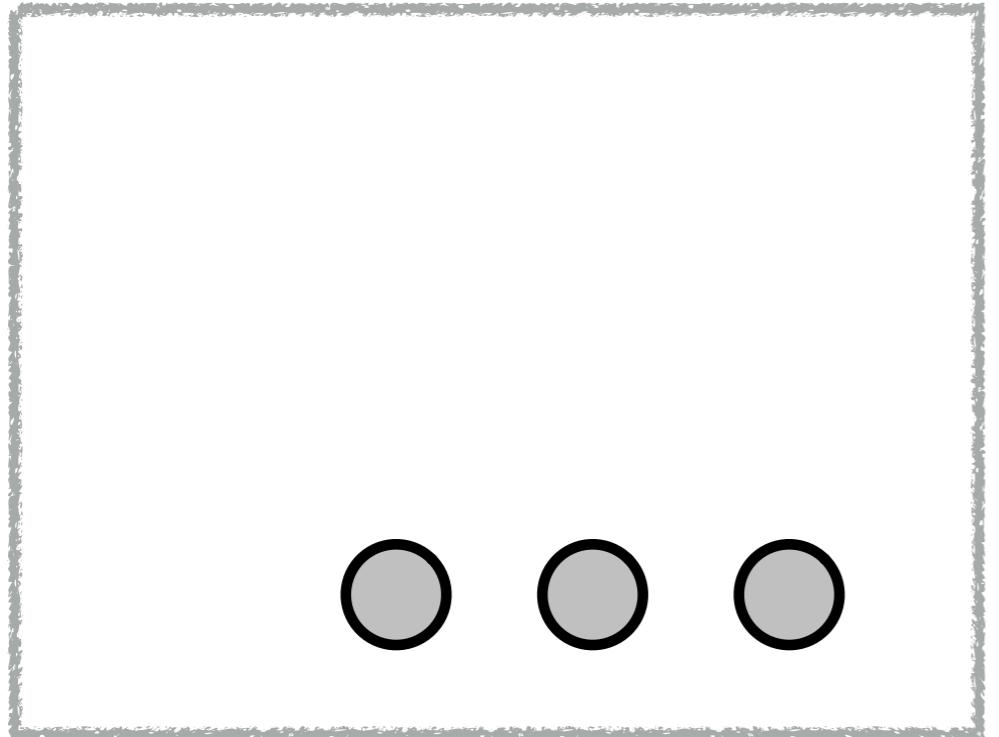


Inference?

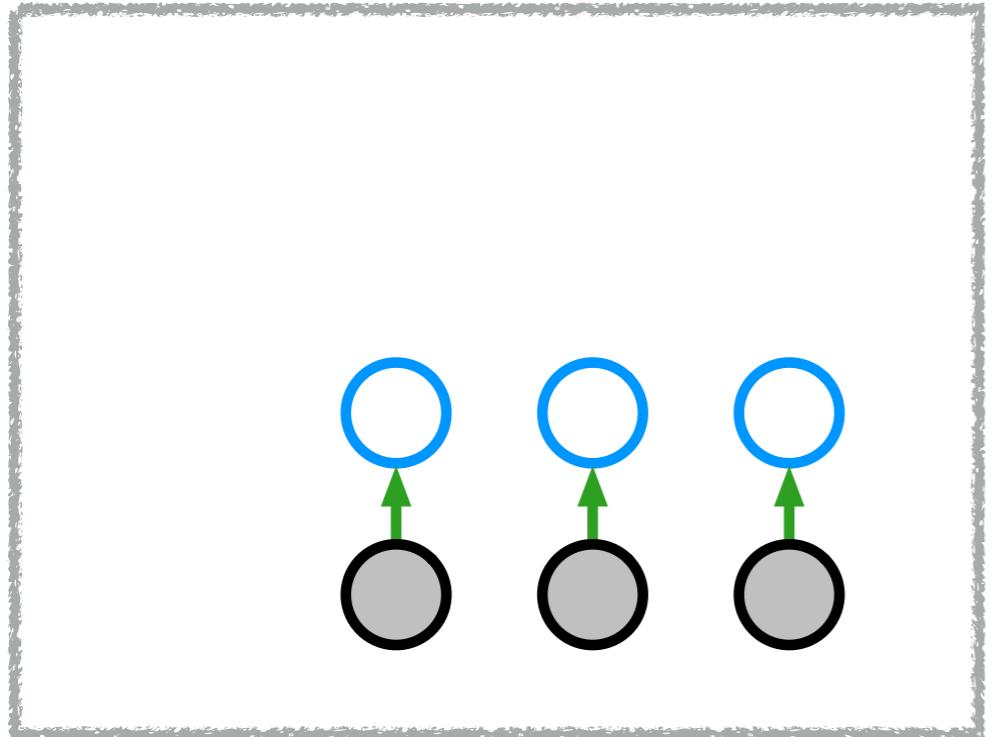
**SVAEs:** recognition networks output conjugate potentials,  
then apply fast graphical model algorithms



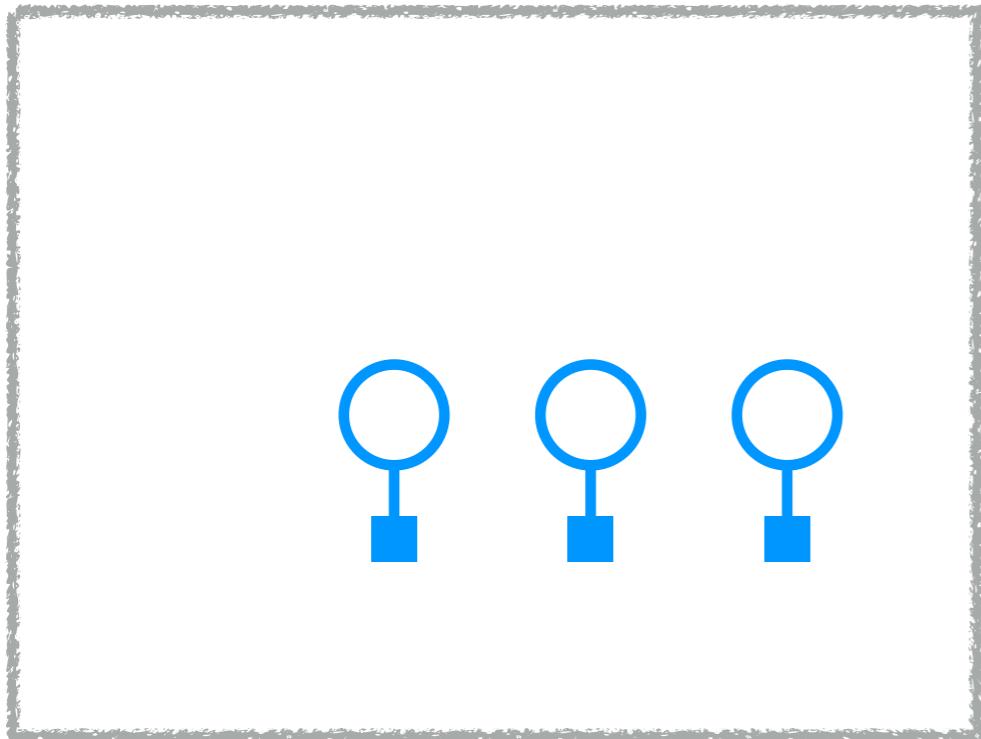
## Step 1: apply recognition network



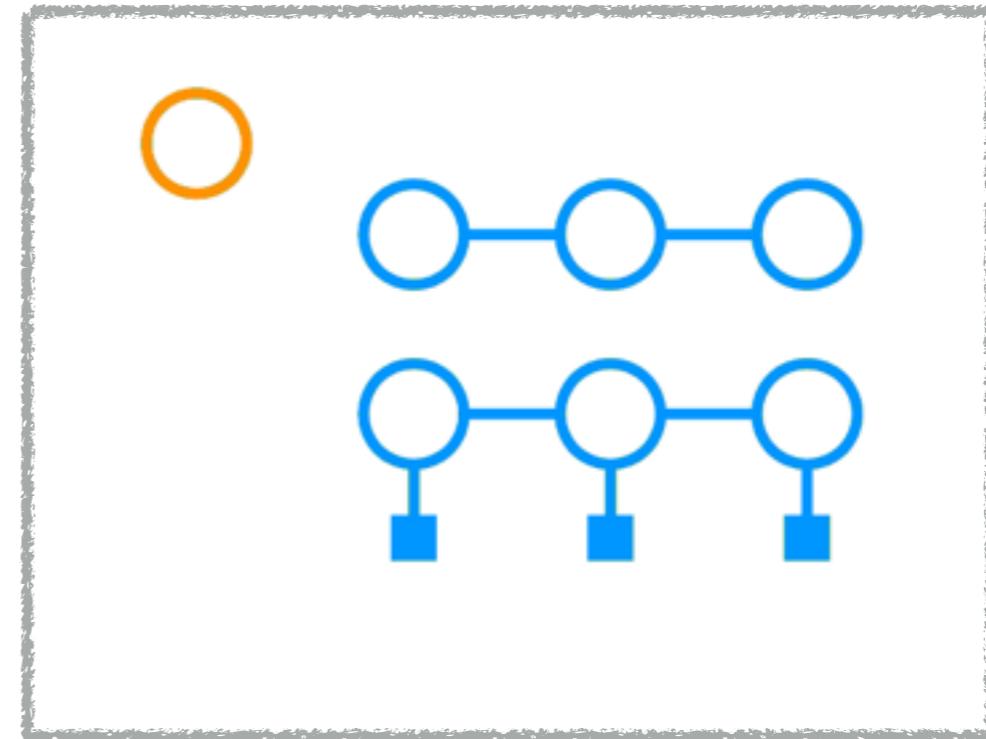
## Step 1: apply recognition network



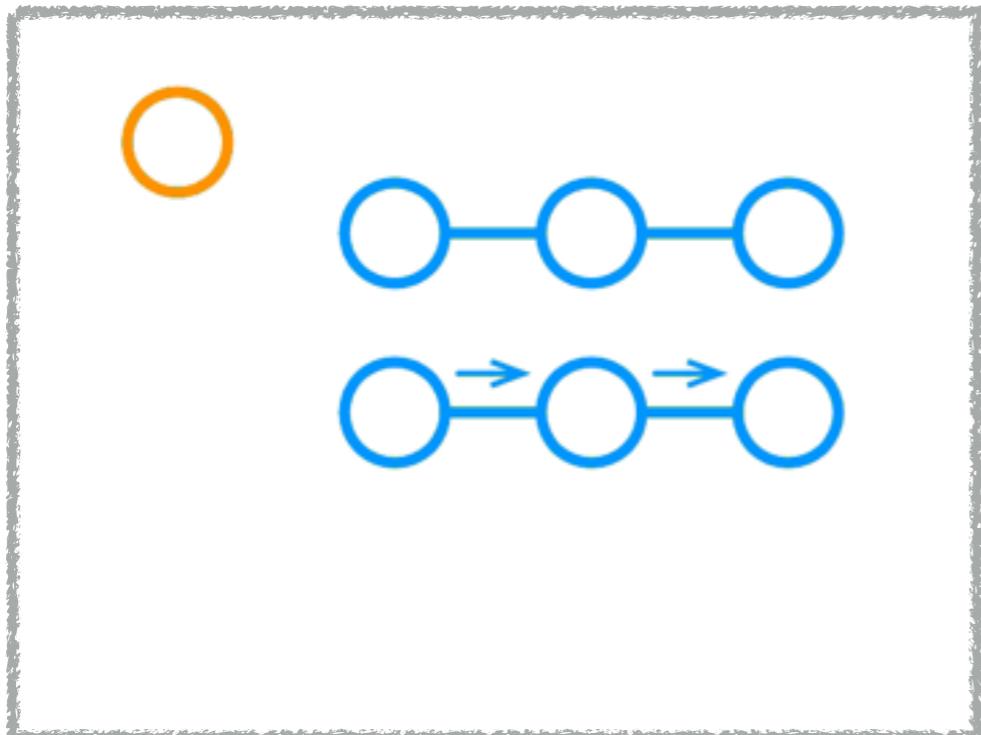
Step 1: apply recognition network



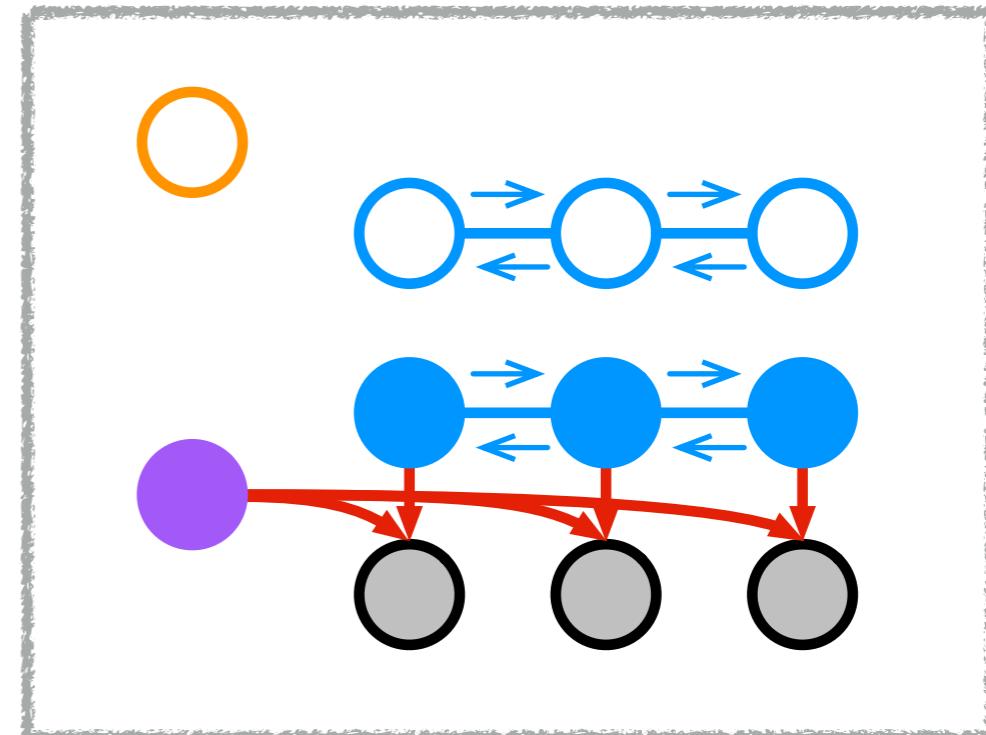
Step 2: run fast PGM algorithms

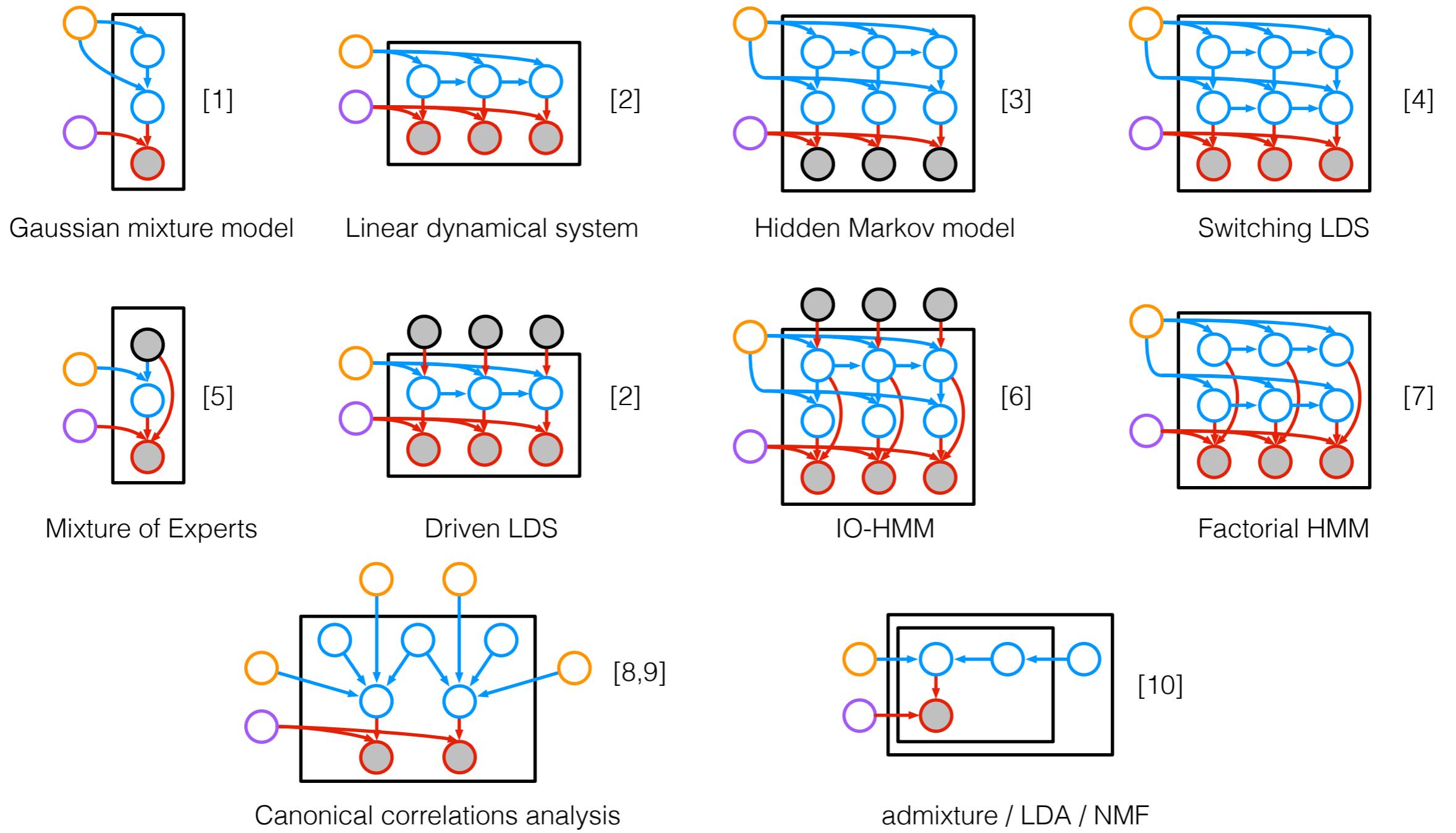


Step 3: sample, compute flat grads



Step 4: compute natural gradient

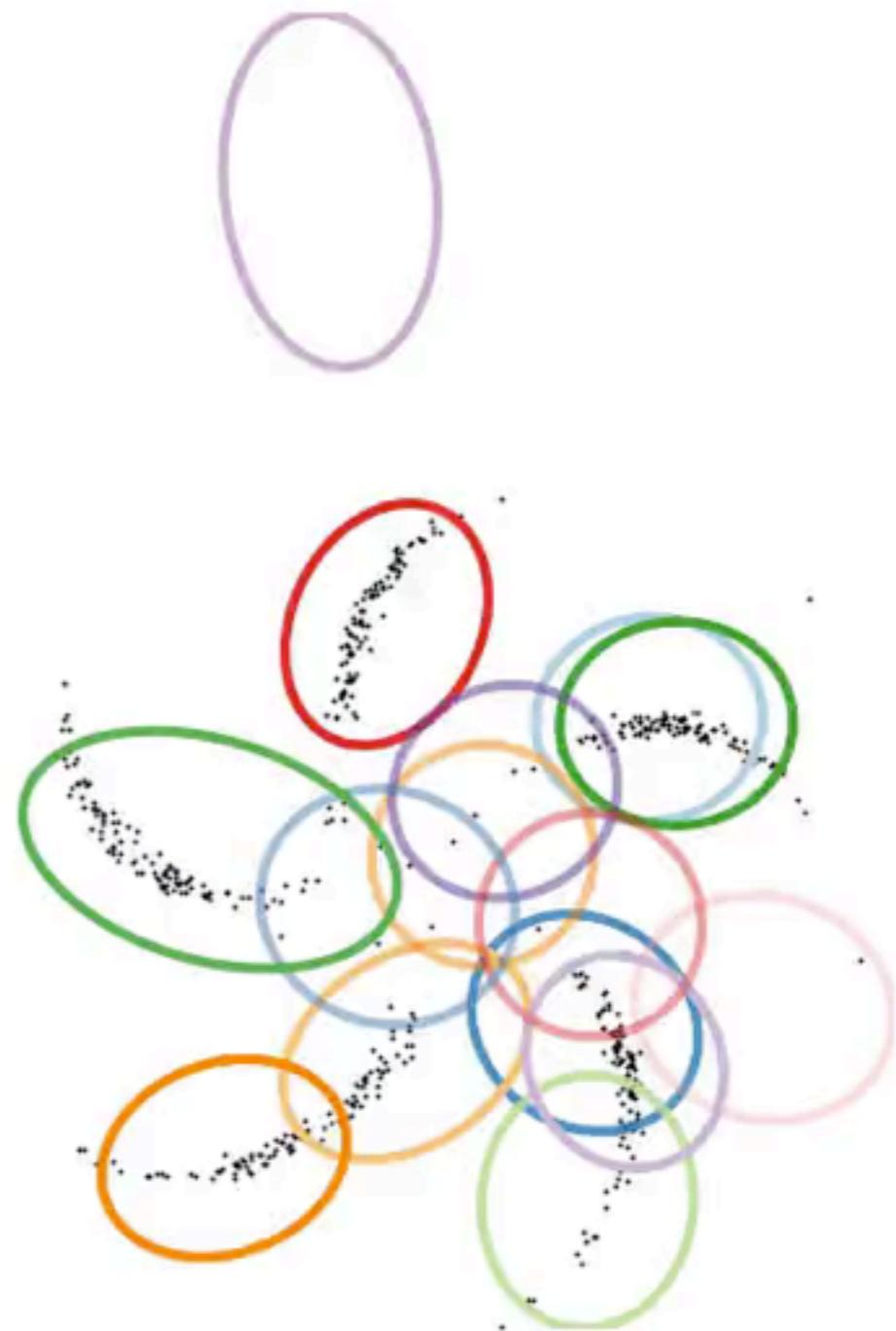




- [1] Palmer, Wipf, Kreutz-Delgado, and Rao. Variational EM algorithms for non-Gaussian latent variable models. NIPS 2005.
- [2] Ghahramani and Beal. Propagation algorithms for variational Bayesian learning. NIPS 2001.
- [3] Beal. Variational algorithms for approximate Bayesian inference, Ch. 3. U of 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.

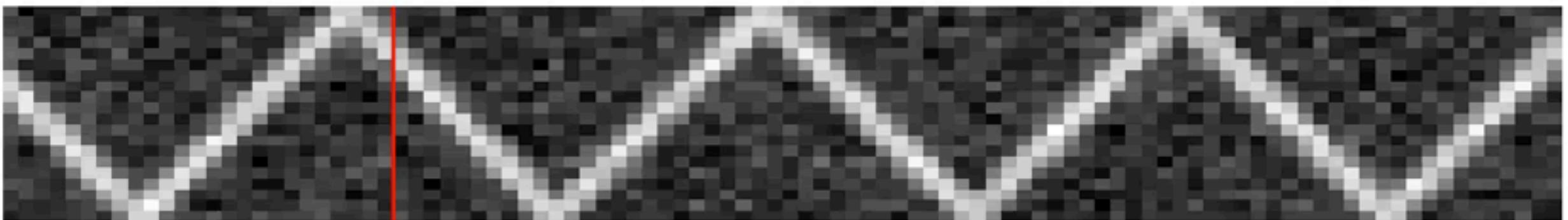


data space



latent space

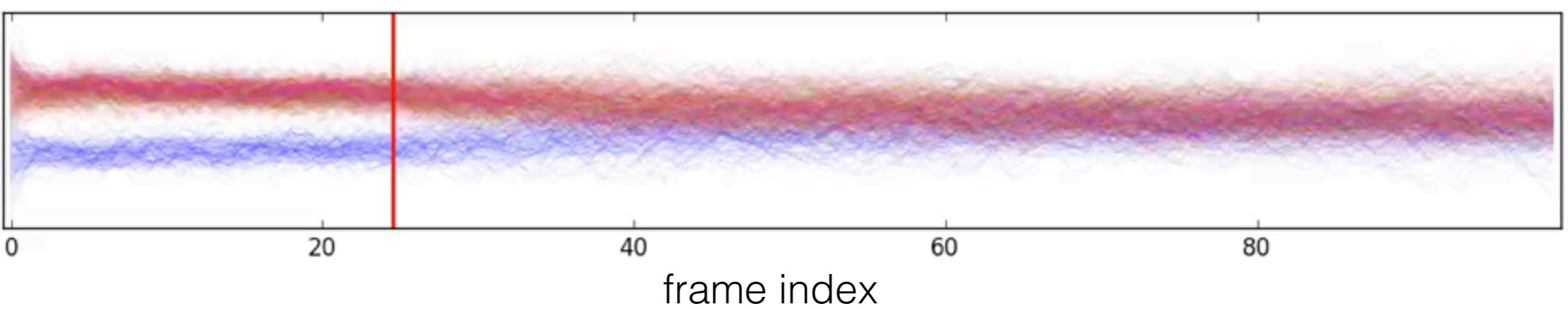
data

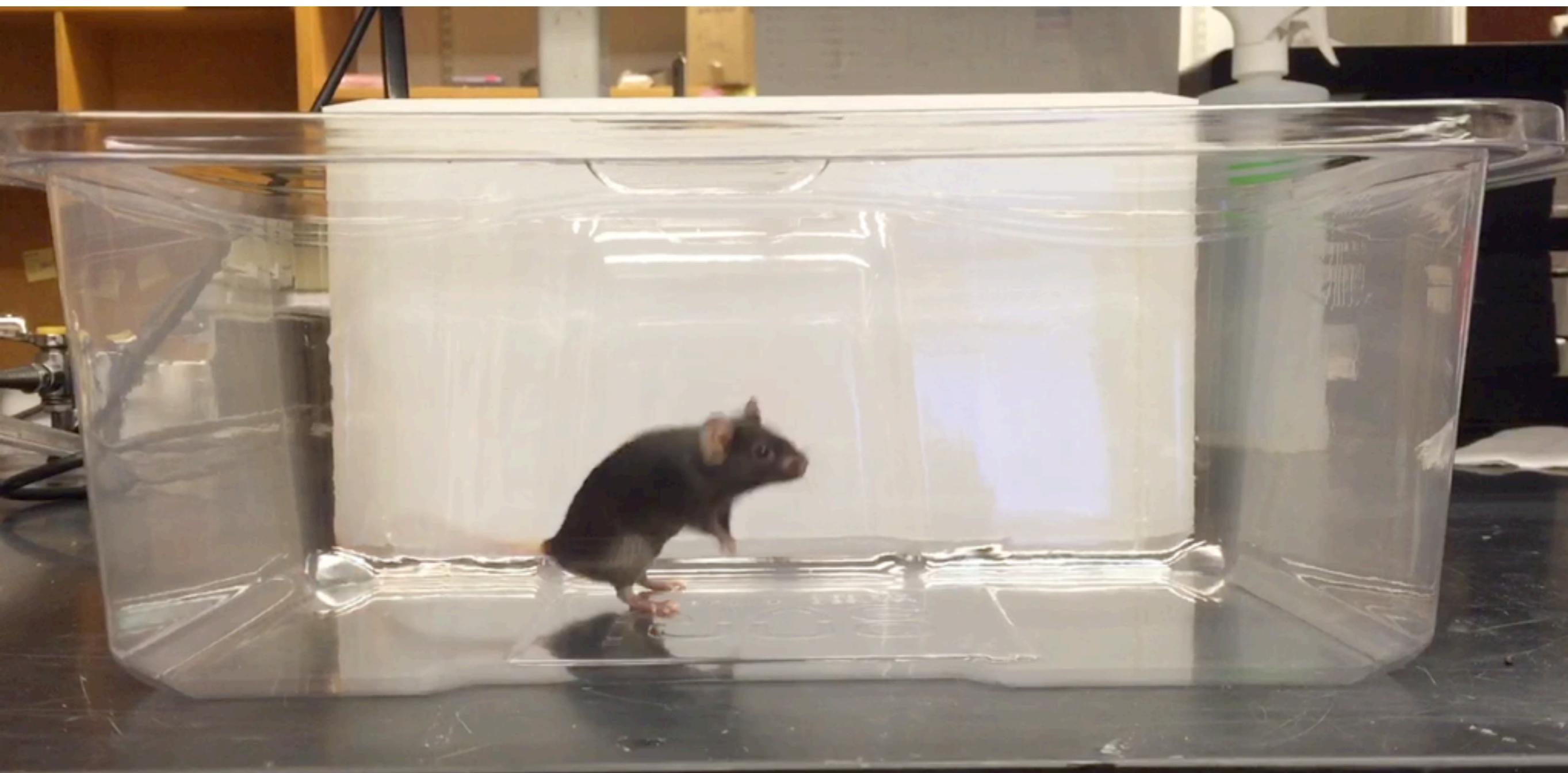


predictions

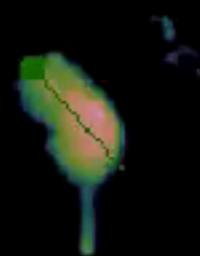
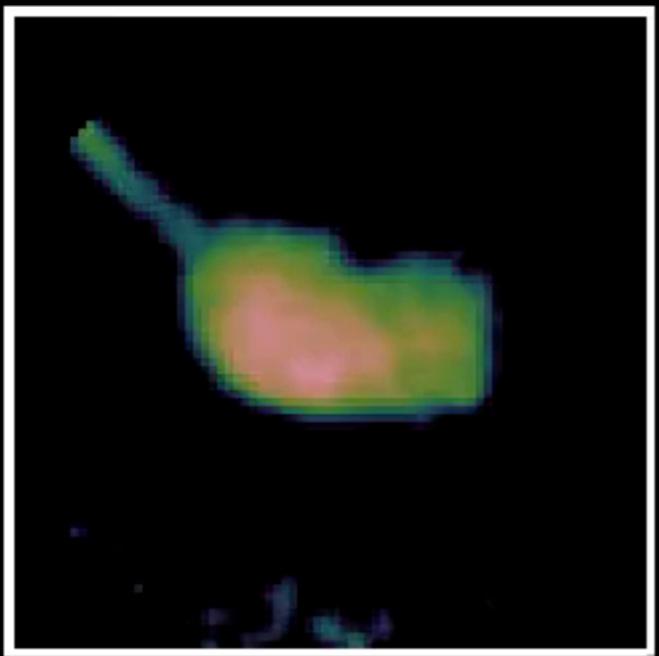


latent states



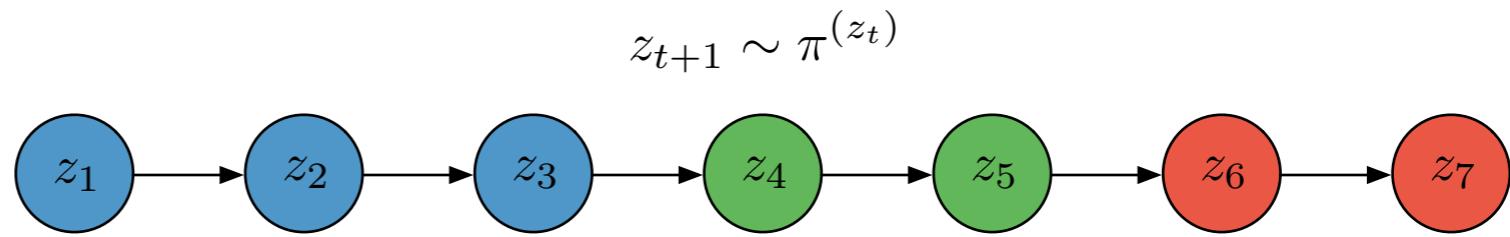


Frame 0



$$\pi = \begin{bmatrix} \textcolor{blue}{\blacksquare} & \textcolor{red}{\blacksquare} & \textcolor{green}{\blacksquare} \\ \hline \textcolor{blue}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \textcolor{red}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \textcolor{green}{\blacksquare} & \textcolor{white}{\rule{0pt}{10pt}} & \textcolor{white}{\rule{0pt}{10pt}} \\ \hline \end{bmatrix} \quad \begin{array}{c} \textcolor{blue}{\blacksquare} \\ \textcolor{red}{\blacksquare} \\ \textcolor{green}{\blacksquare} \end{array} \quad \begin{array}{c} \textcolor{blue}{\blacksquare} \\ \textcolor{red}{\blacksquare} \\ \textcolor{green}{\blacksquare} \end{array}$$

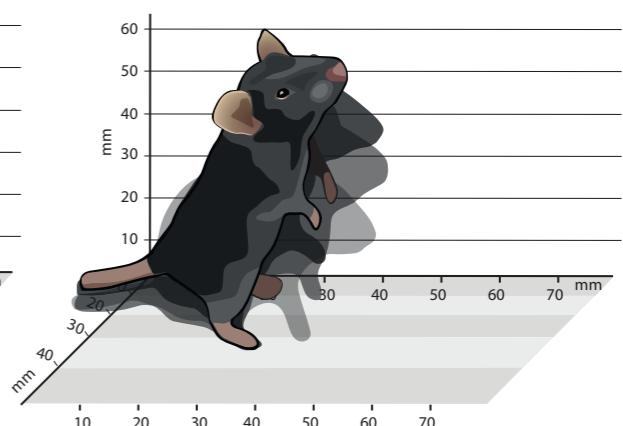
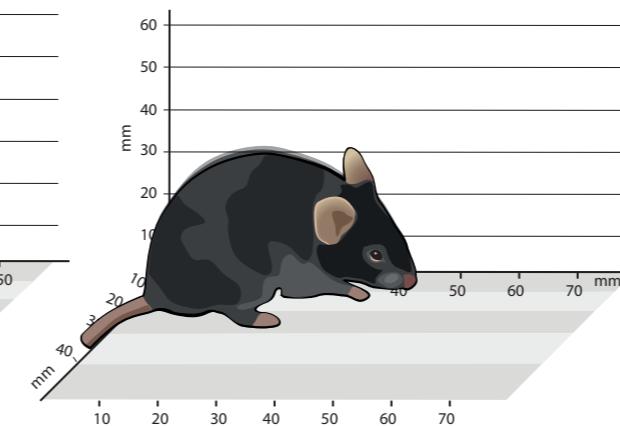
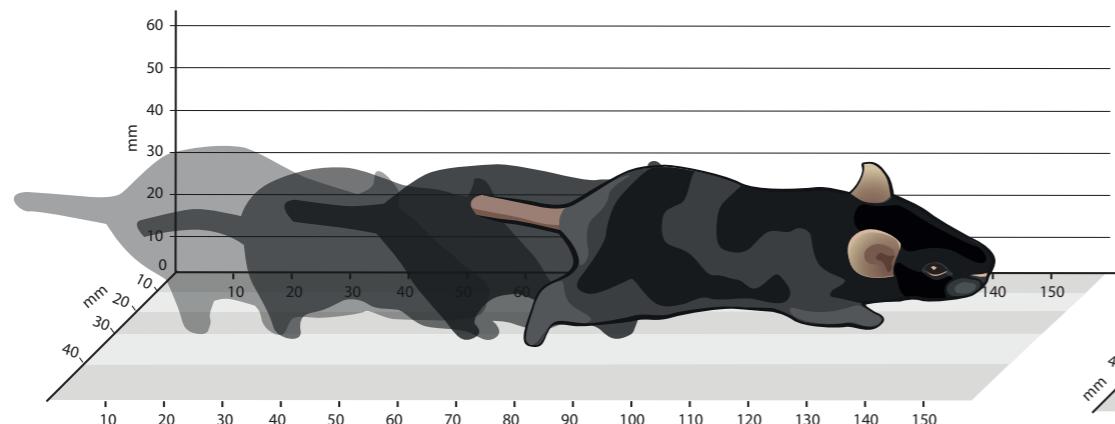
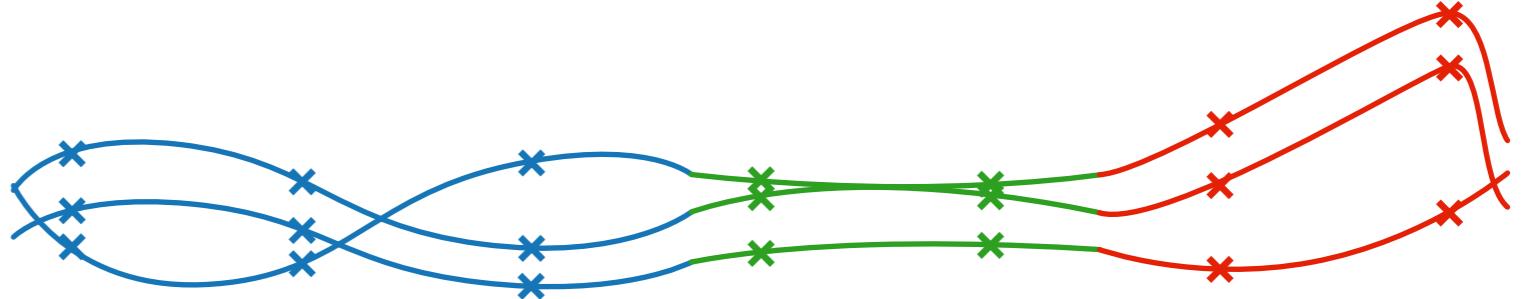
$$\pi^{(1)} \quad \pi^{(2)} \quad \pi^{(3)}$$

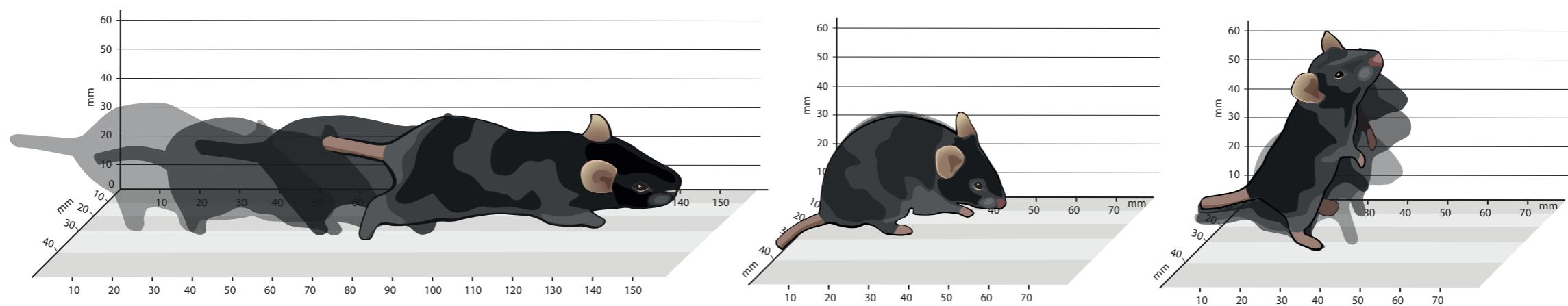
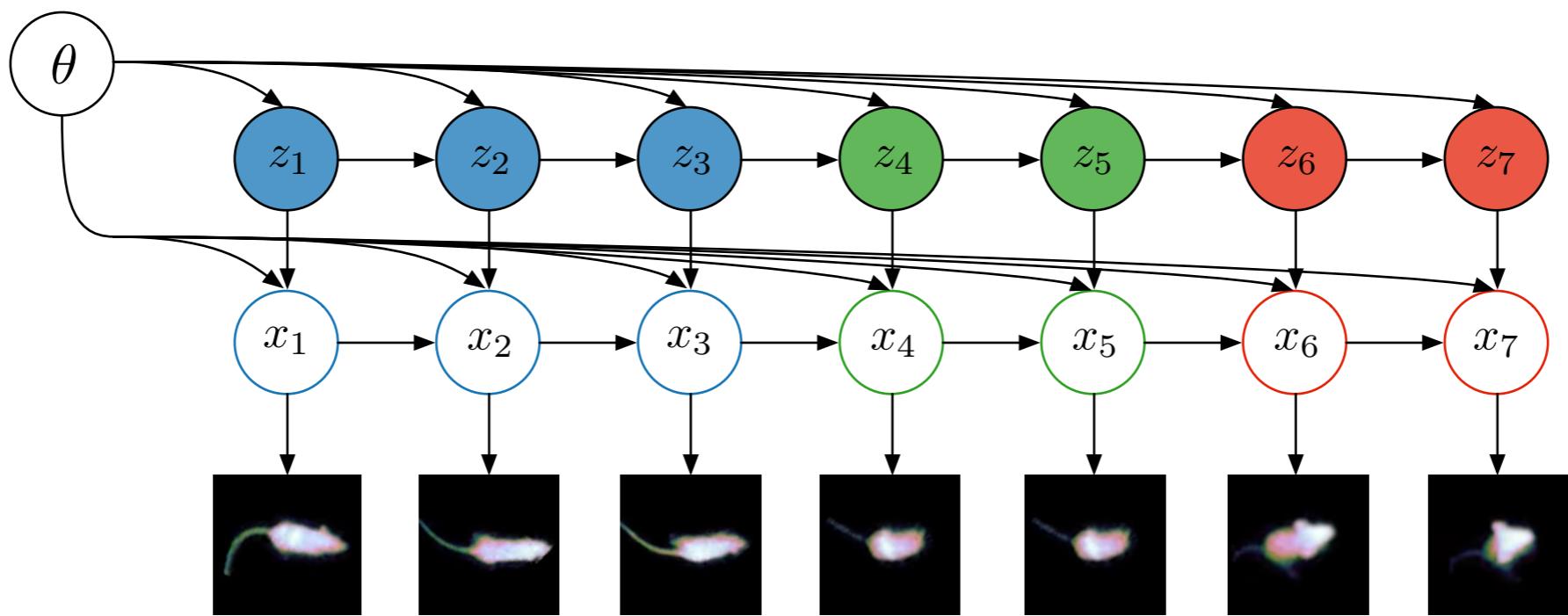


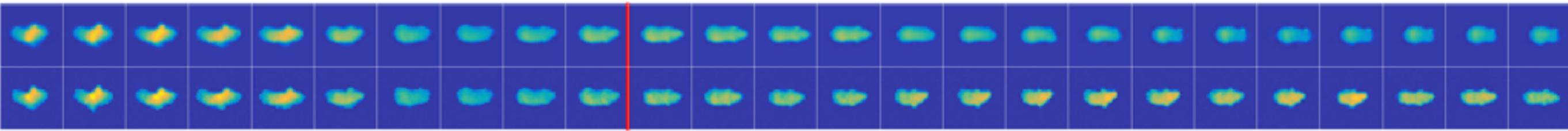
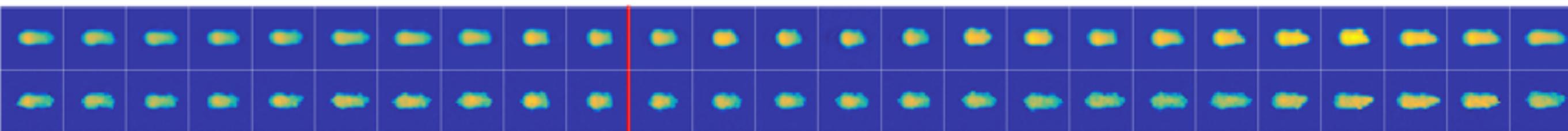
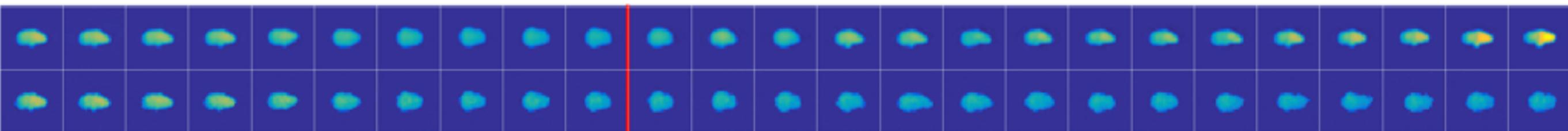
$$A^{(1)} \quad A^{(2)} \quad A^{(3)}$$

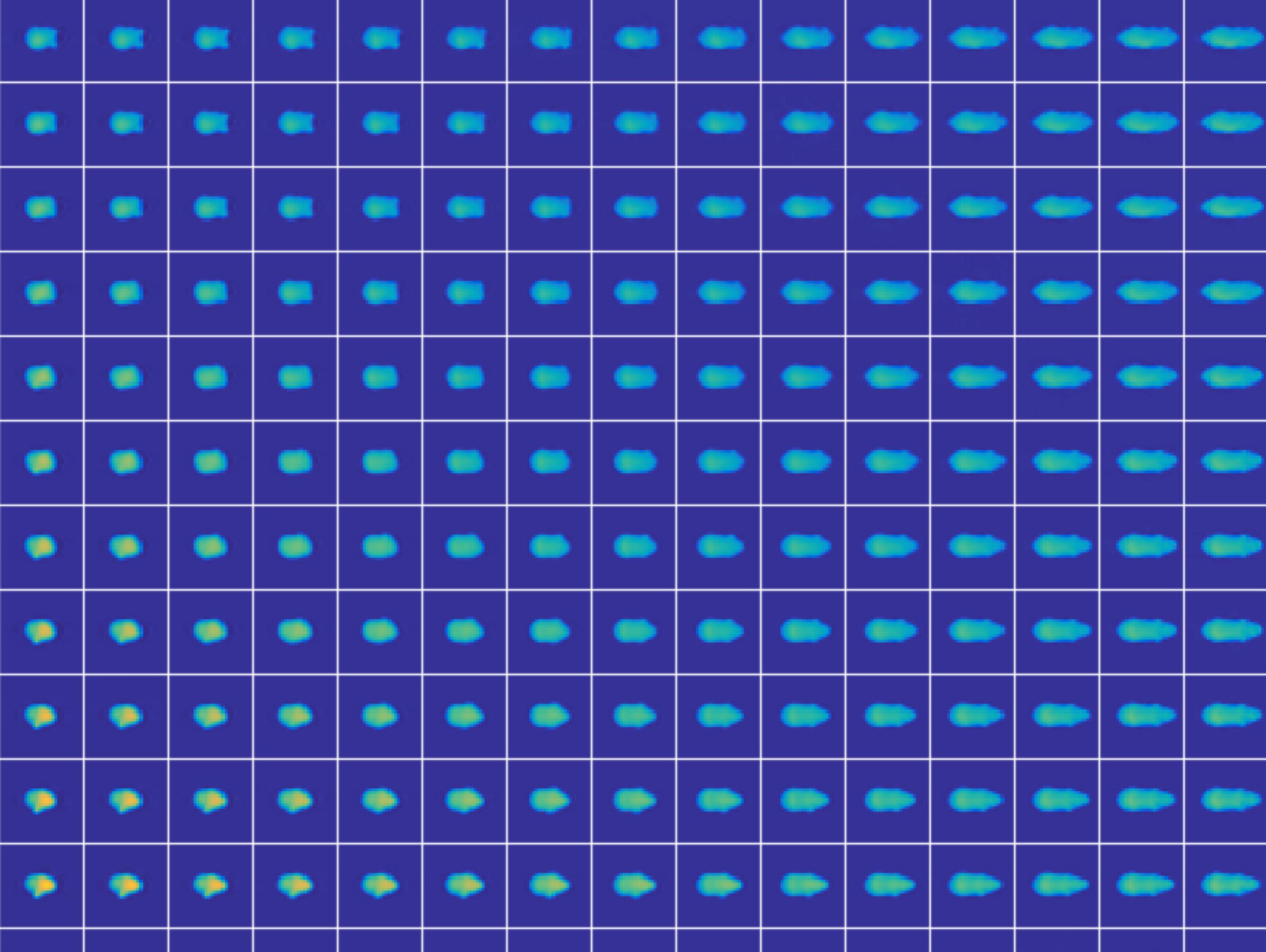
$$B^{(1)} \quad B^{(2)} \quad B^{(3)}$$

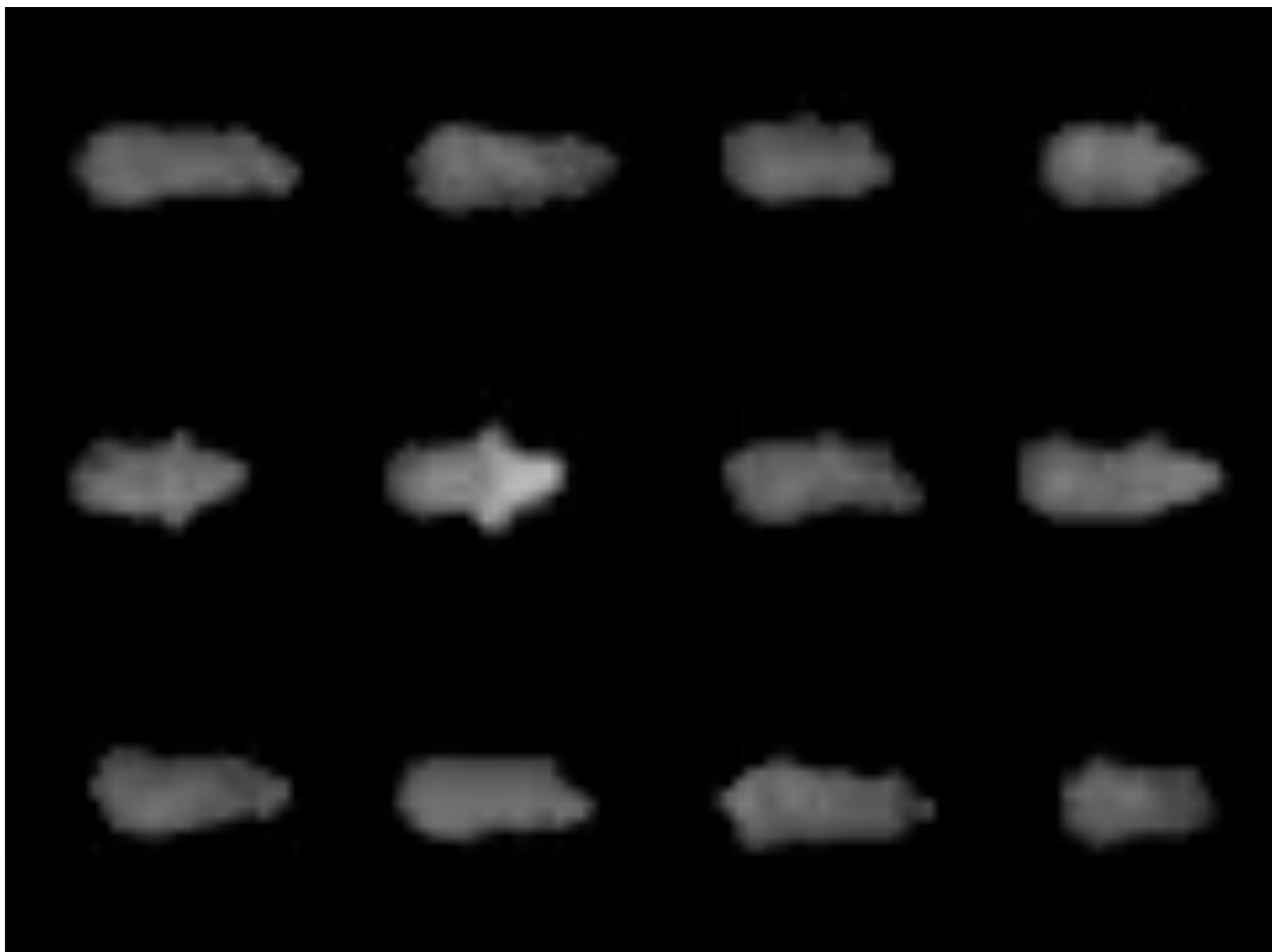
$$x_{t+1} = A^{(z_t)} x_t + B^{(z_t)} u_t \quad u_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I)$$



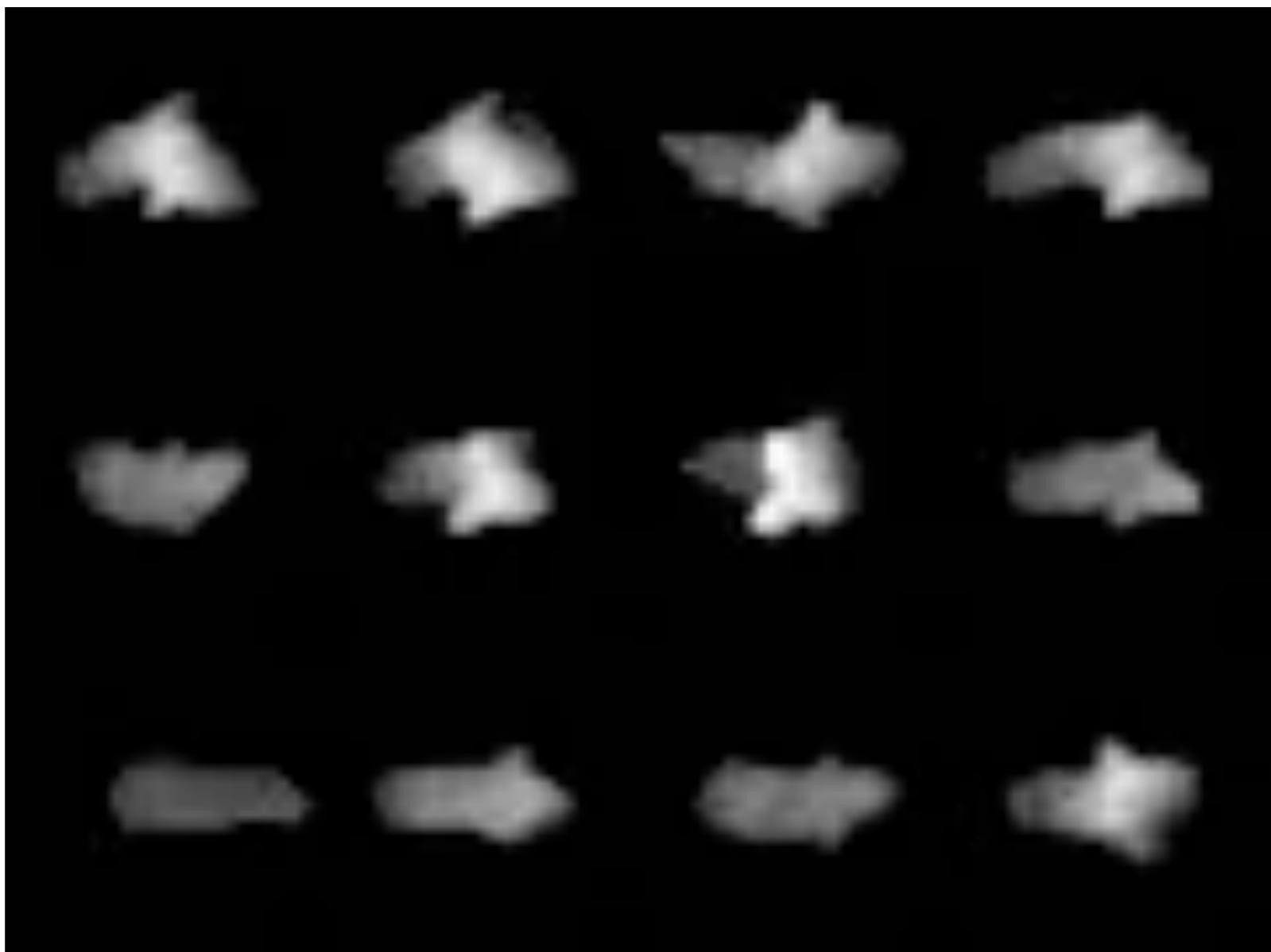




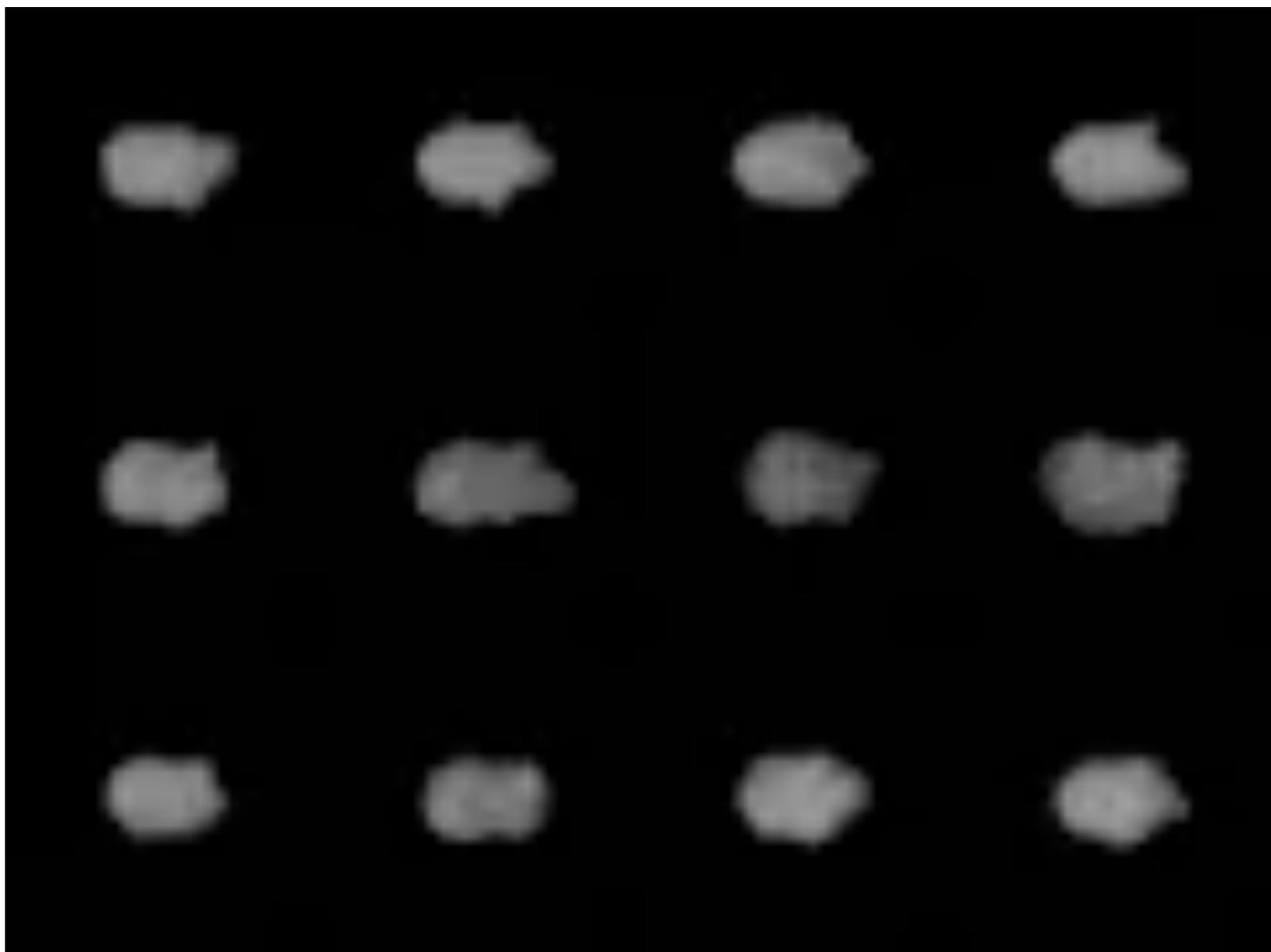




rearing up



fall from rear



grooming

# Limitations and future work

capacity

- How expressive is latent linear structure?
  - word embeddings [1], analogical reasoning in image models
  - SVAE can use nonlinear latent structure

complexity

- PGMs get complicated
  - SVAE keeps complexity modular

future work

- model-based reinforcement learning
- automatic structure search [2,3]
- semi-supervised applications

[1] Hashimoto, Alvarez-Melis, and Jaakkola, Word, graph and manifold embedding from Markov processes, Preprint 2015.

[2] Grosse et al., Exploiting compositionality to explore a large space of model structures, UAI 2012.

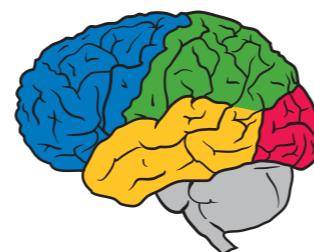
[3] Duvenaud et al., Structure discovery in nonparametric regression through compositional kernel search, ICML 2013.

Matt Johnson, David Duvenaud, Alex Wiltschko, Bob Datta, Ryan Adams

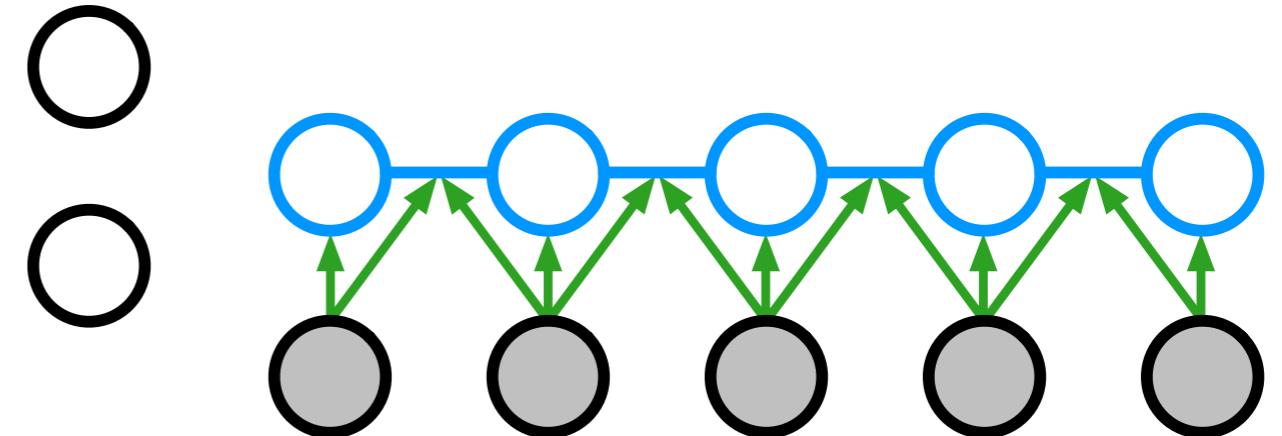


# Thanks!

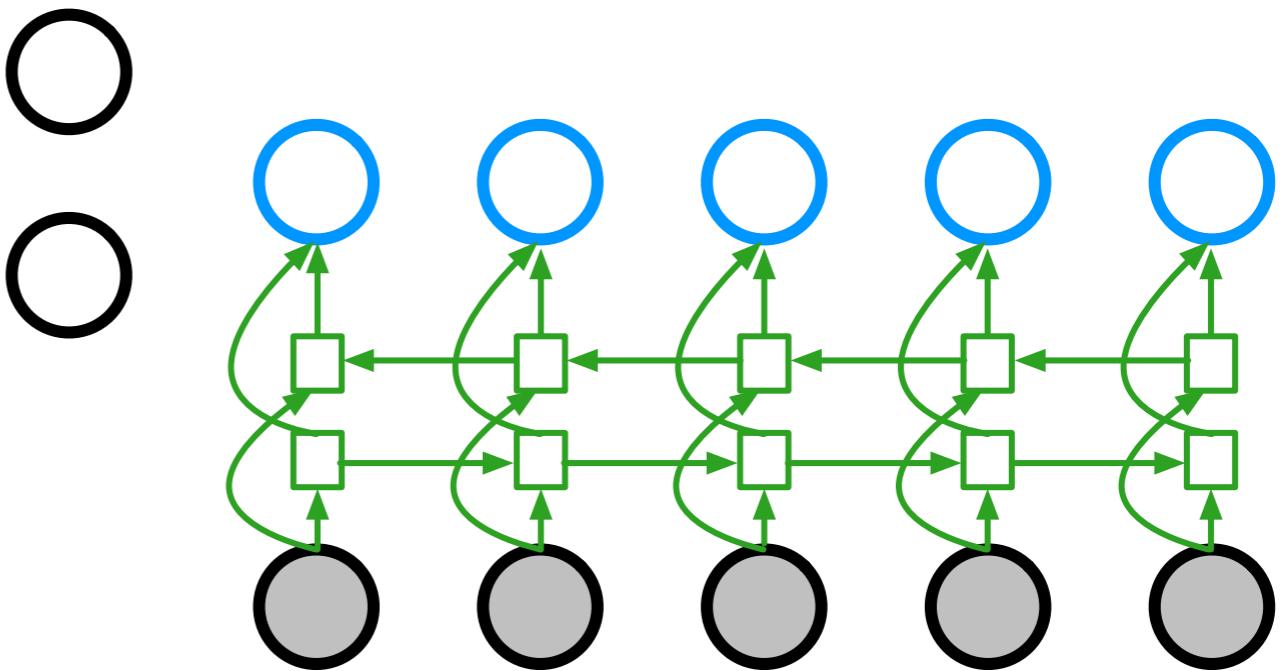
[github.com/mattjj/svae](https://github.com/mattjj/svae)



$$\begin{aligned} & \mu_t(y_t; \phi_\mu) \\ [1,2] \quad & J_{t,t}(y_t; \phi_D) \\ & J_{t,t+1}(y_t, y_{t+1}; \phi_B) \end{aligned}$$

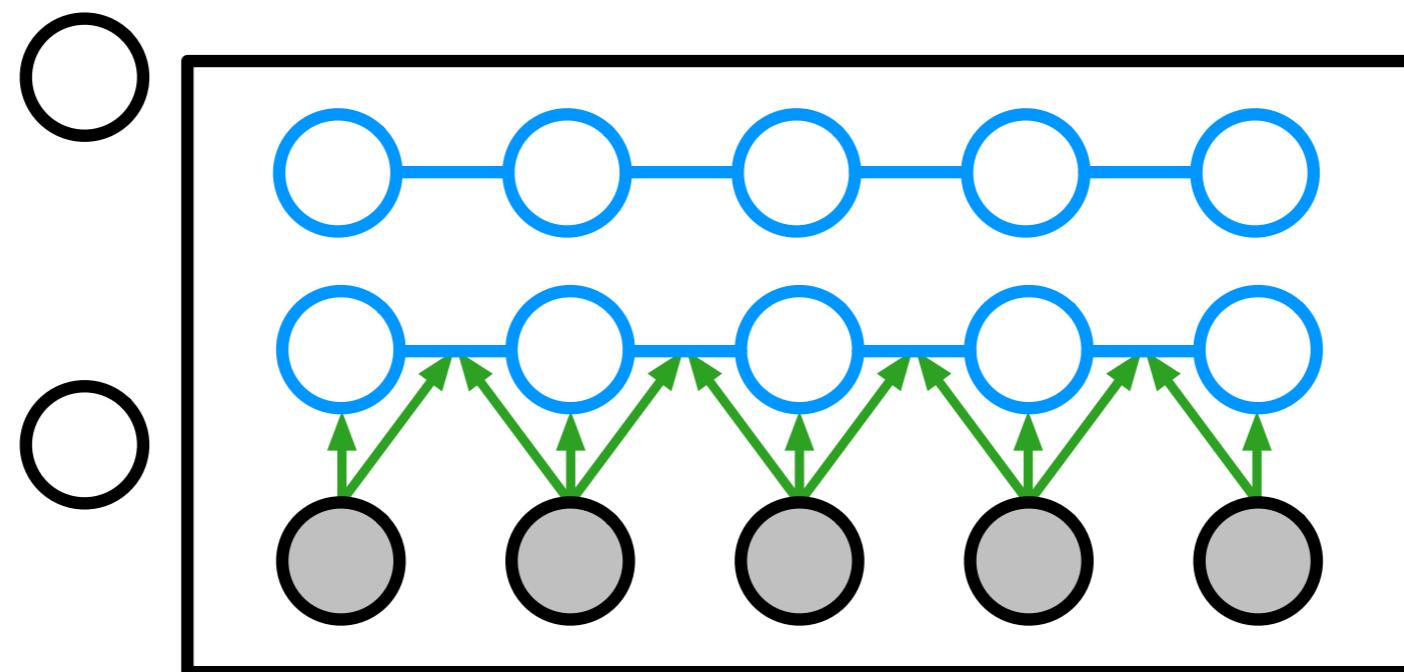


$$\begin{aligned} & \mu_t(y_{1:T}, \hat{x}_{t-1}; \phi) \\ [3] \quad & \Sigma_t(y_{1:T}, \hat{x}_{t-1}; \phi) \end{aligned}$$



- [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.

SVAEs can use any inference network architecture



- [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.

Per-variable recognition nets allow arbitrary inference queries

