



Compositional Methods for Learning and Inference in Deep Probabilistic Programs

Jan-Willem van de Meent



Hao Wu



Heiko Zimmermann



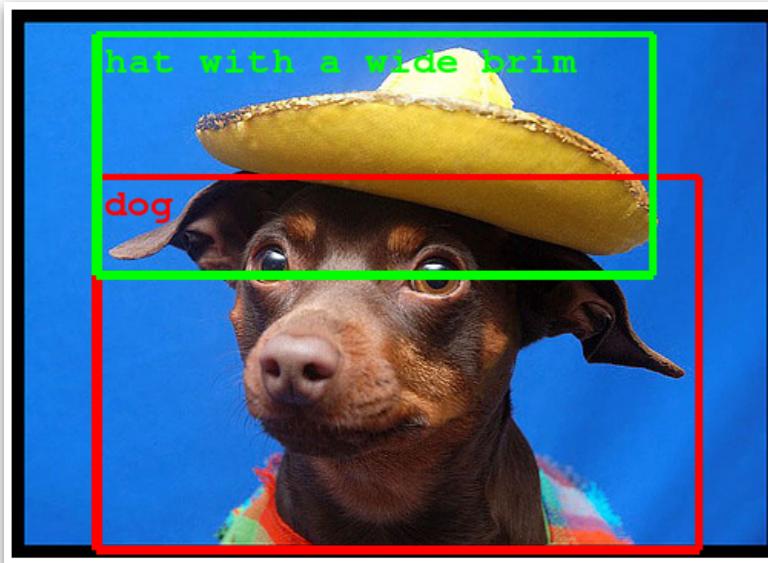
Eli Sennesh



Sam Stites

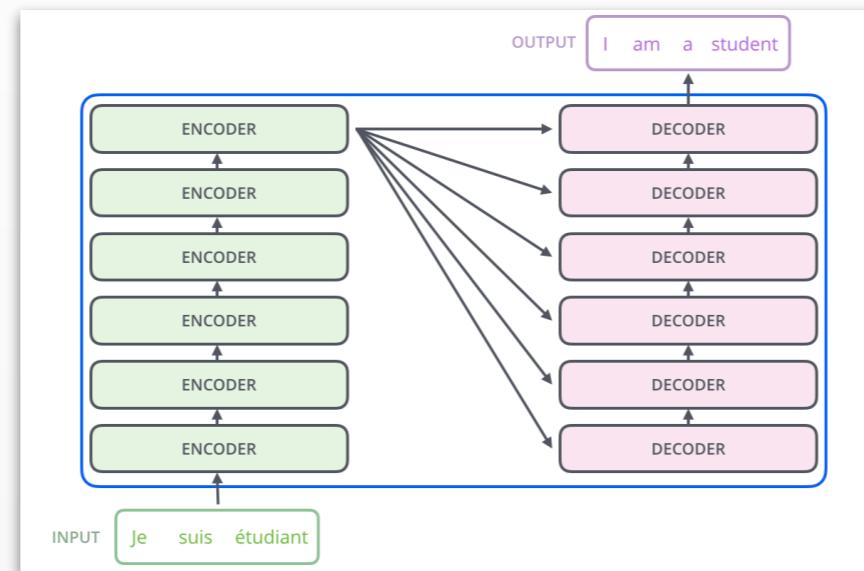
Deep Learning Success Stories

Computer Vision



14M images (ImageNet)
Annotations available

Natural Language



Very large corpora of text
(can self-supervise)

Reinforcement Learning



4.9M games (Self-play)
Clear definition of success

Ingredients for success

1. Abundance of (labeled) data and compute
2. A well-defined general notion of utility

Do we still need models?

The Bitter Lesson

Rich Sutton, March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin.

[http://www.incompleteideas.net/
IncompleteIdeas/BitterLesson.html](http://www.incompleteideas.net/IncompleteIdeas/BitterLesson.html)

Do we still need models or just more data and compute?

Max Welling, April 20, 2019

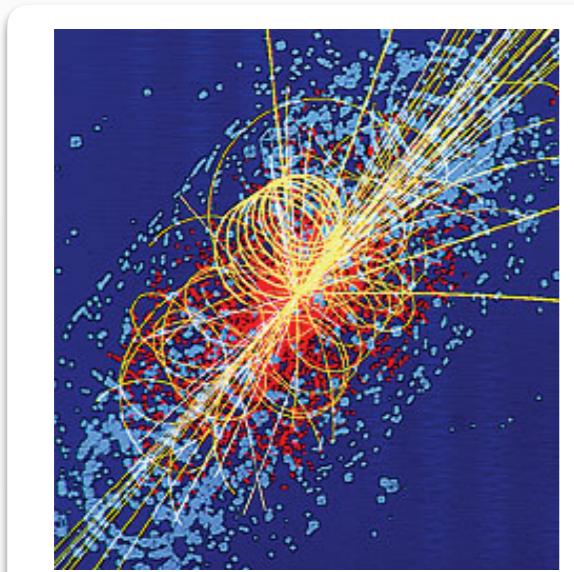
When you need to generalize to new domains, i.e. extrapolate away from the data, you will need a generative model

[https://staff.fnwi.uva.nl/m.welling/
wp-content/uploads/Model-versus-
Data-AI-1.pdf](https://staff.fnwi.uva.nl/m.welling/wp-content/uploads/Model-versus-Data-AI-1.pdf)

Do we ~~still need~~ models?

When are models useful?

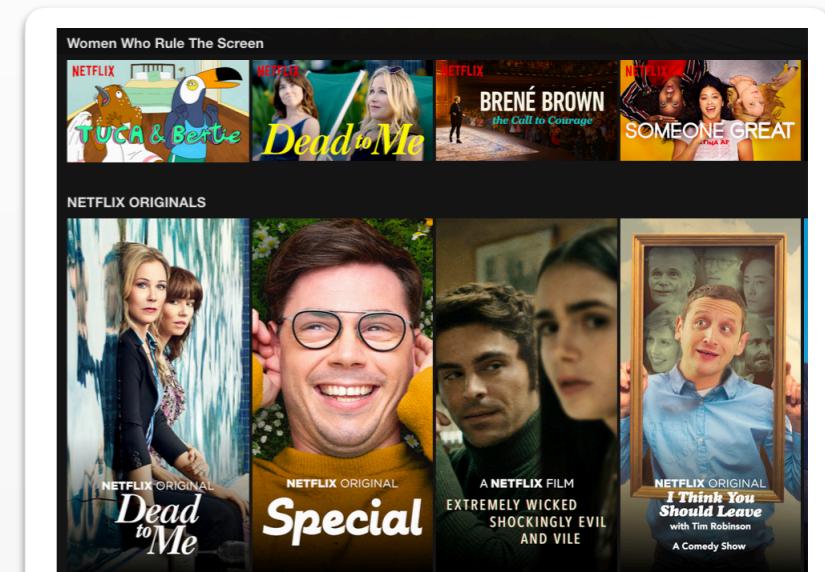
Science & Engineering



Autonomous Vehicles



Recommendation



High quality models
and/or limited data

Generalization to
long tail events

Large collection of
small-data problems

We need *inductive biases* that

1. Improve generalization
2. Safe-guard against overconfident predictions

Deep Probabilistic Models

Deep Learning

- High-capacity models
- Scalable to large datasets
- Easy to try new models



Probabilistic Programming

- Programs as inductive biases
- Structured, interpretable
- Also easy to try new models

SGD + AutoDiff
(very general)

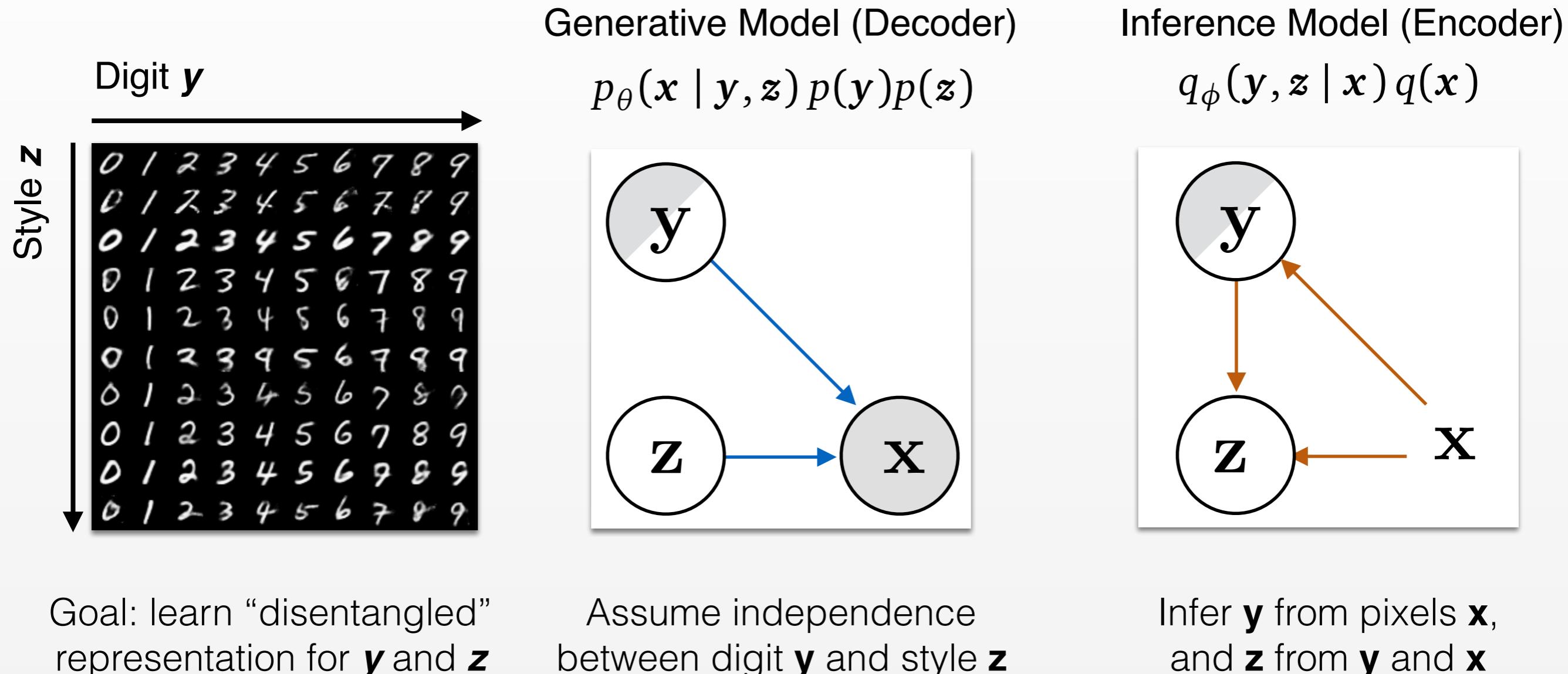


Monte Carlo Methods
(more model specific)



Stochastic Variational Inference
(learn proposals using neural networks)

Structured Variational Autoencoders



Deep Probabilistic Programs

Generative Model (Decoder)

```
class Decoder(torch.nn.Module):
    def __init__(self, x_sz, h_sz, y_sz, z_sz):
        # initializes layers: h, x_mean, ...
        ...

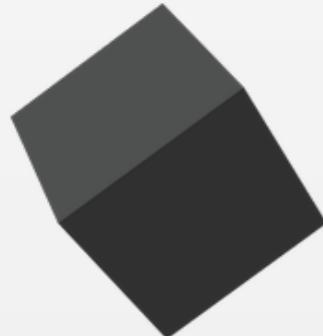
    def forward(self, x, q):
        p = probtorch.Trace()
        y = p.concrete(self.y_log_weights, 0.66,
                        value=q['y'], name='y')
        z = p.normal(0.0, 1.0,
                     value=q['z'], name='z')
        h = self.h(torch.cat([y, z], -1))
        x = p.loss(self.bce,
                   self.x_mean(h), x,
                   name='x')
        return p
```

Inference Model (Encoder)

```
class Encoder(torch.nn.Module):
    def __init__(self, x_sz, h_sz, y_sz, z_sz):
        # initializes layers: h, y_log_weights, ...
        ...

    def forward(self, x, y_values=None):
        q = probtorch.Trace()
        h = self.h(x)
        y = q.concrete(
            self.y_log_weights(h), 0.66,
            value=y_values, name='y')
        hy = torch.cat([h, y], -1)
        z = q.normal(self.z_mean(hy),
                     self.z_std(hy),
                     name='z')
        return q
```

Edward



Probabilistic Torch

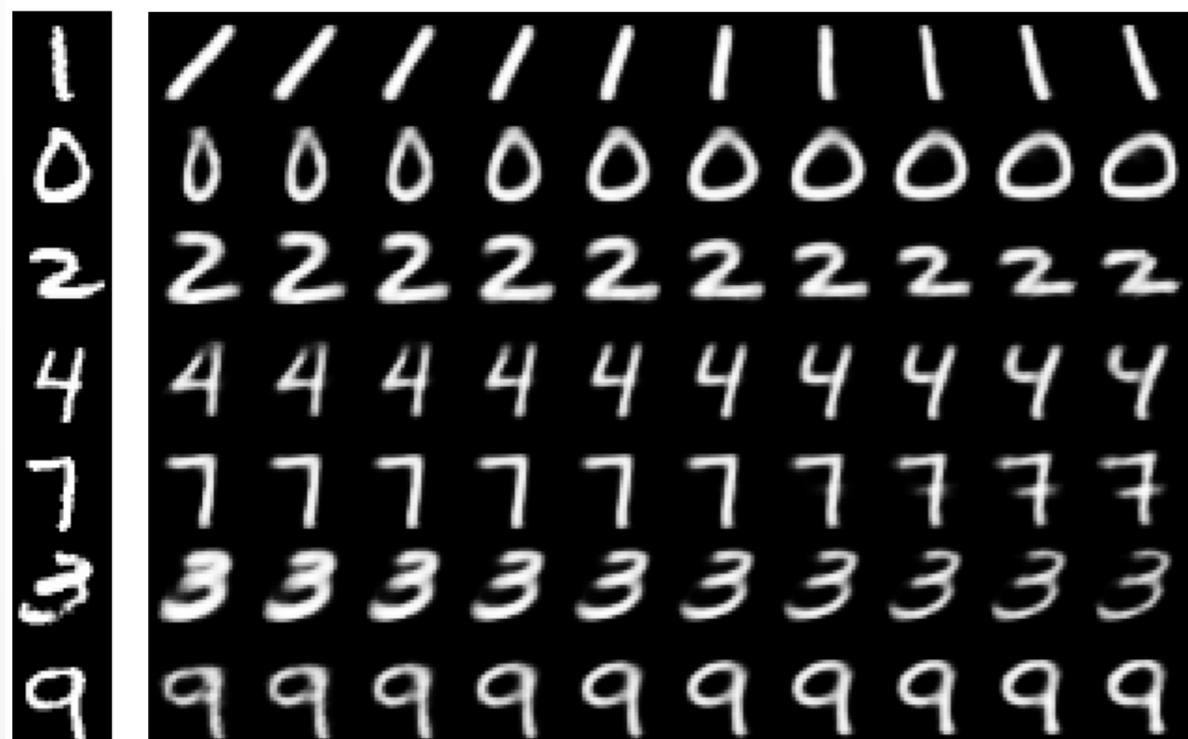


Pyro



Learned Representations (Unsupervised)

Style Variables



Slant

Width

Height

Style 1

Style 2

Thickness

Style 3

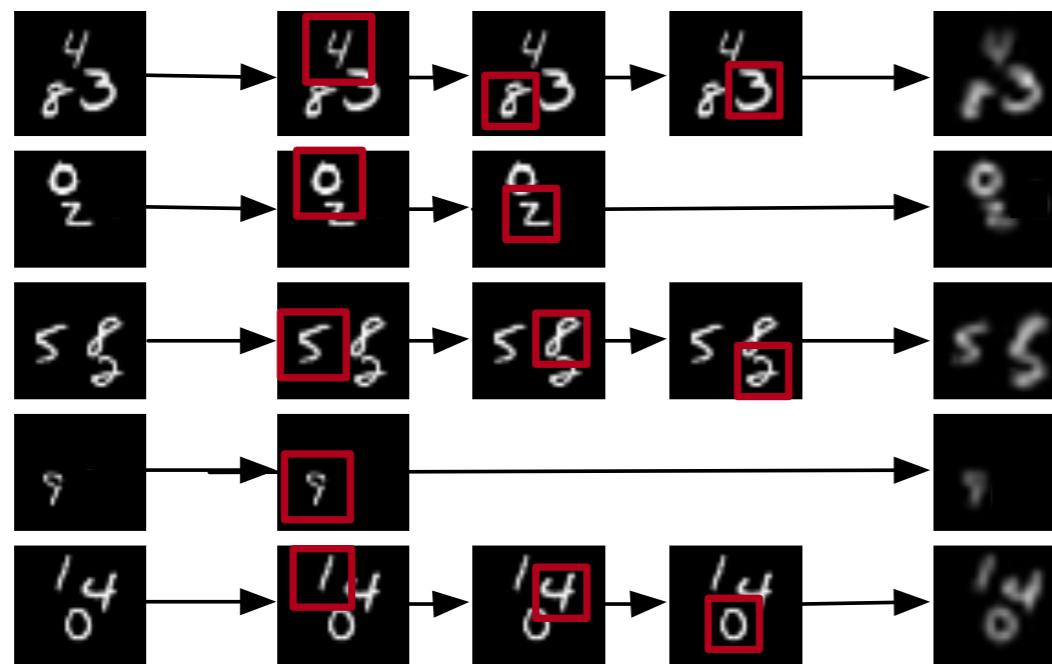
Generalization



Inductive Bias: Style features are uncorrelated with digit label, as well as with other features.

Model Composition

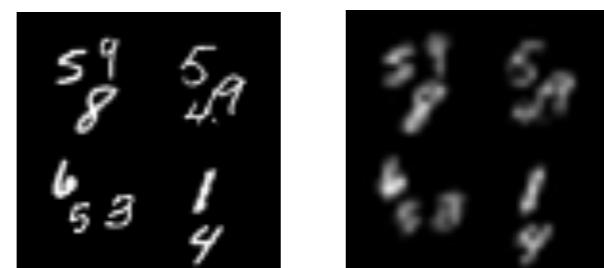
Recurrent Recognition Loop



Decomposition

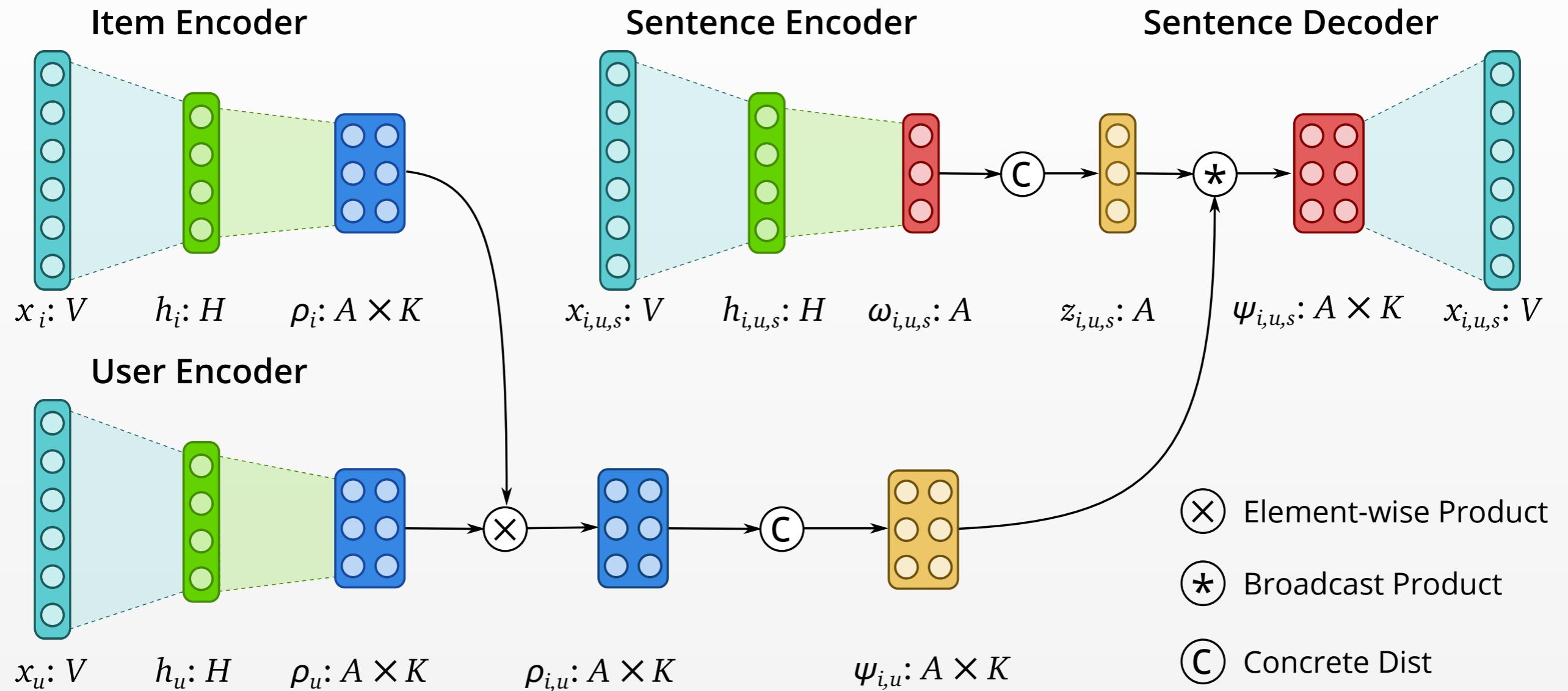


Reconstructions



Idea: Embed model for individual MNIST digits
in a recurrent model for multiple object detection

Example: Modeling Aspects in Reviews



Learn aspect-based representations of users, items, and reviews (*fully unsupervised*)

Example: Modeling Aspects in Reviews

Data: Beer reviews

Amber brown in color with very little head but a nice ring. Nicely carbonated. Smells like a camp fire, malts have a good sweet character with an abundance of smoke. Taste is quite good with smokiness being pungent but not overwhelming. A sweet tasting bock with smokiness coming through around mid drink with a smooth mellow finish. A good warming smoky beer.

Aspects: Look, Mouthfeel, Aroma, Taste, Overall

Example: Modeling Aspects in Reviews

Appearance		Aroma-Taste		Palate		Semantic	
golden	black	roasted	citrus	mouthfeel	mouthfeel	lagers	try
yellow	tans	coffee	grapefruit	bodied	watery	heineken	hype
white	glass	vanilla	pine	smooth	rjt	macro	recommend
orange	pour	chocolate	hops	carbonation	bodied	import	overall
hazy	head	bourbon	lemon	medium	refreshing	euro	founders
color	pitch	oak	floral	drinkability	carbonation	lager	favorite
gold	lacing	malts	clove	drinkable	crisp	bmc	stouts
copper	color	sweet	malt	alcohol	dry	worse	stout
straw	brown	aroma	aroma	finish	finish	bad	ipa
amber	ginger	malt	grass	mouth	thirst	skunk	cheers

Model learns topics that group into aspects
(fully unsupervised)

Where are we today?

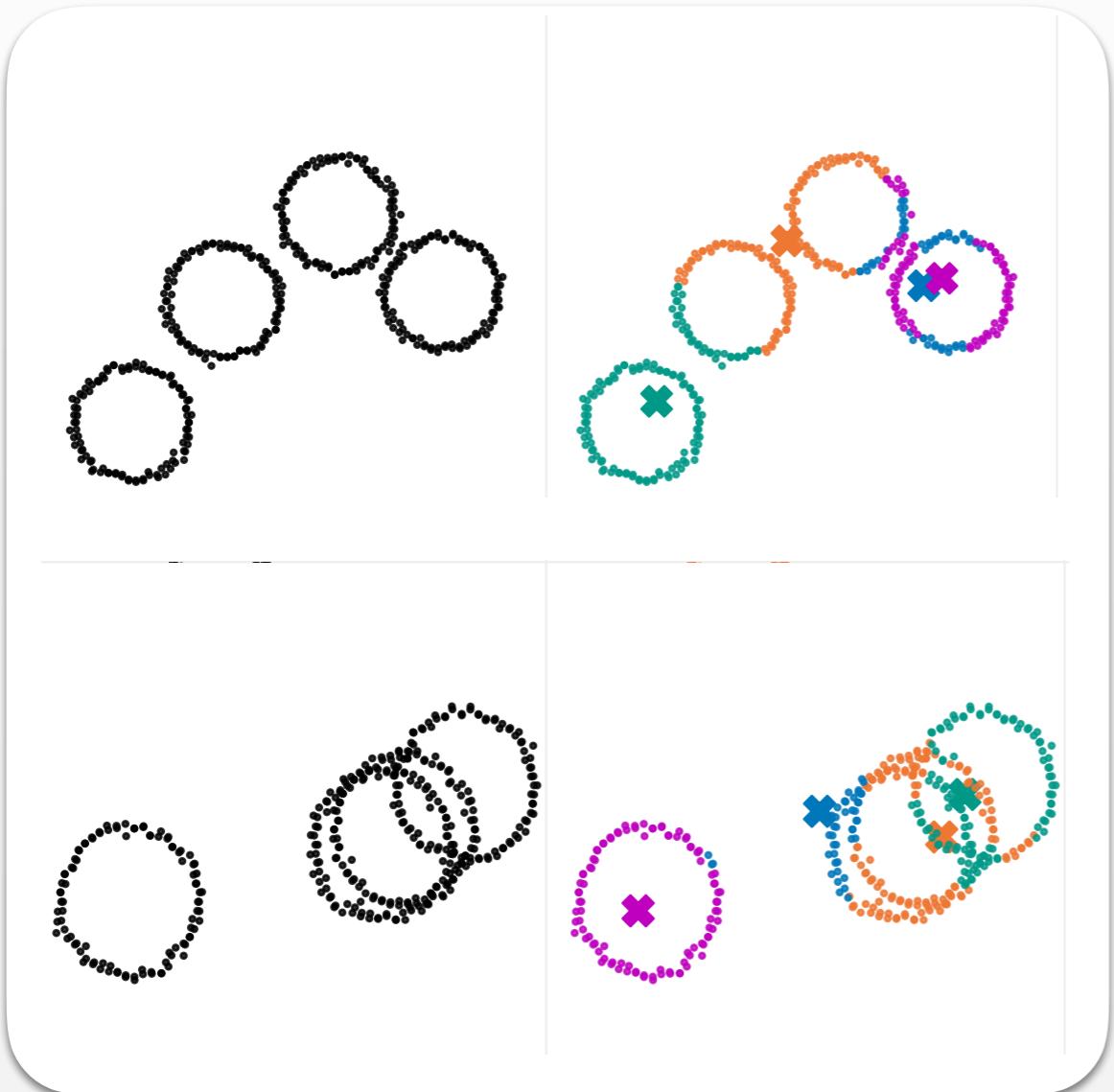
Structured VAEs work well for $O(10)$ variables

- Object recognition (~3 objects)
[Eslami et al, NeurIPS 2016]
- Object tracking (~3 objects, ~10 frames)
[Kosiorek et al, NeurIPS 2018]
- Product Reviews (~10 sentences, ~10 aspects)
[Esmaeli et al, AISTATS 2019]

Still Hard: Clustering

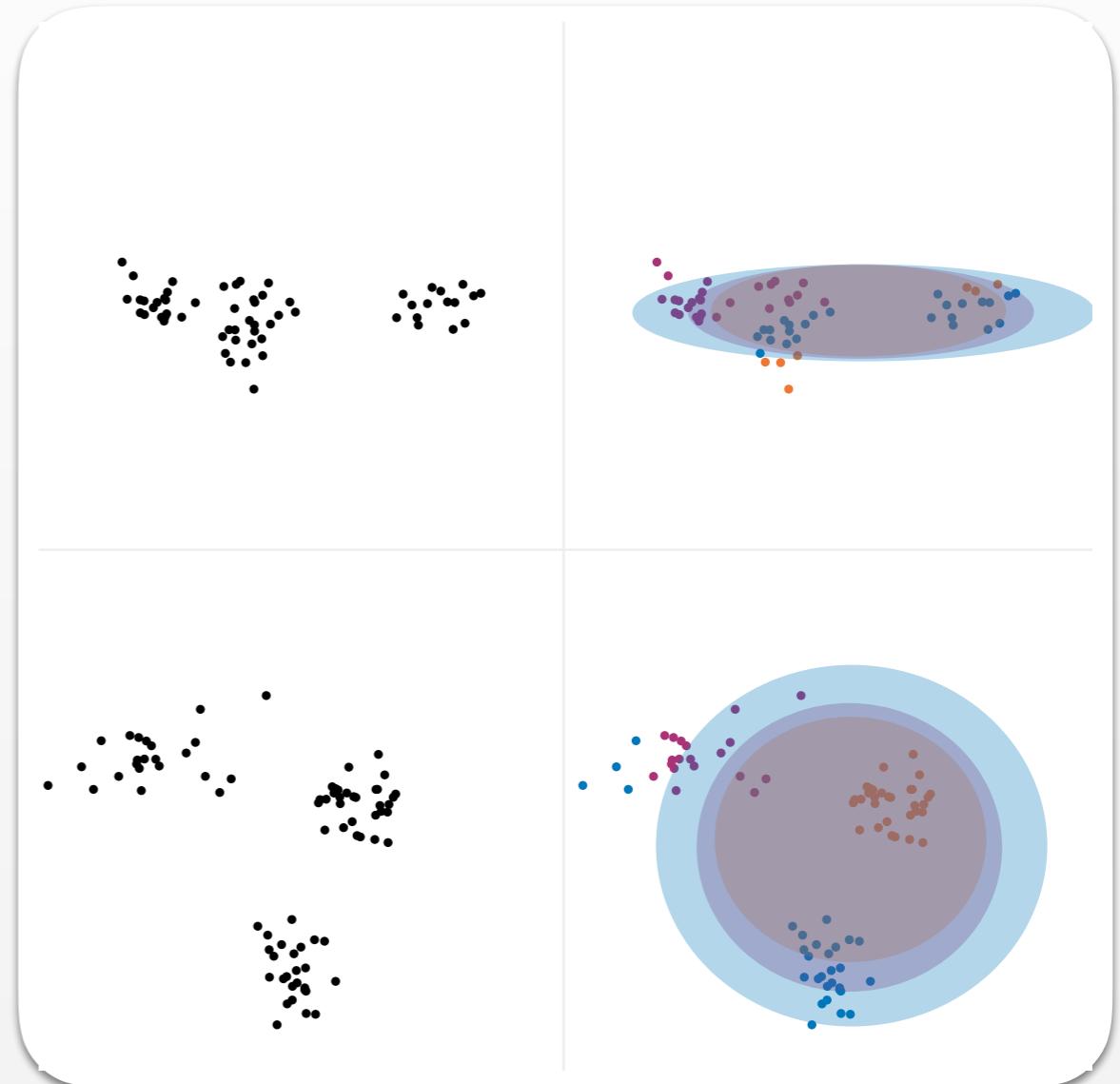
Data

Inference



Data

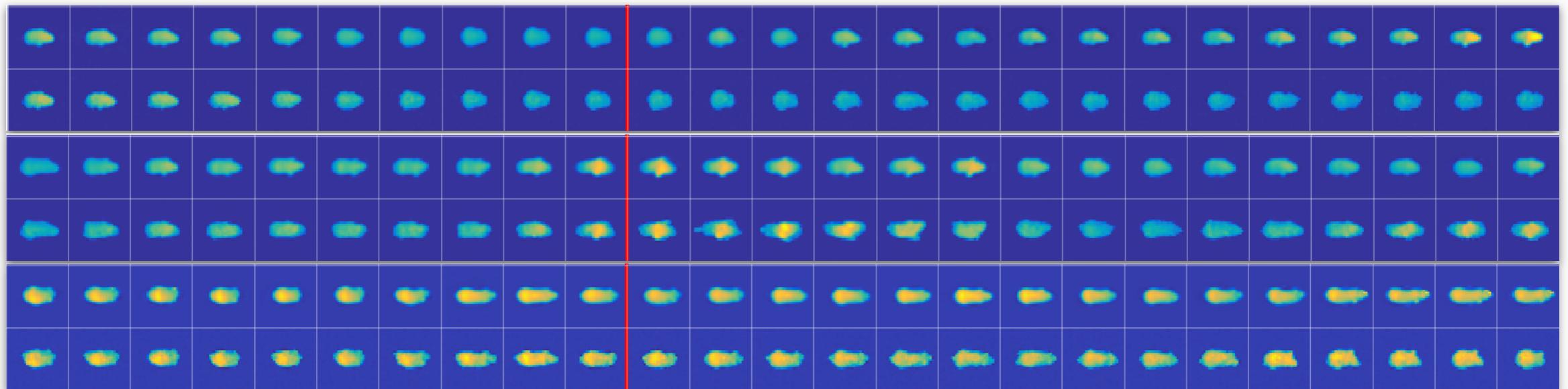
Inference



We need methods that scale to $> O(100)$ variables

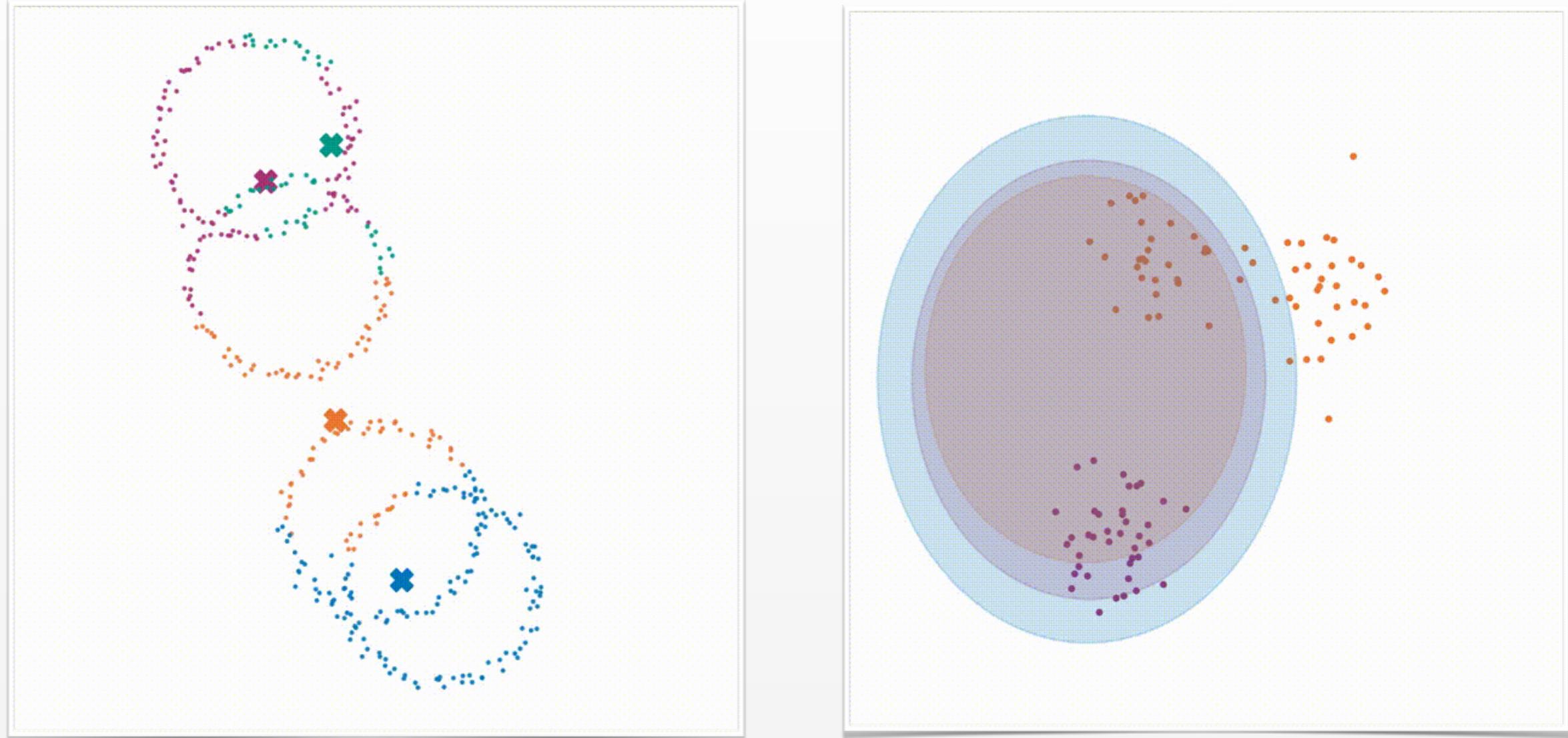
Paths towards Scaling up Inference

Example: Deep State Space Models for Activity Recognition in Mice



1. Generative Model: Exponential family + Neural likelihood
2. Inference: Variational EM + Gradient descent

Amortized Gibbs Samplers



Intuition: Same idea as Gibbs sampling or Variational EM, but learn proposals for conditional updates

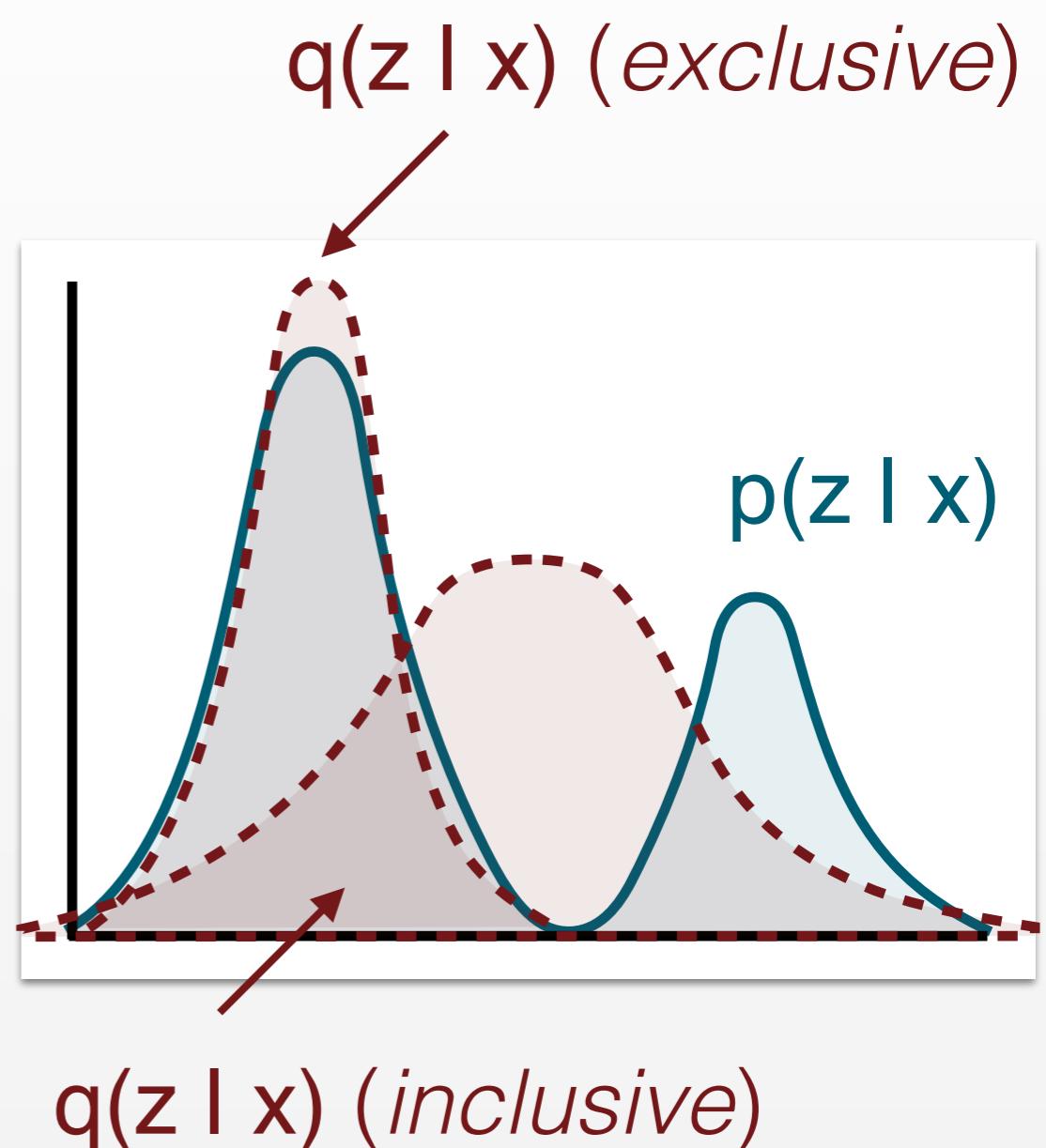
Minimizing the Inclusive KL

Variational Autoencoders

$$\operatorname{argmin}_{\phi} \text{KL}(q_{\phi}(z | x) \| p_{\theta}(z | x))$$

Wake-Sleep / EP Methods

$$\operatorname{argmin}_{\phi} \text{KL}(p_{\theta}(z | x) \| q_{\phi}(z | x))$$



Minimizing the Inclusive KL

Variational Autoencoders

$$\begin{aligned} -\nabla_{\phi} \text{KL}(q_{\phi}(z | x) || p_{\theta}(z | x)) \\ = -\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{q_{\phi}(z | x)}{p_{\theta}(z | x)} \right] \end{aligned}$$

Difficult: Need to approximate expectation that depends on ϕ

Solution: reparameterization or REINFORCE-style estimators

Wake-Sleep / EP Methods

$$\begin{aligned} -\nabla_{\phi} \text{KL}(p_{\theta}(z | x) || q_{\phi}(z | x)) \\ = -\nabla_{\phi} \mathbb{E}_{p_{\theta}(z|x)} \left[\log \frac{p_{\theta}(z | x)}{q_{\phi}(z | x)} \right] \\ = \mathbb{E}_{p_{\theta}(z|x)} [\nabla_{\phi} \log q_{\phi}(z | x)] \end{aligned}$$

Easier: Expectation depends on θ (generative parameters)

Solution: sample from $p_{\theta}(z | x)$

Reweighted Wake-sleep Style Methods

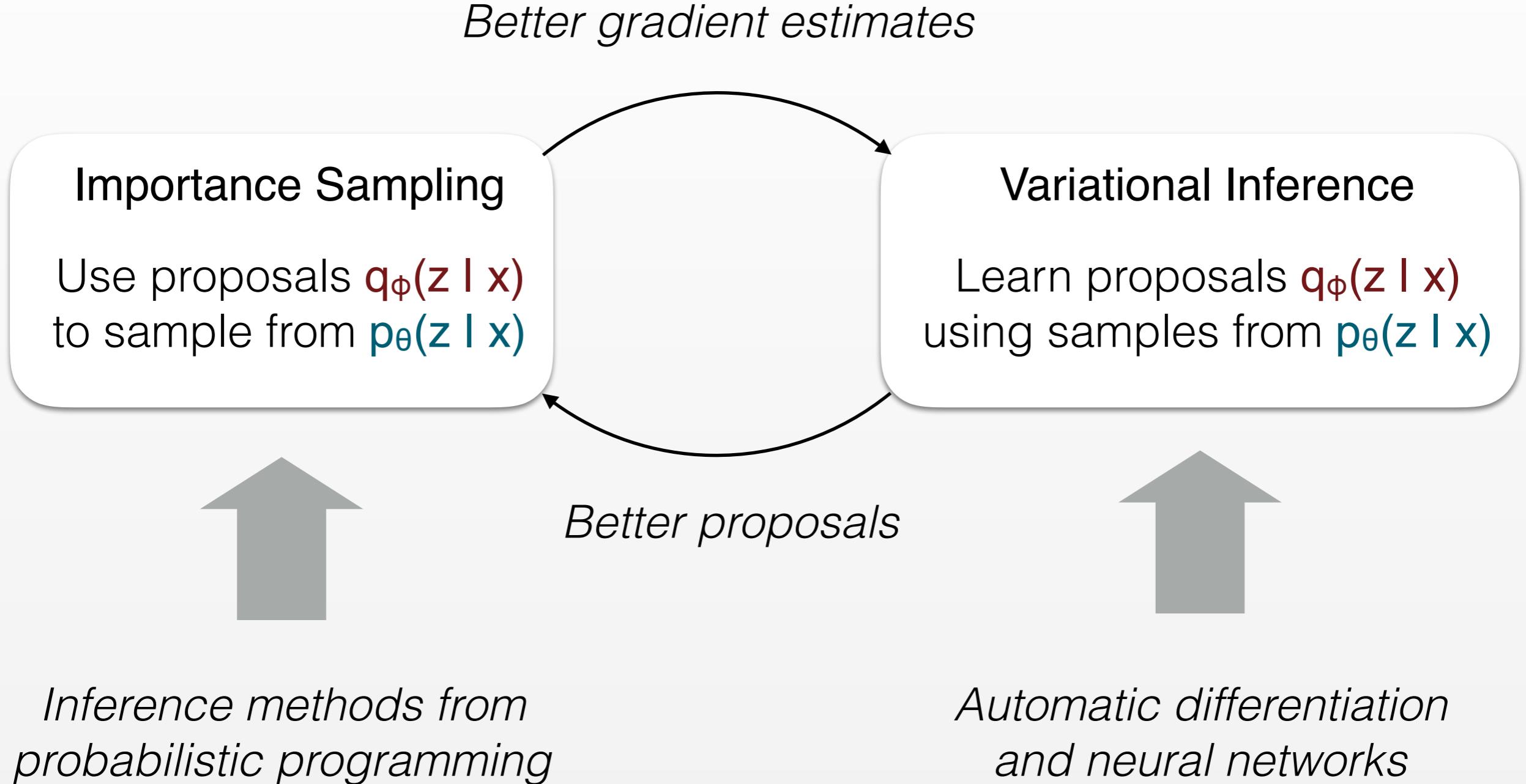
$$-\nabla_{\phi} \text{KL}(p_{\theta}(z|x) \parallel q_{\phi}(z|x)) = \mathbb{E}_{p_{\theta}(z|x)} \left[\nabla_{\phi} \log q_{\phi}(z|x) \right]$$

$$\nabla_{\theta} \log p_{\theta}(x) = \mathbb{E}_{p_{\theta}(z|x)} \left[\nabla_{\theta} \log p_{\theta}(x, z) \right]$$



Approximate with *any* importance sampler
(lots of probabilistic programming methods available,
can use learned $q_{\phi}(z|x)$ as proposals)

Reweighted Wake-sleep Style Methods



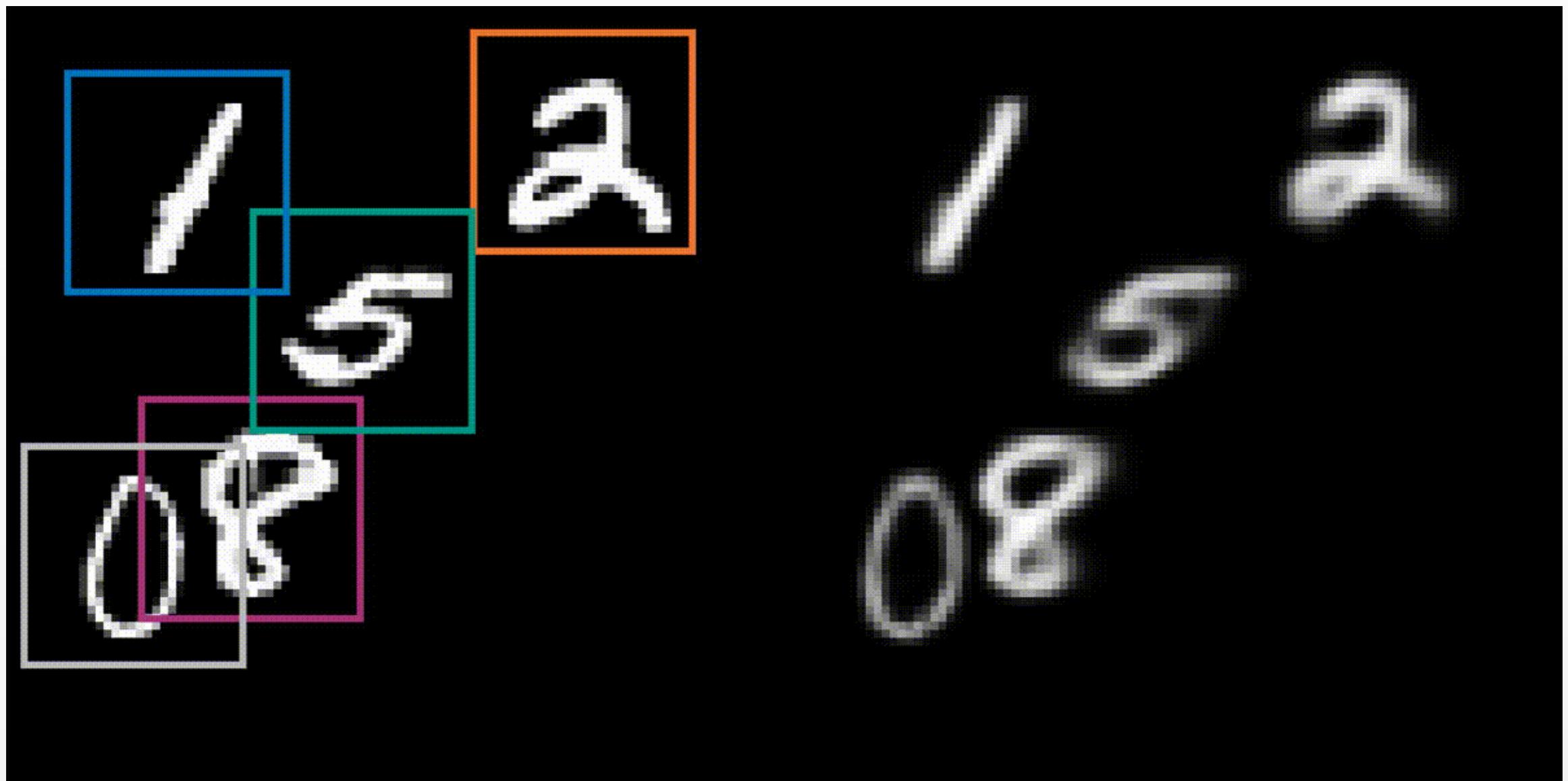
[Bornstein and Bengio, ICLR 2015]

[Le, Kosiorek, Siddarth, Teh, Wood, UAI 2019]

Example: Unsupervised Tracking

Inferred Locations

Reconstructions



Fully unsupervised (learns sub-model for MNIST)
Scales to $O(100)$ frames and ~5 digits (possibly more)

Thinking Compositionally

What (inference) DSL could define this sampler?

Algorithm 1 Amortized Population Gibbs Sampling

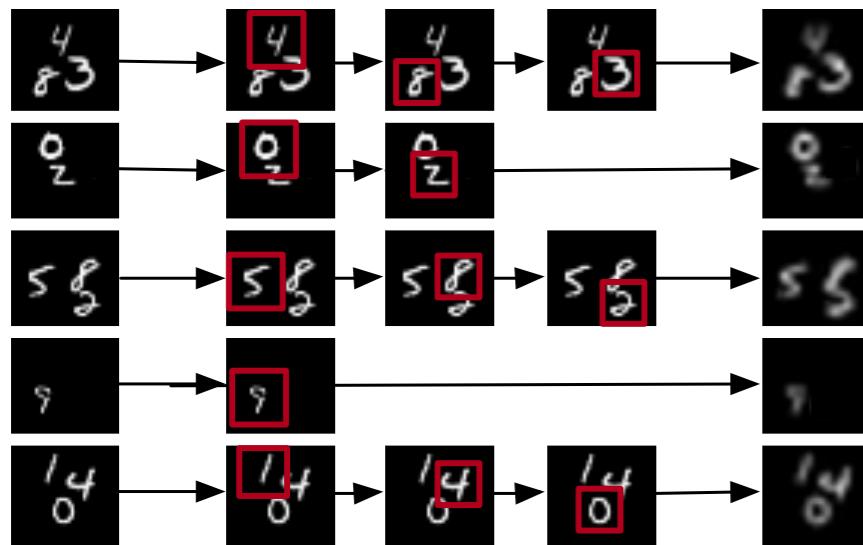
```
1: for  $n$  in  $1, \dots, N$  do
2:    $G_\phi = 0$ 
3:    $x^n \sim p^{\text{DATA}}(x)$ 
4:   for  $l$  in  $1, \dots, L$  do
5:      $z^{n,1,l} \sim q_\phi(z | x^n)$ 
6:      $w^{n,1,l} \leftarrow p_\theta(x^n, z^{n,1,l}) / q_\phi(z^{n,1,l})$ 
7:   for  $k$  in  $2, \dots, K$  do
8:      $\tilde{z}, \tilde{w} = z^{n,k-1}, w^{n,k-1}$ 
9:     for  $b$  in  $1, \dots, B$  do
10:       $\tilde{z}, \tilde{w} = \text{RESAMPLE}(\tilde{z}, \tilde{w})$ 
11:      for  $l$  in  $1, \dots, L$  do
12:         $\tilde{z}'^l_b \sim q_\phi(\cdot | x^n, \tilde{z}_{-b}^l)$ 
13:         $\tilde{w}^l = \frac{p_\theta(x^n, \tilde{z}'^l_b, \tilde{z}_{-b}^l)}{p_\theta(x^n, \tilde{z}_b^l, \tilde{z}_{-b}^l)} \frac{q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)}{q_\phi(\tilde{z}'^l_b | x^n, \tilde{z}_{-b}^l)} \tilde{w}^l$ 
14:         $\tilde{z}_b^l = \tilde{z}'^l_b$ 
15:         $G_\phi = G_\phi + \sum_{l=1}^L \frac{\tilde{w}^l}{\sum_{l'} \tilde{w}^{l'}} \frac{d}{d\phi} \log q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)$ 
16:       $z^{n,k}, w^{n,k} = \tilde{z}, \tilde{w}$ 
17: return  $G_\phi, z, w$ 
```

▷ Output: Grac

- APG combines inference from probabilistic programming with SGD-based methods
- Known building blocks (SMC and RWS), but not trivial to combine correctly
- Can we define compositional methods for importance sampling and gradient estimation?

Static vs Dynamic Models

Recurrent Recognition Loop

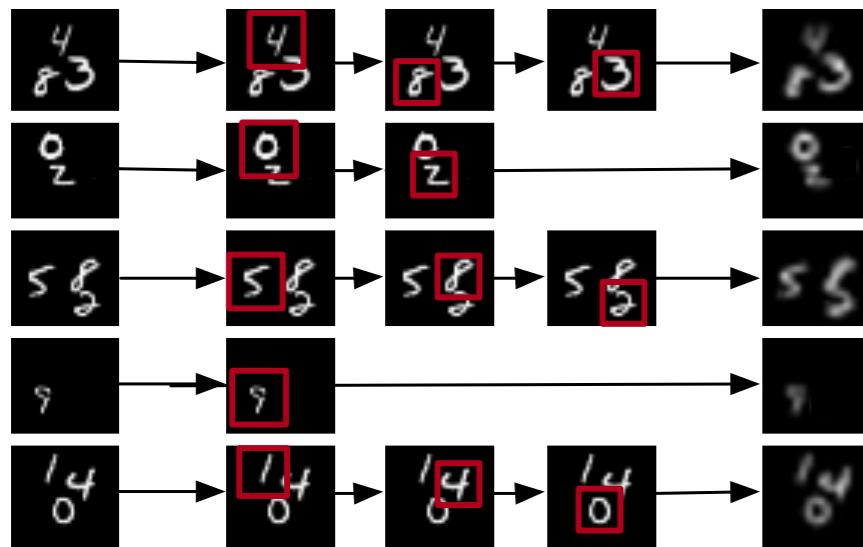


```
def next_object(image, objects):
    area = detect(image, objects)
    features = recognize(image, area)
    return area, features

def detect_objects(image, canvas):
    objects = []
    while not similar(image, objects):
        objects.append(
            next_object(image, canvas))
    return objects, canvas
```

Static vs Dynamic Models

Recurrent Recognition Loop



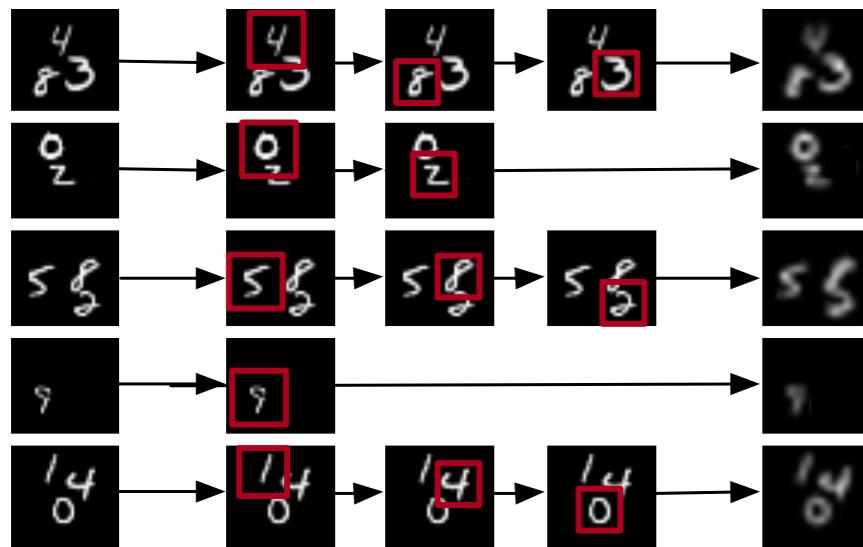
```
def next_object(image, objects):
    area = detect(image, objects)
    features = recognize(image, area)
    return area, features

def detect_objects(image, canvas):
    objects = []
    while not similar(image, objects):
        objects.append(
            next_object(image, canvas))
    return objects, canvas
```

Dynamic Computation Graphs: Number of variable nodes is data-dependent and/or stochastic

Static vs Dynamic Models

Recurrent Recognition Loop



```
def next_object(image, objects):
    area = detect(image, objects)
    features = recognize(image, area)
    return area, features

def detect_objects(image):
    objects = []
    for k in range(3):
        object =
            next_object(image, canvas)
        if not similar(image, objects):
            objects.append(object)
    return objects
```

Static Computation Graphs: Set of nodes and dependency graph determinable from static analysis.

Static vs Dynamic Models

Prob. Prog.

Static Graphs



Infer.NET



Stan*

Dynamic Graphs



Anglican

Venture
Turing Gen
WebPPL

Birch

(*dynamic graphs, static support)

Deep Learning



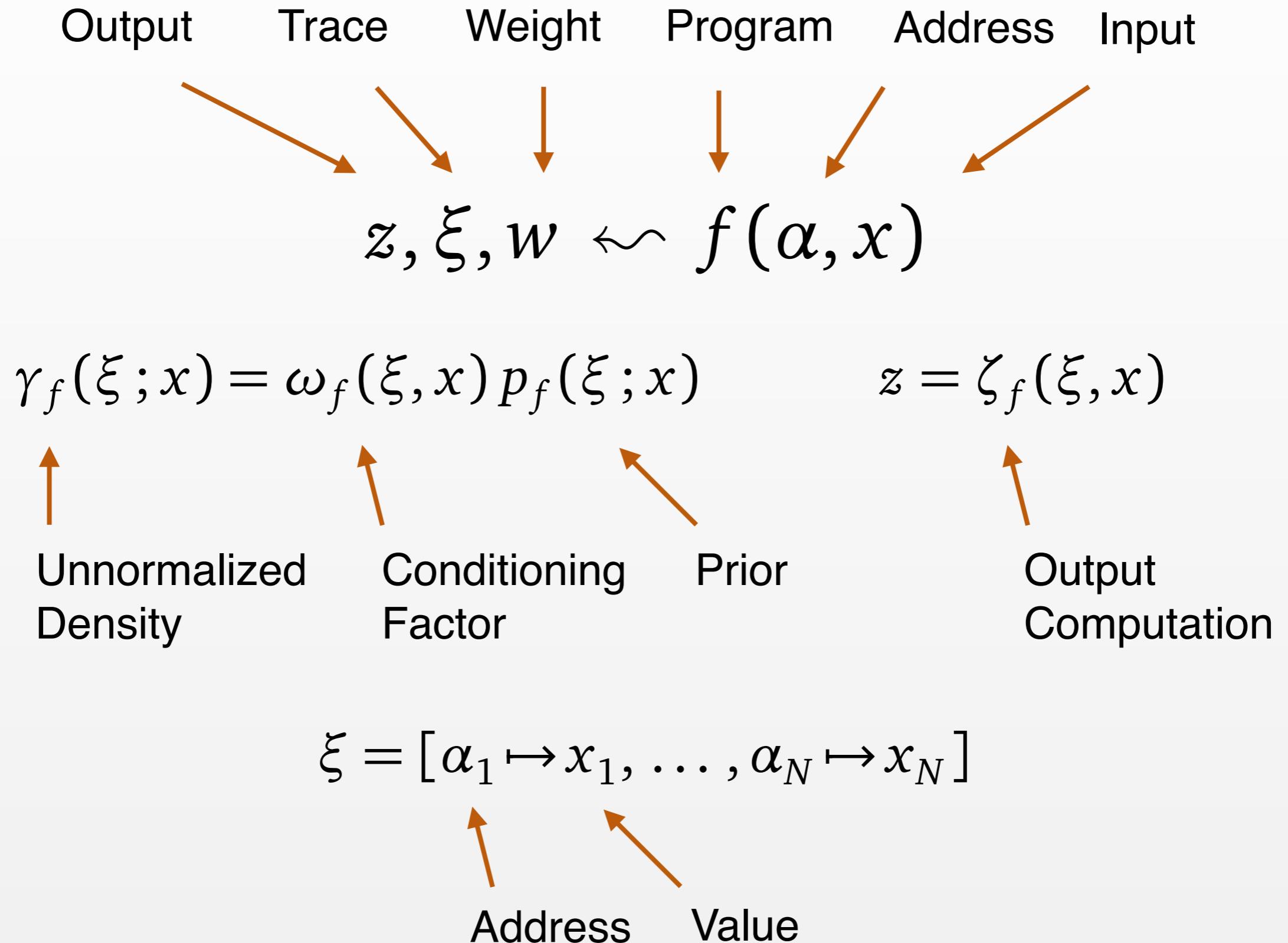
TensorFlow^{**}

(**in non-eager mode)

PyTorch

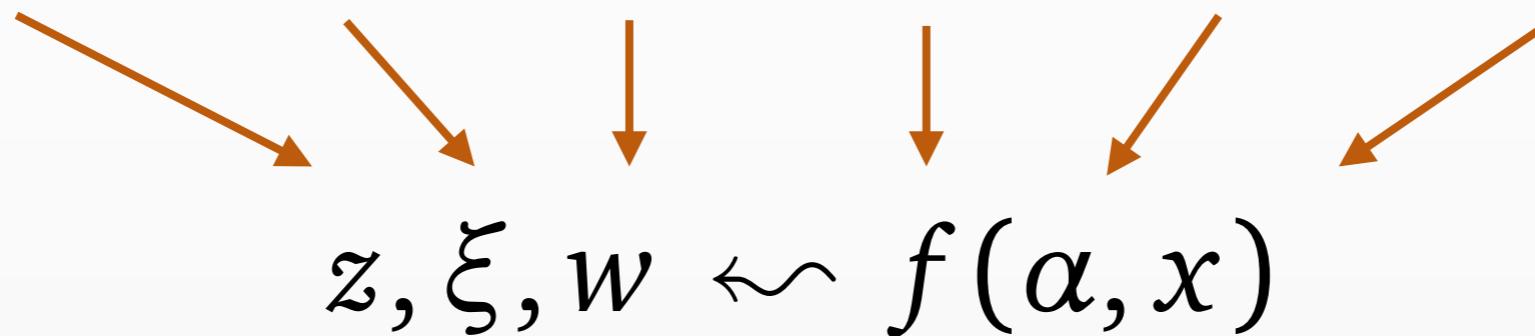
mxnet

Programs as Importance Samplers



Programs as Importance Samplers

Output Trace Weight Program Address Input



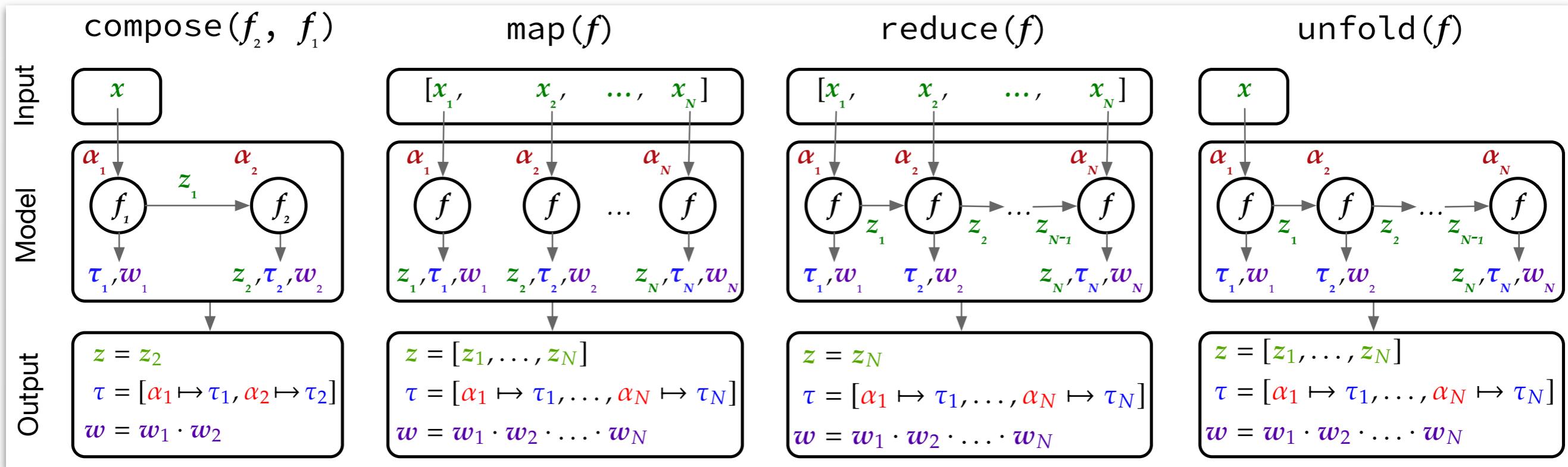
$$\gamma_f(\xi; x) = \omega_f(\xi, x) p_f(\xi; x) \quad z = \zeta_f(\xi, x)$$

(measure semantics)

$$\xi \sim p_f(\cdot; x) \quad w = \omega_f(\xi; x)$$

(likelihood weighting semantics)

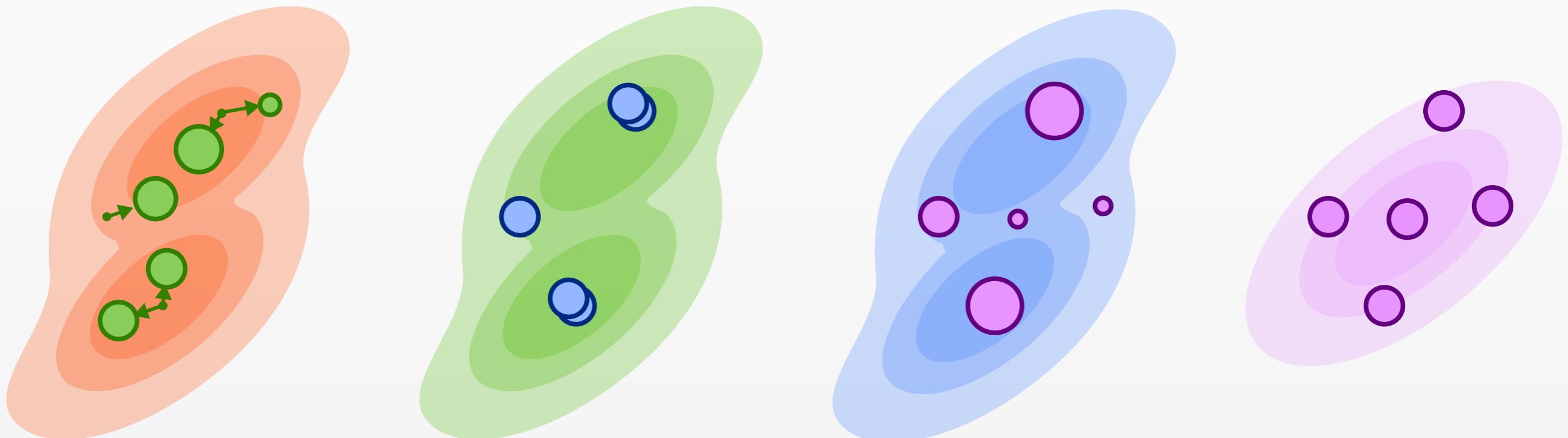
Model Combinators



Idea: *Static graphs with dynamic complexity*

Sampling Primitives

3 Building Blocks for Importance Samplers

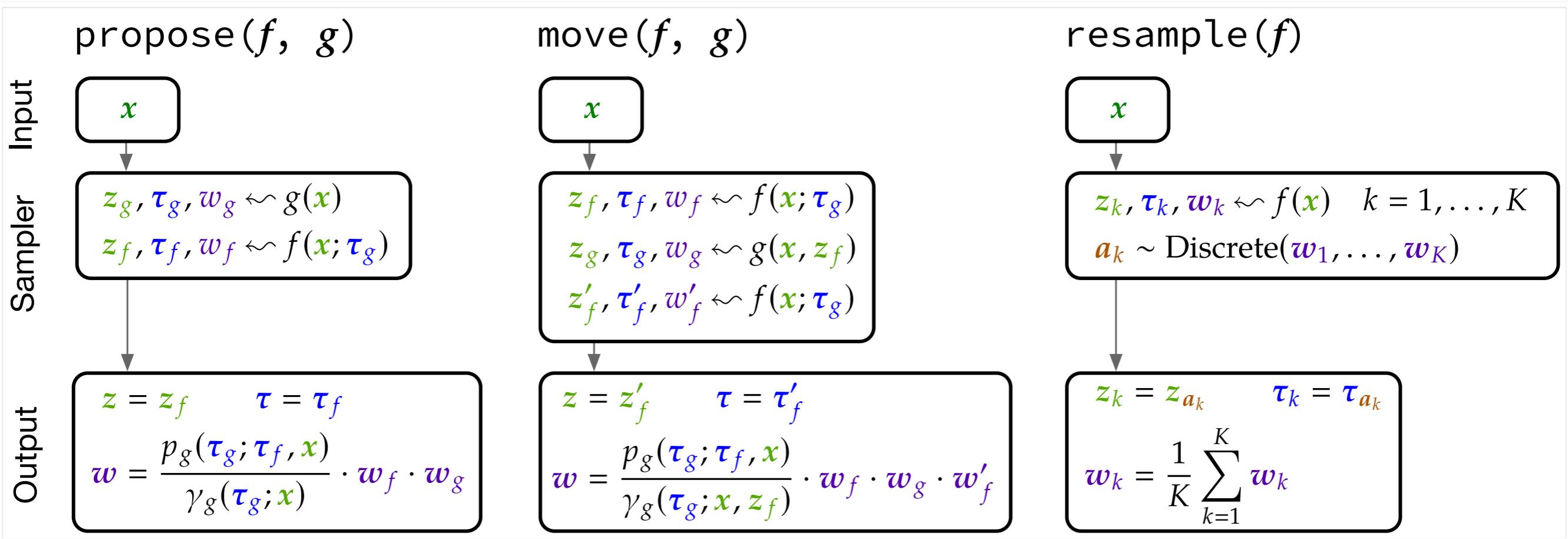


Move

Resample

Propose

Inference Combinators



Idea: Change the evaluation (sampling) strategy, whilst leaving the target density unaffected.

Dealing with Missing Variables

Target (variables: v, x1, x2)

```
def f(x0):  
    x1 = sample(categorical(prior_1), name='x1')  
    x2 = sample(normal(prior_2), name='x2')  
    v = sample(normal(dec_v(x1, x2)), name='v')  
    observe(normal(dec_0(v)), x0, name='x0')
```

$$\gamma_f(x_1, x_2, \nu; x_0) = p_f(x_0 | \nu) p_f(\nu | x_1, x_2) p_f(x_1, x_2)$$

Proposal (variables: u, x1, x2)

```
def g(x0):  
    u = sample(normal(enc_u(x0)), name='u')  
    x1 = sample(categorical(enc_1(u)), name='x1')  
    x2 = sample(normal(enc_2(u)), name='x2')
```

$$p_g(x_1, x_2, u; x_0) = p_g(x_1, x_2 | u) p_g(u; x_0)$$

Dealing with Missing Variables

Extended Target (variables: v, x1, x2, u)

```
def f(x0):  
    x1 = sample(categorical(prior_1), name='x1')  
    x2 = sample(normal(prior_2), name='x2')  
    v = sample(normal(dec_v(x1, x2)), name='v')  
    observe(normal(dec_0(v)), x0, name='x0')
```

$$\tilde{\gamma}_f(x_1, x_2, \nu; x_0) = p_f(x_0 | \nu) p_f(\nu | x_1, x_2) p_f(x_1, x_2) p_g(u; x_0)$$

Extended Proposal (variables: u, x1, x2, v)

```
def g(x0):  
    u = sample(normal(enc_u(x0)), name='u')  
    x1 = sample(categorical(enc_1(u)), name='x1')  
    x2 = sample(normal(enc_2(u)), name='x2')
```

$$\tilde{p}_g(x_1, x_2, u; x_0) = p_g(x_1, x_2 | u) p_g(u; x_0) p_f(\nu | x_1, x_2)$$

Dealing with Missing Variables

Extended Target (variables: v, x1, x2, u)

```
def f(x₀):  
    x₁ = sample(categorical(prior_1), name='x₁')  
    x₂ = sample(normal(prior_2), name='x₂')  
    v = sample(normal(dec_v(x₁, x₂)), name='v')  
    observe(normal(dec_θ(v)), x₀, name='x₀')
```

$$\tilde{\gamma}_f(x_1, x_2, v; x_0) = p_f(x_0 | v) p_f(v | x_1, x_2) p_f(x_1, x_2) p_g(u; x_0)$$

Extended Proposal (variables: u, x1, x2, v)

```
def g(x₀):  
    u = sample(normal(enc_u(x₀)), name='u')  
    x₁ = sample(categorical(enc_1(u)), name='x₁')  
    x₂ = sample(normal(enc_2(u)), name='x₂')
```

$$\tilde{p}_g(x_1, x_2, u; x_0) = p_g(x_1, x_2 | u) p_g(u; x_0) p_f(v | x_1, x_2)$$

Dealing with Missing Variables

$$w_f = \frac{\tilde{r}_f(x_1, x_2, v; x_0)}{\tilde{p}_g(x_1, x_2, u; x_0)} = \frac{p_f(x_0 | v) p_f(x_1, x_2)}{p_g(x_1, x_2 | u)}$$

(ratio for variables common to both models)

Proposal

$$z_1, \xi_1, w_1 \leftarrow g(x_0)$$

Conditioned Target

$$z_2, \xi_2, w_2 \leftarrow f(x_0 ; \xi_1)$$

$$z_2, \xi_2, w_1 \cdot w_2 \cdot \frac{p_g(\xi_1 ; \xi_2, x_0)}{p_g(\xi_1 ; x_0)} \leftarrow \text{propose}(f, g)(x_0)$$

Conditioned reverse

Combinator

Model Combinators

define compositional model structure

Function Composition

$$\frac{z_1, \tau_1, w_1 \leftarrow g(x_0 ; \tau_0(\alpha_1)) \quad z_2, \tau_2, w_2 \leftarrow f(z_1 ; \tau_0(\alpha_2))}{z_2, [\alpha_1 \mapsto \tau_1, \alpha_2 \mapsto \tau_2], \quad w_1 \cdot w_2 \leftarrow \text{compose}(f, g)(x_0 ; \tau_0)}$$

Map over Sequence

$$\frac{z_n, \tau_n, w_n \leftarrow f(x_n ; \tau_0(\alpha_n)) \quad \text{for } n = 1, \dots, N}{(z_1, \dots, z_N), [\alpha_1 \mapsto \tau_1, \dots, \alpha_N \mapsto \tau_N], \prod_{n=1}^N w_n \leftarrow \text{map}(f)((x_1, \dots, x_N) ; \tau_0)}$$

Reduce / Fold over Sequence

$$\frac{z_1, \tau_1, w_1 \leftarrow g(x_0 ; \tau_0(\alpha_1)) \quad z_n, \tau_n, w_n \leftarrow f(z_{n-1} ; \tau_0(\alpha_n)) \quad \text{for } n = 2, \dots, N}{z_N, [\alpha_1 \mapsto \tau_1, \dots, \alpha_N \mapsto \tau_N], \prod_{n=1}^N w_n \leftarrow \text{reduce}(f, g)((x_1, \dots, x_N) ; \tau_0)}$$

Inference Combinators

define a compositional sampling strategy

Nested Importance Sampling

$$\frac{z_1, \tau_1, w_1 \leftarrow g(x_0) \quad z_2, \tau_2, w_2 \leftarrow f(x_0; \tau_1)}{z_2, \tau_2, \frac{w_1 \cdot w_2}{\omega_g(\tau_1; \tau_2, x_0)} \leftarrow \text{propose}(f, g)(x_0)}$$

Importance Resampling

$$\frac{z^k, \tau^k, w^k \leftarrow f(\alpha, x) \quad \text{for } k = 1, \dots, K \quad a^1, \dots, a^K \sim \text{Discrete}(w^1, \dots, w^K)}{z^{a^k}, \tau^{a^k}, \frac{1}{K} \sum_k w^k \leftarrow \text{resample}(f, K)(\alpha, x) \quad \text{for } k = 1, \dots, K}$$

Application of a Transition Kernel

$$\frac{z_1, \tau_1, w_1 \leftarrow f(x_0) \quad z_2, \tau_2, w_2 \leftarrow g(z_1) \quad z_3, \tau_3, w_3 \leftarrow f(x_0; \tau_2)}{z_3, \tau_3, w_1 \cdot \frac{w_2 \cdot w_3}{\omega_g(\tau_2; \tau_3, z_1)} \leftarrow \text{move}(f, g)(x_0)}$$

Scaling up Amortized Inference

Deep Generative Model
program $p_\theta(x, z)$

Inference Model
program $q_\phi(z | x)$

Better gradient estimates

Importance Sampling

Use proposals $q_\phi(z | x)$
to sample from $p_\theta(z | x)$

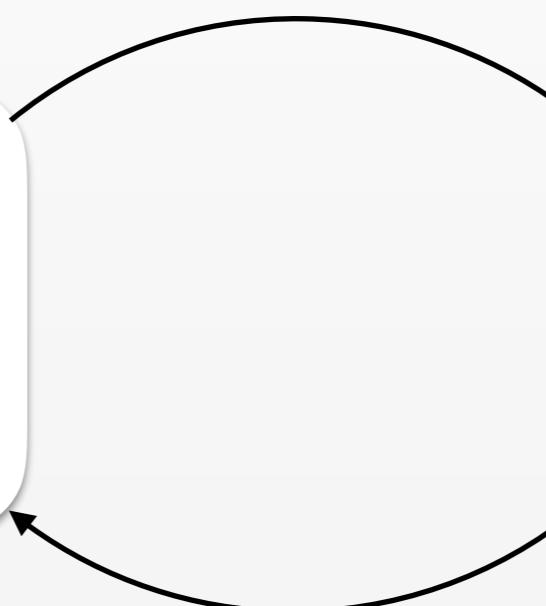
Variational Inference

Learn proposals $q_\phi(z | x)$
using samples from $p_\theta(z | x)$

Better proposals

User-specified
importance sampler
(inference combinator)

User-specified
variational objective
(next talk)



Thank You!



Hao
Wu



Eli
Sennesh



Sam Stites



Heiko
Zimmermann



Babak
Esmaeli



Alican
Bozkurt

Composing Modeling and Inference Operations

with Probabilistic Program Combinators

E. Sennesh, A. Scibior, H. Wu, J.-W. van de Meent

ArXiv 2018 [<https://arxiv.org/abs/1811.05965>]

Amortized Gibbs Samplers with Neural Sufficient Statistics

H. Wu, H. Zimmermann, E. Sennesh, J.-W. Van de Meent

ArXiv 2019 [<https://arxiv.org/abs/1911.01382>]

Structured Neural Topic Models for Reviews

B. Esmaeli, H. Huang, B. Wallace, J.-W. van de Meent

AISTATS 2019 [<http://proceedings.mlr.press/v89/esmaeli19b.html>]

An Introduction to Probabilistic Programming

J-W van de Meent, B. Paige, H. Yang, F. Wood

ArXiv 2018 [<https://arxiv.org/abs/1809.10756>]

Evaluating Combinatorial Generalization in VAEs

A. Bozkurt, B. Esmaeli, D. H. Brooks,

J. Dy, J.-W. van de Meent

ArXiv 2019 [<https://arxiv.org/abs/1911.04594>]

Source Code (Apache 2.0)

<https://github.com/probtorch/probtorch>

