# SEQUENTIAL LATENT VARIABLE MODELS & FILTERING

*JOSEPH MARINO*

**CALTECH**

# OUTLINE

- sequence models

- amortized variational filtering

- sequential autoregressive flows
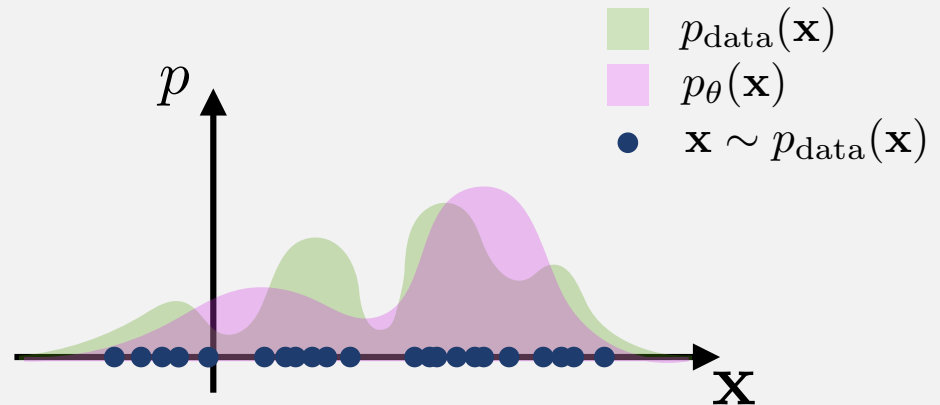
# SEQUENCE MODELS

# GENERATIVE MODEL

*a model of the density of observed data*

# MAXIMUM LIKELIHOOD



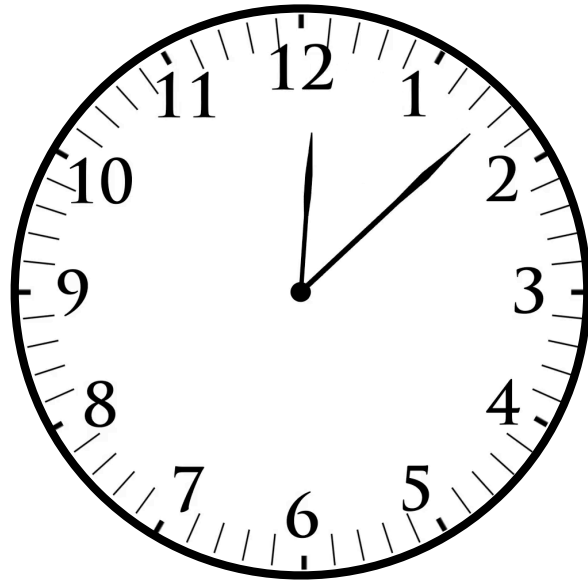data: $p_{\text{data}}(\mathbf{x})$
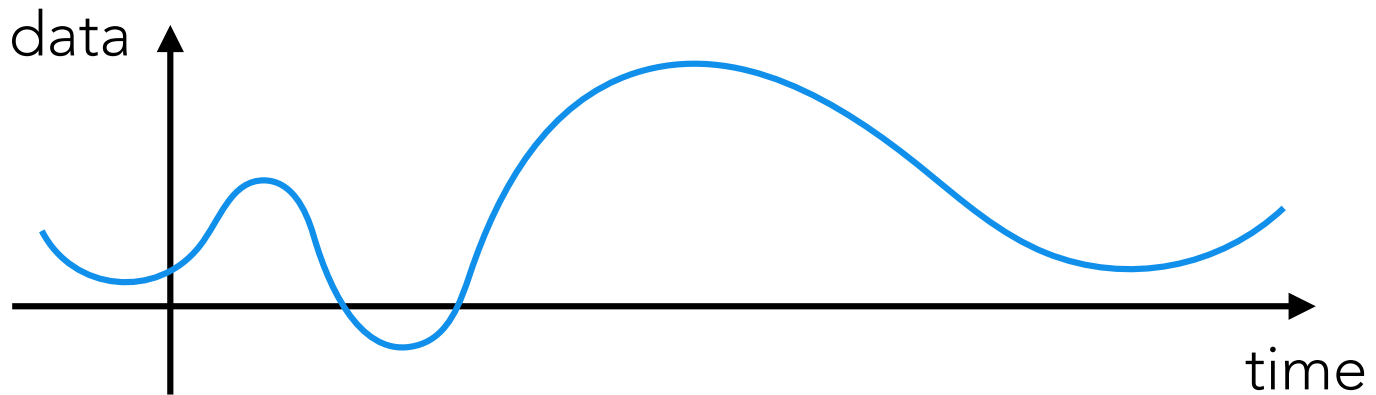
model: $p_\theta(\mathbf{x})$

parameters: $\theta$

Legend:
- $p_{\text{data}}(\mathbf{x})$
- $p_\theta(\mathbf{x})$
- $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$

**maximum likelihood estimation**

find the model that assigns the *maximum likelihood* to the data

$$\theta^* = \arg\max_\theta \; \mathbb{E}_{p_{\text{data}}(\mathbf{x})}\left[\log p_\theta(\mathbf{x})\right] \approx \frac{1}{N}\sum_{i=1}^{N}\log p_\theta(\mathbf{x}^{(i)})$$
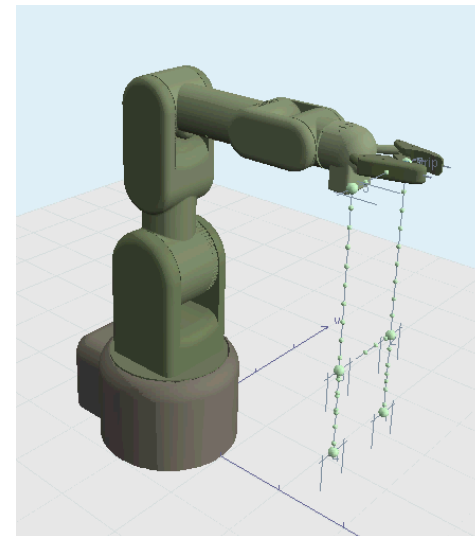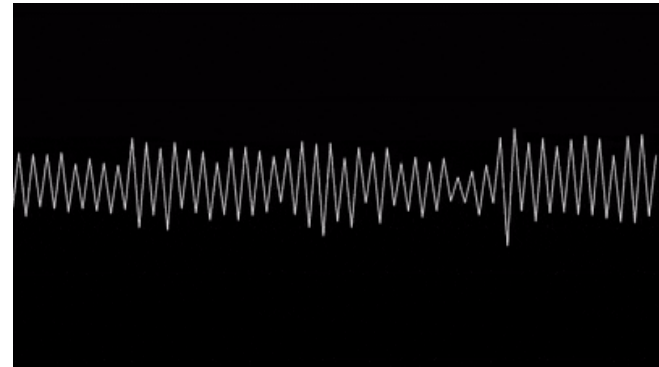
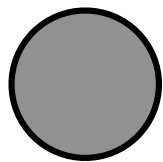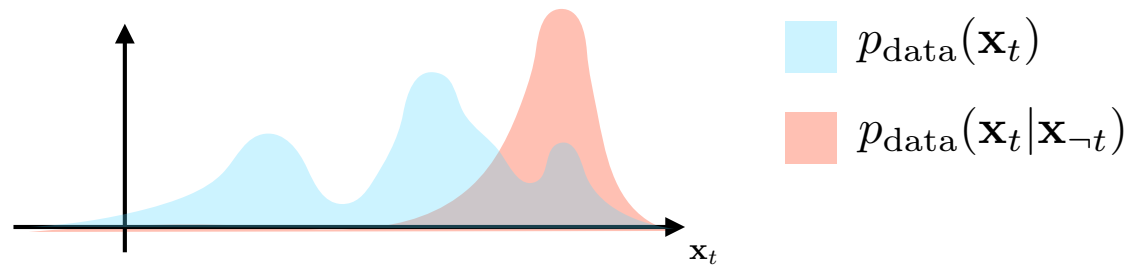**observed data are often sequential**
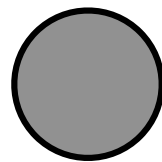
vision

audio



joint angles

# *dynamics*: dependence in time

*multi-information:*  $\mathcal{I}(\mathbf{x}_{1:T}) = \sum_{t} \mathcal{H}(\mathbf{x}_t) - \mathcal{H}(\mathbf{x}_{1:T}) \geq 0$

observing $\mathbf{x}_{\neg t}$ *reduces uncertainty* in $\mathbf{x}_t$



$p_{\mathrm{data}}(\mathbf{x}_t)$

$p_{\mathrm{data}}(\mathbf{x}_t | \mathbf{x}_{\neg t})$

$t - 1$          $t$          $t + 1$

*model temporal dependencies*



$$t-1 \qquad\qquad t \qquad\qquad t+1$$

*model temporal dependencies*

$t-1$       $t$       $t+1$
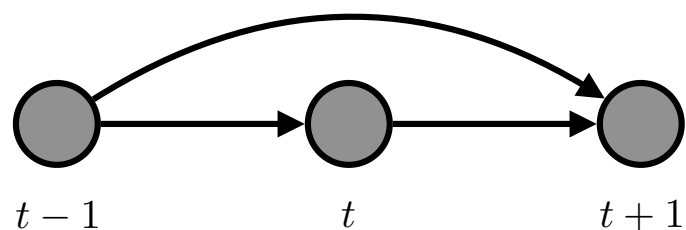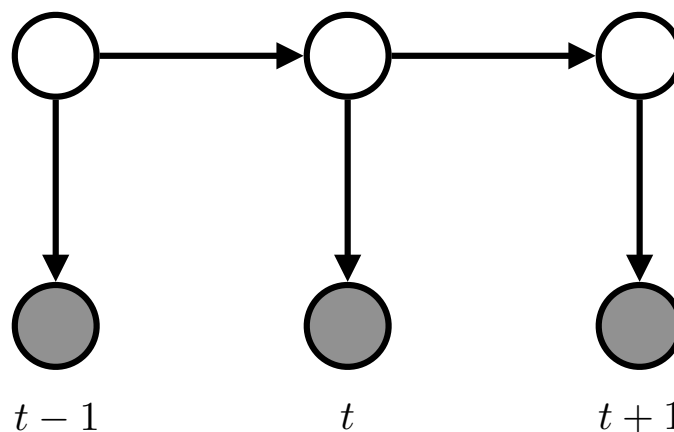
*model temporal dependencies*

$t-1$          $t$          $t+1$

# MODELING DYNAMICS



fully-observed

latent

$$p_\theta(\mathbf{x}_t|\mathbf{x}_{<t}) = \int p_\theta(\mathbf{x}_t|\mathbf{z}_t)p_\theta(\mathbf{z}_t|\mathbf{x}_{<t})d\mathbf{z}_t$$

mixture component

mixture probability

may be more flexible than a fixed-form $p_\theta(\mathbf{x}_t|\mathbf{x}_{<t})$

12

# SEQUENTIAL LATENT VARIABLE MODELS

general form:

$$p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^{T} \underbrace{p_\theta(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t})}_{\text{likelihood}} \underbrace{p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})}_{\text{prior}}$$

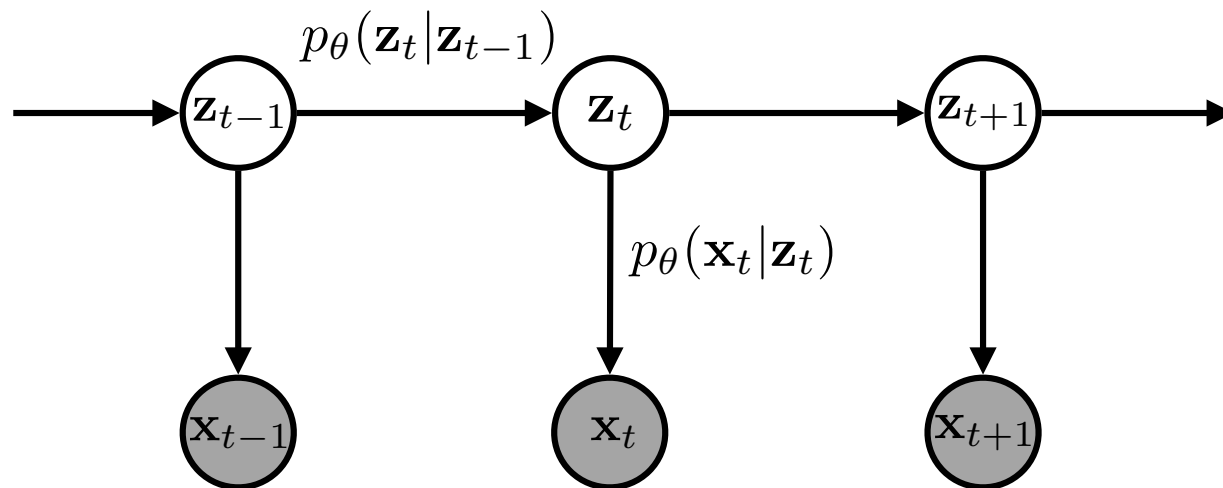where $\quad$ $\mathbf{x}_{\leq T}$ is a sequence of $T$ observed variables

$\mathbf{z}_{\leq T}$ is a sequence of $T$ latent variables
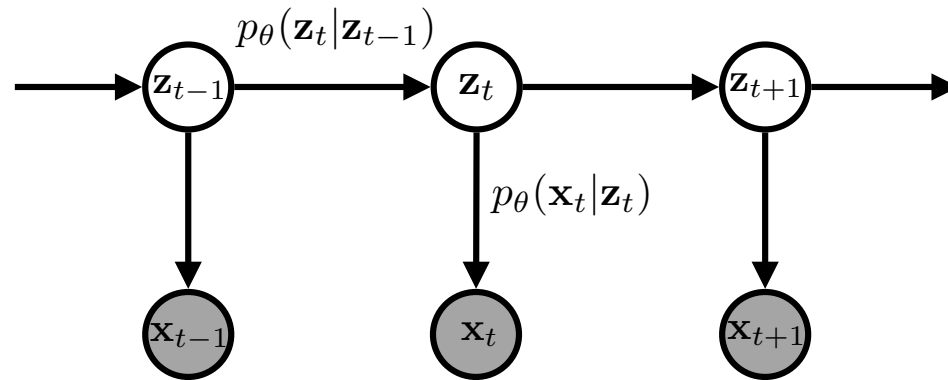
# SEQUENTIAL LATENT VARIABLE MODELS

general form:

$$p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^{T} \underbrace{p_\theta(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t})}_{\text{likelihood}} \underbrace{p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})}_{\text{prior}}$$

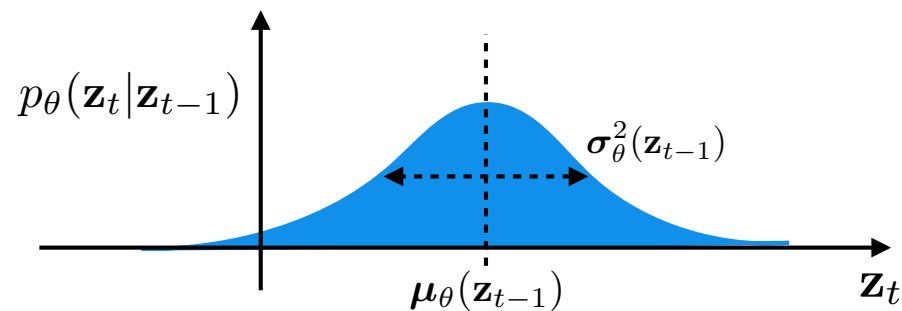simplified case (hidden Markov model):

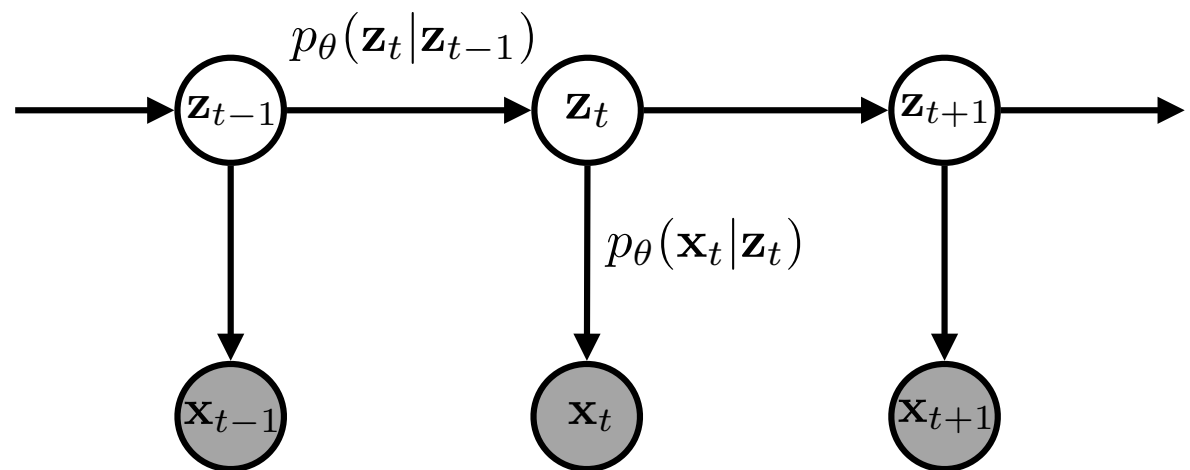# SEQUENTIAL LATENT VARIABLE MODELS

Markov model:

$$p_\theta(\mathbf{z}_t|\mathbf{z}_{t-1})$$

$$\mathbf{z}_{t-1} \longrightarrow \mathbf{z}_t \longrightarrow \mathbf{z}_{t+1}$$

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

$$\mathbf{x}_{t-1} \quad \mathbf{x}_t \quad \mathbf{x}_{t+1}$$

Parameterization:

$p_\theta(\mathbf{z}_t|\mathbf{z}_{t-1})$ is typically an analytical distribution

for example, $p_\theta(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_\theta(\mathbf{z}_{t-1}), \mathrm{diag}(\boldsymbol{\sigma}_\theta^2(\mathbf{z}_{t-1})))$

$p_\theta(\mathbf{z}_t|\mathbf{z}_{t-1})$

$\boldsymbol{\sigma}_\theta^2(\mathbf{z}_{t-1})$

$\boldsymbol{\mu}_\theta(\mathbf{z}_{t-1})$

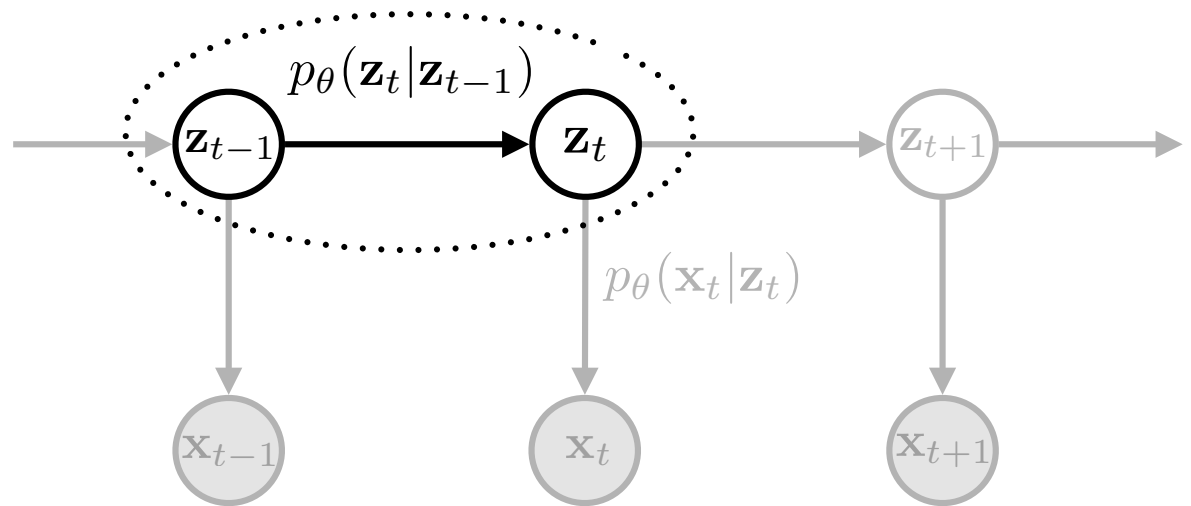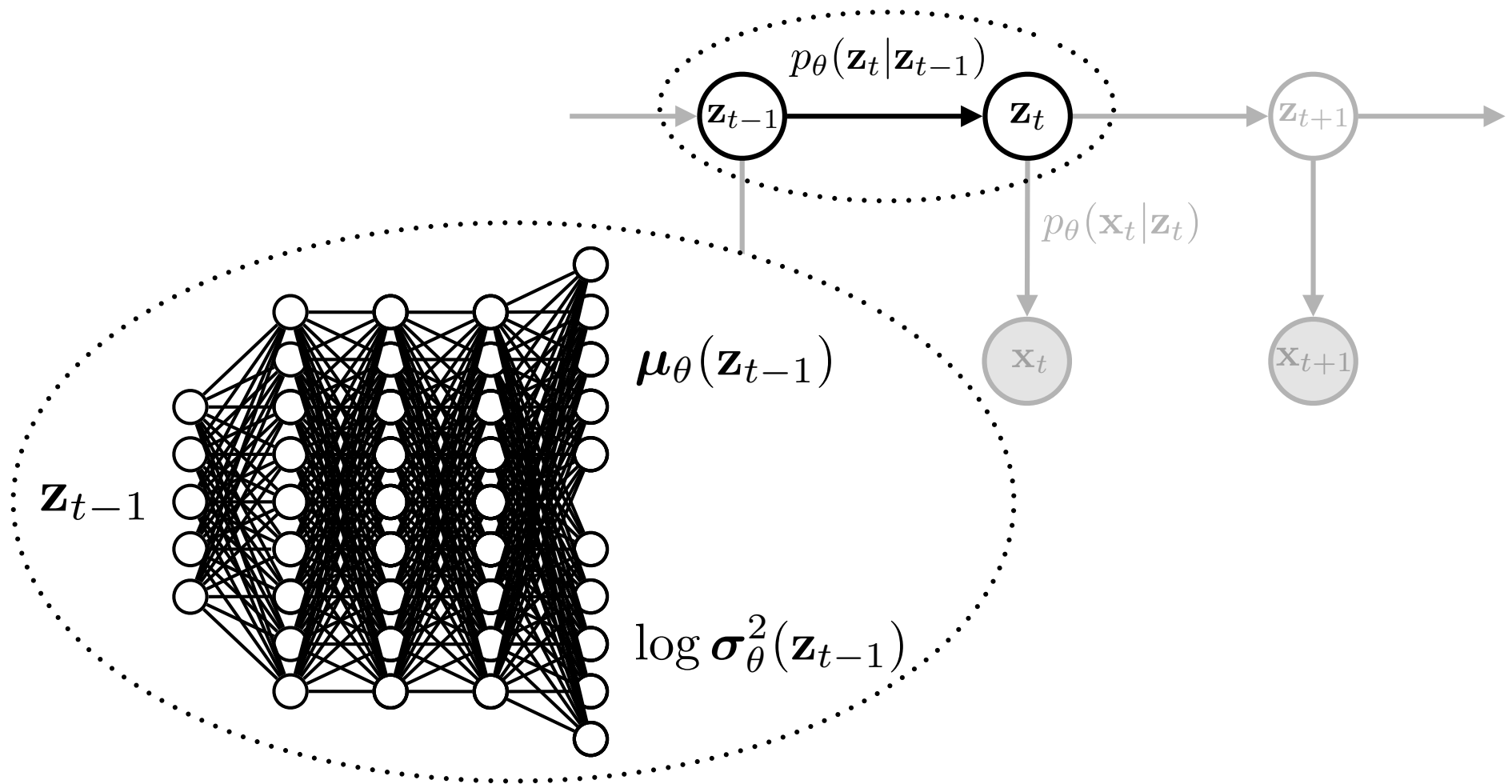$\mathbf{z}_t$

# SEQUENTIAL LATENT VARIABLE MODELS

the parameters of these analytical distributions are
functions, often *deep networks*
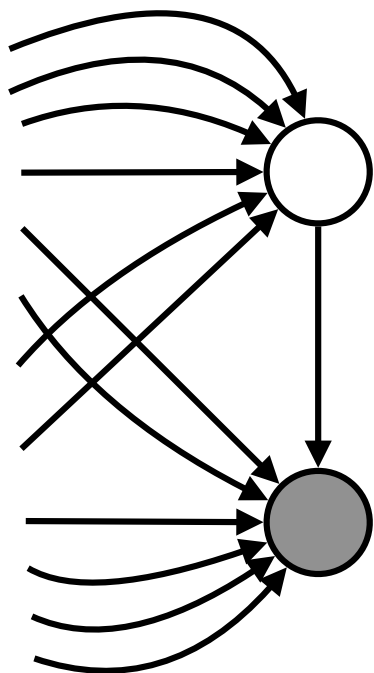
# SEQUENTIAL LATENT VARIABLE MODELS

the parameters of these analytical distributions are functions, often *deep networks*
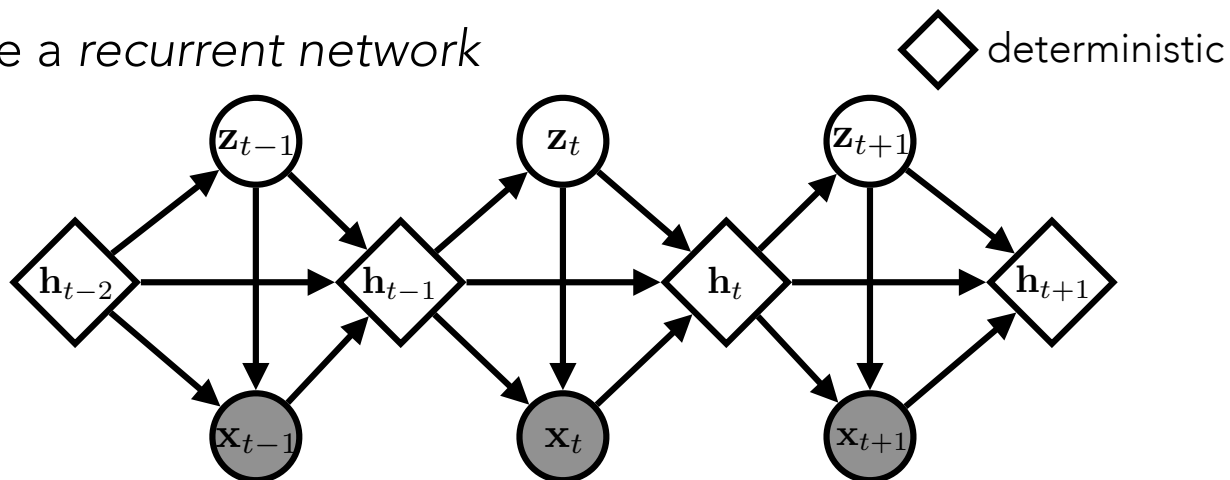
# SEQUENTIAL LATENT VARIABLE MODELS

the parameters of these analytical distributions are functions, often *deep networks*

# LONG-TERM DEPENDENCIES

general model form $\quad p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \displaystyle\prod_{t=1}^{T} p_\theta(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t}) p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$

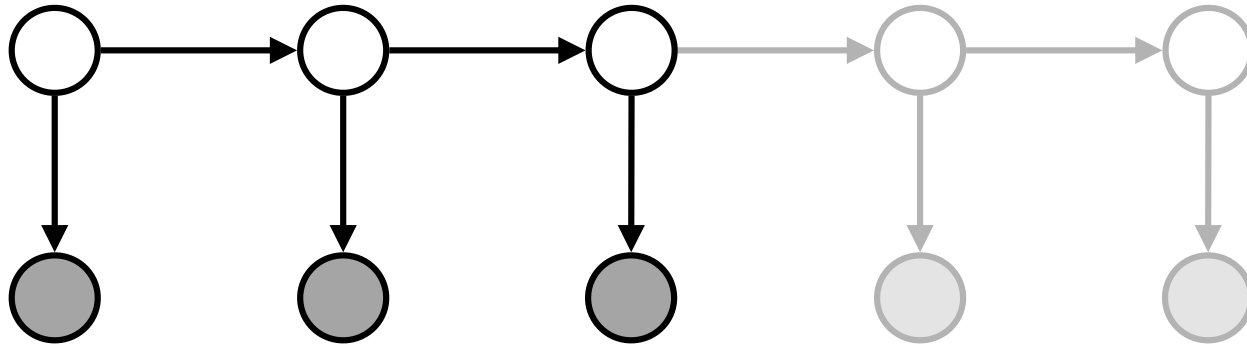how do we model long-term dependencies?

use a *recurrent network* $\qquad\qquad\qquad\qquad$ ◇ deterministic

$$p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^{T} p_\theta(\mathbf{x}_t | \mathbf{h}_{t-1}, \mathbf{z}_t) p_\theta(\mathbf{z}_t | \mathbf{h}_{t-1})$$
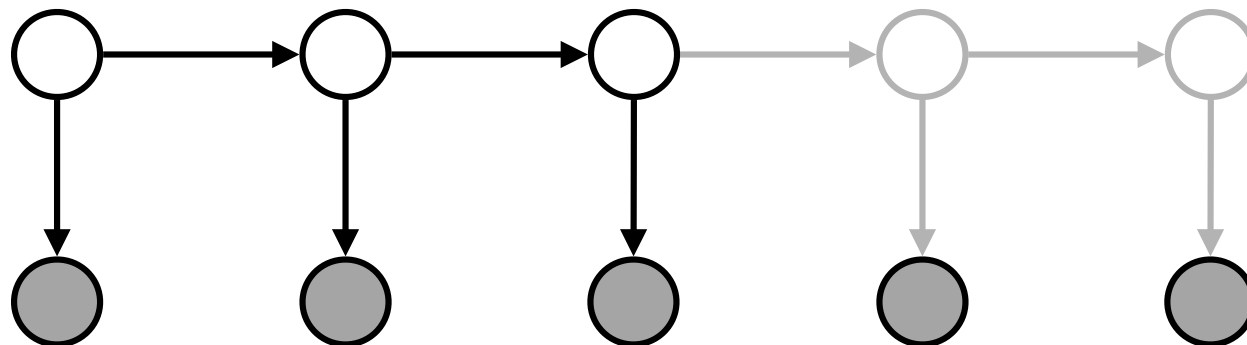
# INFERENCE

given a sequence of observations, $\mathbf{x}_{\leq T}$, infer $p_\theta(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})$
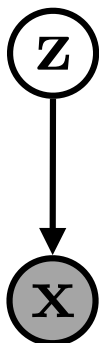


filtering inference



smoothing inference

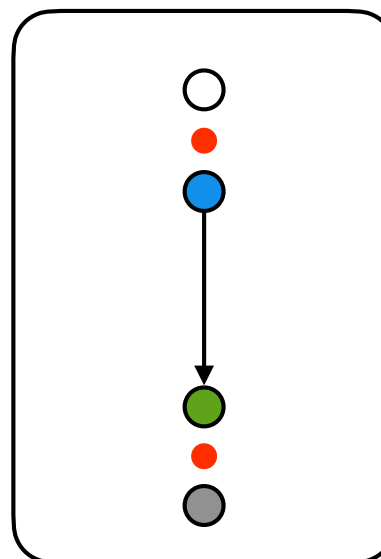# *ASIDE*: VARIATIONAL INFERENCE

# VARIATIONAL INFERENCE

graphical model

computation graph



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

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

*intractable*

# VARIATIONAL INFERENCE
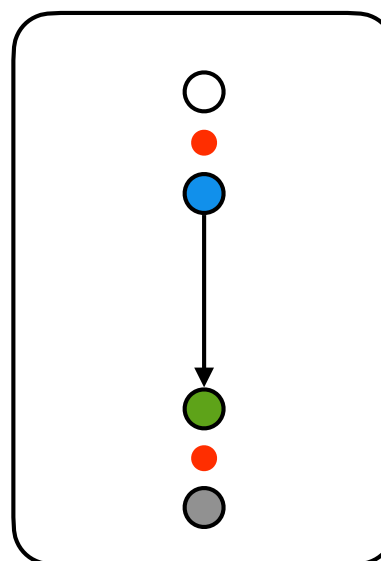
approximate posterior $q(\mathbf{z}|\mathbf{x})$

variational lower bound

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; q)$$

where

$$\mathcal{L}(\mathbf{x}; q) = \mathbb{E}_q \left[ \underbrace{\log p_\theta(\mathbf{x}|\mathbf{z})}_{\text{"reconstruction"}} - \underbrace{\log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})}}_{\text{"regularization"}} \right]$$

# VARIATIONAL INFERENCE

approximate posterior $q(\mathbf{z}|\mathbf{x})$

computation graph

*latent space*

variational lower bound

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; q)$$

$p_\theta(\mathbf{z})$

$-\mathbb{E}_q \left[ \log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} \right]$

$q(\mathbf{z}|\mathbf{x})$

where

$$\mathcal{L}(\mathbf{x}; q) = \mathbb{E}_q \left[ \underbrace{\log p_\theta(\mathbf{x}|\mathbf{z})}_{\text{"reconstruction"}} - \underbrace{\log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})}}_{\text{"regularization"}} \right]$$

# VARIATIONAL INFERENCE

approximate posterior $q(\mathbf{z}|\mathbf{x})$
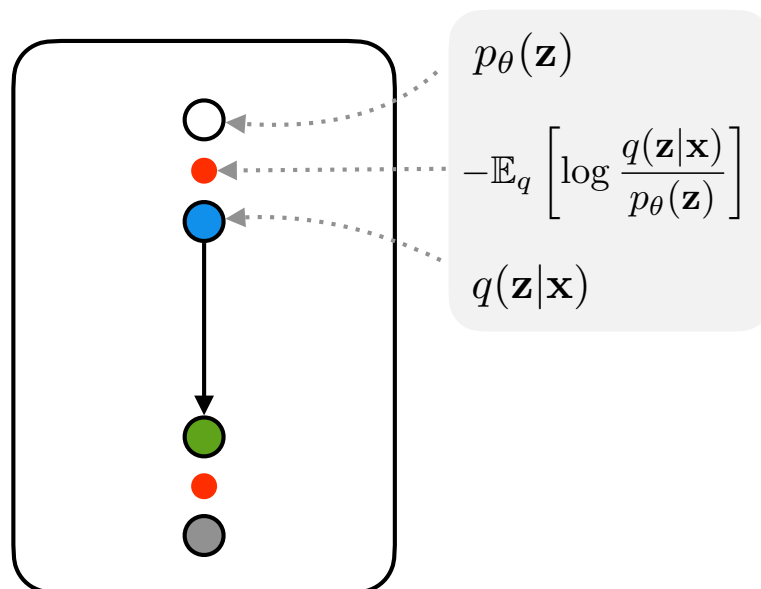
computation graph

*latent space*

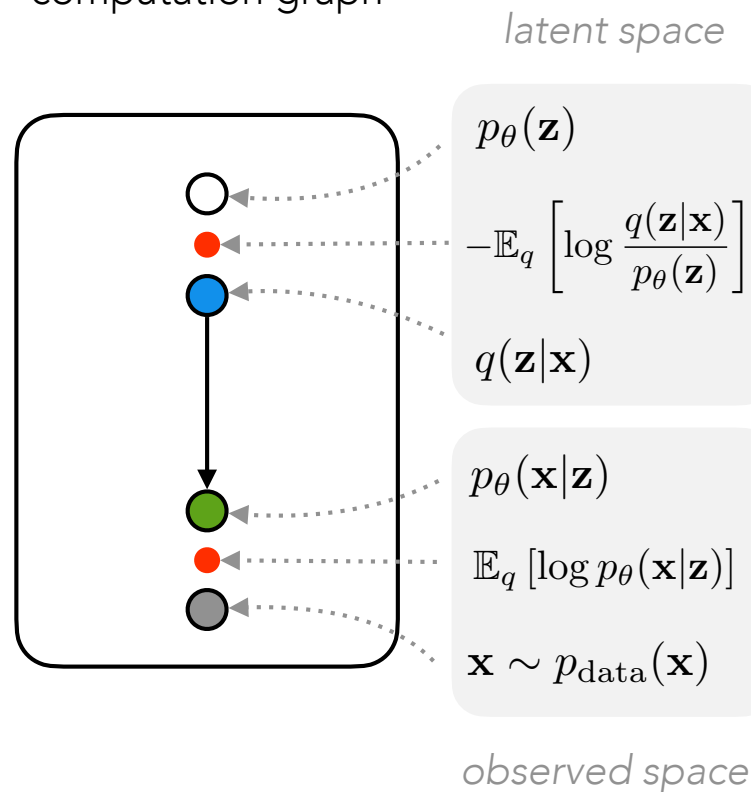variational lower bound

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; q)$$

where

$$\mathcal{L}(\mathbf{x}; q) = \mathbb{E}_q \left[ \underbrace{\log p_\theta(\mathbf{x}|\mathbf{z})}_{\text{"reconstruction"}} - \underbrace{\log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})}}_{\text{"regularization"}} \right]$$
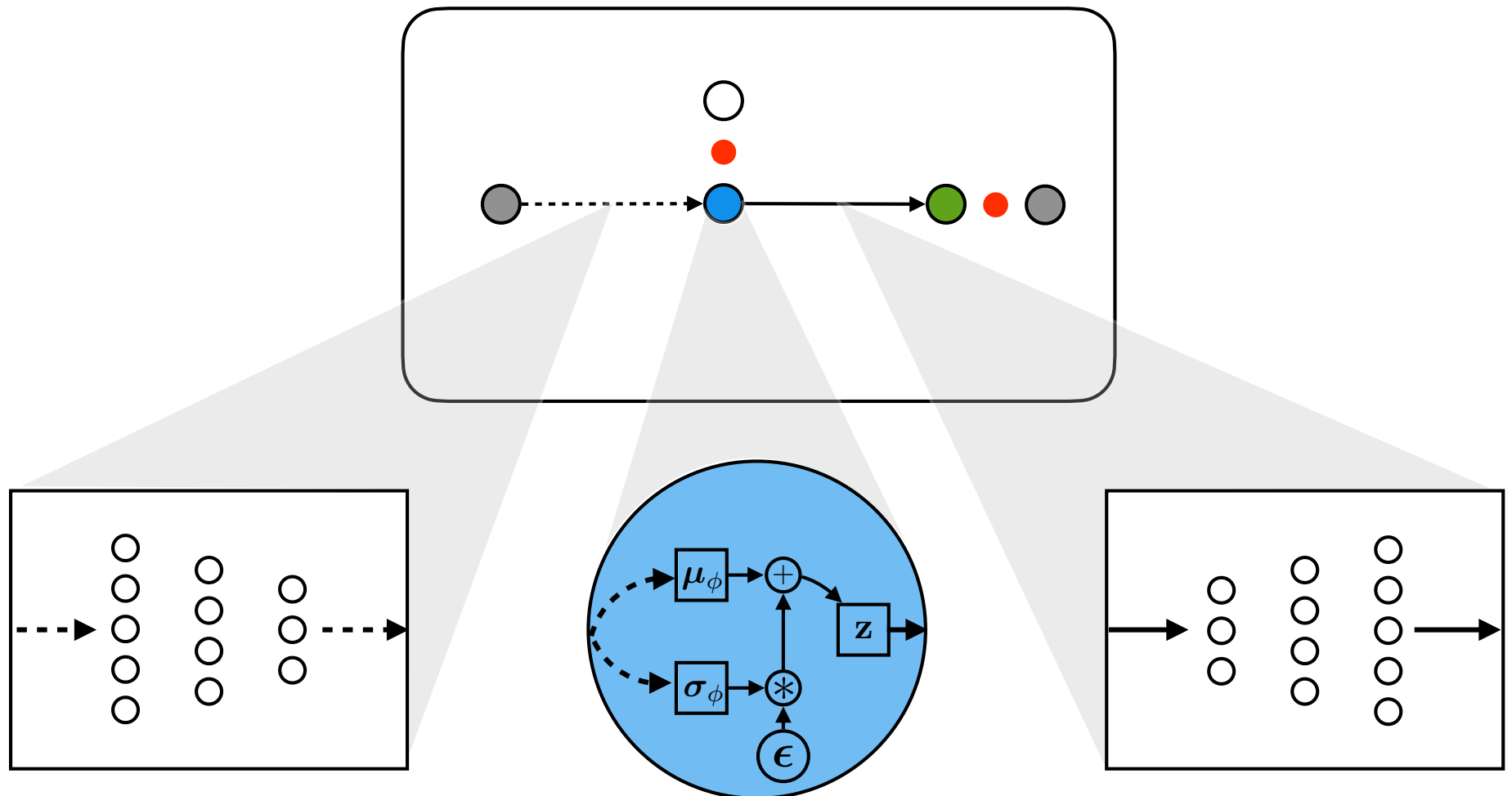
$p_\theta(\mathbf{z})$

$-\mathbb{E}_q \left[ \log \frac{q(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} \right]$

$q(\mathbf{z}|\mathbf{x})$

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

$\mathbb{E}_q \left[ \log p_\theta(\mathbf{x}|\mathbf{z}) \right]$

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

*observed space*

# AMORTIZED INFERENCE

**Variational Autoencoder (VAE):**

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



Kingma & Welling, 2014
Rezende et al., 2014

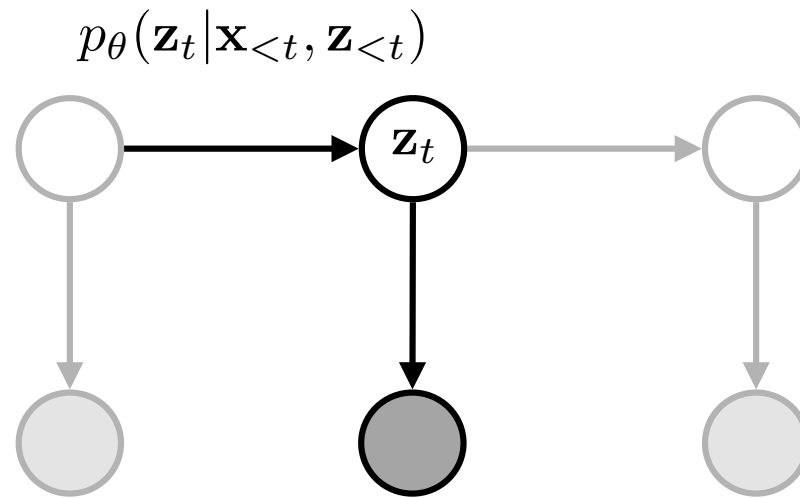introduce an approximate posterior $q(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})$

$$\text{ELBO:} \quad \mathcal{L}(\mathbf{x}_{\leq T}, q) = \mathbb{E}_q \left[ \log \frac{p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T})}{q(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})} \right]$$

choices about the form of $q(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})$ determine how we evaluate $\mathcal{L}$

→ often $q(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T})$ is *structured*
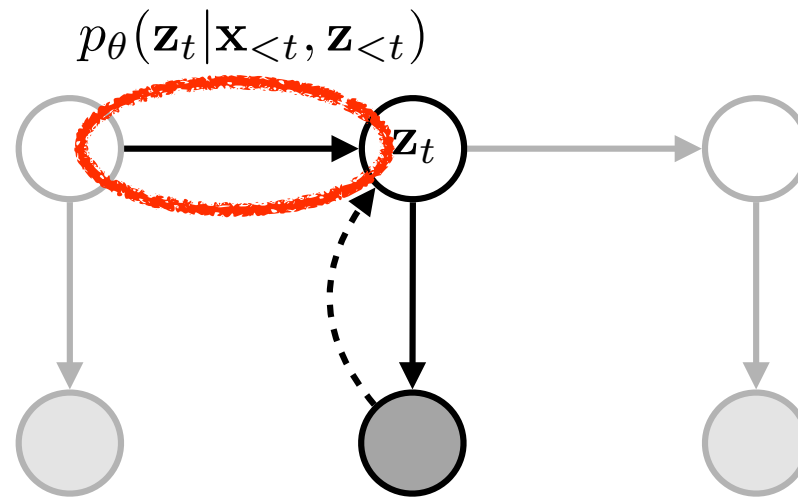
# STRUCTURED VARIATIONAL INFERENCE

the model contains temporal dependencies

$$p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$$



the approximate posterior should account for these dependencies

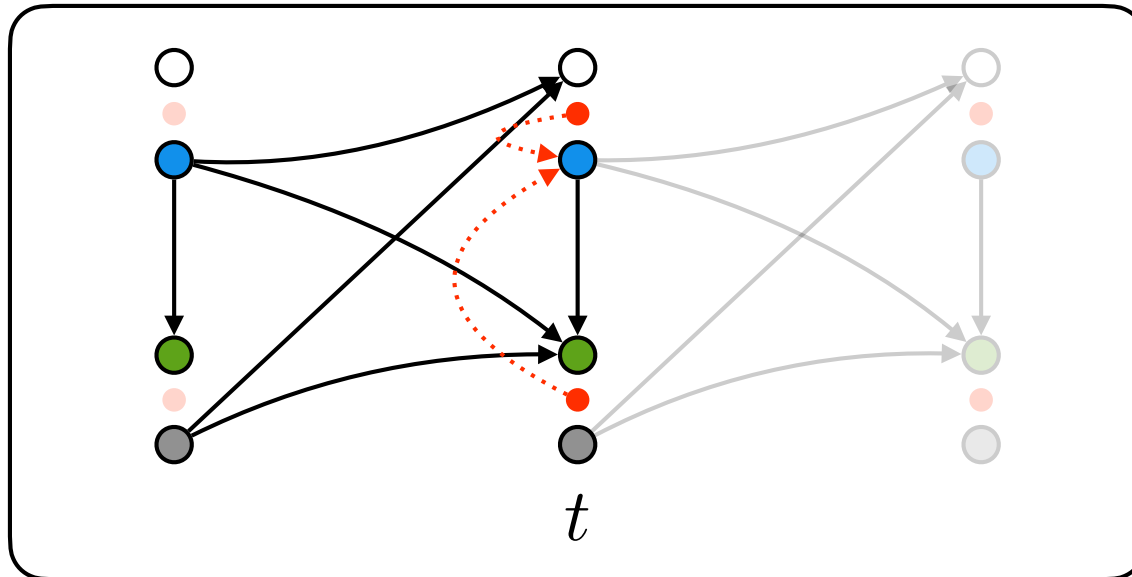# STRUCTURED VARIATIONAL INFERENCE

the model contains temporal dependencies

$$p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$$



the approximate posterior should account for these dependencies

$\longrightarrow$ if we use $q(\mathbf{z}_t | \mathbf{x}_t)$, we cannot account for $\mathbf{x}_{<t}$ and $\mathbf{z}_{<t}$

# FILTERING INFERENCE

filtering approximate posterior

$$q(\mathbf{z}_{\leq T}|\mathbf{x}_{\leq T}) = \prod_{t=1}^{T} q(\mathbf{z}_t|\mathbf{x}_{\leq t}, \mathbf{z}_{<t})$$



condition on observations at past and present time steps
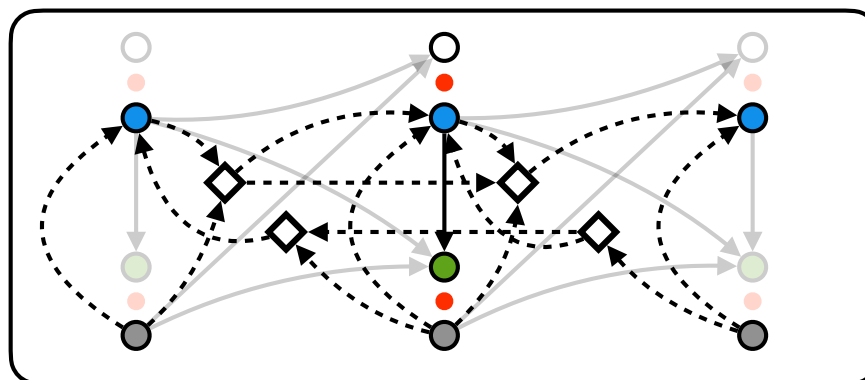
# AMORTIZED VARIATIONAL INFERENCE

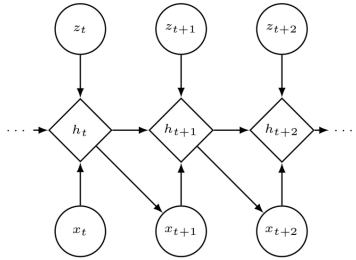*how do we amortize inference in sequential models?*

typical approach:

**filtering**: use a recurrent network
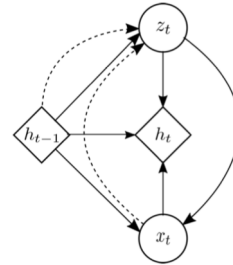


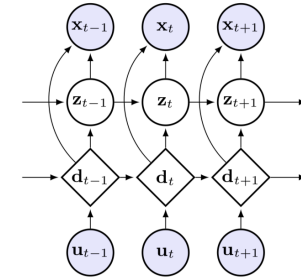**smoothing**: use a bi-directional recurrent network

# MODELS



**STORN**
Bayer & Osendorfer, 2014

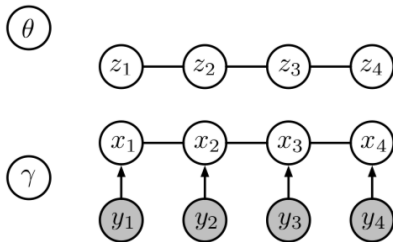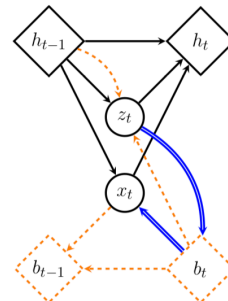**VRNN**
Chung *et al.*, 2015

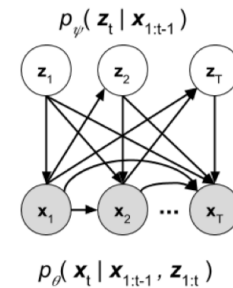**SRNN**
Fraccaro *et al.*, 2016

**Structured VAE**
Johnson *et al.*, 2016
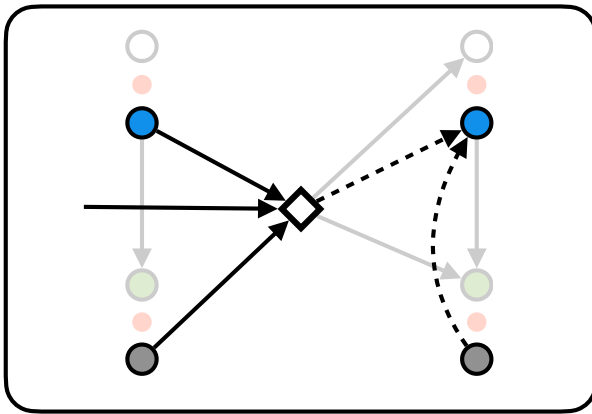
**Z-Forcing**
Goyal *et al.*, 2017

**SVG**
Denton & Fergus, 2018

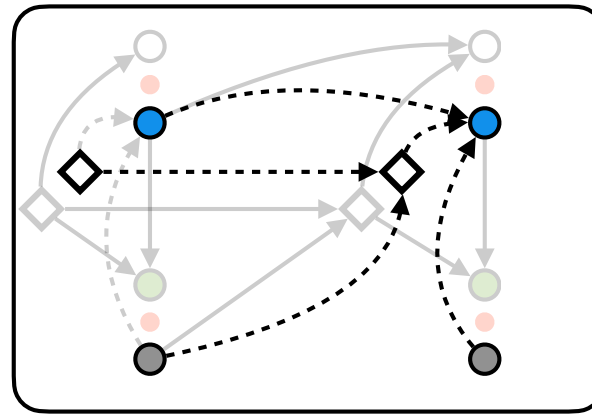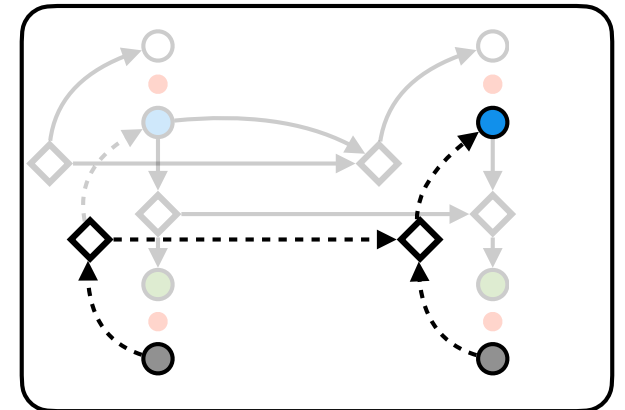# FILTERING INFERENCE MODELS



**VRNN**
Chung et al., 2015

**SRNN**
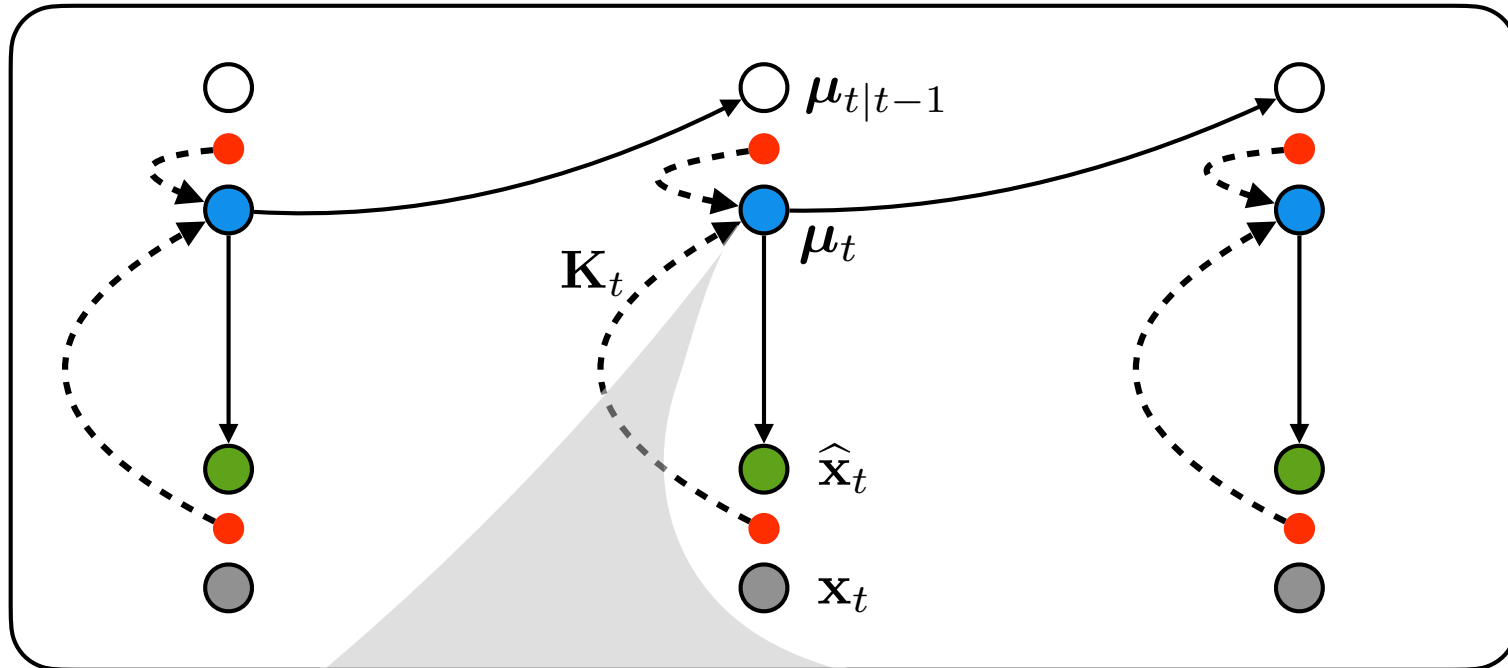Fraccaro et al., 2016

**SVG**
Denton & Fergus, 2018

*custom-designed inference models*

# AMORTIZED VARIATIONAL FILTERING

# KALMAN FILTERING

**Kalman filtering**: exact Bayesian inference in linear-Gaussian model



$$\boldsymbol{\mu}_t \leftarrow \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t(\mathbf{x}_t - \widehat{\mathbf{x}}_t)$$

$$\text{estimate} \leftarrow \text{prediction} + \text{gain} \cdot \text{prediction error}$$

# ITERATIVE AMORTIZED INFERENCE

let $\boldsymbol{\lambda}$ be the distribution parameters of $q(\mathbf{z}|\mathbf{x})$, for example, $\boldsymbol{\lambda} = \{\boldsymbol{\mu}, \boldsymbol{\sigma}^2\}$

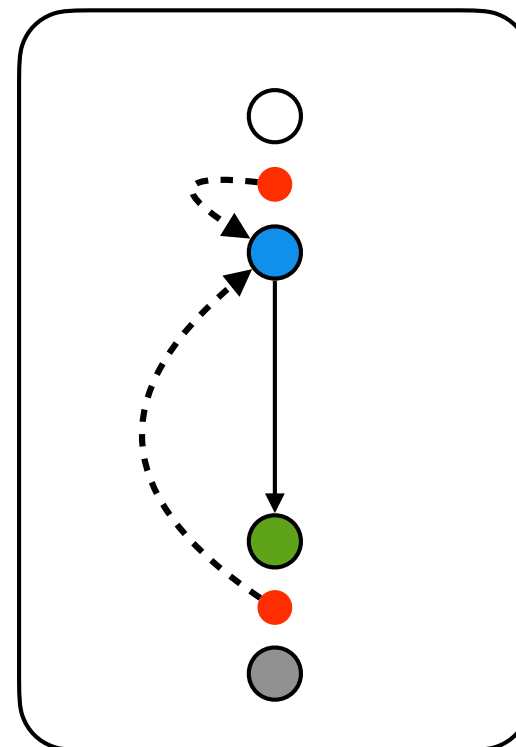$$\text{inference optimization:} \quad q(\mathbf{z}|\mathbf{x}) \leftarrow \arg\max_{q} \mathcal{L}(\mathbf{x}; q)$$

ITERATIVE AMORTIZED INFERENCE

learn an iterative mapping

$$\boldsymbol{\lambda} \leftarrow f_\phi(\boldsymbol{\lambda}, \nabla_{\boldsymbol{\lambda}} \mathcal{L})$$
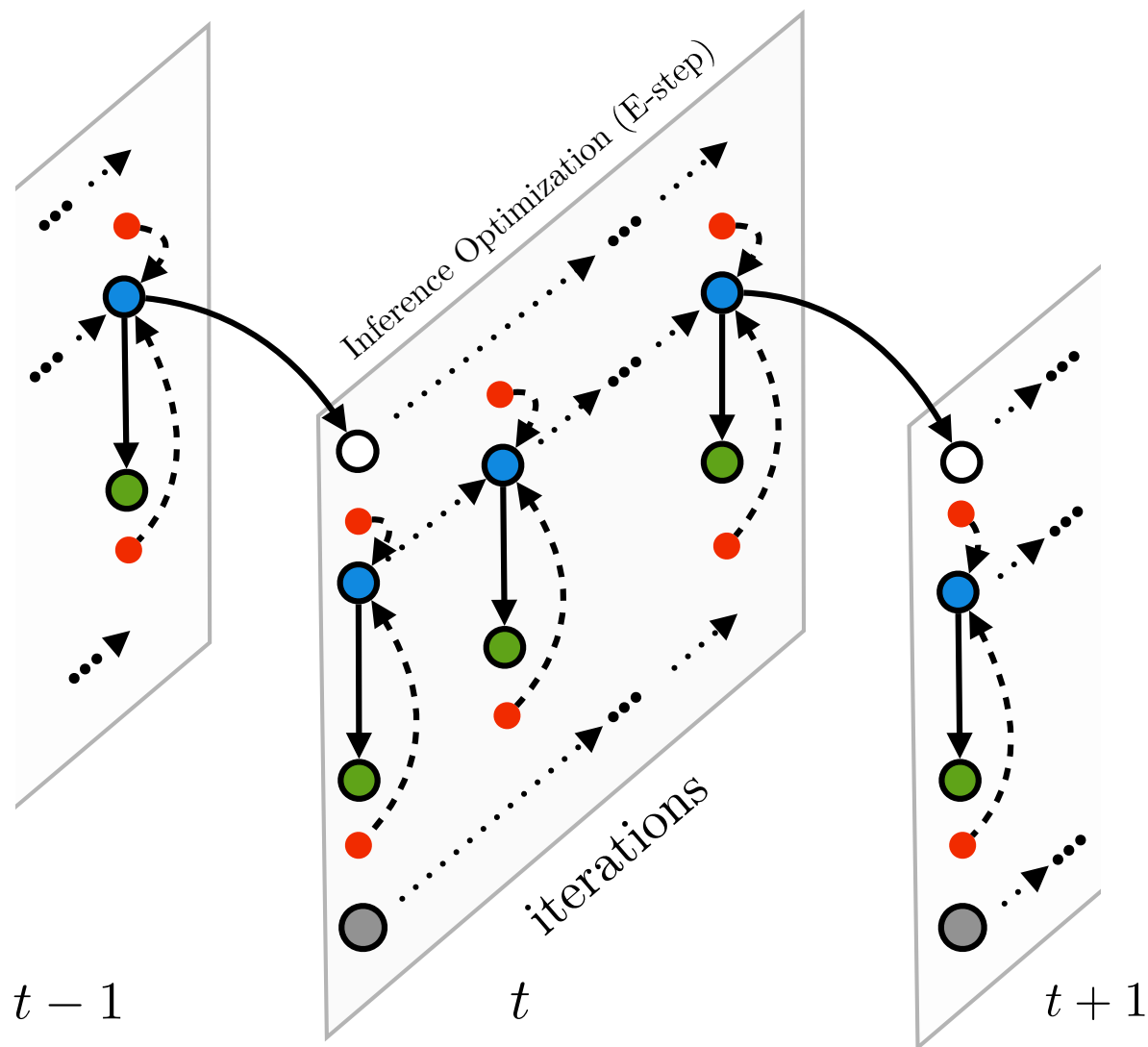
contains prediction errors

$$\mathbf{x} - \widehat{\mathbf{x}}$$

Marino *et al.*, 2018a

36

# AMORTIZED VARIATIONAL FILTERING



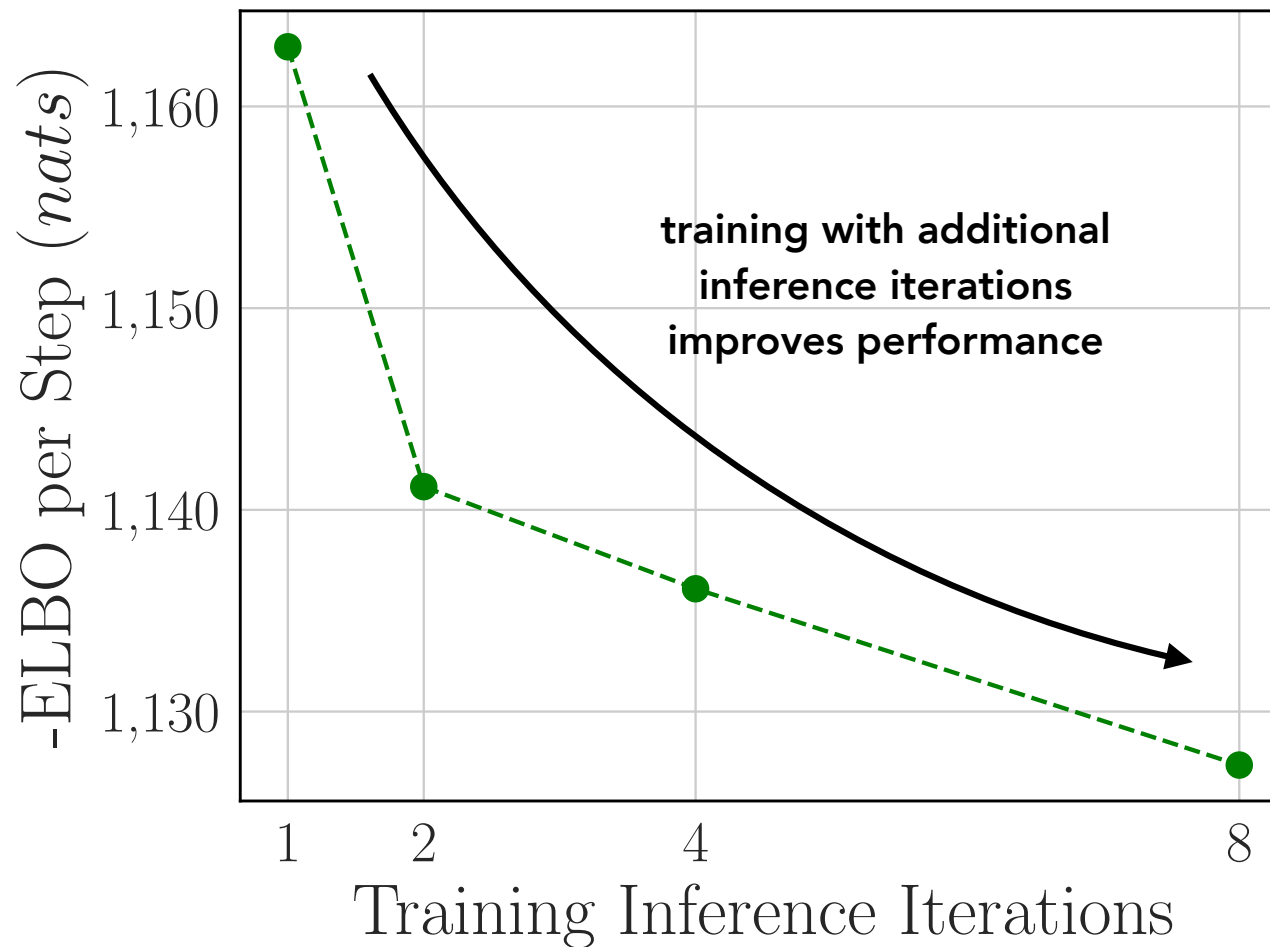perform iterative amortized inference at each time step

Marino *et al.*, 2018b

# INFERENCE IMPROVEMENT

TIMIT audio waveforms



observation

model output,
iteration 0

model output,
iteration 1

Marino *et al.*, 2018b

# INFERENCE ITERATIONS



ON TIMIT VAL SET

training with additional inference iterations improves performance

-ELBO per Step (*nats*)

Training Inference Iterations

Marino *et al.*, 2018b

# INFERENCE ITERATIONS



ON TIMIT VAL SET

each inference iteration yields
diminishing relative improvement

Marino *et al.*, 2018b

# PERFORMANCE

*one inference method*, consistent improvement across models & domains

## AUDIO

|          | TIMIT     |
|----------|-----------|
| VRNN     |           |
| baseline | 1,082     |
| AVF (1 Iter.) | 1,105 |
| AVF (2 Iter.) | **1,071** |
| SRNN     |           |
| baseline | 1,026     |
| AVF (1 Iter.) | **1,024** |

## VIDEO

|          | KTH Actions |
|----------|-------------|
| SVG      |             |
| baseline | 3.69        |
| AVF (1 Iter.) | **2.86** |

## MIDI MUSIC

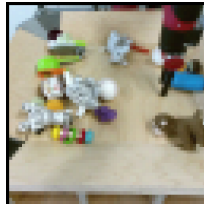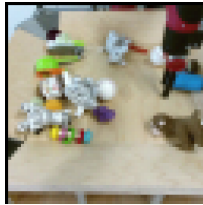|                              | Piano-midi.de | MuseData | JSB Chorales | Nottingham |
|------------------------------|---------------|----------|--------------|------------|
| SRNN                         |               |          |              |            |
| baseline (Fraccaro et al., 2016) | 8.20      | 6.28     | 4.74         | 2.94       |
| baseline                     | 8.19          | 6.27     | 6.92         | 3.19       |
| AVF (1 Iter.)                | **8.12**      | **5.99** | 6.97         | **3.13**   |
| AVF (5 Iter.)                | –             | –        | **6.77**     | –          |

Marino *et al.*, 2018b

# SEQUENTIAL AUTOREGRESSIVE FLOWS

$\mathbf{x}_{t-3}$ $\qquad$ $\mathbf{x}_{t-2}$ $\qquad$ $\mathbf{x}_{t-1}$ $\qquad$ $\mathbf{x}_{t}$

$$\mathbf{y}_t$$

$$\boldsymbol{\sigma}_\theta(\mathbf{x}_{<t})$$

$$\theta$$

$$\div$$

$$-$$

$$\boldsymbol{\mu}_\theta(\mathbf{x}_{<t})$$

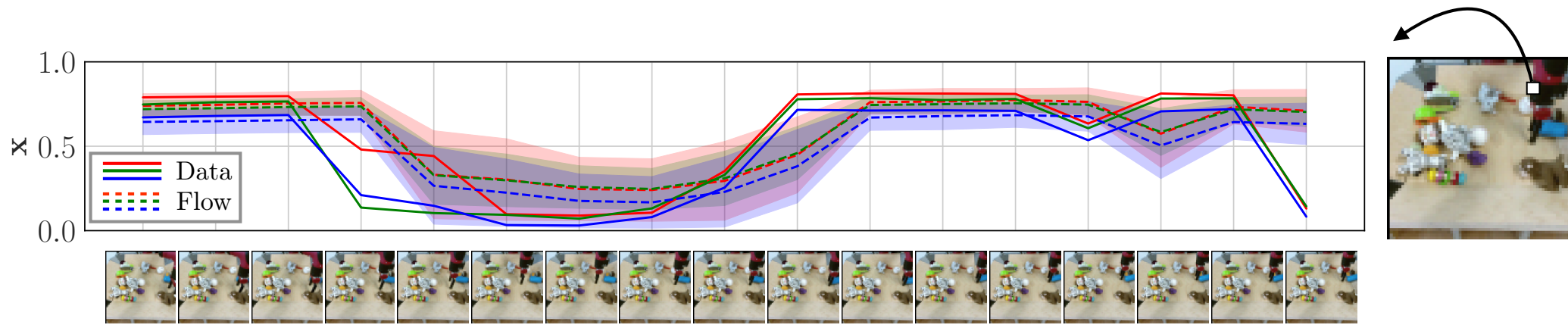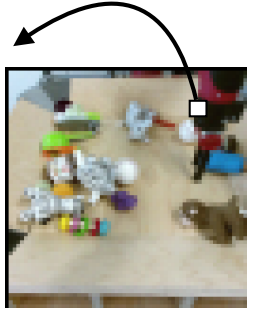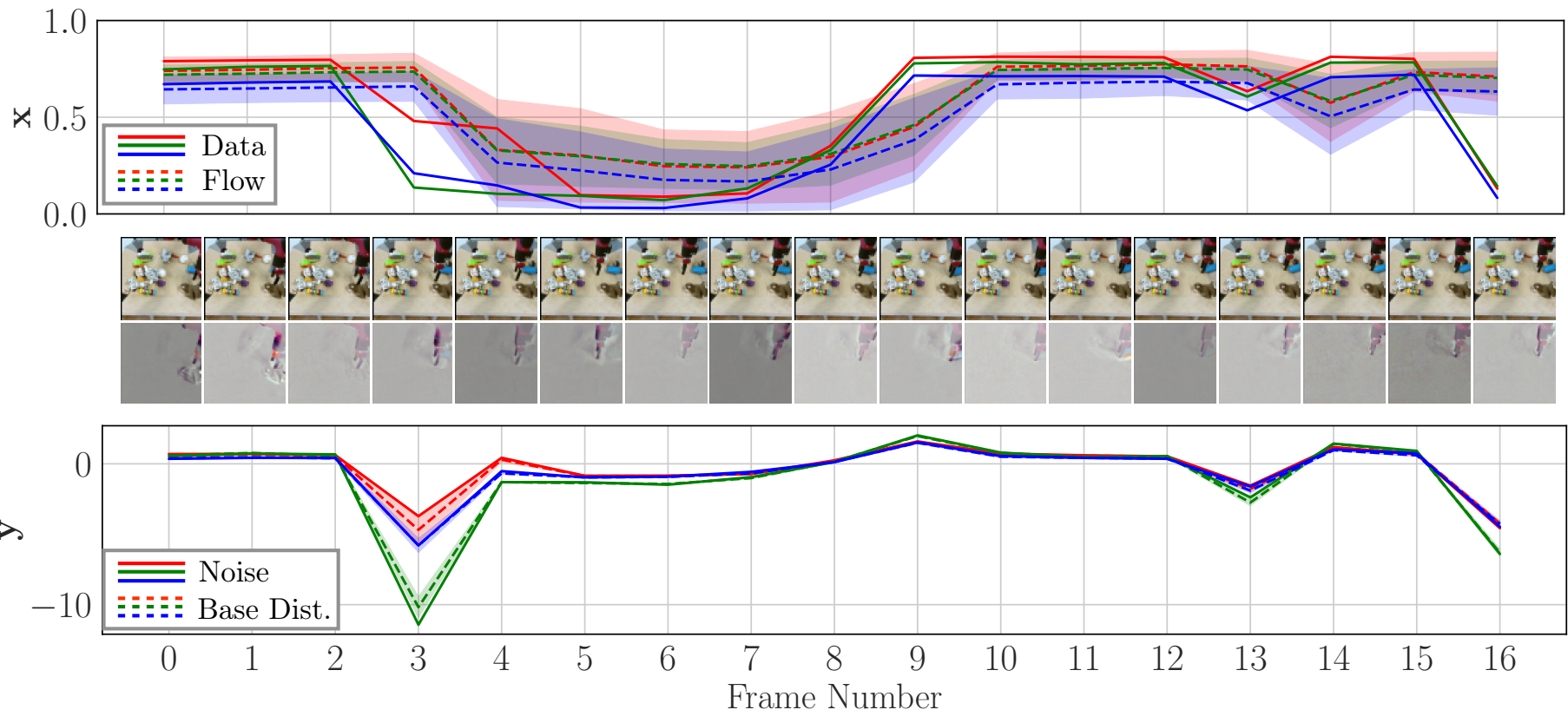$$\mathbf{x}_{t-3} \qquad \mathbf{x}_{t-2} \qquad \mathbf{x}_{t-1} \qquad\qquad \mathbf{x}_t$$

44

# SEQUENTIAL AUTOREGRESSIVE FLOWS



remove low-level temporal dependencies using an autoregressive flow

Marino *et al.*, 2020
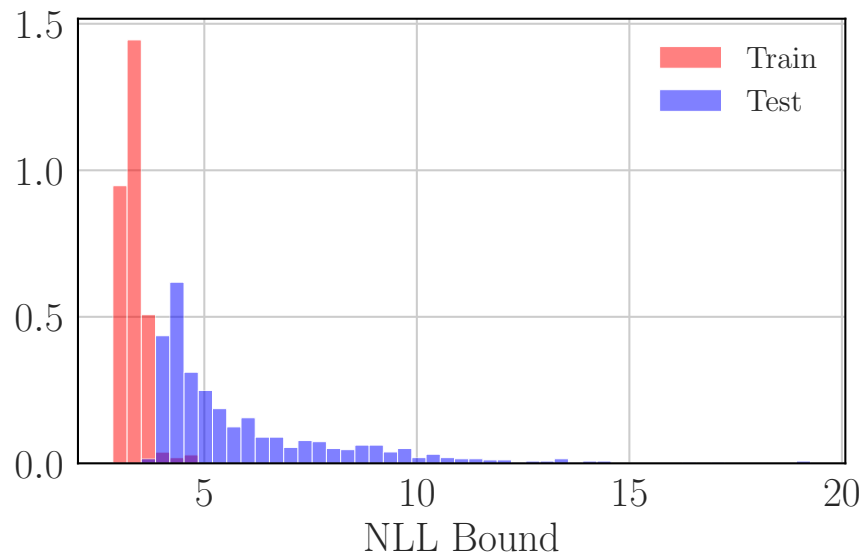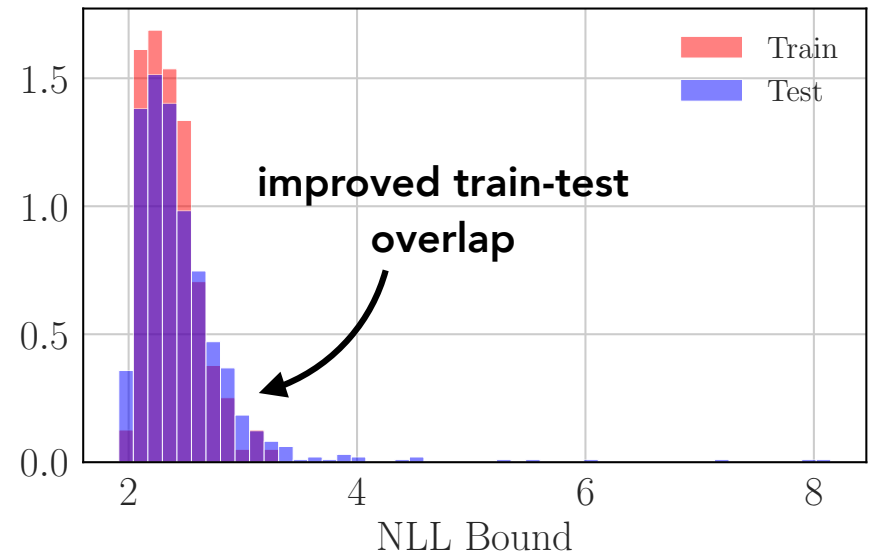
# REDUCED TEMPORAL CORRELATION



during training

Marino *et al.*, 2020

# IMPROVED GENERALIZATION

**KTH Actions**

SLVM

SLVM + AF



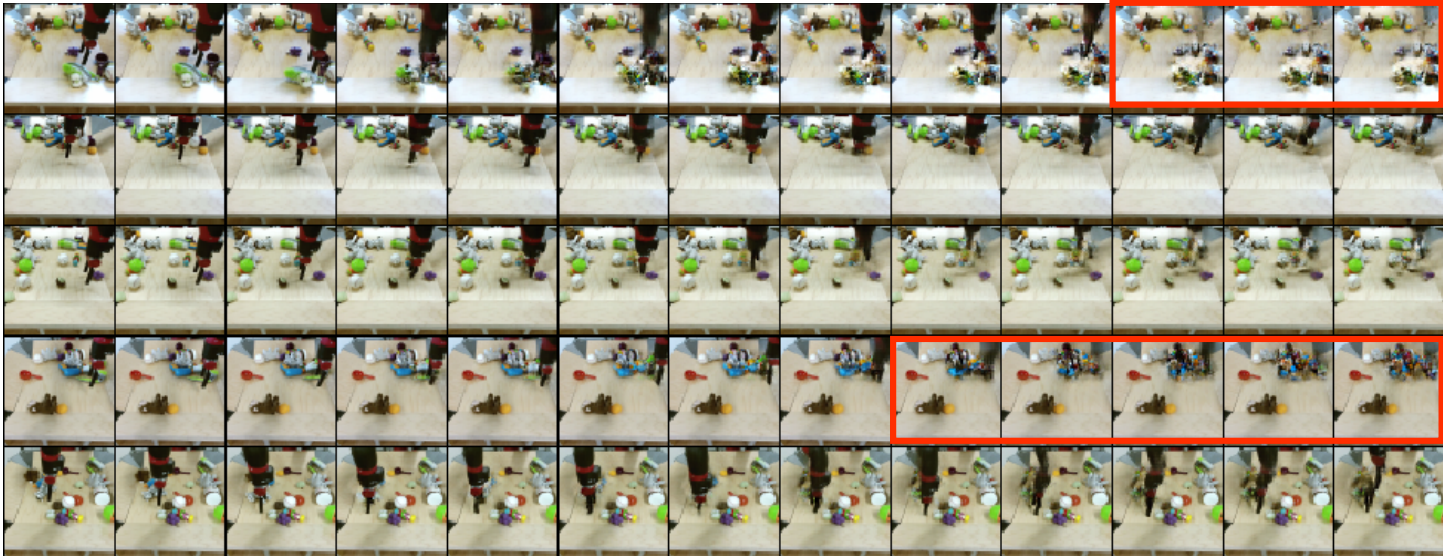improved train-test overlap

Marino *et al.*, 2020

# IMPROVED SAMPLES

VideoFlow



VideoFlow + AF



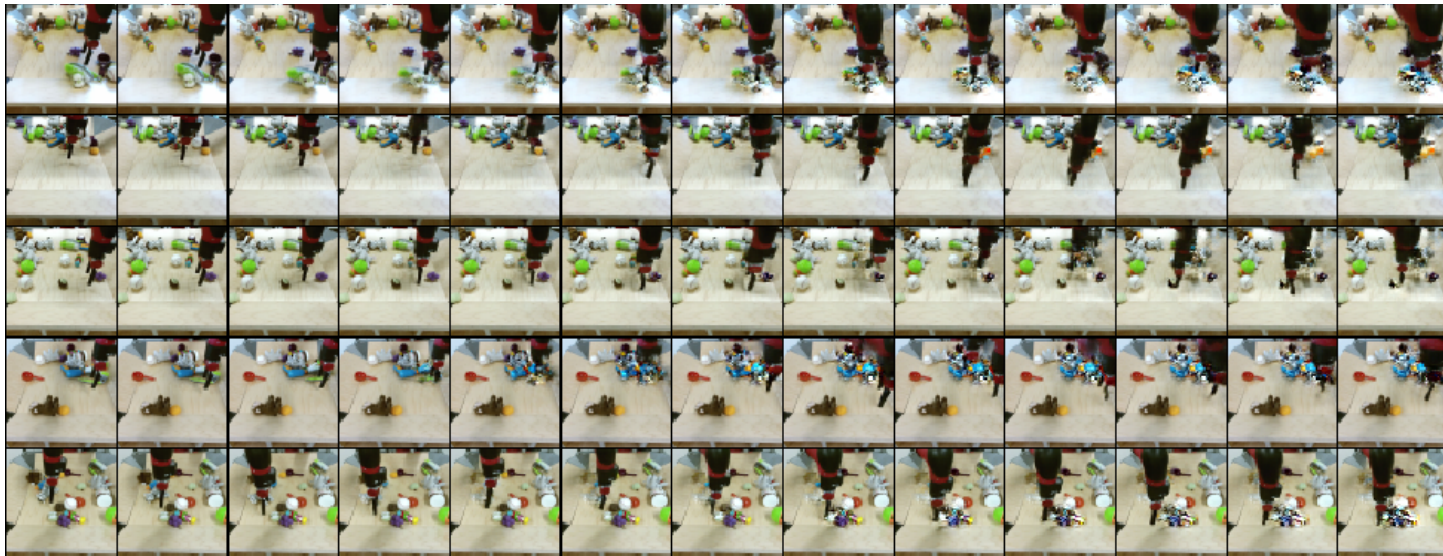Marino *et al.*, 2020

# RECAP

- sequence models

- amortized variational filtering

- sequential autoregressive flows