

Machine Learning Methods for Neural Data Analysis

Lecture 13: Switching linear dynamical systems

Scott Linderman

STATS 220/320 (*NBIO220, CS339N*). Winter 2021.

Announcements

- **Lab 7 Errata:**
 - Don't change `log_likes` in-place in Problem 1a. The test function reuses it for computing the reference answers.
 - Typo in the commented functions for downloading crowd movies.
 - How to fill in the feature vector for first couple time steps?
- **Updated proposals** due this Friday.
 - Looking for progress on downloading and extracting the data.

Agenda

- Switching linear dynamical systems (SLDS)
- Hardness of exact EM for SLDS
- Approximate message passing

Recap: Gaussian HMM

Generative Model:

$$\begin{aligned} z_1 &\sim \text{Cat}(\pi), \\ z_t \mid z_{t-1} &\sim \text{Cat}(P_{z_{t-1}}), \quad \text{for } t = 2, \dots, T. \\ x_t \mid z_t &\sim \mathcal{N}(b_{z_t}, Q_{z_t}) \quad \text{for } t = 1, \dots, T \end{aligned}$$

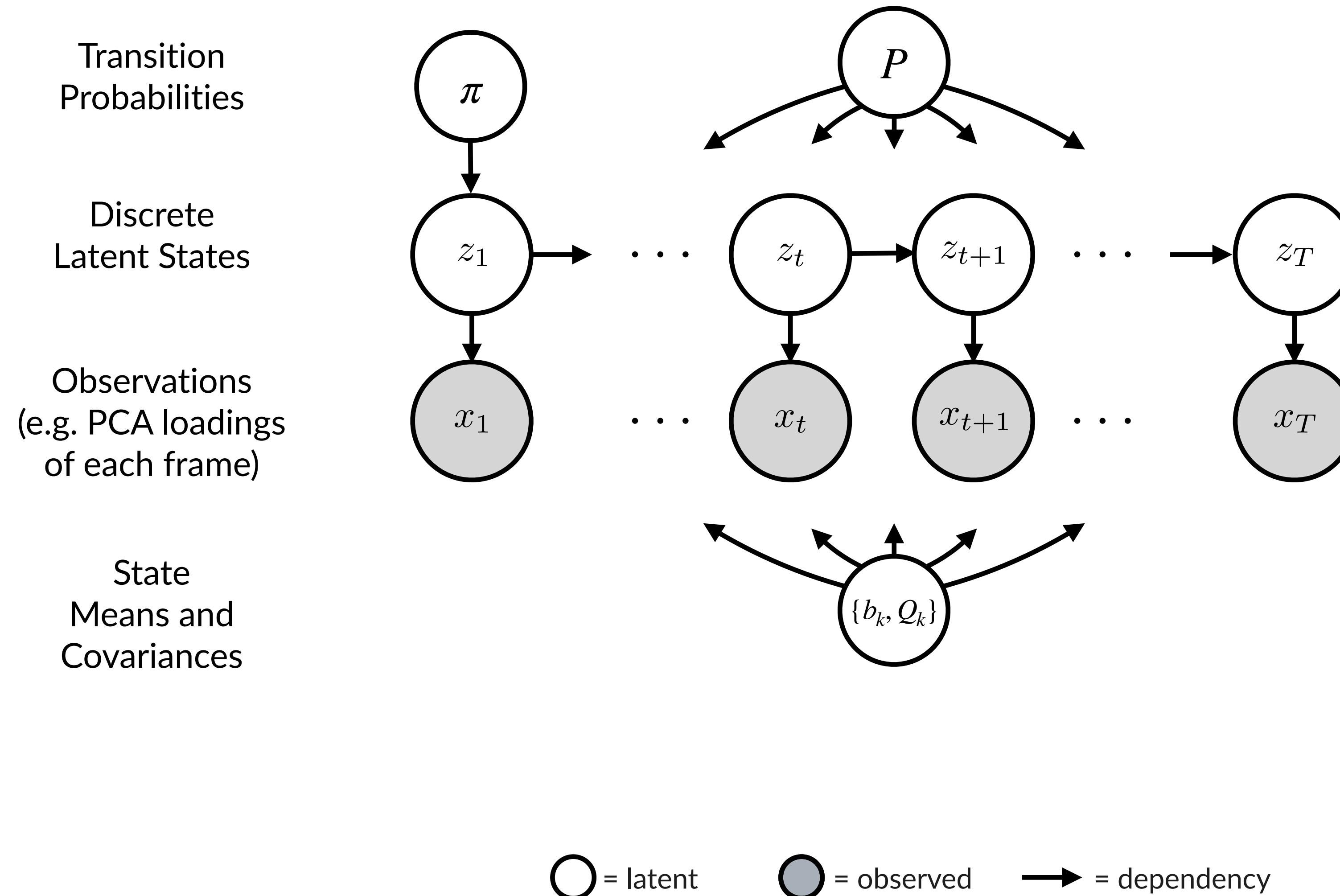
Parameters:

$$\Theta = \pi, P, \{b_k, Q_k\}_{k=1}^K$$

Joint probability:

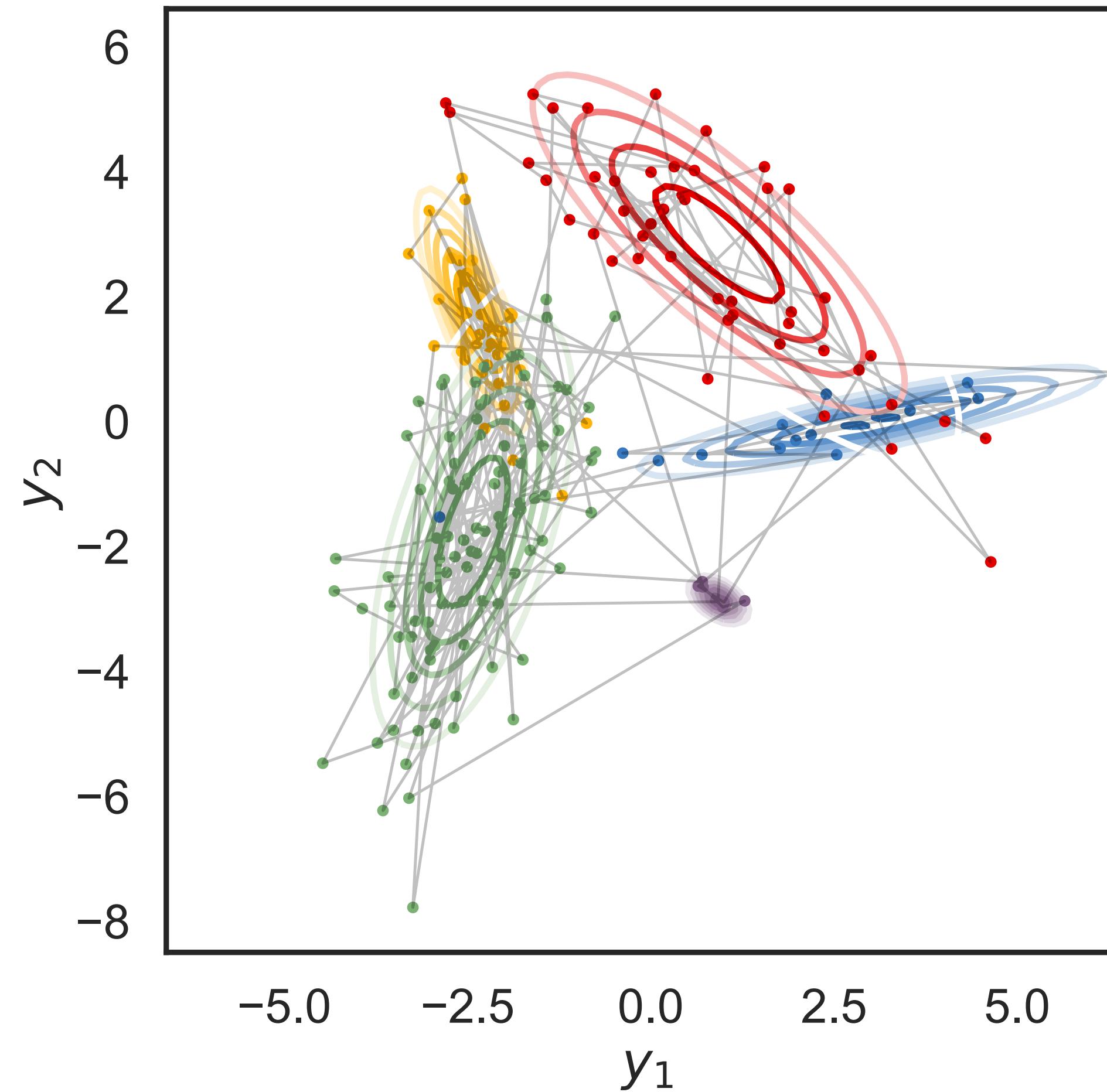
$$p(x, z \mid \Theta) = p(z_1) \prod_{t=1}^{T-1} p(z_{t+1} \mid z_t) \prod_{t=1}^T p(x_t \mid z_t)$$

Recap: The Gaussian HMM Graphical Model

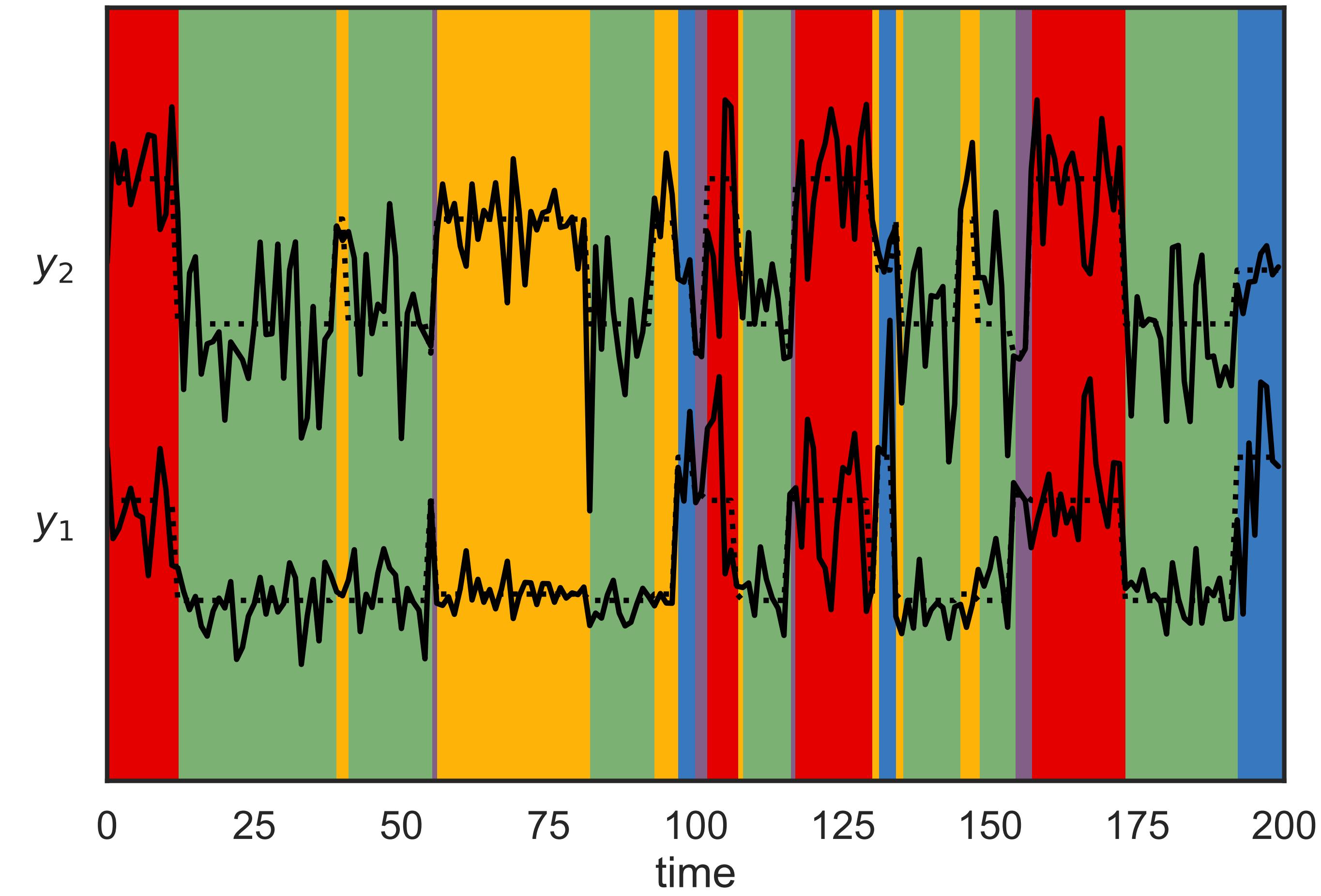


Simulated data from a Gaussian HMM

Observation Distributions



Simulated data from an HMM



Recap: Autoregressive (AR) HMM

Generative Model:

$$\begin{aligned} z_1 &\sim \text{Cat}(\pi), \\ z_t \mid z_{t-1} &\sim \text{Cat}(P_{z_{t-1}}), & \text{for } t = 2, \dots, T \\ x_1 \mid z_1 &\sim \mathcal{N}(b_{z_1}, Q_{z_1}) \\ x_t \mid x_{t-1}, z_t &\sim \mathcal{N}(A_{z_t}x_{t-1} + b_{z_t}, Q_{z_t}) & \text{for } t = 1, \dots, T \end{aligned}$$

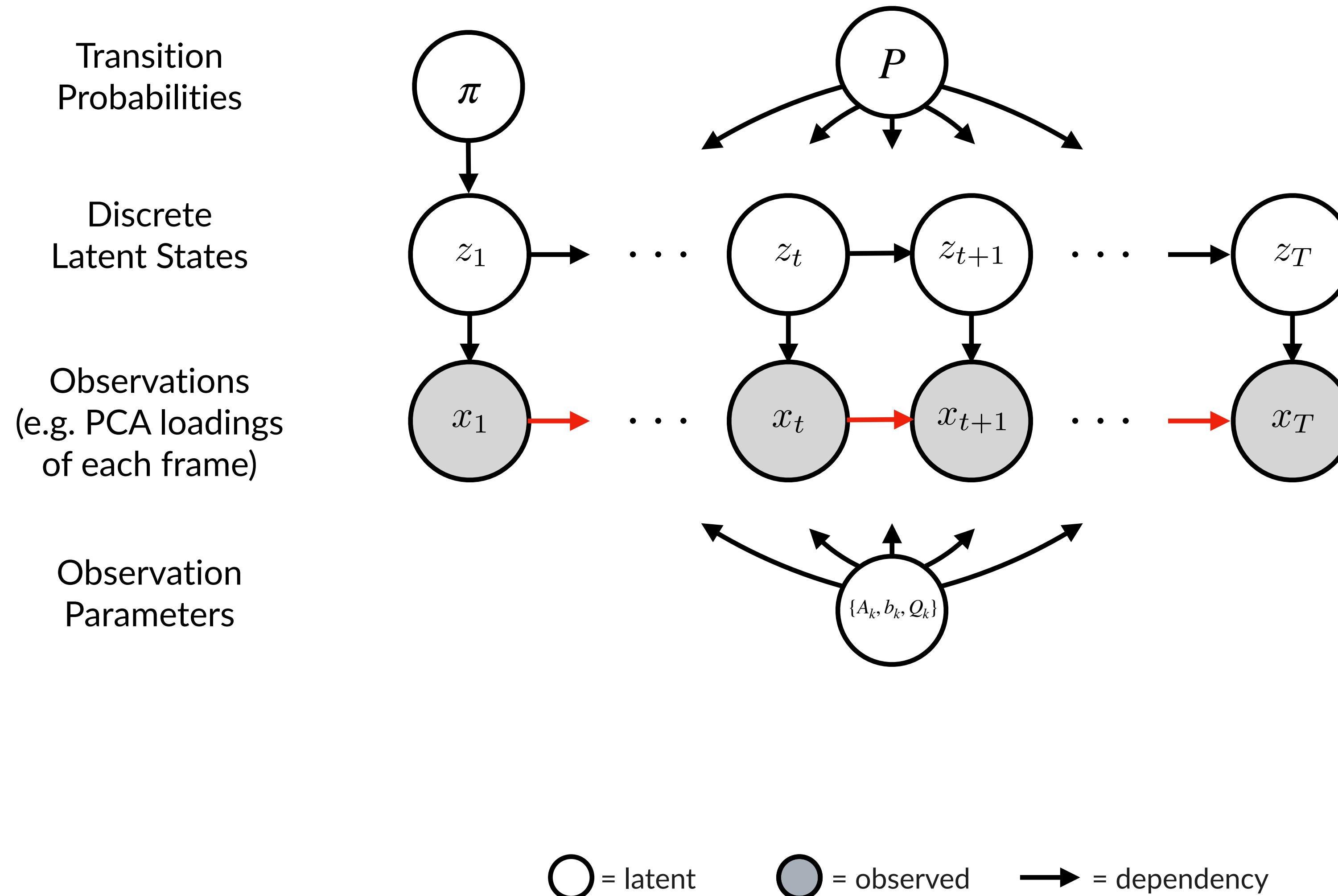
Parameters:

$$\Theta = \pi, P, \{A_k, b_k, Q_k\}_{k=1}^K$$

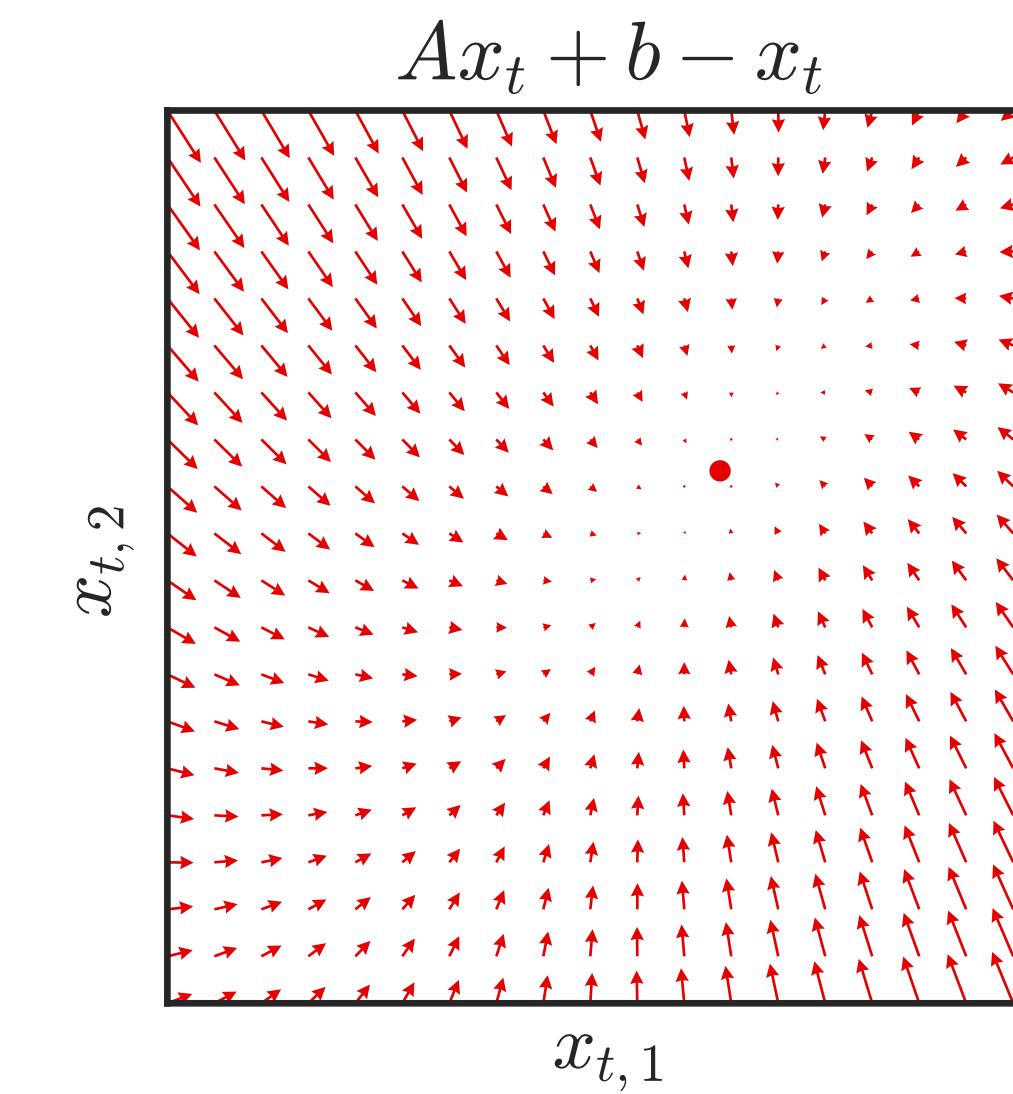
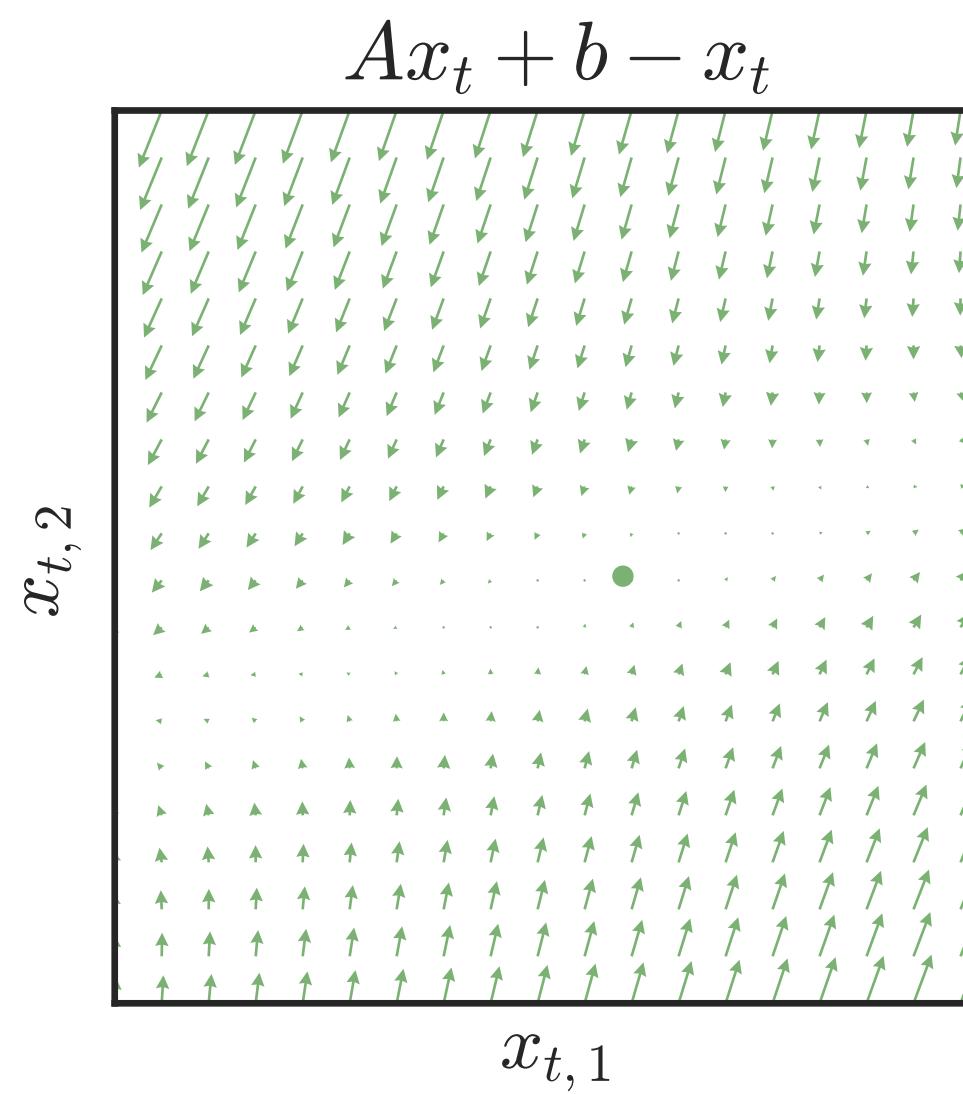
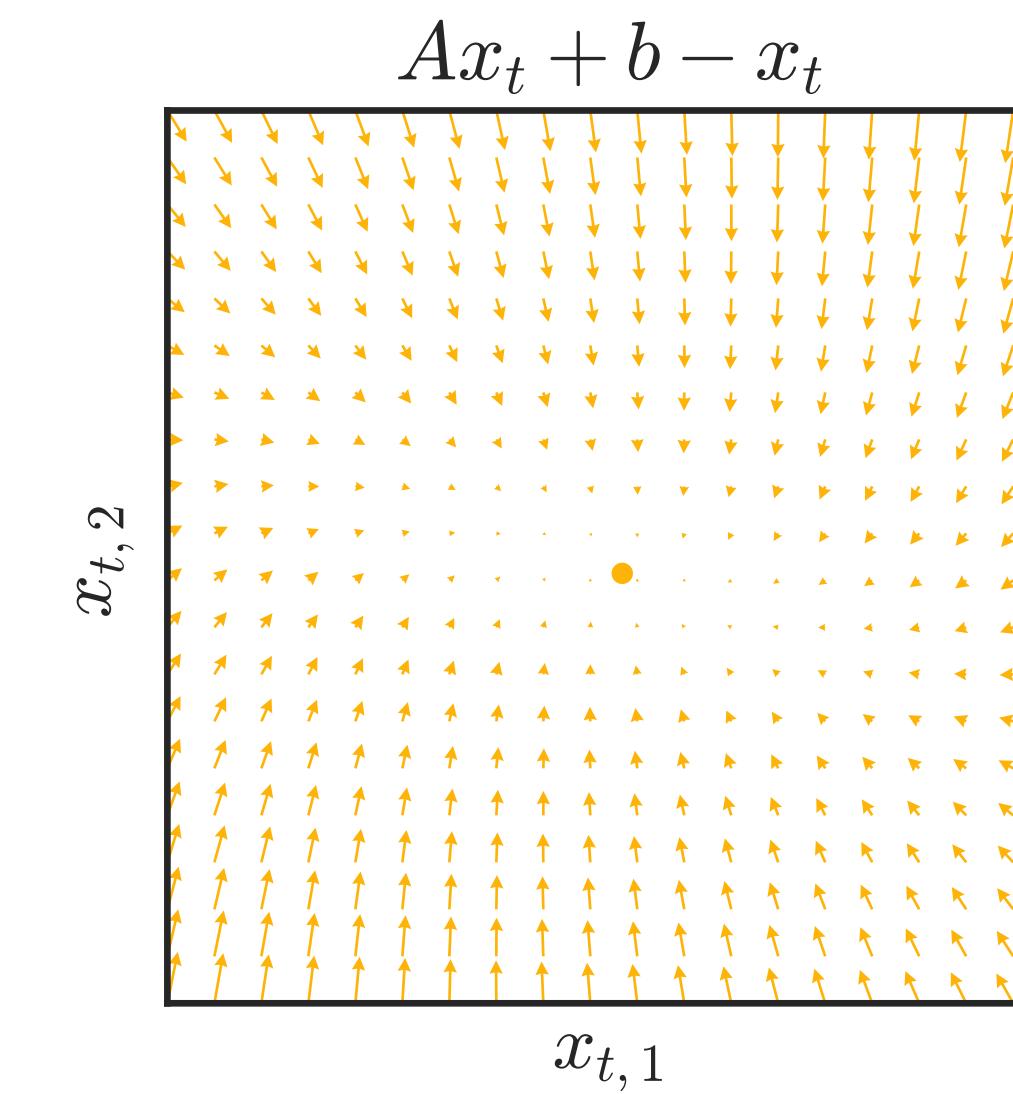
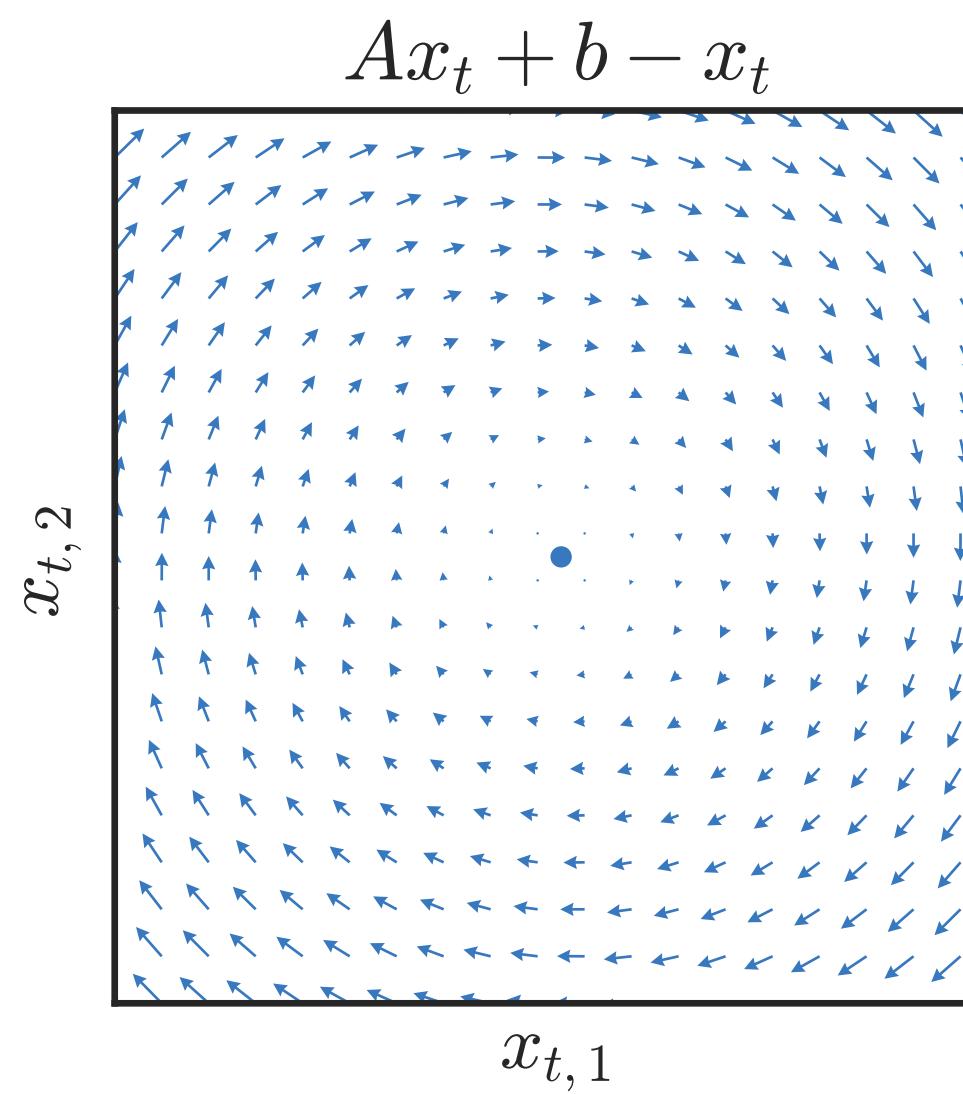
Joint probability:

$$p(x, z \mid \Theta) = p(z_1) p(x_1 \mid z_1) \prod_{t=1}^{T-1} p(z_{t+1} \mid z_t) \prod_{t=1}^T p(x_t \mid x_{t-1}, z_t)$$

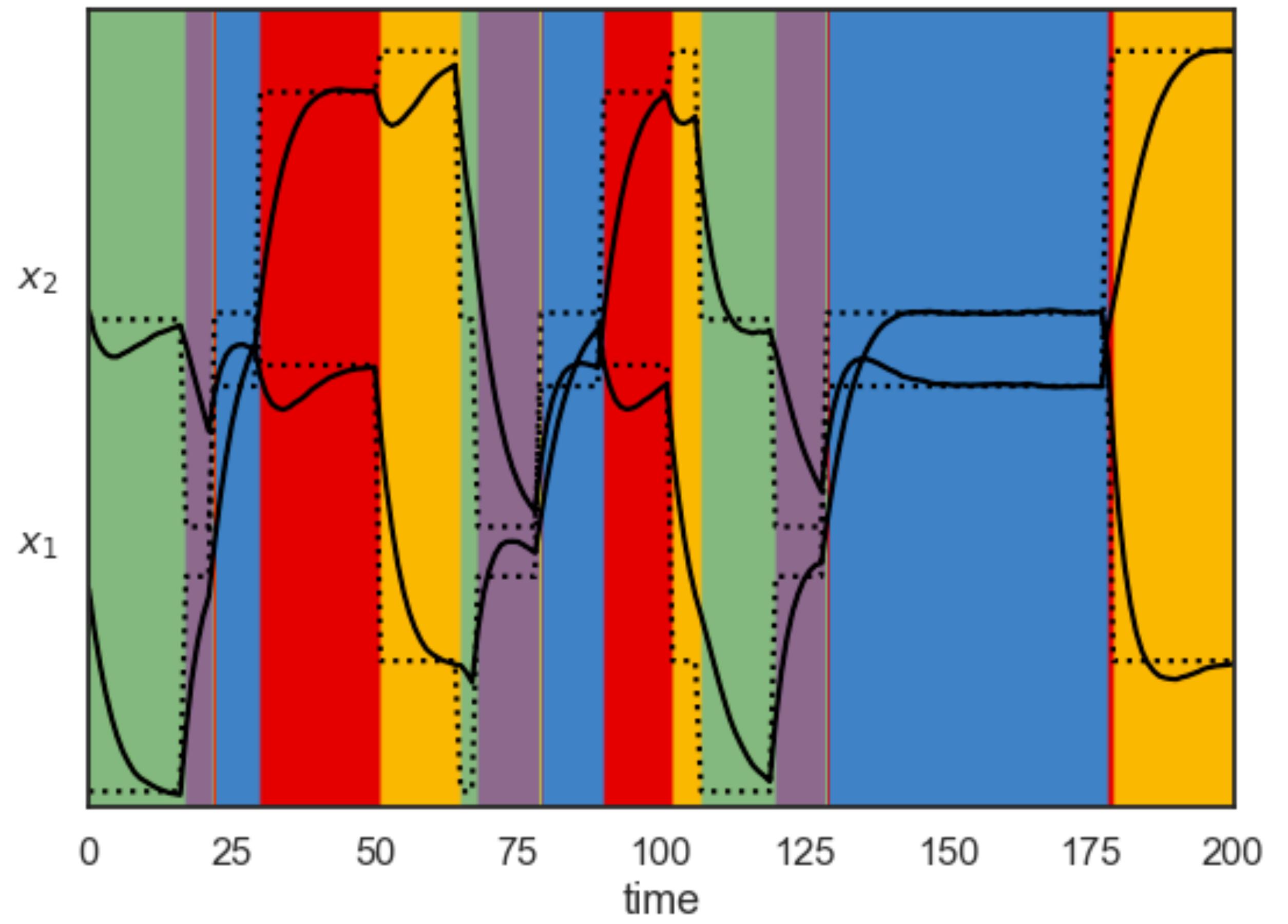
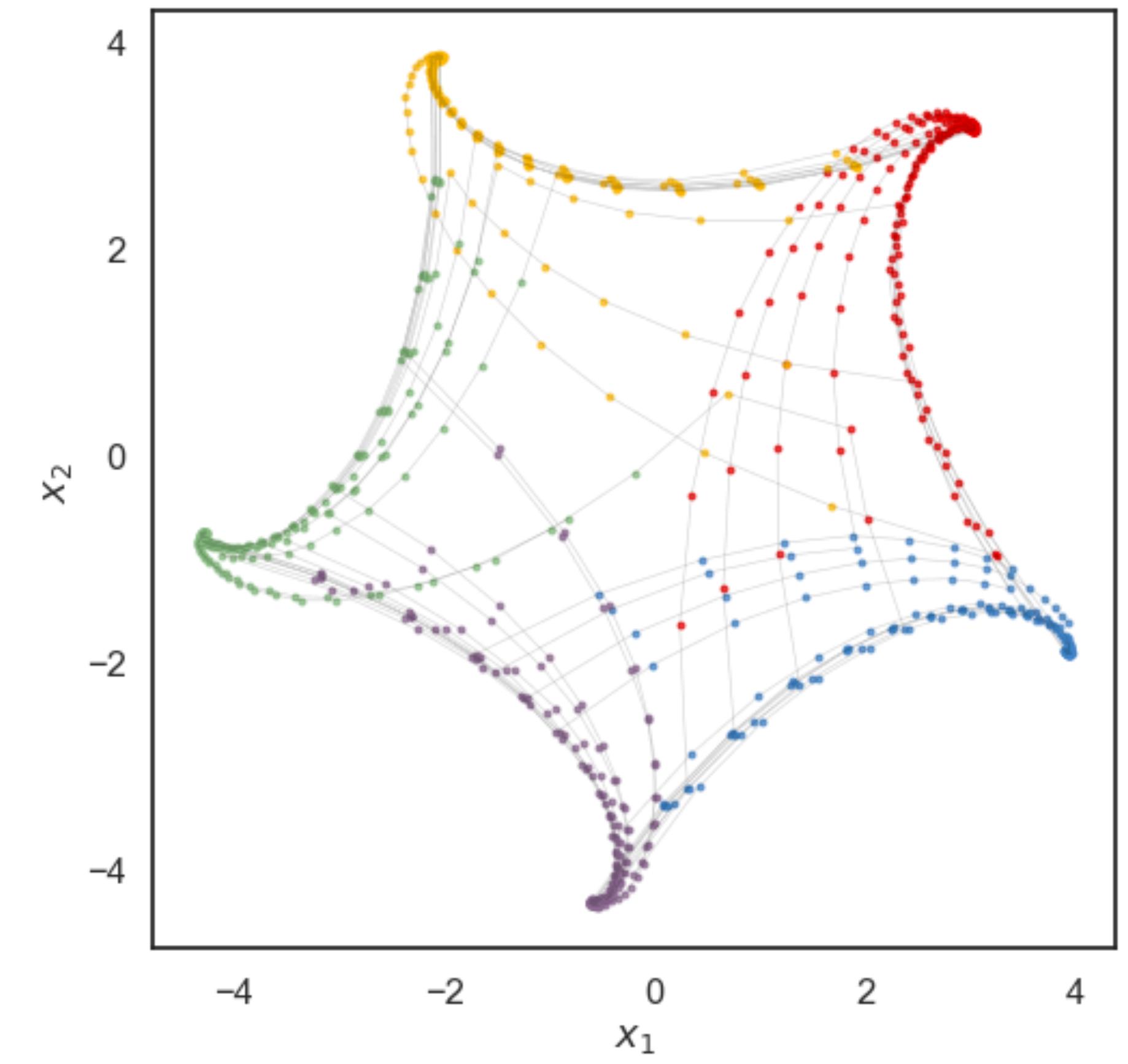
Recap: Autoregressive (AR) HMM Graphical Model



Visualizing Linear Dynamics



Simulated data from an ARHMM



Recap: Gaussian Linear Dynamical Systems

Generative Model:

$$\begin{aligned}x_1 &\sim \mathcal{N}(b, Q), \\x_t \mid x_{t-1} &\sim \mathcal{N}(Ax_{t-1} + b, Q), & \text{for } t = 2, \dots, T. \\y_t \mid x_t &\sim \mathcal{N}(Cx_t + d, R) & \text{for } t = 1, \dots, T\end{aligned}$$

Parameters:

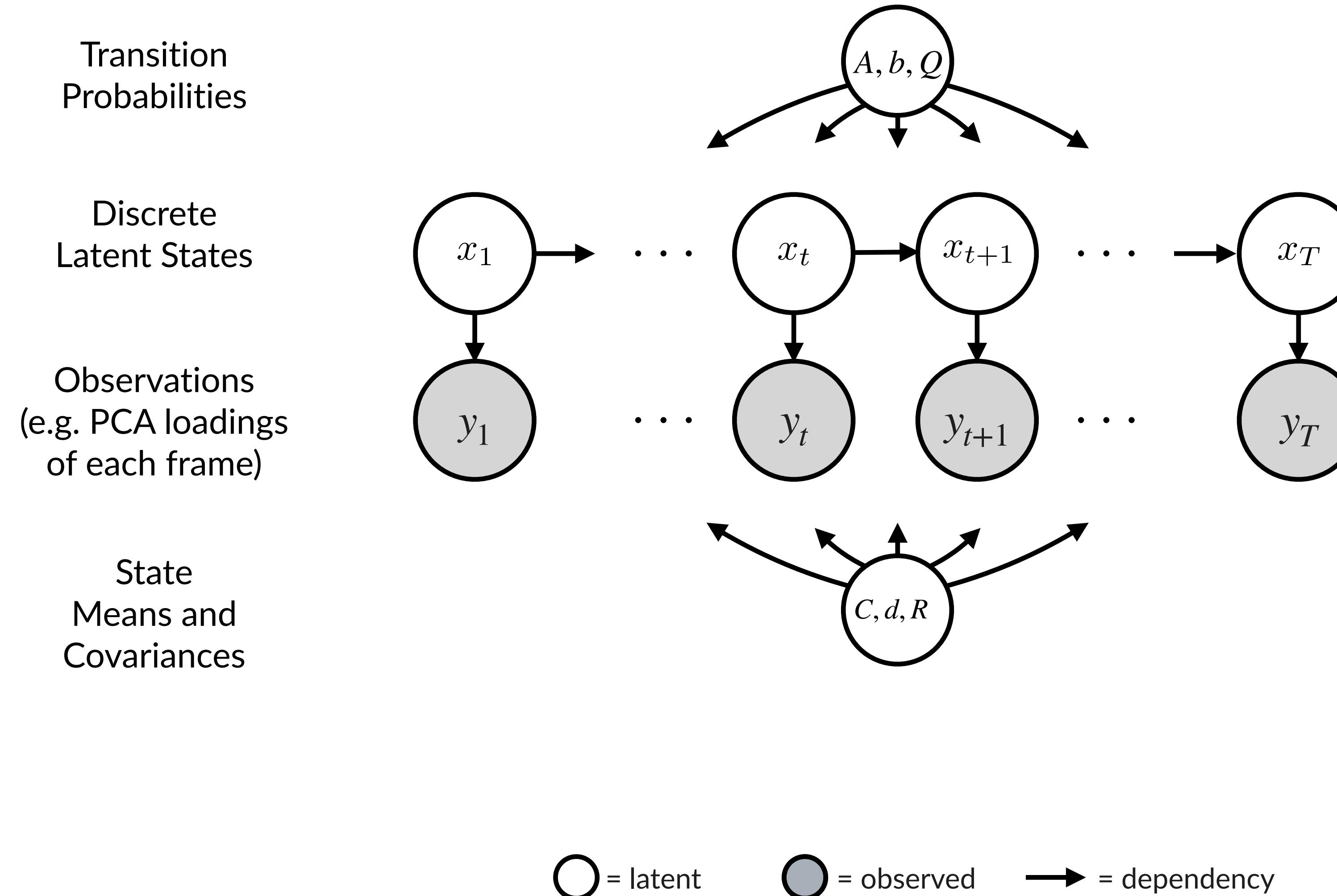
$$\Theta = A, b, Q, C, d, R$$

Joint probability:

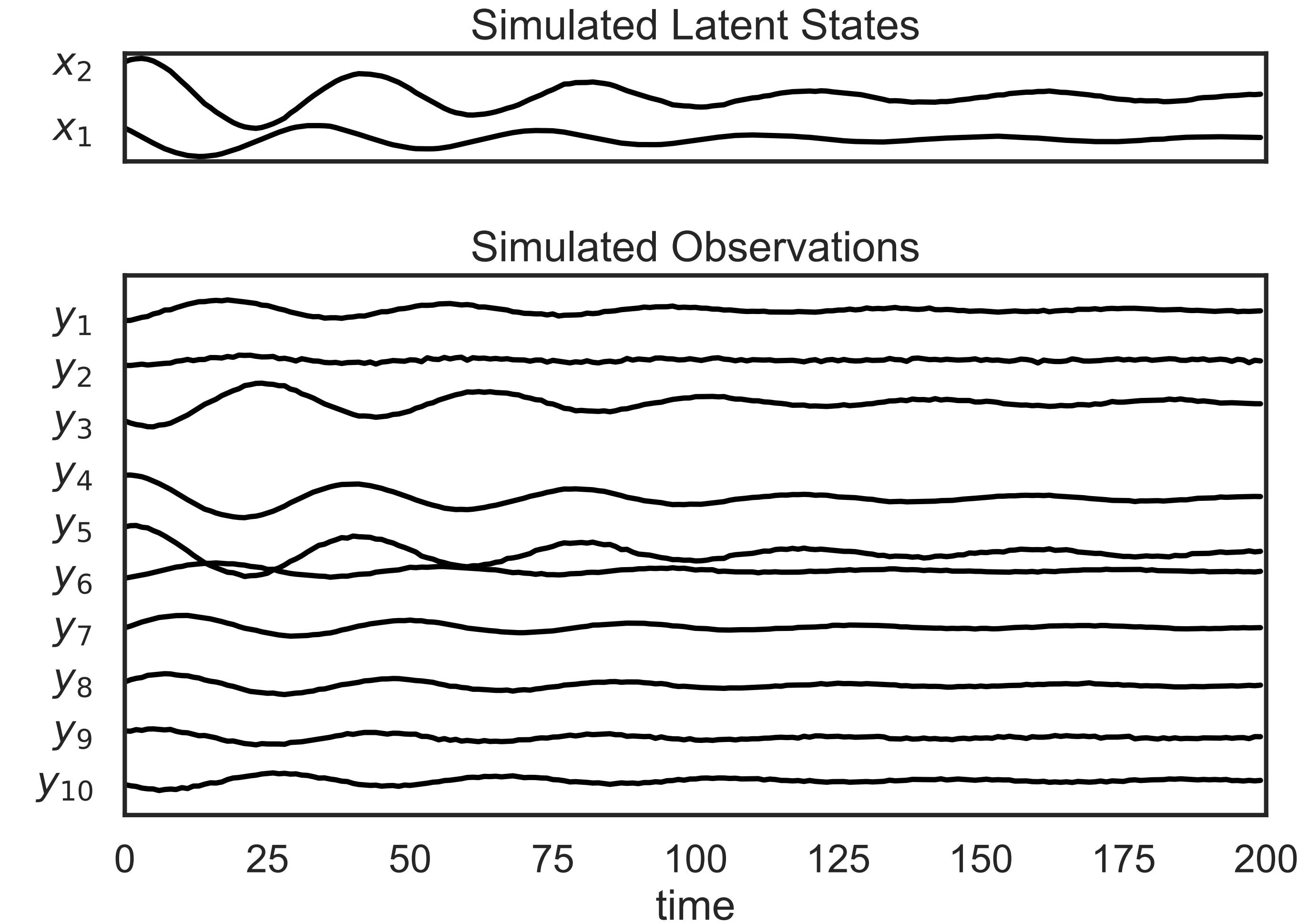
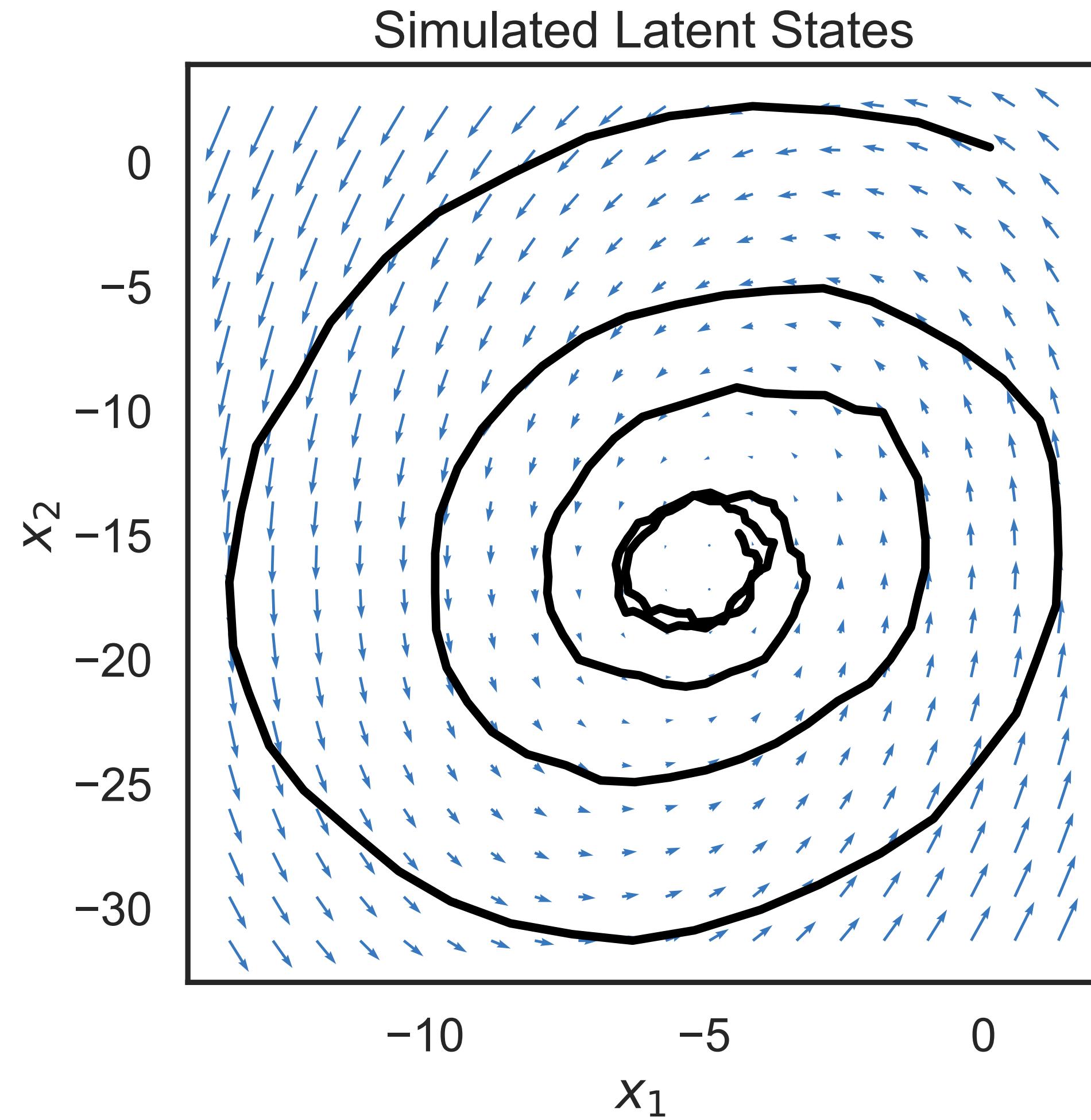
$$p(y, x \mid \Theta) = p(x_1) \prod_{t=1}^{T-1} p(x_{t+1} \mid x_t) \prod_{t=1}^T p(y_t \mid x_t)$$

Recap: Gaussian Linear Dynamical Systems

Graphical Model



Simulated data from an LDS



Switching LDS: Best of Both Worlds

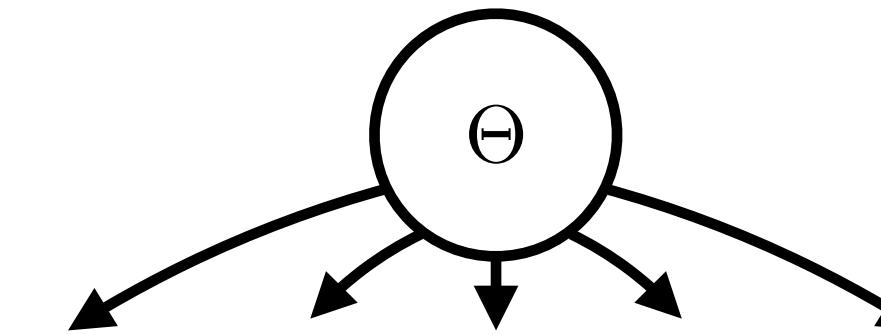
○ = latent

● = observed

→ = dependency

Switching LDS: Best of Both Worlds

Global
Parameters

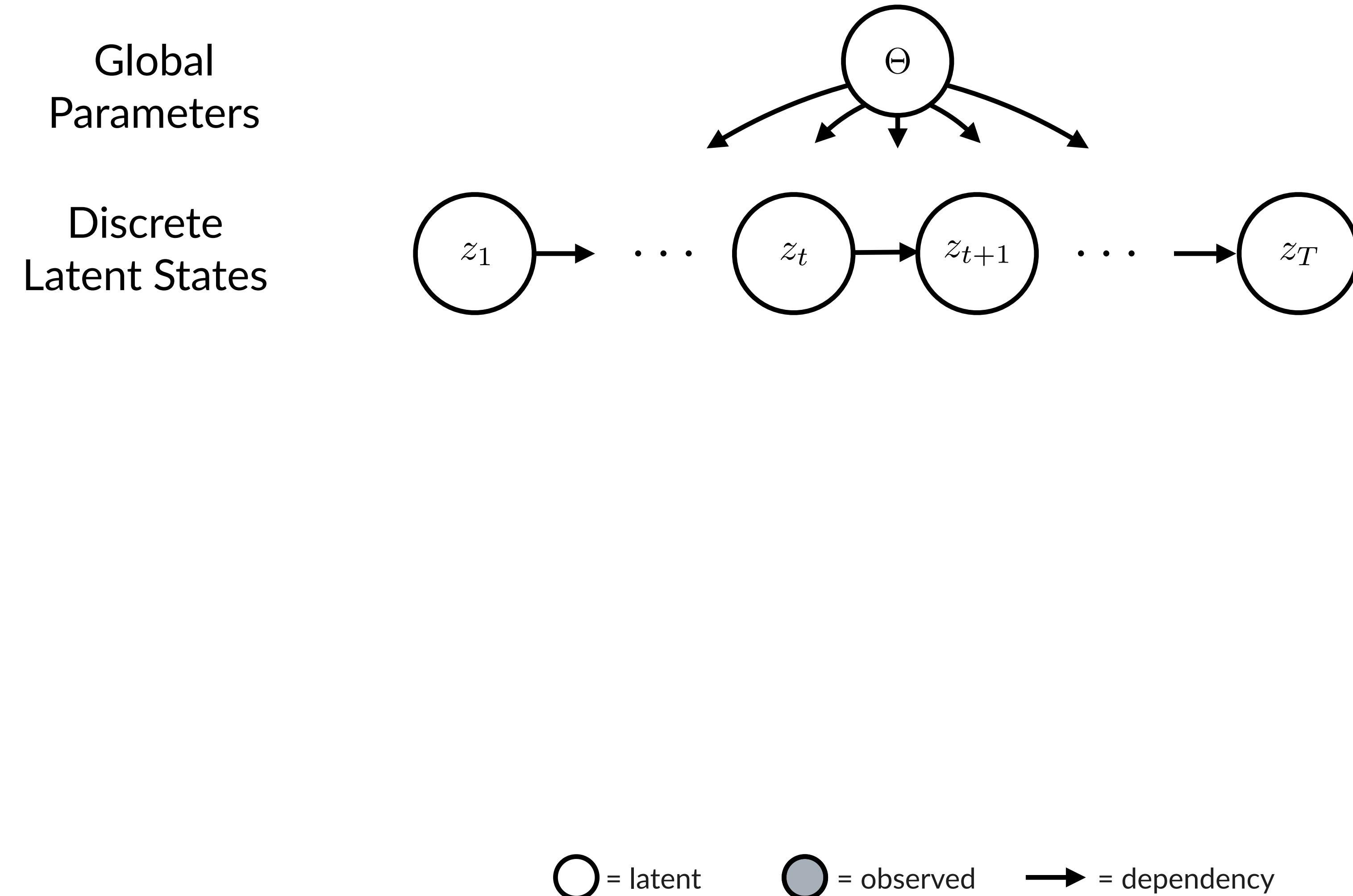


○ = latent

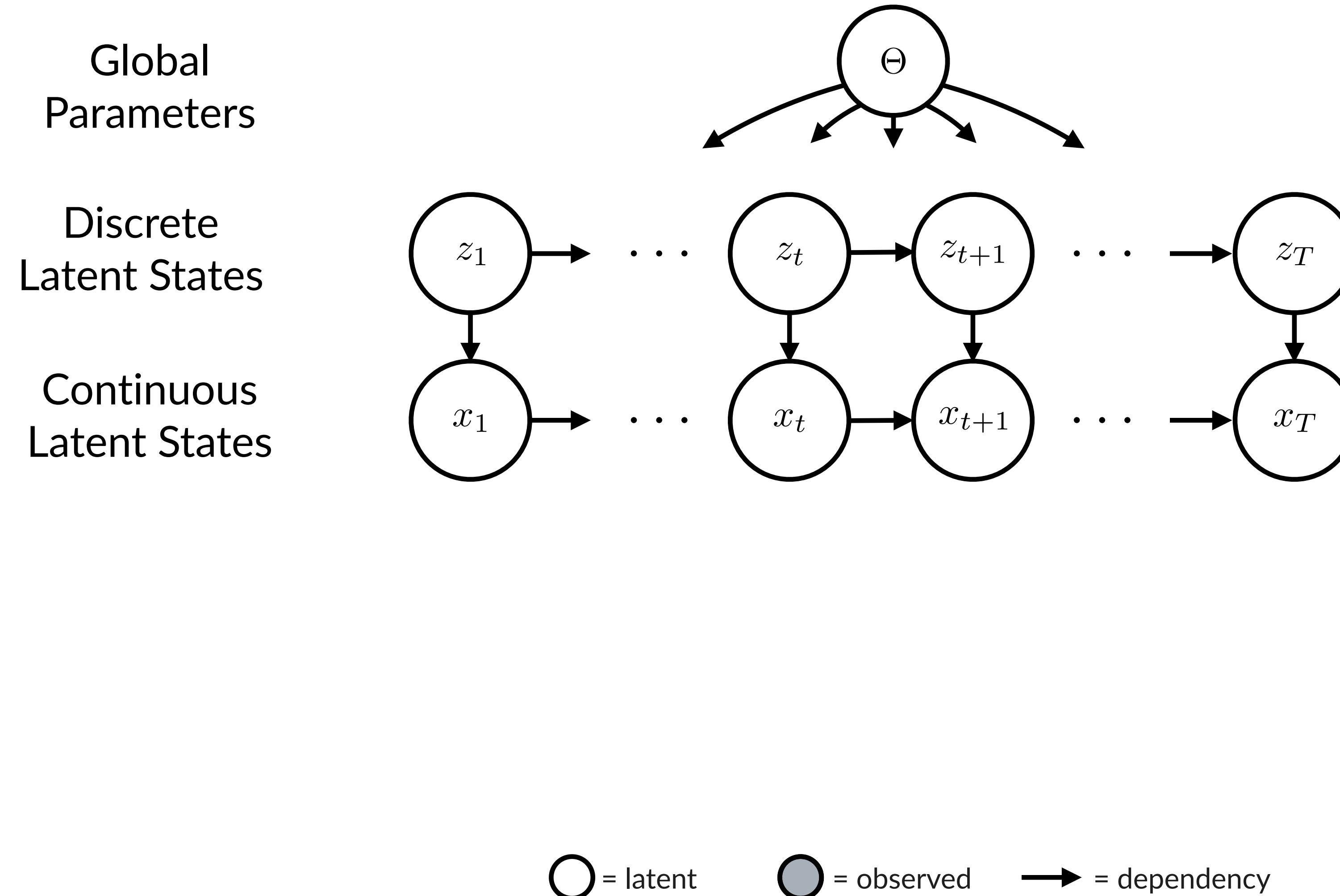
● = observed

→ = dependency

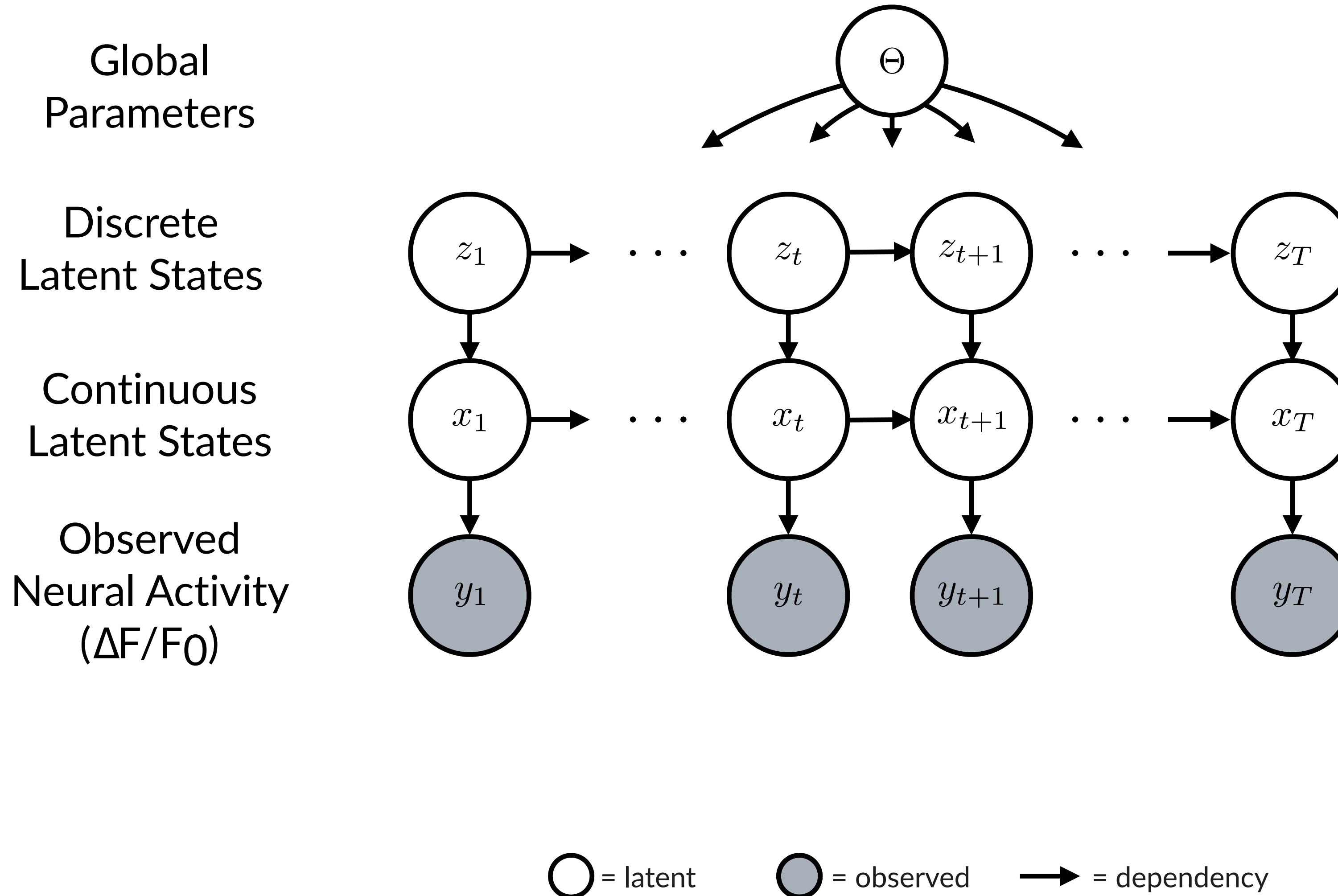
Switching LDS: Best of Both Worlds



Switching LDS: Best of Both Worlds



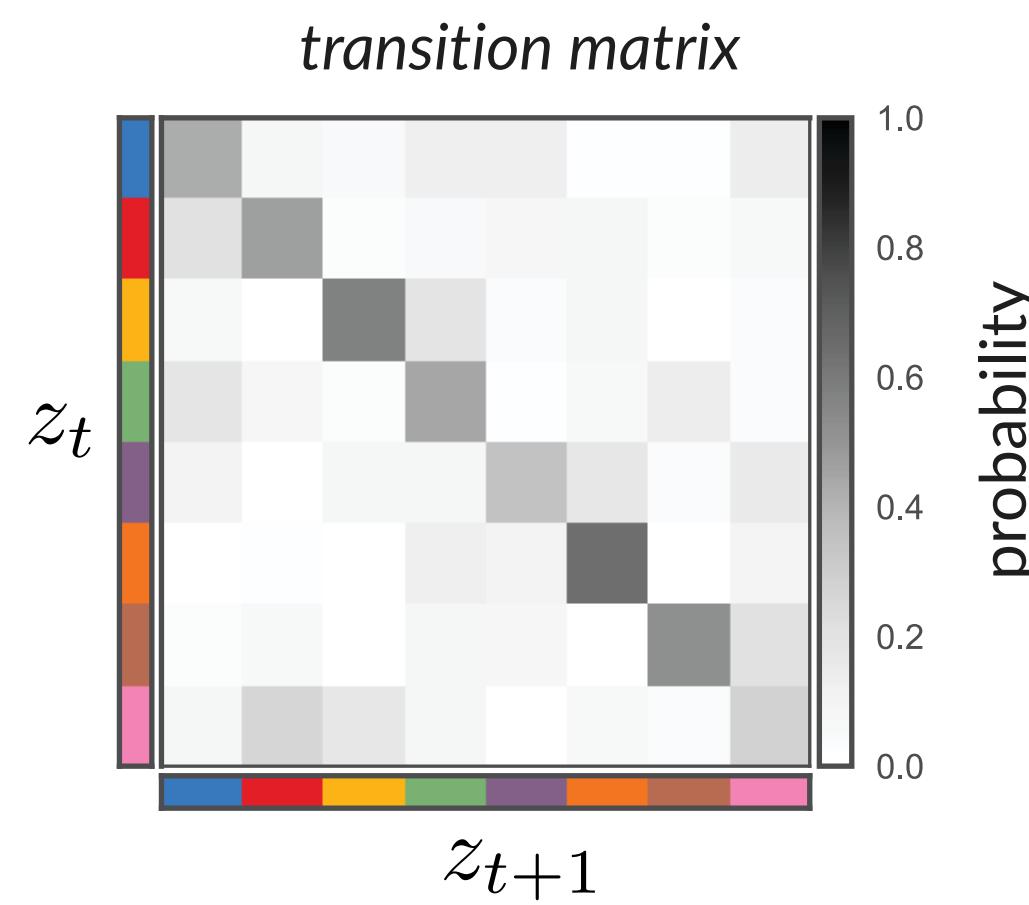
Switching LDS: Best of Both Worlds



Specifying the form of the dependencies

Specifying the form of the dependencies

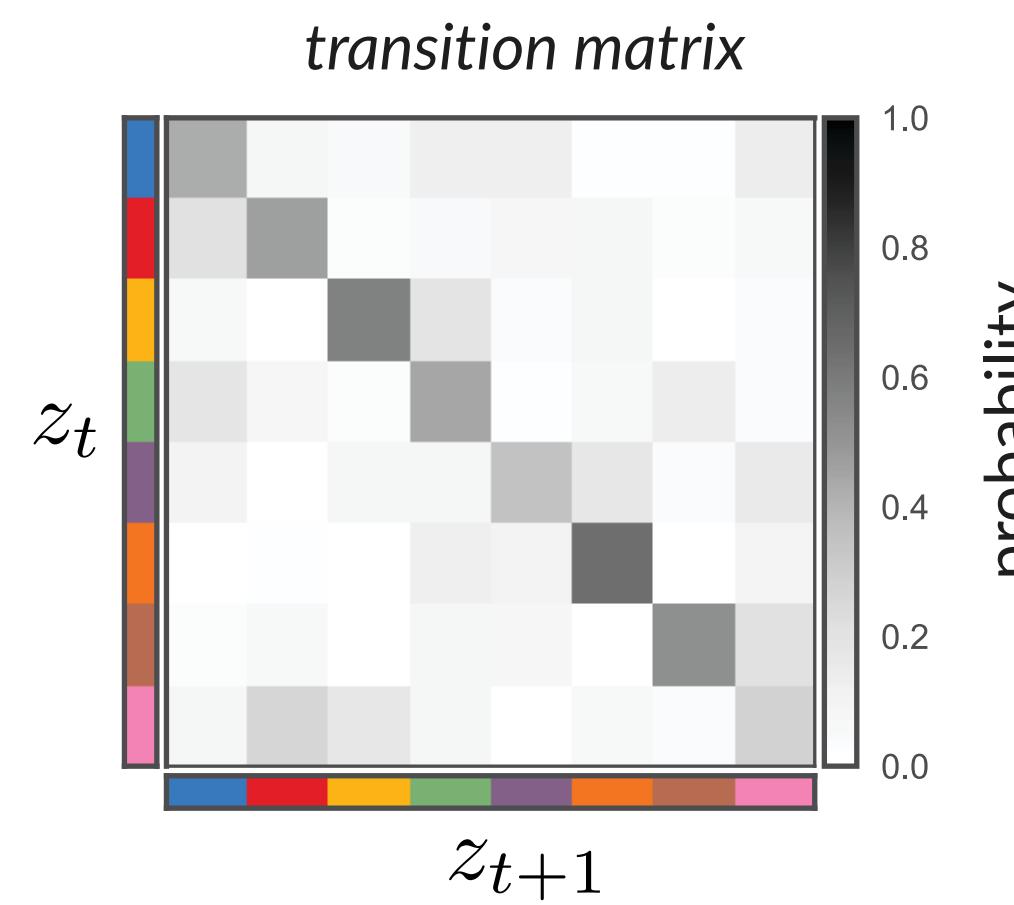
State-dependent
switching probabilities



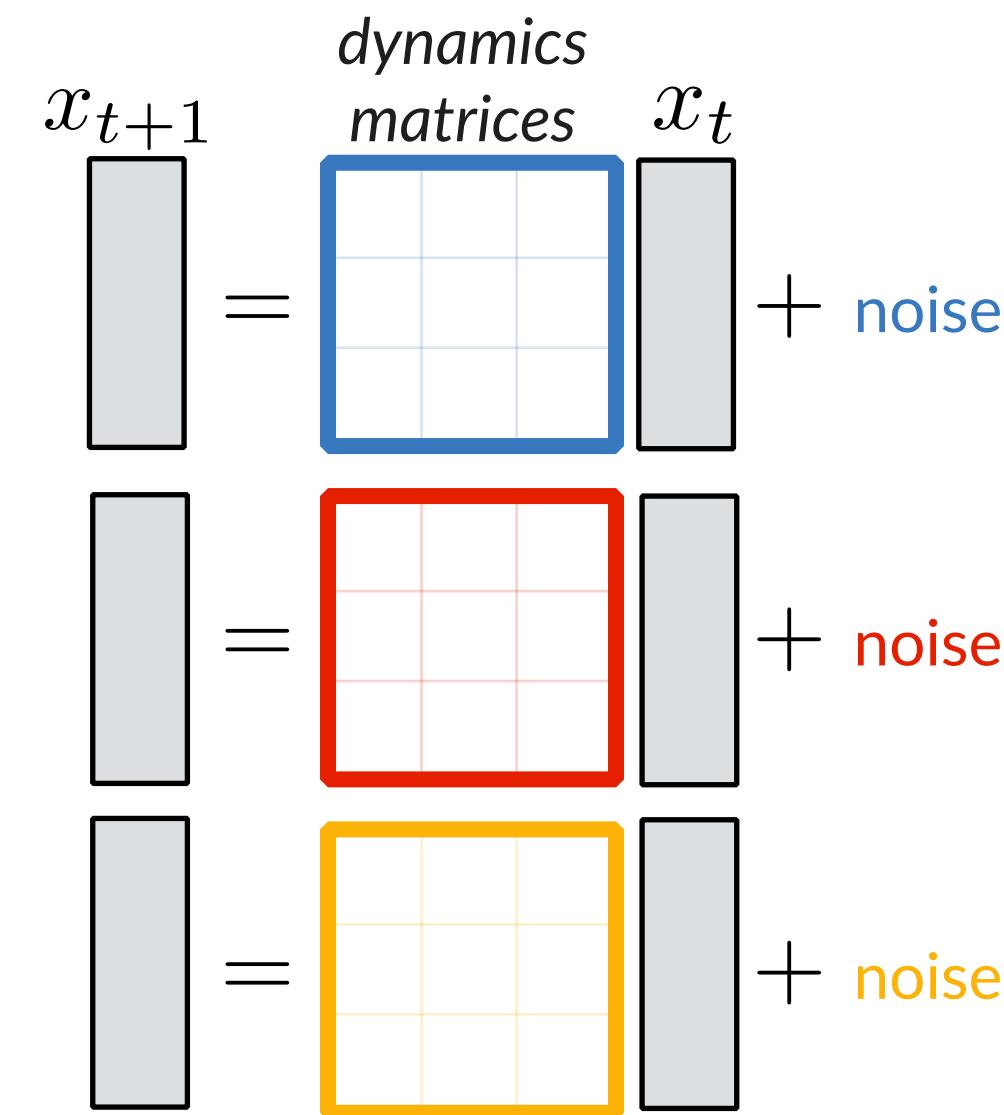
$$\Pr(z_{t+1} = j \mid z_t = i) = P_{ij}$$

Specifying the form of the dependencies

State-dependent
switching probabilities



Different linear dynamics
in each discrete state

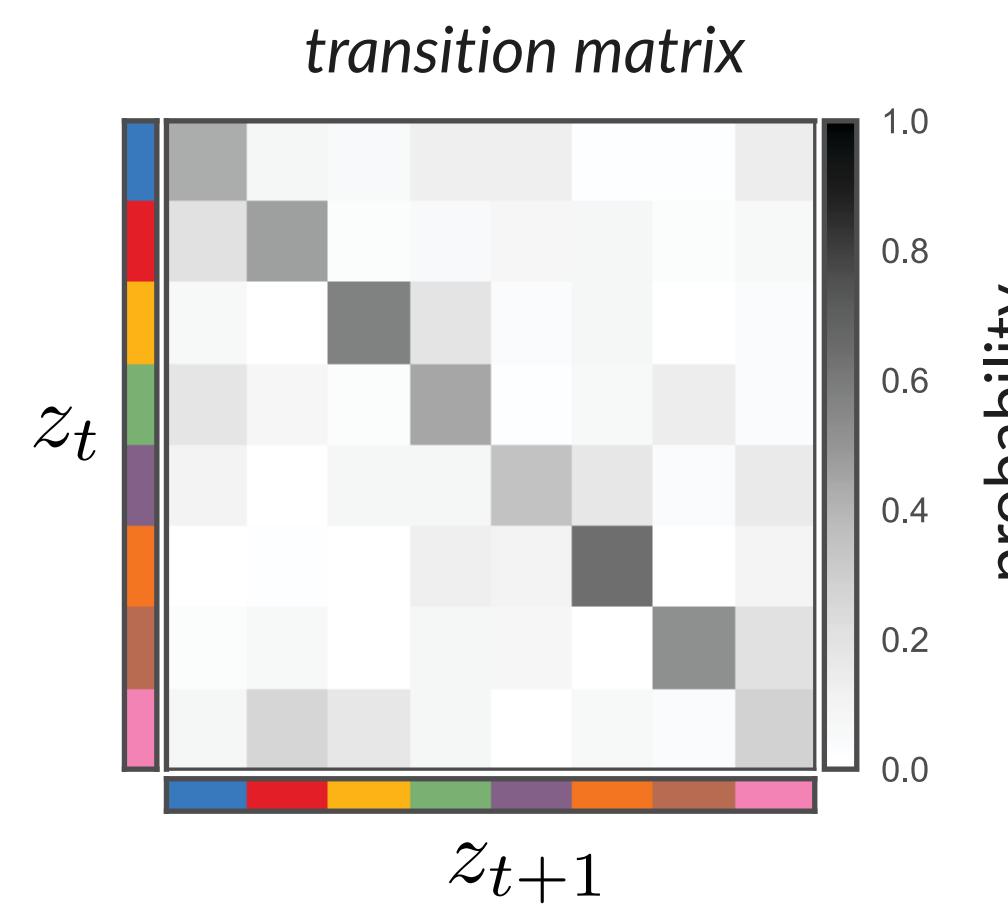


$$\Pr(z_{t+1} = j \mid z_t = i) = P_{ij}$$

$$x_{t+1} = A_{z_{t+1}} x_t + b_{z_{t+1}} + \epsilon_{t+1}$$

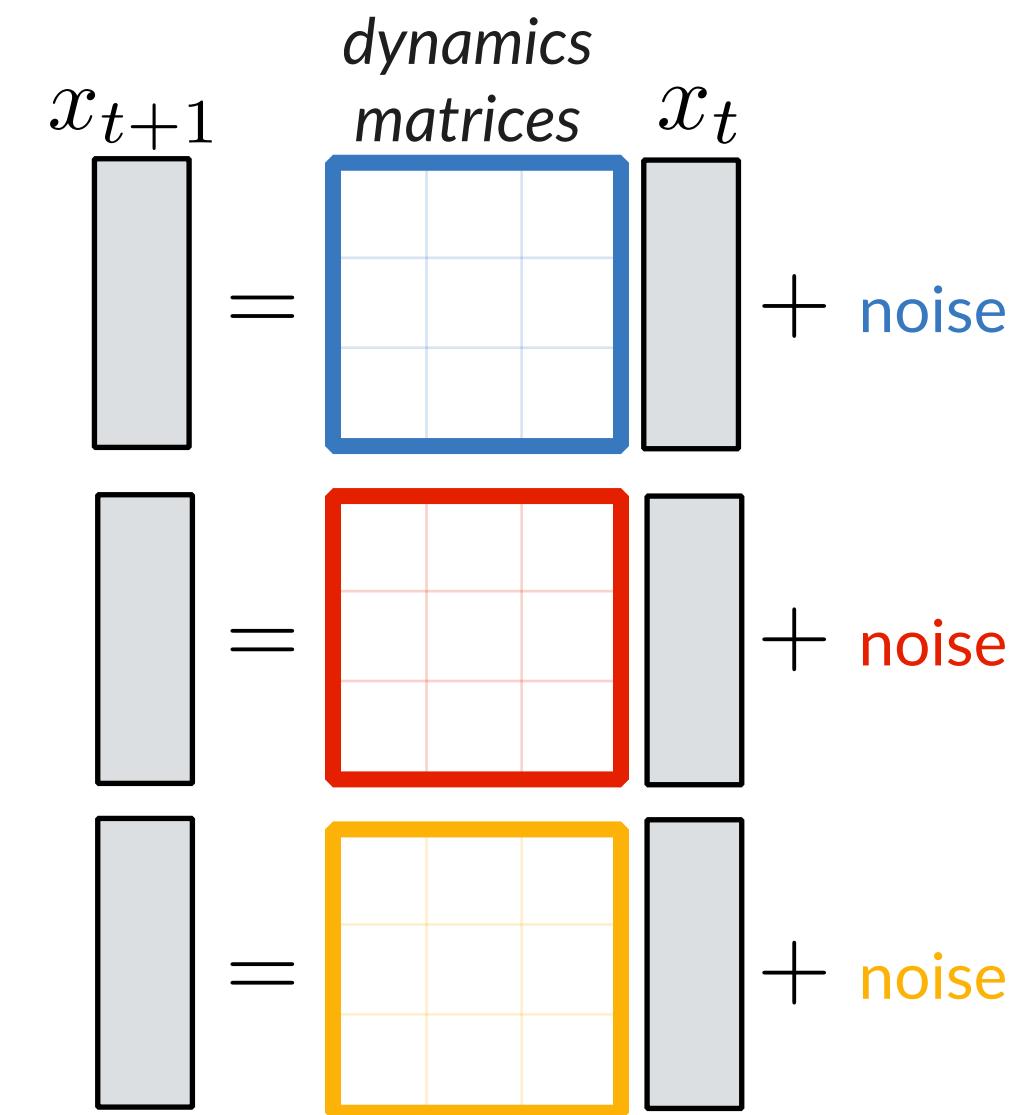
Specifying the form of the dependencies

State-dependent
switching probabilities



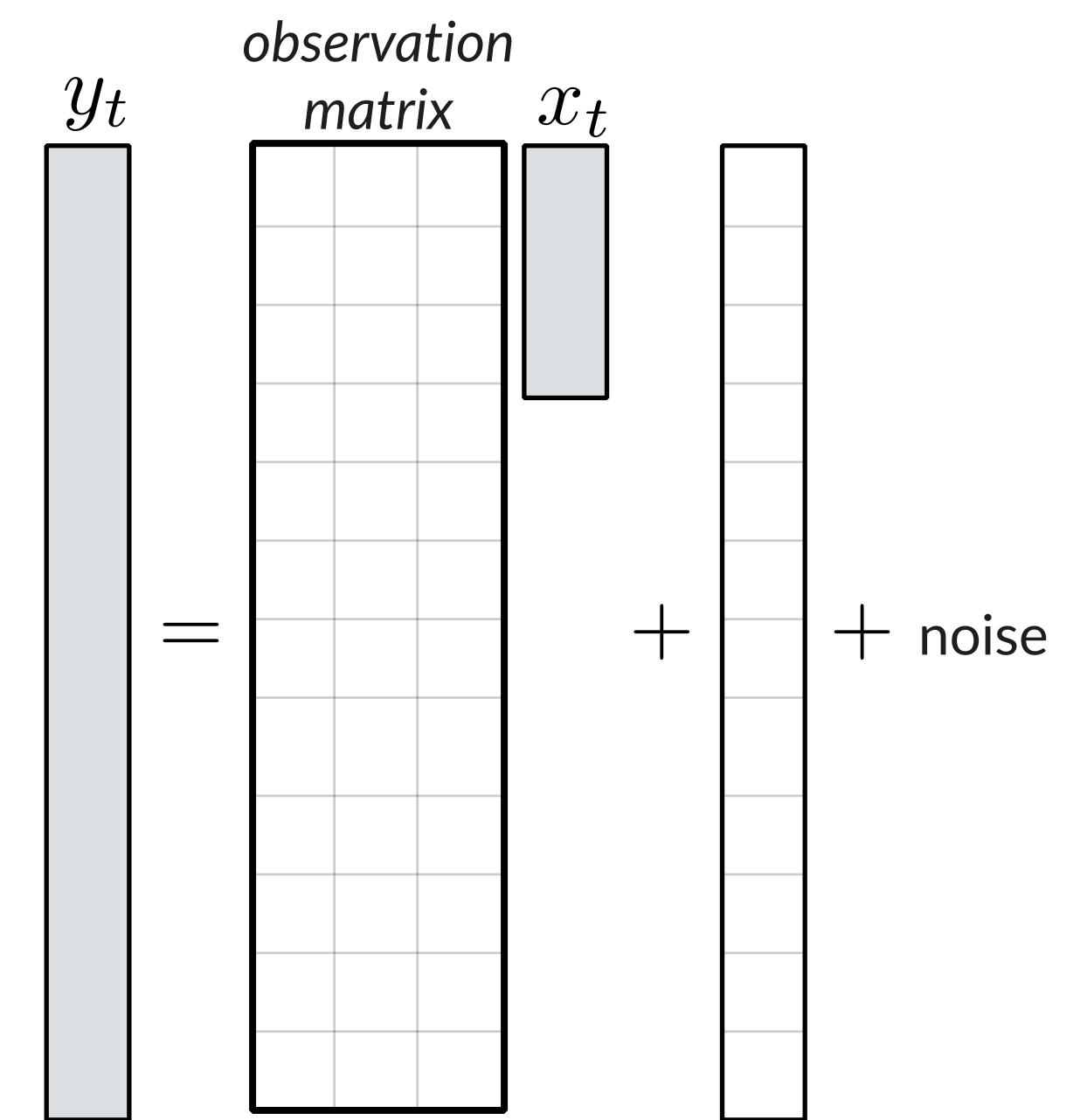
$$\Pr(z_{t+1} = j \mid z_t = i) = P_{ij}$$

Different linear dynamics
in each discrete state



$$x_{t+1} = A_{z_{t+1}} x_t + b_{z_{t+1}} + \epsilon_{t+1}$$

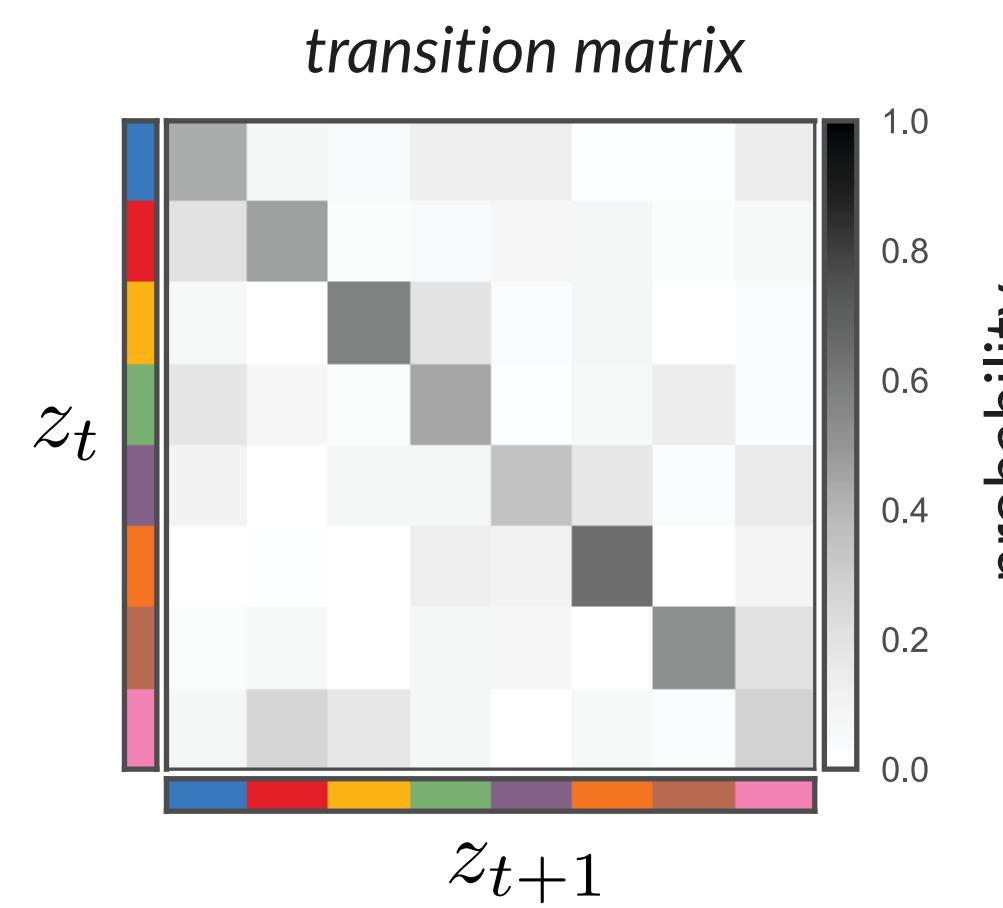
Linear mapping from continuous latent
states to observed neural activity



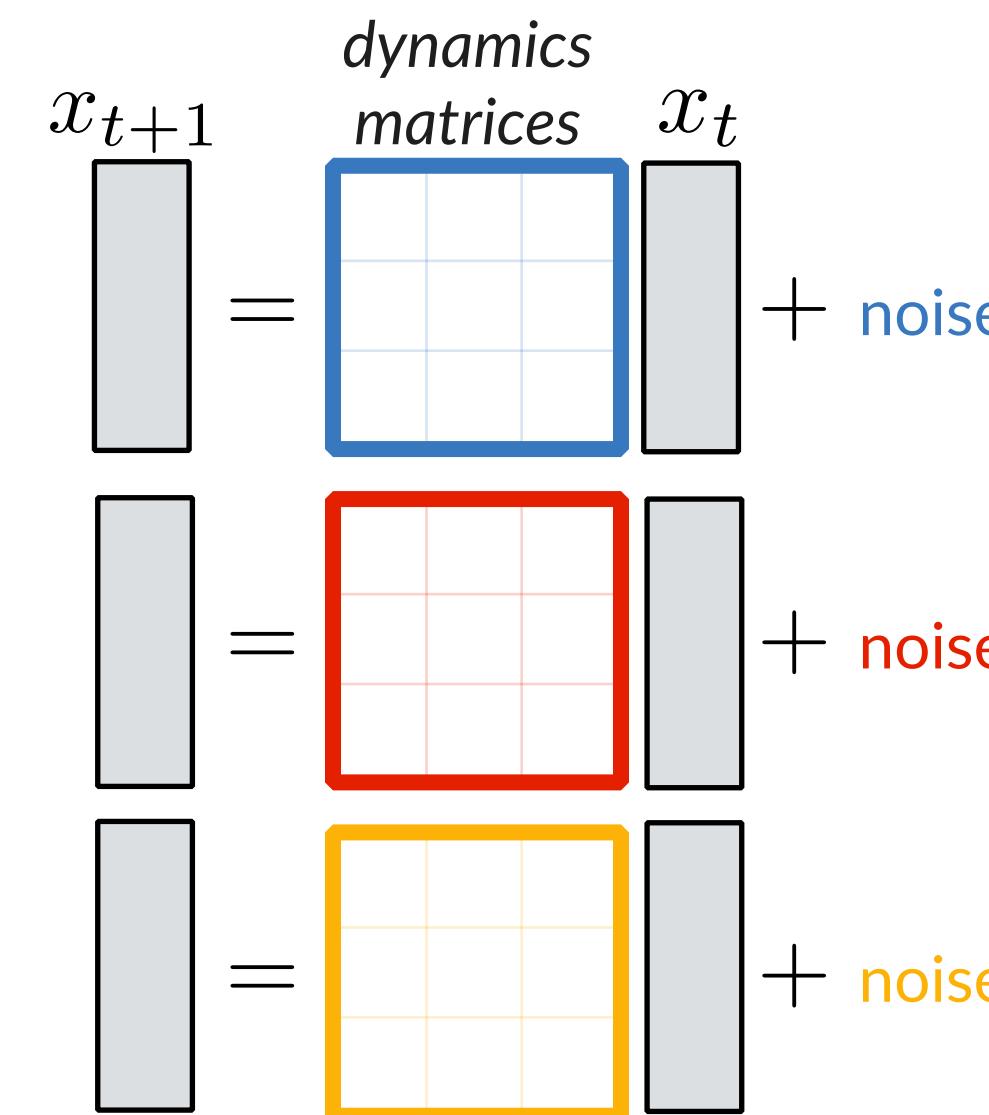
$$y_t = C x_t + d + \delta_t$$

Specifying the form of the dependencies

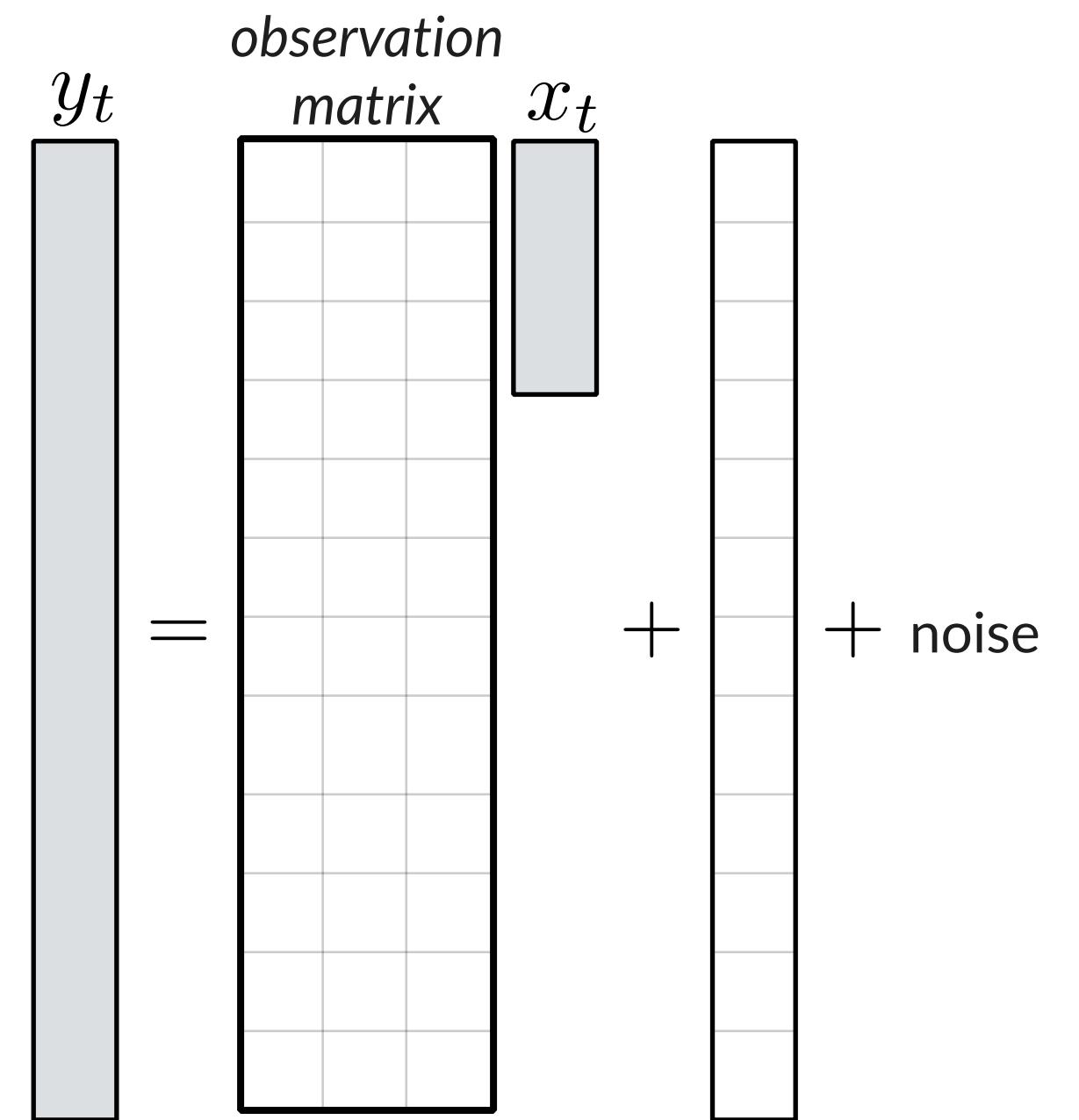
State-dependent
switching probabilities



Different linear dynamics
in each discrete state



Linear mapping from continuous latent
states to observed neural activity



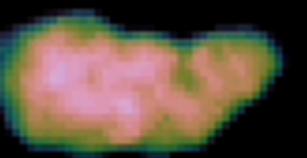
$$\Pr(z_{t+1} = j \mid z_t = i) = P_{ij}$$

$$x_{t+1} = A_{z_{t+1}} x_t + b_{z_{t+1}} + \epsilon_{t+1}$$

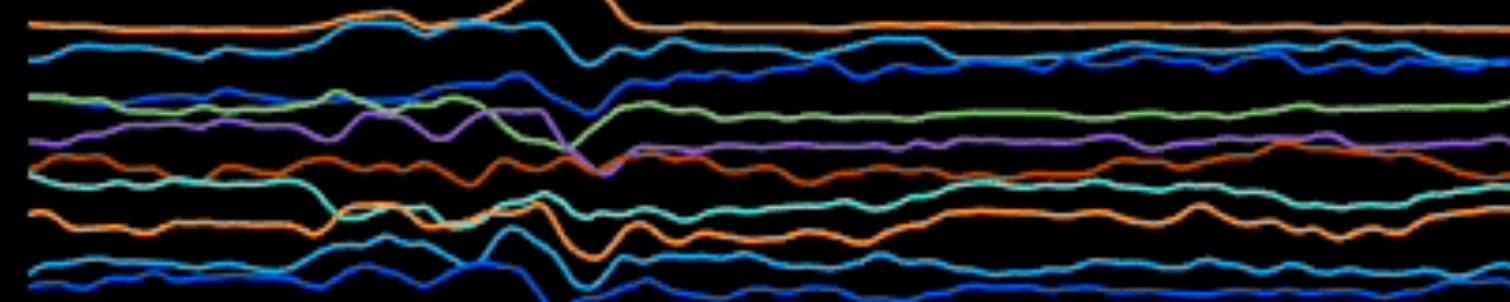
$$y_t = Cx_t + d + \delta_t$$

Switching Linear Dynamical System (SLDS)
Combines LDS (Kalman, 1960) and HMMs (Rabiner, 1989).

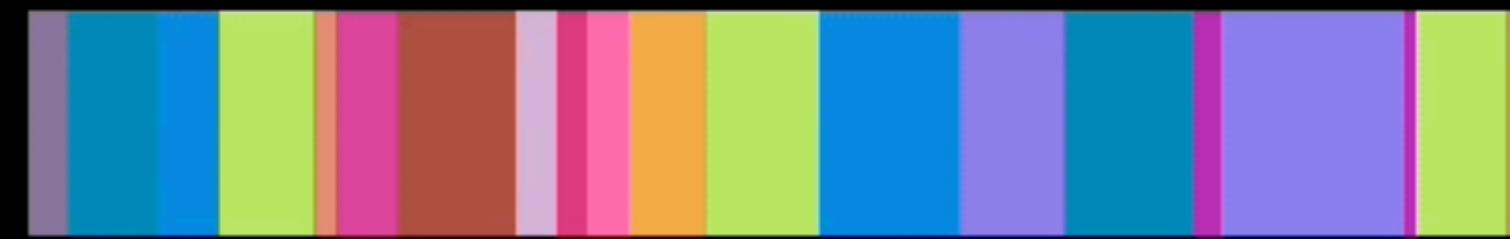
Raw data



Cropped
and
Rotated



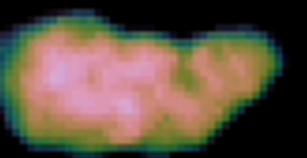
Continuous
Latents x



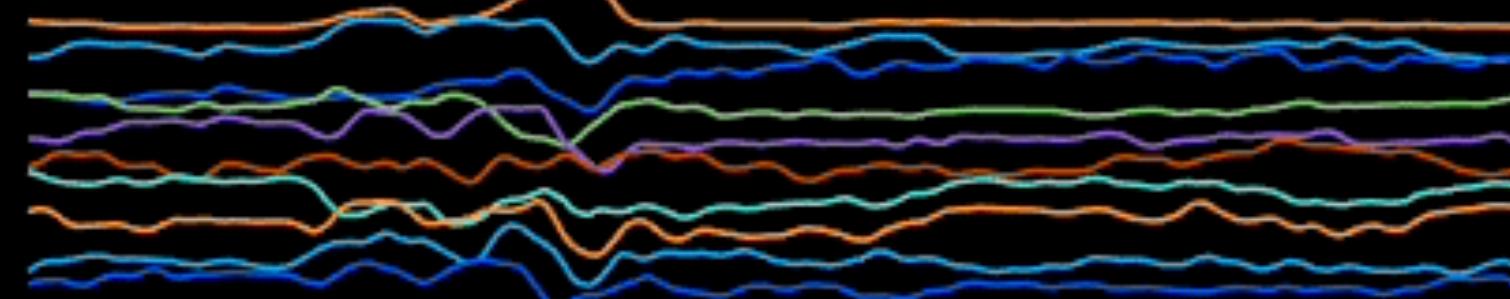
Discrete
States z

Time →
Time

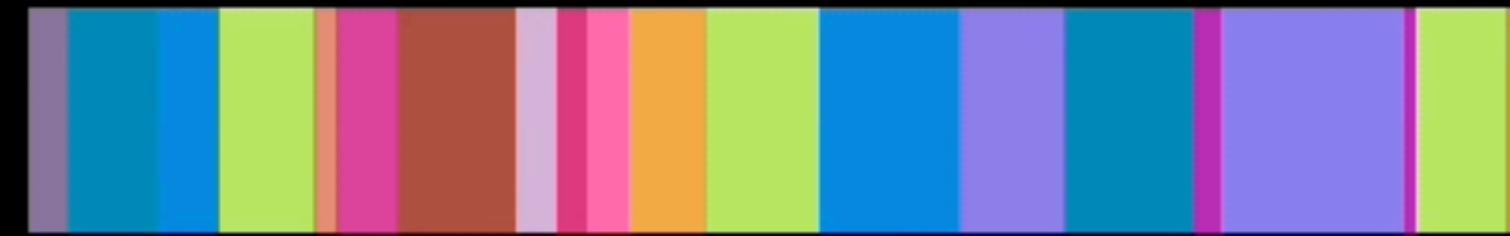
Raw data



Cropped
and
Rotated



Continuous
Latents x



Discrete
States z

Time →
Time

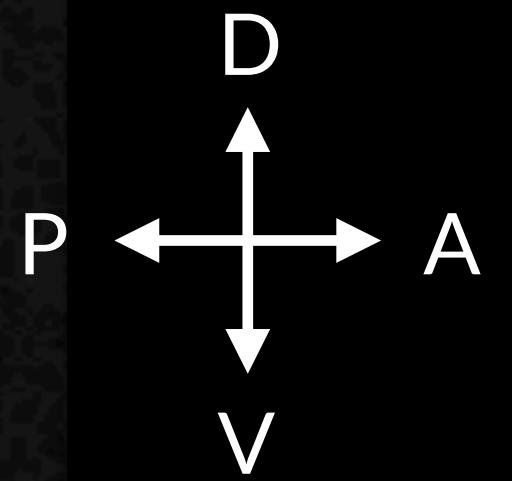
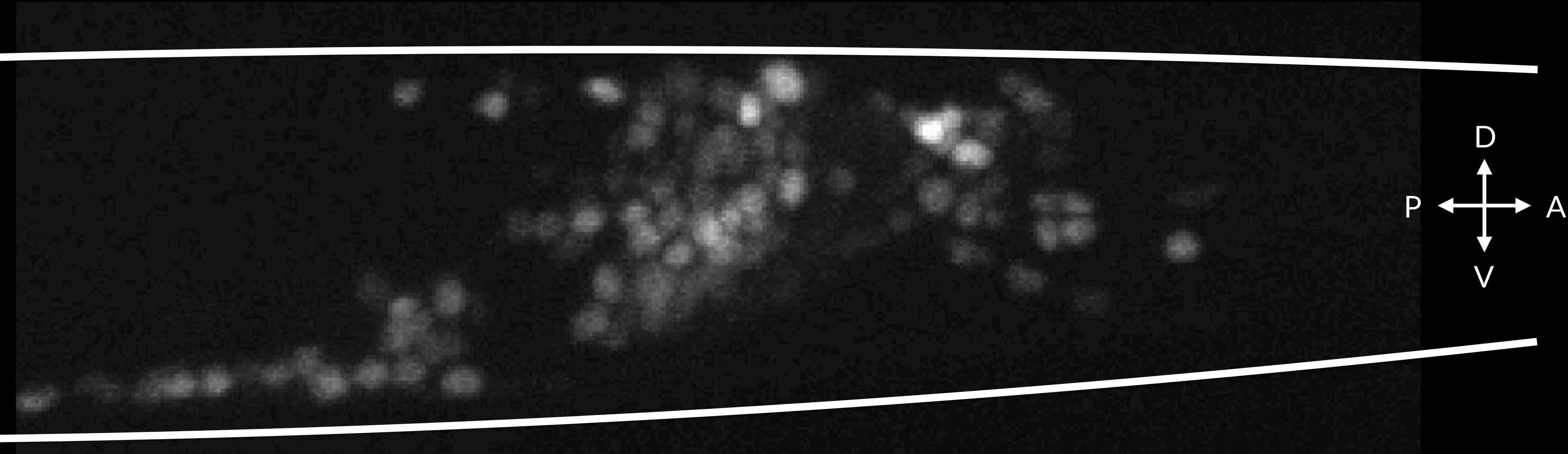


A. E. X. Brown, Imperial College

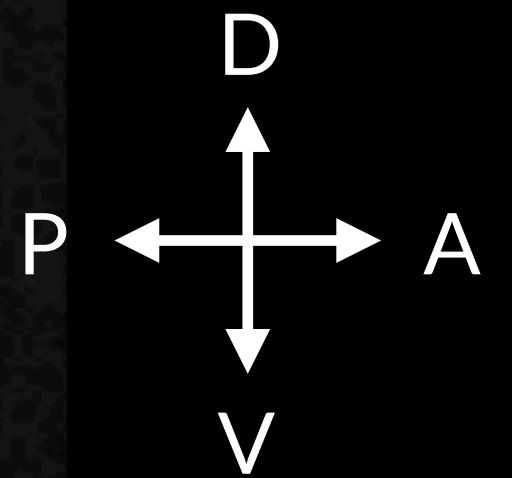
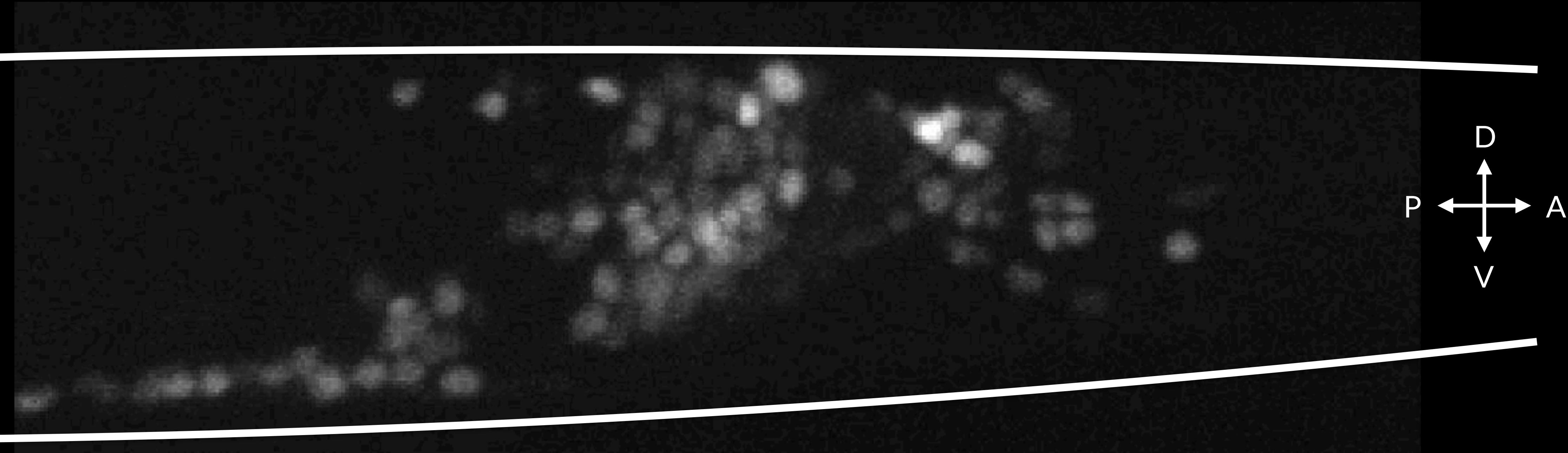


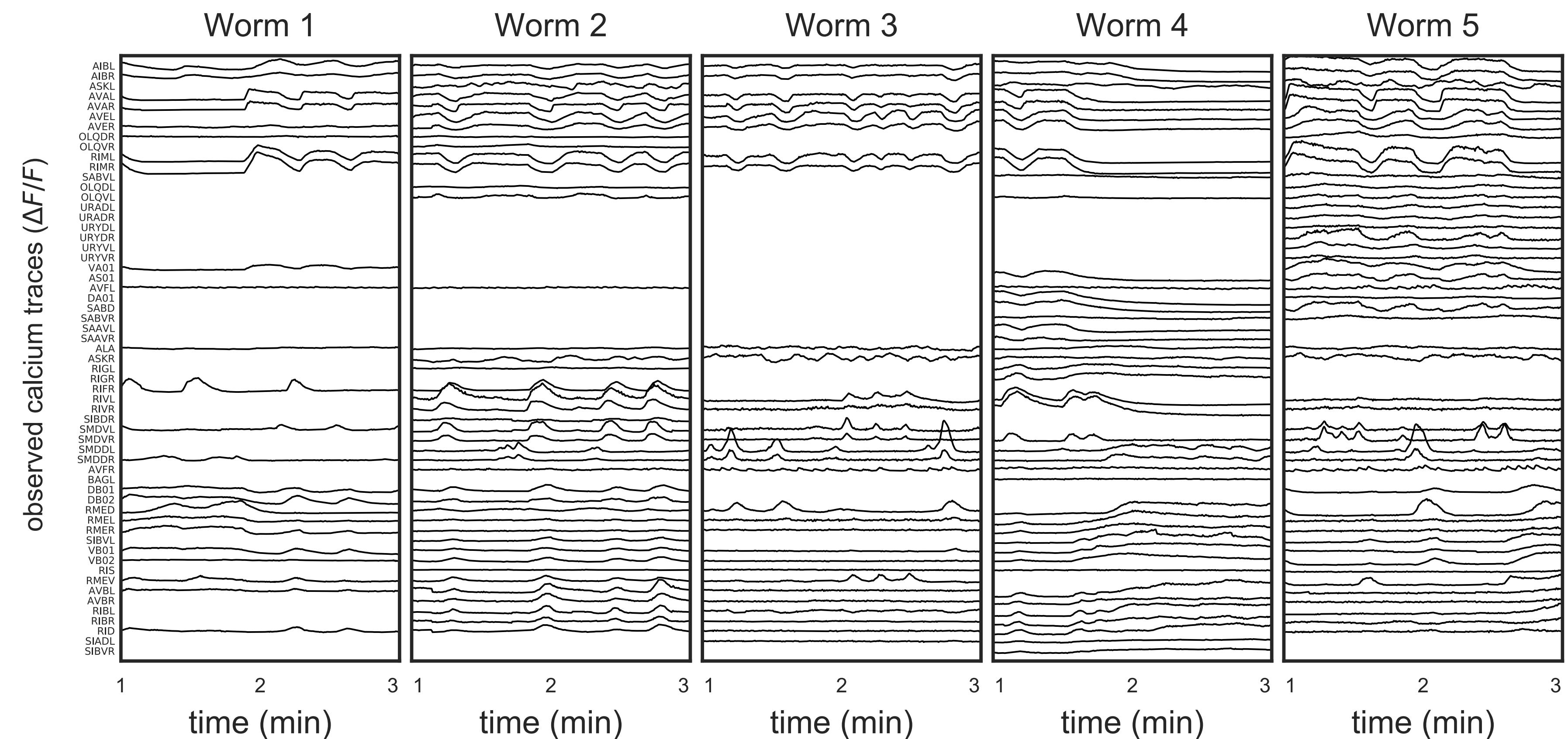
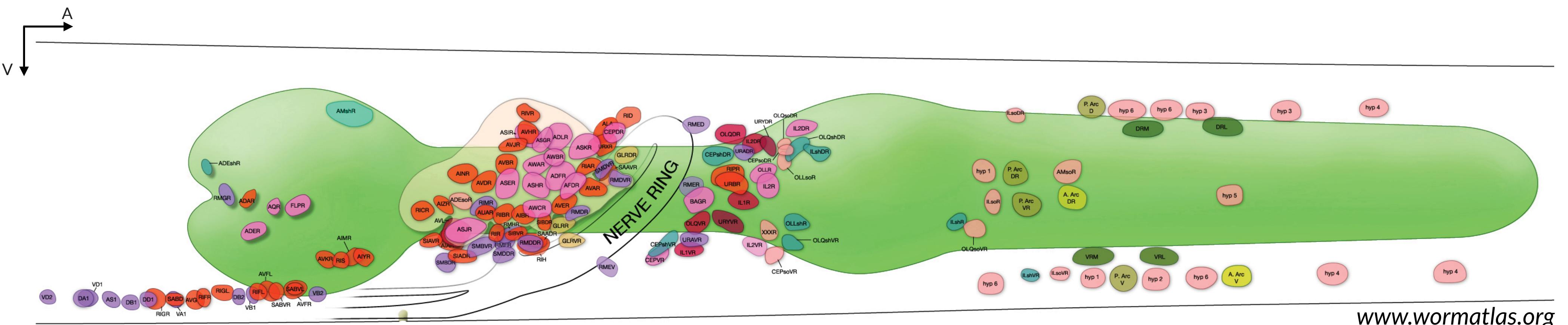
A. E. X. Brown, Imperial College

Calcium imaging of ~100 head ganglia neurons in immobilized *C. elegans*

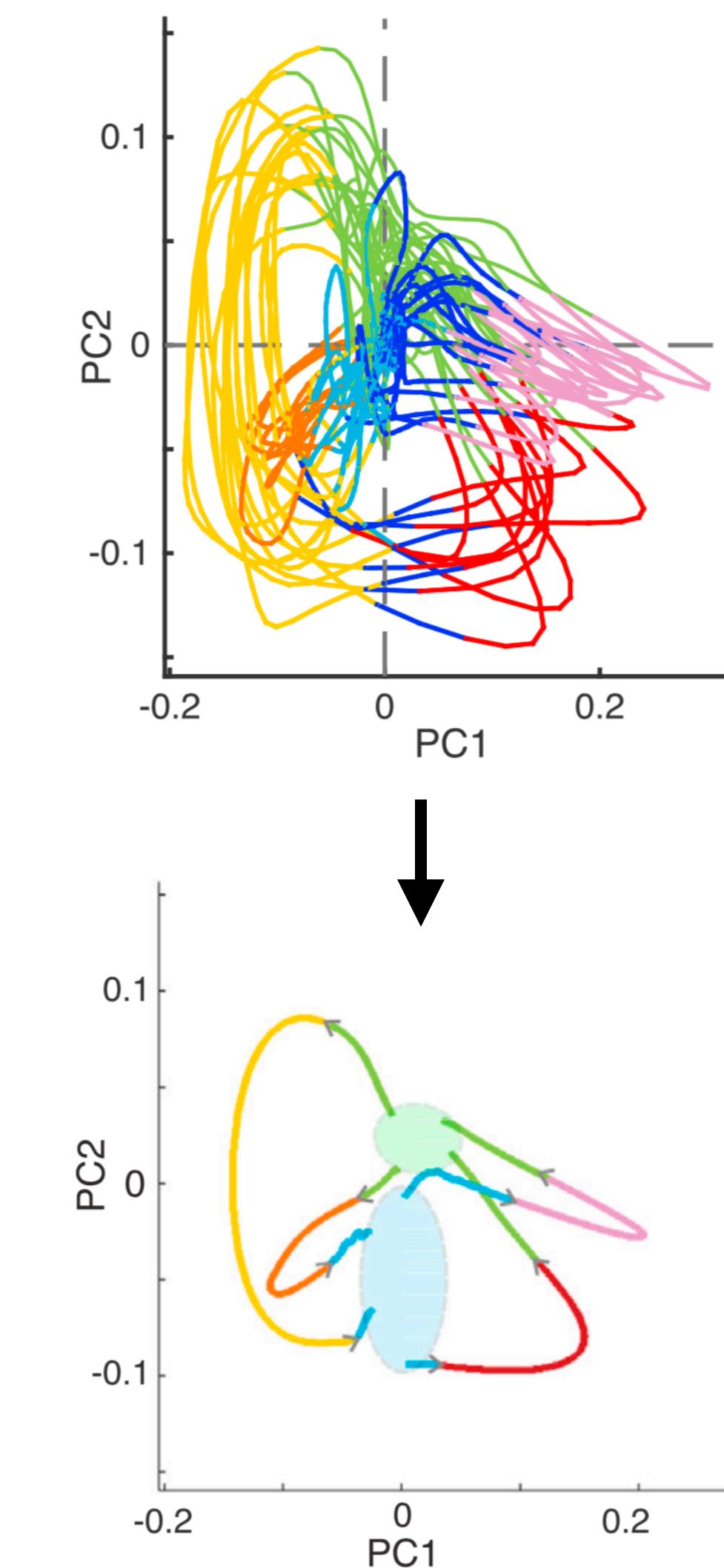
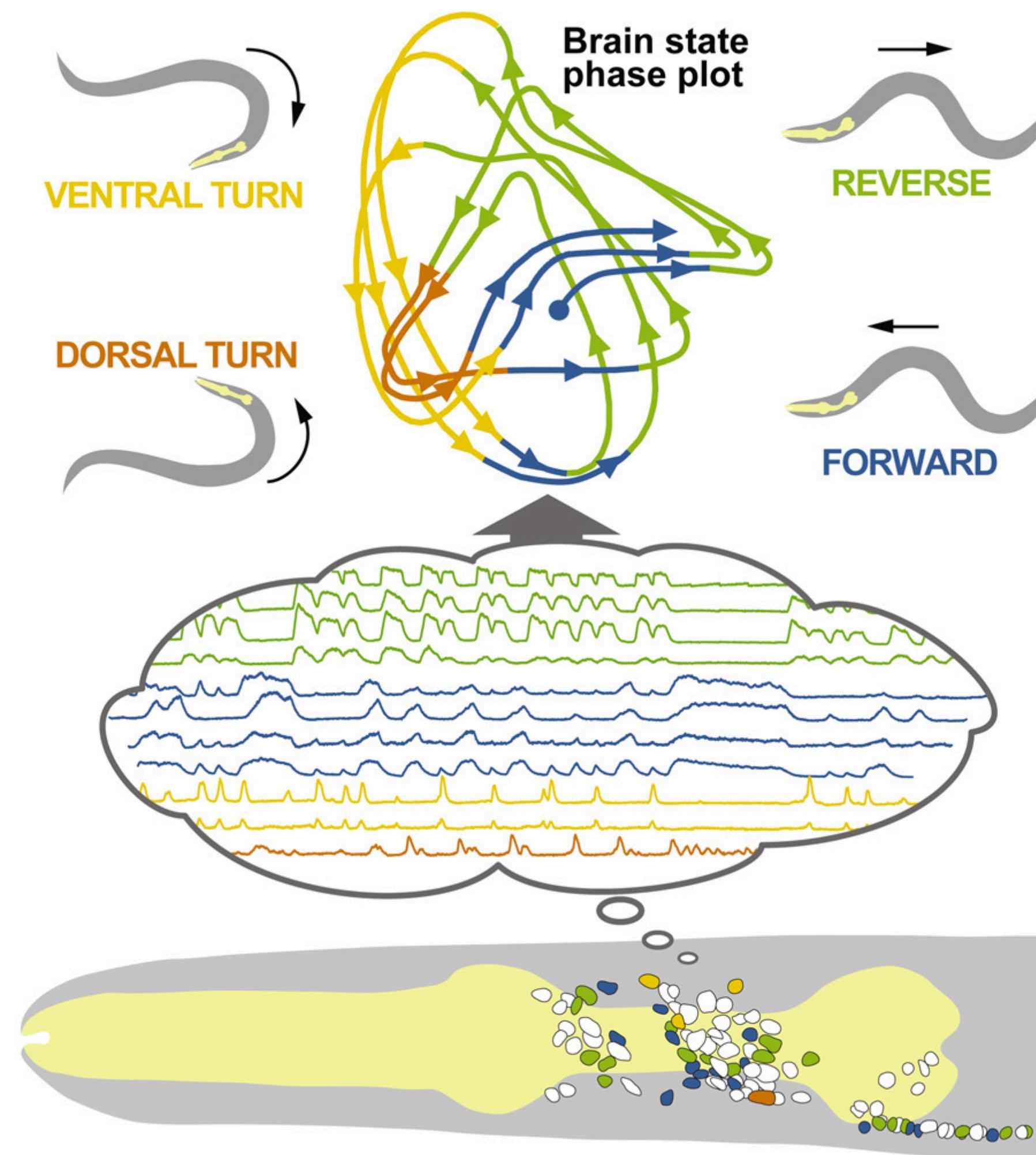


Calcium imaging of ~100 head ganglia neurons in immobilized *C. elegans*



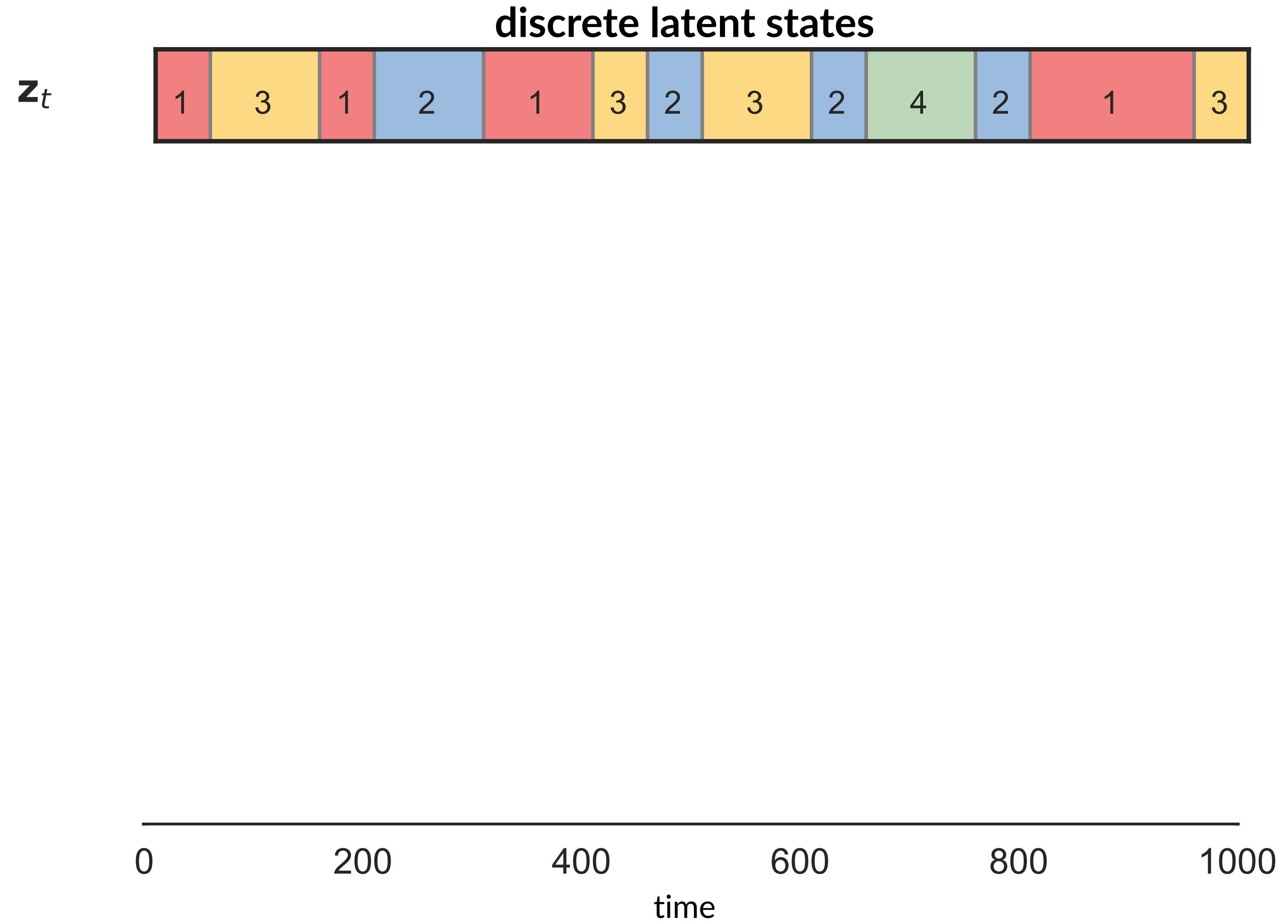
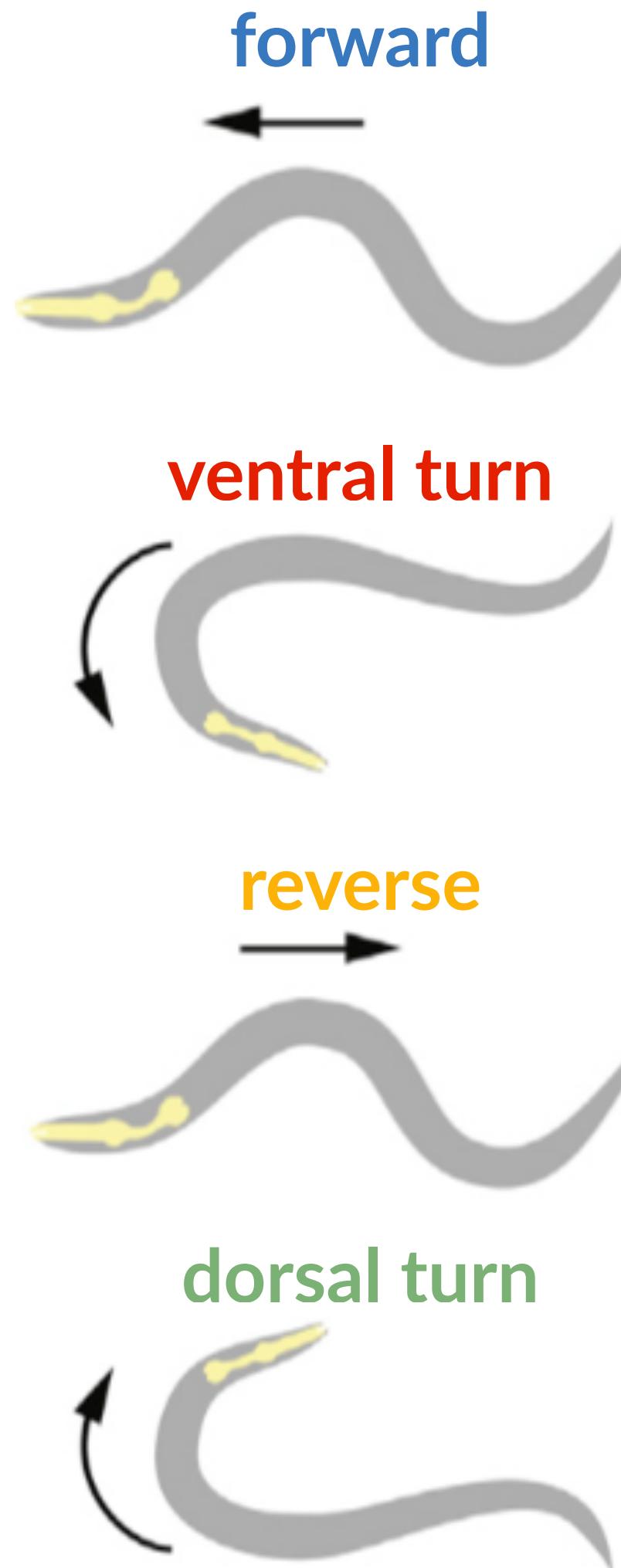


Previous work suggests that this neural activity lies on a low dimensional manifold partitioned by behavior

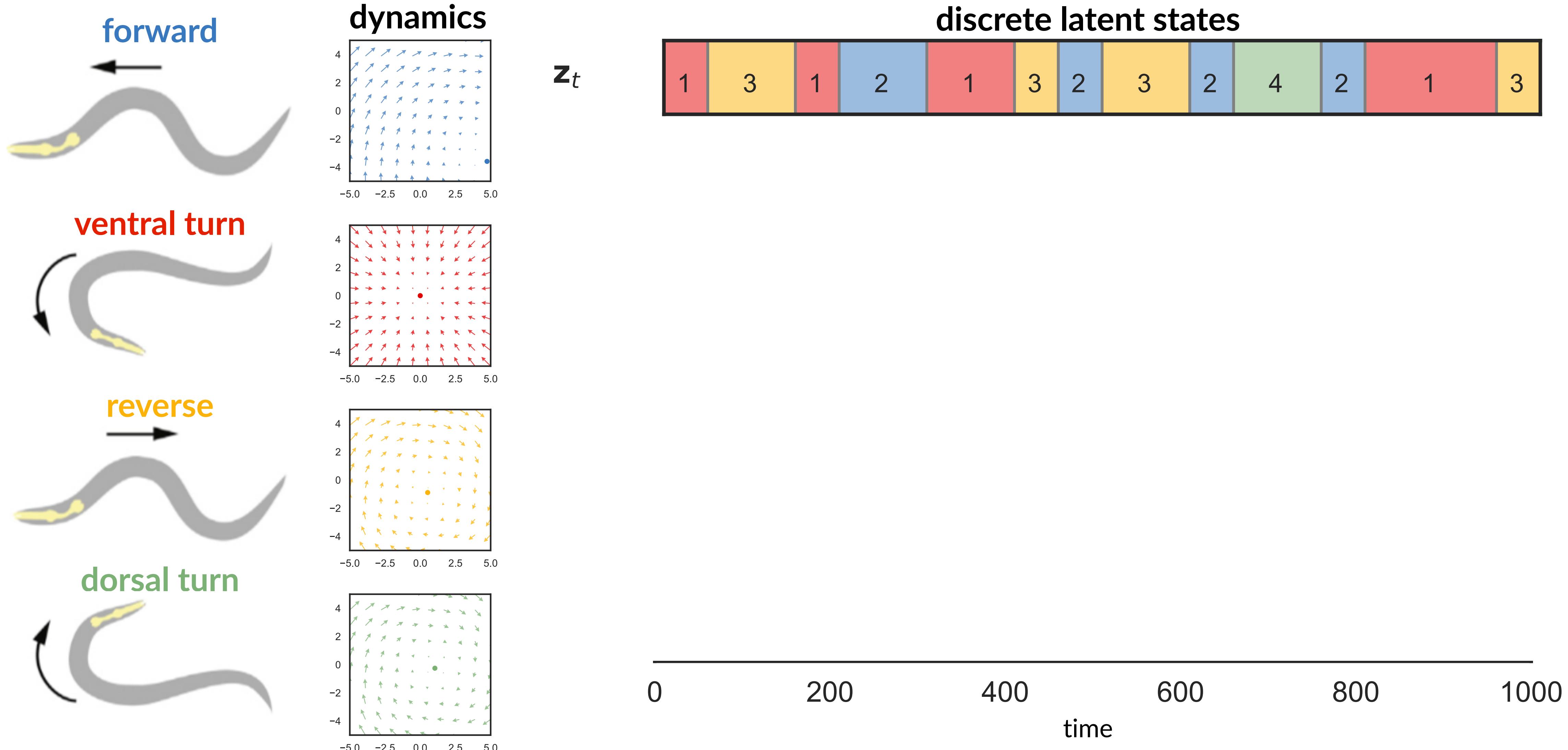


Kato et al (2015)

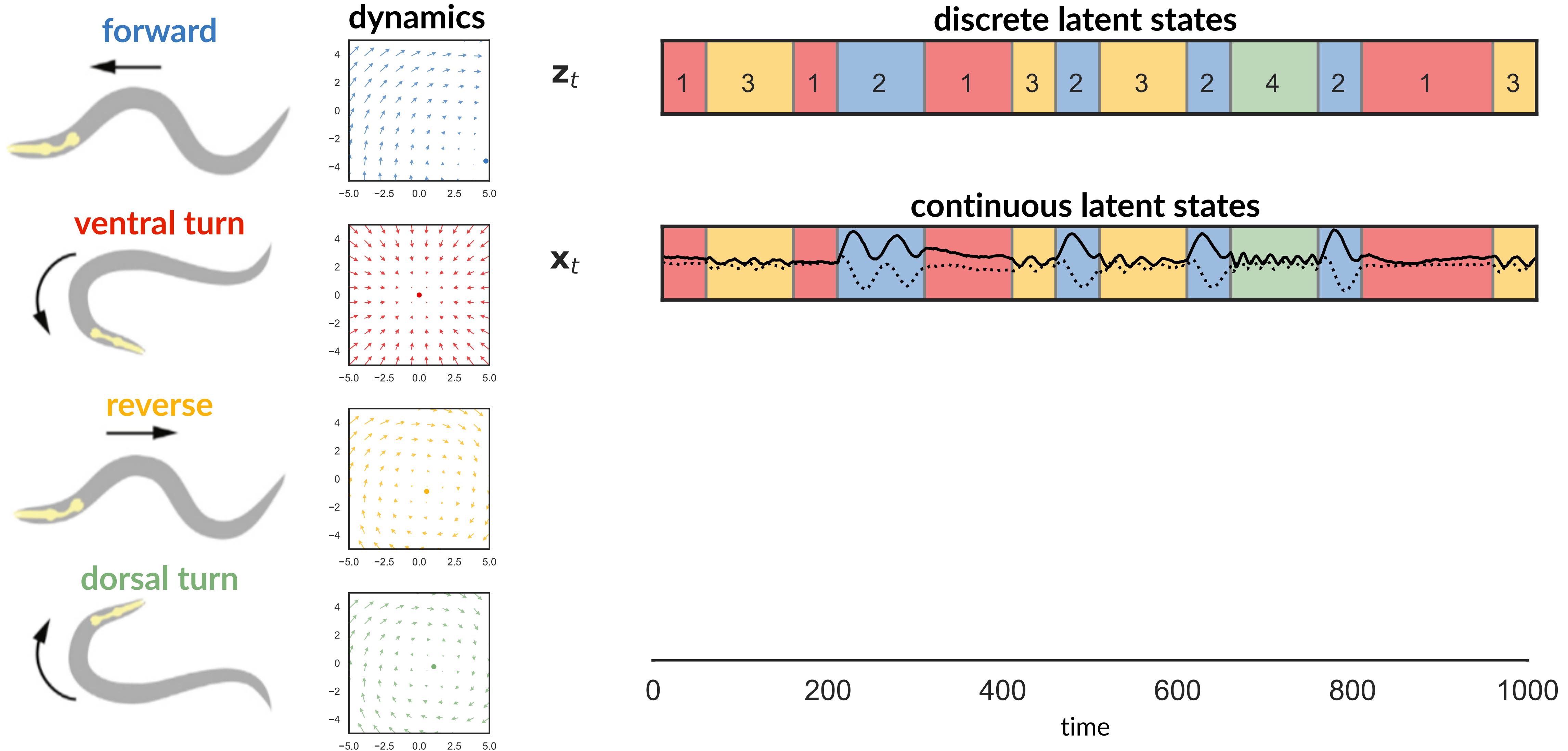
Building a probabilistic model of neural data



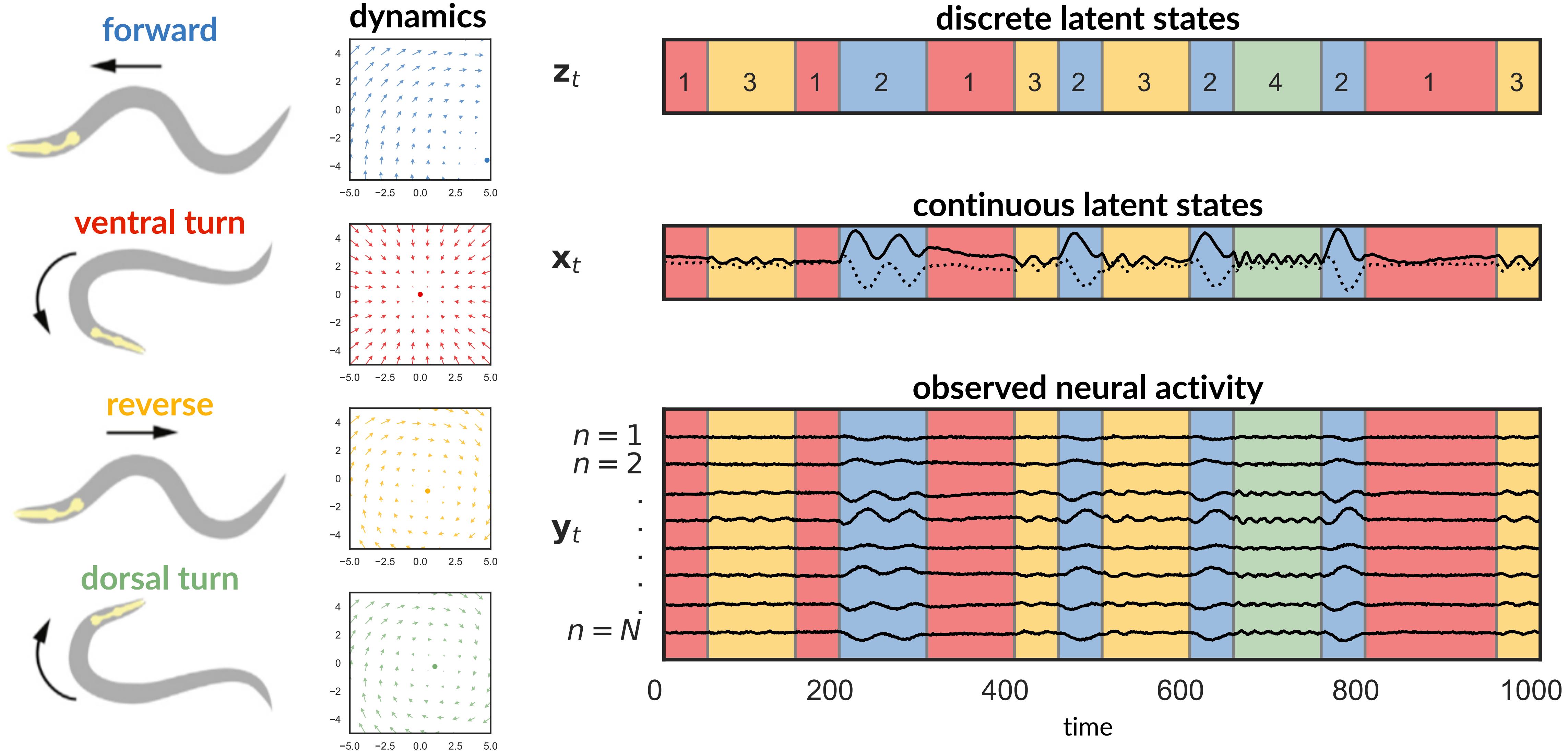
Building a probabilistic model of neural data



Building a probabilistic model of neural data



Building a probabilistic model of neural data



Hardness of exact EM for SLDS

Exact EM for the SLDS

- **E-step:** Update the posterior over latent variables,

$$q(z, x) \leftarrow p(z, x | y, \Theta) = \frac{p(z, x, y | \Theta)}{p(y | \Theta)}$$

- As before, we only need certain expectations under q ,

$$\mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t x_t^\top], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t x_{t+1}^\top],$$

- **M-step:** Update the parameters,

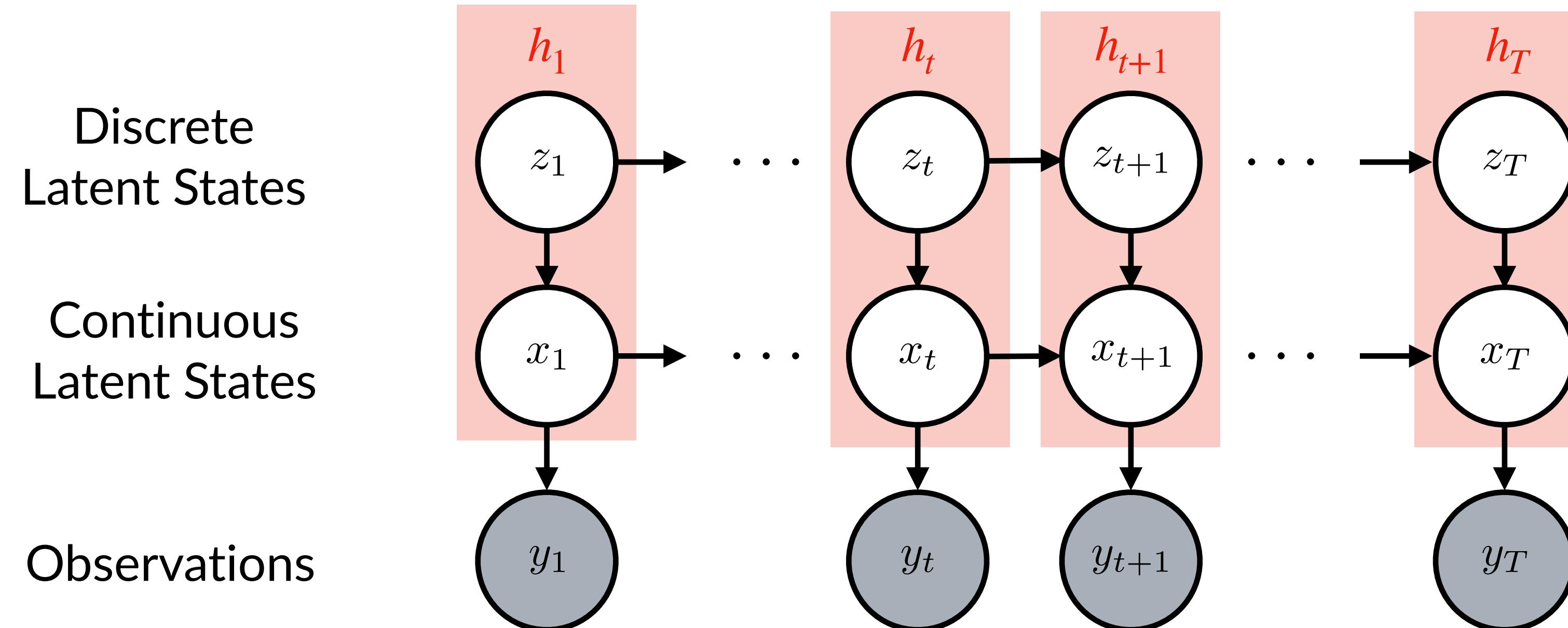
$$\Theta \leftarrow \arg \max \mathbb{E}_{q(z,x)} [\log p(z, x, y | \Theta)]$$

- Unfortunately, computing the necessary expectations is a lot harder now!

Combining the latent states

SLDS as a “hybrid” state space model

- Let $h_t = (z_t, x_t)$ denote the hybrid discrete & continuous latent state



○ = latent

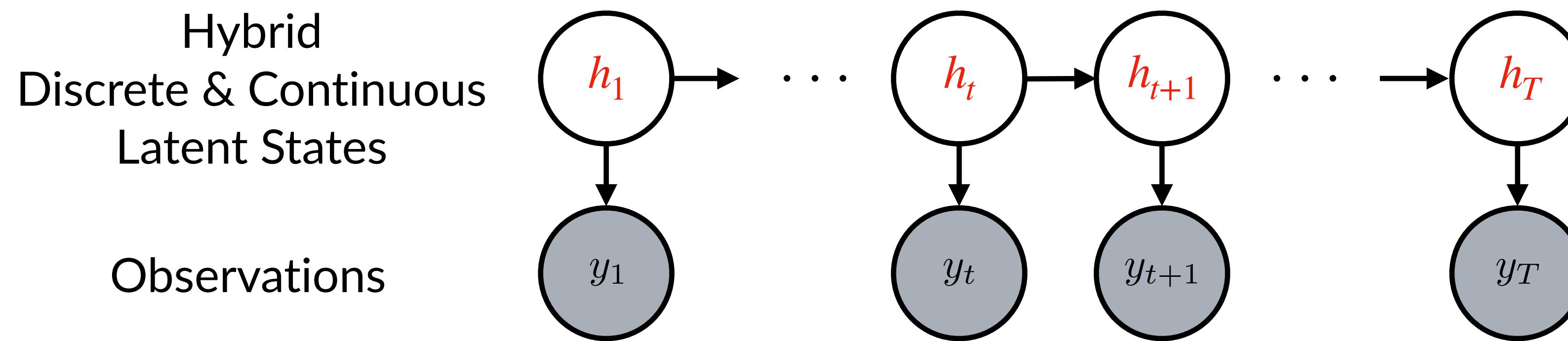
● = observed

→ = dependency

Combining the latent states

SLDS as a “hybrid” state space model

- Let $h_t = (z_t, x_t)$ denote the hybrid discrete & continuous latent state



○ = latent

● = observed

→ = dependency

Exact EM for SLDS

Computing the marginal distributions

- Consider the marginal probability of the latent states at time t :

$$q(h_t) = \int dh_1 \cdots \int dh_{t-1} \int dh_{t+1} \cdots \int dh_T q(h_1, \dots, h_{t-1}, h_t, h_{t+1}, \dots, h_T)$$

Exact EM for SLDS

Computing the marginal distributions

- Consider the marginal probability of the latent states at time t :

$$\begin{aligned} q(h_t) &= \int dh_1 \cdots \int dh_{t-1} \int dh_{t+1} \cdots \int dh_T q(h_1, \dots, h_{t-1}, h_t, h_{t+1}, \dots, h_T) \\ &\propto \left[\int dh_1 \cdots \int dh_{t-1} p(h_1) \prod_{s=1}^{t-1} p(h_s | h_s) p(h_{s+1} | h_s) \right] \times \left[p(y_t | h_t) \right] \\ &\quad \times \left[\int dh_{t+1} \cdots \int dh_T \prod_{u=t+1}^T p(h_u | h_{u-1}) p(y_u | h_u) \right] \end{aligned}$$

Exact EM for SLDS

Computing the marginal distributions

- Consider the marginal probability of the latent states at time t :

$$\begin{aligned} q(h_t) &= \int dh_1 \cdots \int dh_{t-1} \int dh_{t+1} \cdots \int dh_T q(h_1, \dots, h_{t-1}, h_t, h_{t+1}, \dots, h_T) \\ &\propto \left[\int dh_1 \cdots \int dh_{t-1} p(h_1) \prod_{s=1}^{t-1} p(h_s | h_s) p(h_{s+1} | h_s) \right] \times \left[p(y_t | h_t) \right] \\ &\quad \times \left[\int dh_{t+1} \cdots \int dh_T \prod_{u=t+1}^T p(h_u | h_{u-1}) p(y_u | h_u) \right] \\ &\triangleq \alpha_t(h_t) \times p(y_t | h_t) \times \beta_t(h_t) \end{aligned}$$

Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- Consider the “forward messages”:

$$\alpha_t(h_t) \triangleq \int dh_1 \cdots \int dh_{t-1} p(h_1) \prod_{s=1}^{t-1} p(h_s | h_s) p(h_{s+1} | h_s)$$

Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- Consider the “forward messages”:

$$\begin{aligned}\alpha_t(h_t) &\triangleq \int dh_1 \cdots \int dh_{t-1} p(h_1) \prod_{s=1}^{t-1} p(h_s | h_s) p(h_{s+1} | h_s) \\ &= \int dh_{t-1} \left[\left(\int dh_1 \cdots \int dh_{t-2} p(h_1) \prod_{s=1}^{t-2} p(y_s | h_s) p(h_{s+1} | h_s) \right) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1}) \right]\end{aligned}$$

Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- Consider the “forward messages”:

$$\begin{aligned}\alpha_t(h_t) &\triangleq \int dh_1 \cdots \int dh_{t-1} p(h_1) \prod_{s=1}^{t-1} p(h_s | h_s) p(h_{s+1} | h_s) \\ &= \int dh_{t-1} \left[\left(\int dh_1 \cdots \int dh_{t-2} p(h_1) \prod_{s=1}^{t-2} p(y_s | h_s) p(h_{s+1} | h_s) \right) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1}) \right] \\ &= \int dh_{t-1} \alpha_{t-1}(h_{t-1}) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1})\end{aligned}$$

- We can compute these messages **recursively!**

Hardness of Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- Now substitute $h_t = (z_t, x_t)$.
- Base case:

$$\begin{aligned}\alpha_1(z_1, x_1) &= p(z_1) p(x_1 \mid z_1) \\ &= \sum_{k=1}^K [\text{Cat}(z_1 \mid \pi) \mathcal{N}(x_1 \mid b_k, Q_k)]^{\mathbb{I}[z_1=k]}\end{aligned}$$

Hardness of Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- Now substitute $h_t = (z_t, x_t)$.
- Base case:

$$\begin{aligned}\alpha_1(z_1, x_1) &= p(z_1) p(x_1 \mid z_1) \\ &= \sum_{k=1}^K \left[\text{Cat}(z_1 \mid \pi) \mathcal{N}(x_1 \mid b_k, Q_k) \right]^{\mathbb{I}[z_1=k]}\end{aligned}$$

- Second time step:

$$\begin{aligned}\alpha_2(z_2, x_2) &= \sum_{z_1=1}^K \int dx_1 \alpha_1(z_1, x_1) p(y_1 \mid x_1) p(z_2 \mid z_1) p(x_2 \mid x_1, z_2) \\ &= \sum_{z_1=1}^K \sum_{k=1}^K \left[\text{Cat}(z_2 \mid \rho(z_1, y_1)) \mathcal{N}(x_2 \mid \mu(z_1, z_2, y_1), \Sigma(z_1, z_2, y_1)) \right]^{\mathbb{I}[z_2=k]}\end{aligned}$$

Hardness of Exact EM for SLDS

Computing the forward messages $\alpha_t(h_t)$

- By the t -th time step,

$$\alpha_t(z_t, x_t) = \sum_{z_1=1}^K \cdots \sum_{z_{t-1}=1}^K \sum_{k=1}^K \left[\text{Cat}\left(z_t \mid \rho(z_{1:t-1}, y_{1:t-1})\right) \mathcal{N}\left(x_t \mid \mu(z_t, z_{1:t-1}, y_{1:t-1}), \Sigma(z_t, z_{1:t-1}, y_{1:t-1})\right) \right]^{\mathbb{I}[z_t=k]}$$

- This is still a mixture of Gaussians.
- **Question:** How many components does it have?

Approximate Message Passing

Approximating the forward messages

Assumed density filtering (ADF)

Approximating the forward messages

Assumed density filtering (ADF)

- **Idea:** Approximate the forward messages with a tractable family of distributions \mathcal{A} .
 - For example, assume $\alpha_t(h_t) \approx \tilde{\alpha}_t(h_t) \in \mathcal{A}$ where \mathcal{A} is the set of Gaussian mixtures with at most M components.

Approximating the forward messages

Assumed density filtering (ADF)

- **Idea:** Approximate the forward messages with a tractable family of distributions \mathcal{A} .
 - For example, assume $\alpha_t(h_t) \approx \tilde{\alpha}_t(h_t) \in \mathcal{A}$ where \mathcal{A} is the set of Gaussian mixtures with at most M components.
 - Suppose we have $\tilde{\alpha}_{t-1}(h_{t-1}) \in \mathcal{A}$. Our **target for** $\alpha_t(h_t)$ is,

$$\hat{\alpha}_t(h_t) \triangleq \int dh_{t-1} \tilde{\alpha}_{t-1}(h_{t-1}) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1})$$

This may not be in \mathcal{A} !

Approximating the forward messages

Assumed density filtering (ADF)

- **Idea:** Approximate the forward messages with a tractable family of distributions \mathcal{A} .
 - For example, assume $\alpha_t(h_t) \approx \tilde{\alpha}_t(h_t) \in \mathcal{A}$ where \mathcal{A} is the set of Gaussian mixtures with at most M components.
 - Suppose we have $\tilde{\alpha}_{t-1}(h_{t-1}) \in \mathcal{A}$. Our **target for** $\alpha_t(h_t)$ is,

$$\hat{\alpha}_t(h_t) \triangleq \int dh_{t-1} \tilde{\alpha}_{t-1}(h_{t-1}) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1})$$

This may not be in \mathcal{A} !

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(h_t)$:

$$\tilde{\alpha}_t(h_t) \leftarrow \arg \min_{\mathcal{A}} D\left(\tilde{\alpha}_t(h_t) \parallel \hat{\alpha}_t(h_t)\right)$$

where $D(\cdot \parallel \cdot)$ is a measure of divergence between two densities; e.g. reverse Kullback-Leibler (KL) divergence.

Approximating the forward messages

Assumed density filtering (ADF)

- **Idea:** Approximate the forward messages with a tractable family of distributions \mathcal{A} .
 - For example, assume $\alpha_t(h_t) \approx \tilde{\alpha}_t(h_t) \in \mathcal{A}$ where \mathcal{A} is the set of Gaussian mixtures with at most M components.
 - Suppose we have $\tilde{\alpha}_{t-1}(h_{t-1}) \in \mathcal{A}$. Our **target for** $\alpha_t(h_t)$ is,

$$\hat{\alpha}_t(h_t) \triangleq \int dh_{t-1} \tilde{\alpha}_{t-1}(h_{t-1}) p(y_{t-1} | h_{t-1}) p(h_t | h_{t-1})$$

This may not be in \mathcal{A} !

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(h_t)$:

$$\tilde{\alpha}_t(h_t) \leftarrow \arg \min_{\mathcal{A}} D\left(\tilde{\alpha}_t(h_t) \parallel \hat{\alpha}_t(h_t)\right)$$

where $D(\cdot \parallel \cdot)$ is a measure of divergence between two densities; e.g. reverse Kullback-Leibler (KL) divergence.

- This is called **assumed density filtering (ADF)**, and it is closely related to **expectation-propagation (EP)** and the **unscented/extended Kalman filter**.

Approximating the forward messages

Assumed density filtering (ADF) with GMM messages

- For example, let \mathcal{A} be the set of Gaussian mixtures with parameters $\rho_t, \{\mu_{t,k}, \Sigma_{t,k}\}_{k=1}^K$:

$$\tilde{\alpha}_t(z_t, x_t) = \tilde{\alpha}_t(z_t) \tilde{\alpha}_t(x_t | z_t) = \text{Cat}(z_t | \rho_t) \mathcal{N}(x_t | \mu_{t,z_t}, \Sigma_{t,z_t}).$$

Approximating the forward messages

Assumed density filtering (ADF) with GMM messages

- For example, let \mathcal{A} be the set of Gaussian mixtures with parameters $\rho_t, \{\mu_{t,k}, \Sigma_{t,k}\}_{k=1}^K$:

$$\tilde{\alpha}_t(z_t, x_t) = \tilde{\alpha}_t(z_t) \tilde{\alpha}_t(x_t | z_t) = \text{Cat}(z_t | \rho_t) \mathcal{N}(x_t | \mu_{t,z_t}, \Sigma_{t,z_t}).$$

- The target is

$$\hat{\alpha}_t(z_t, x_t) \triangleq \sum_{z_{t-1}=1}^K \int dx_{t-1} \tilde{\alpha}_{t-1}(z_{t-1}, x_{t-1}) p(y_{t-1} | x_{t-1}) p(z_t | z_{t-1}) p(x_t | x_{t-1}, z_{t-1})$$

Approximating the forward messages

Assumed density filtering (ADF) with GMM messages

- For example, let \mathcal{A} be the set of Gaussian mixtures with parameters $\rho_t, \{\mu_{t,k}, \Sigma_{t,k}\}_{k=1}^K$:

$$\tilde{\alpha}_t(z_t, x_t) = \tilde{\alpha}_t(z_t) \tilde{\alpha}_t(x_t | z_t) = \text{Cat}(z_t | \rho_t) \mathcal{N}(x_t | \mu_{t,z_t}, \Sigma_{t,z_t}).$$

- The target is

$$\begin{aligned}\hat{\alpha}_t(z_t, x_t) &\triangleq \sum_{z_{t-1}=1}^K \int dx_{t-1} \tilde{\alpha}_{t-1}(z_{t-1}, x_{t-1}) p(y_{t-1} | x_{t-1}) p(z_t | z_{t-1}) p(x_t | x_{t-1}, z_{t-1}) \\ &= \sum_{z_{t-1}=1}^K \int dx_{t-1} \rho_{t-1, z_{t-1}} \mathcal{N}(x_{t-1} | \mu_{t-1, z_{t-1}}, \Sigma_{t-1, z_{t-1}}) \mathcal{N}(y_{t-1} | Cx_{t-1} + d, R) P_{z_{t-1}, z_t} \mathcal{N}(x_t | A_{z_t} x_{t-1} + b_{z_t}, Q_{z_t})\end{aligned}$$

Approximating the forward messages

Assumed density filtering (ADF) with GMM messages

- For example, let \mathcal{A} be the set of Gaussian mixtures with parameters $\rho_t, \{\mu_{t,k}, \Sigma_{t,k}\}_{k=1}^K$:

$$\tilde{\alpha}_t(z_t, x_t) = \tilde{\alpha}_t(z_t) \tilde{\alpha}_t(x_t | z_t) = \text{Cat}(z_t | \rho_t) \mathcal{N}(x_t | \mu_{t,z_t}, \Sigma_{t,z_t}).$$

- The target is

$$\begin{aligned}\hat{\alpha}_t(z_t, x_t) &\triangleq \sum_{z_{t-1}=1}^K \int dx_{t-1} \tilde{\alpha}_{t-1}(z_{t-1}, x_{t-1}) p(y_{t-1} | x_{t-1}) p(z_t | z_{t-1}) p(x_t | x_{t-1}, z_{t-1}) \\ &= \sum_{z_{t-1}=1}^K \int dx_{t-1} \rho_{t-1, z_{t-1}} \mathcal{N}(x_{t-1} | \mu_{t-1, z_{t-1}}, \Sigma_{t-1, z_{t-1}}) \mathcal{N}(y_{t-1} | Cx_{t-1} + d, R) P_{z_{t-1}, z_t} \mathcal{N}(x_t | A_{z_t} x_{t-1} + b_{z_t}, Q_{z_t}) \\ &\propto \sum_{z_{t-1}=1}^K \rho(z_t, z_{t-1}) \mathcal{N}(x_t | \mu(z_t, z_{t-1}), \Sigma(z_t, z_{t-1}))\end{aligned}$$

- This is a GMM with K components for each assignment of z_t . We want to approximate it a single Gaussian for each z_t .

Approximating the forward messages

Projecting onto the set \mathcal{A}

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(z_t, h_t)$ by minimizing the reverse Kullback-Leibler divergence,

$$\text{KL}\left(\hat{\alpha}_t(z_t, x_t) \parallel \tilde{\alpha}_t(z_t, x_t)\right) = \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \hat{\alpha}_t(z_t, x_t) - \log \tilde{\alpha}_t(z_t, x_t)]$$

Approximating the forward messages

Projecting onto the set \mathcal{A}

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(z_t, h_t)$ by minimizing the reverse Kullback-Leibler divergence,

$$\begin{aligned}\text{KL}\left(\hat{\alpha}_t(z_t, x_t) \parallel \tilde{\alpha}_t(z_t, x_t)\right) &= \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \hat{\alpha}_t(z_t, x_t) - \log \tilde{\alpha}_t(z_t, x_t)] \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \tilde{\alpha}_t(z_t, x_t)] + c\end{aligned}$$

Approximating the forward messages

Projecting onto the set \mathcal{A}

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(z_t, h_t)$ by minimizing the reverse Kullback-Leibler divergence,

$$\begin{aligned}\text{KL}\left(\hat{\alpha}_t(z_t, x_t) \parallel \tilde{\alpha}_t(z_t, x_t)\right) &= \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \hat{\alpha}_t(z_t, x_t) - \log \tilde{\alpha}_t(z_t, x_t)] \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \tilde{\alpha}_t(z_t, x_t)] + c \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \text{Cat}(z_t \mid \rho_t) + \log \mathcal{N}(x_t \mid \mu_{t,z_t}, \Sigma_{t,z_t})] + c\end{aligned}$$

Approximating the forward messages

Projecting onto the set \mathcal{A}

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(z_t, h_t)$ by minimizing the reverse Kullback-Leibler divergence,

$$\begin{aligned} \text{KL}\left(\hat{\alpha}_t(z_t, x_t) \parallel \tilde{\alpha}_t(z_t, x_t)\right) &= \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \hat{\alpha}_t(z_t, x_t) - \log \tilde{\alpha}_t(z_t, x_t)] \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \tilde{\alpha}_t(z_t, x_t)] + c \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \text{Cat}(z_t \mid \rho_t) + \log \mathcal{N}(x_t \mid \mu_{t,z_t}, \Sigma_{t,z_t})] + c \\ &= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} \left[\sum_{k=1}^K \mathbb{I}[z_t = k] \left(\rho_{tk} - \frac{1}{2} \log |\Sigma_{tk}| - \frac{1}{2} x_t^\top \Sigma_{tk}^{-1} x_t + \mu_{tk}^\top \Sigma_{tk}^{-1} x_t - \frac{1}{2} \mu_{tk}^\top \Sigma_{tk}^{-1} \mu_{tk} \right) \right] + c \end{aligned}$$

Approximating the forward messages

Projecting onto the set \mathcal{A}

- Find the member of \mathcal{A} that best approximates $\hat{\alpha}_t(z_t, h_t)$ by minimizing the reverse Kullback-Leibler divergence,

$$\begin{aligned}
\text{KL}\left(\hat{\alpha}_t(z_t, x_t) \parallel \tilde{\alpha}_t(z_t, x_t)\right) &= \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \hat{\alpha}_t(z_t, x_t) - \log \tilde{\alpha}_t(z_t, x_t)] \\
&= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} [\log \tilde{\alpha}_t(z_t, x_t)] + c \\
&= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} \left[\log \text{Cat}(z_t \mid \rho_t) + \log \mathcal{N}(x_t \mid \mu_{t,z_t}, \Sigma_{t,z_t}) \right] + c \\
&= -\mathbb{E}_{\hat{\alpha}_t(z_t, x_t)} \left[\sum_{k=1}^K \mathbb{I}[z_t = k] \left(\rho_{tk} - \frac{1}{2} \log |\Sigma_{tk}| - \frac{1}{2} x_t^\top \Sigma_{tk}^{-1} x_t + \mu_{tk}^\top \Sigma_{tk}^{-1} x_t - \frac{1}{2} \mu_{tk}^\top \Sigma_{tk}^{-1} \mu_{tk} \right) \right] + c \\
&= -\sum_{k=1}^K \left\langle \rho_{tk} - \frac{1}{2} \log |\Sigma_{tk}|, \bar{N}_k \right\rangle + \left\langle -\frac{1}{2} \Sigma_{tk}^{-1}, \bar{\psi}_{k,1} \right\rangle + \left\langle \Sigma_{tk}^{-1} \mu_{tk}, \bar{\psi}_{k,2} \right\rangle + \left\langle -\frac{1}{2} \mu_{tk}^\top \Sigma_{tk}^{-1} \mu_{tk}, \bar{\psi}_{k,3} \right\rangle + c
\end{aligned}$$

Where $\bar{N}_k = \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)}[\mathbb{I}[z_t = k]]$, $\bar{\psi}_{k,1} = \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)}[\mathbb{I}[z_t = k] x_t x_t^\top]$, $\bar{\psi}_{k,2} = \mathbb{E}_{\hat{\alpha}_t(z_t, x_t)}[\mathbb{I}[z_t = k] x_t]$, $\bar{\psi}_{k,3} = \bar{N}_k$.

Conclusion

- Switching LDS combine ARHMMs and LDS to get the best of both worlds.
- They approximate nonlinear dynamical systems by switching between linear dynamical states.
- However, posterior inference is much harder! The posterior has K^T modes.
- Approximate message passing techniques leverage the Markovianity of the system with recursive updates, but approximate the messages with a simpler class.
- Next time, we'll look at alternative methods of approximate inference for SLDS.

Further Reading

- Barber, David. 2012. Bayesian Reasoning and Machine Learning. Cambridge University Press. **Chapter 25.**
- Linderman, Scott W., Matthew J. Johnson, Andrew C. Miller, Ryan P. Adams, David M. Blei, and Liam Paninski. 2017. “Bayesian Learning and Inference in Recurrent Switching Linear Dynamical Systems.” In Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS).