# Structured Disentangled Representations

**Jan-Willem van de Meent**

Babak Esmaeli    Hao Wu    Sarthak Jain    Alican Bozkurt    Siddharth N.    Brooks Paige

*Northeastern University*    *Oxford*    *Alan Turing Institute*
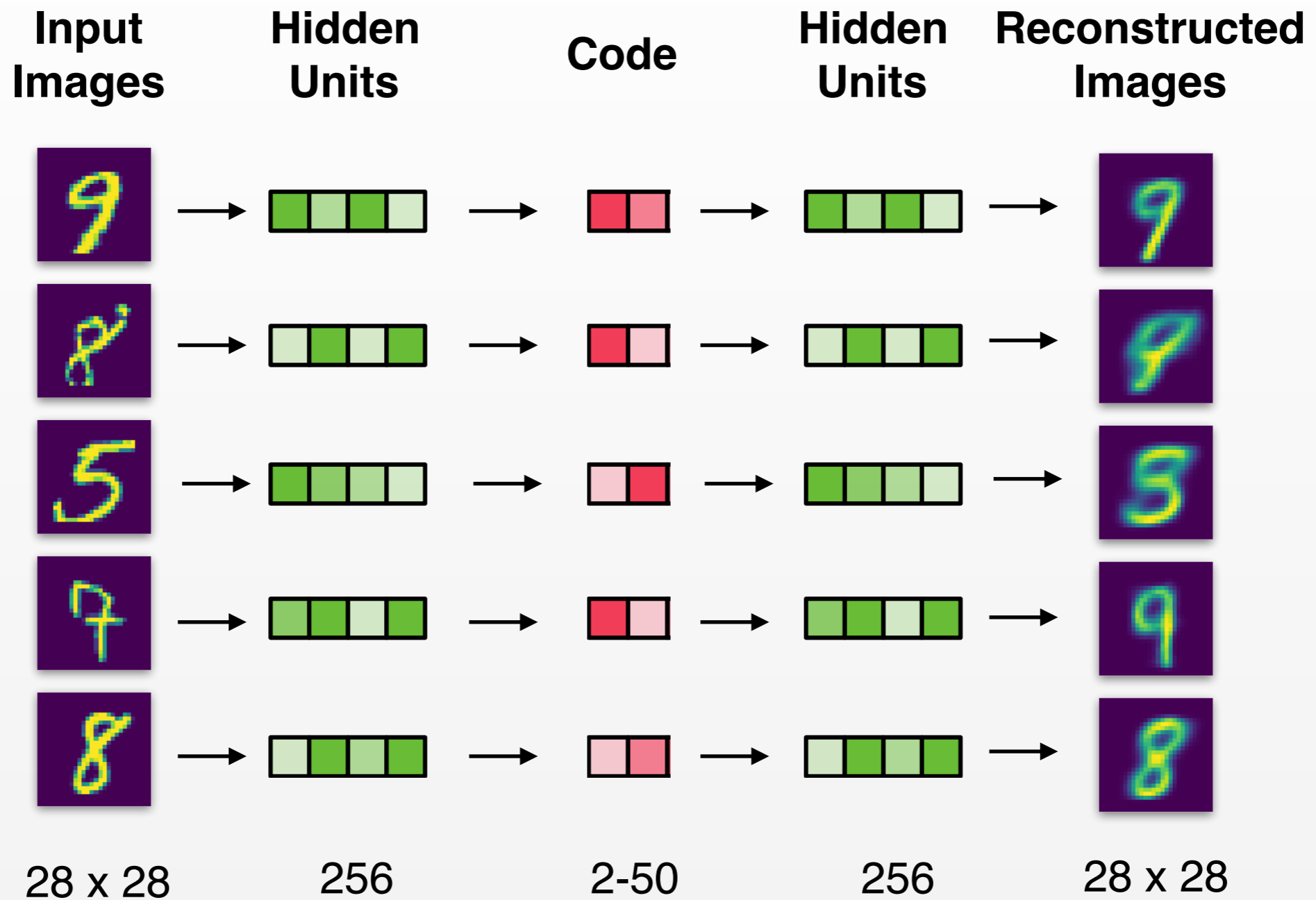
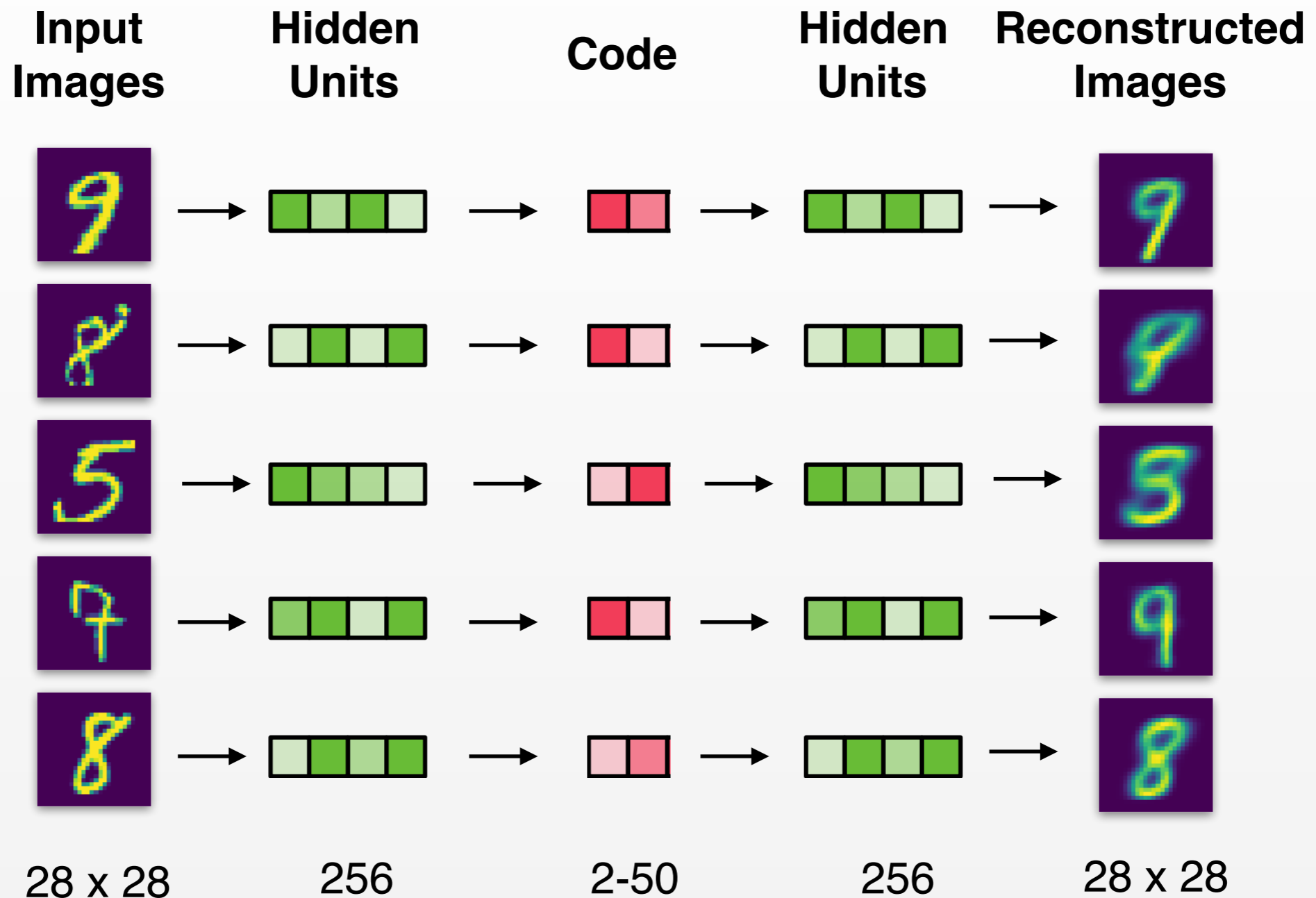# Learning Representations

Everyone's Favorite Example: MNIST



**Goal:** Learn features that capture *similarities* and *dissimilarities*

**Requirement:** Objective that defines notion of *utility* (*task-dependent*)

# Autoencoders



| Input Images | Hidden Units | Code | Hidden Units | Reconstructed Images |
| --- | --- | --- | --- | --- |
| 28 x 28 | 256 | 2-50 | 256 | 28 x 28 |

# Autoencoders



| Input Images | Hidden Units | Code | Hidden Units | Reconstructed Images |
|:---:|:---:|:---:|:---:|:---:|
| 28 x 28 | 256 | 2-50 | 256 | 28 x 28 |

**Notion of Utility:** Ability to reconstruction of pixels from code
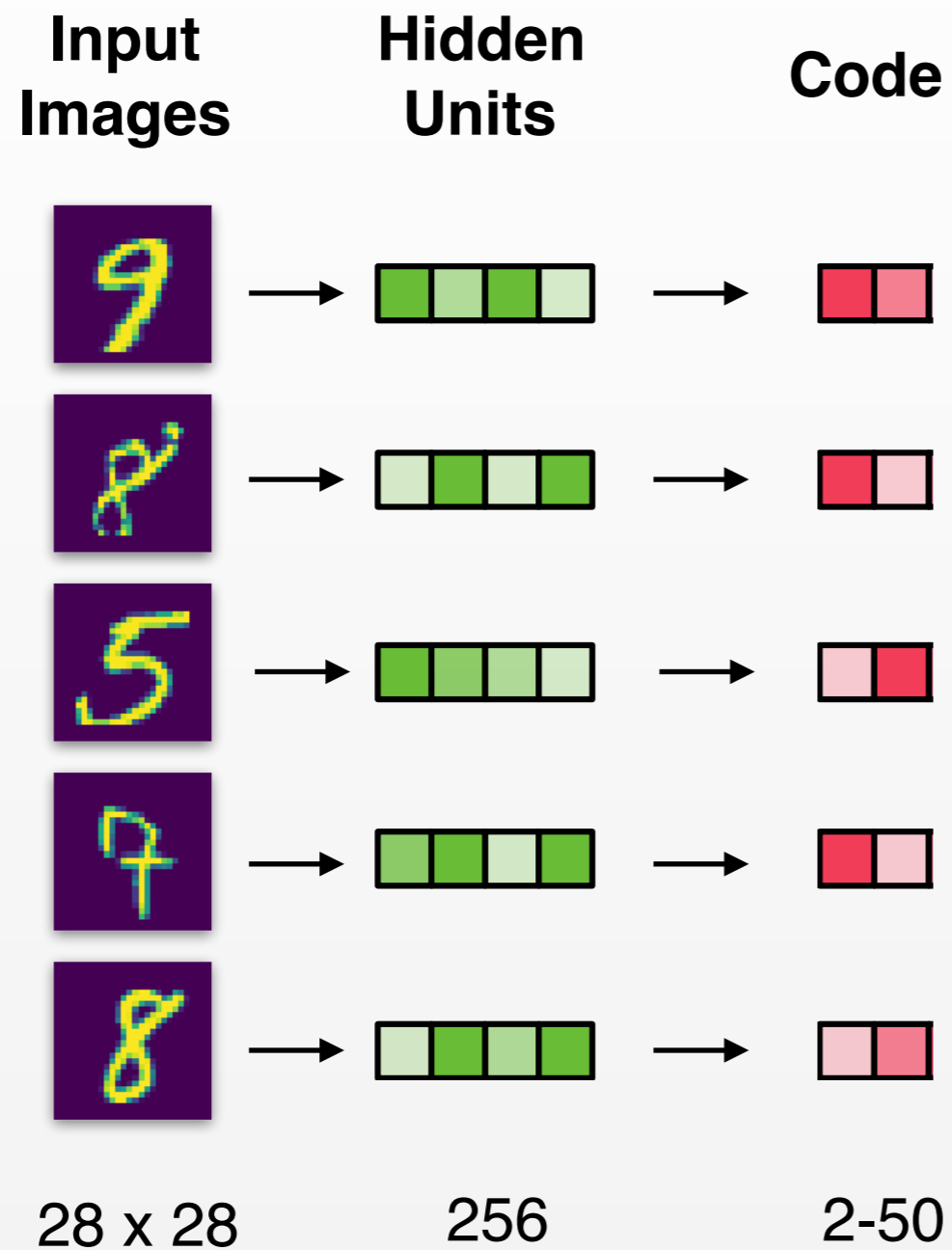(features are a *compressed* representation of original data)

# Autoencoders



**Notion of Utility:** Ability to reconstruction of pixels from code (features are a *compressed* representation of original data)
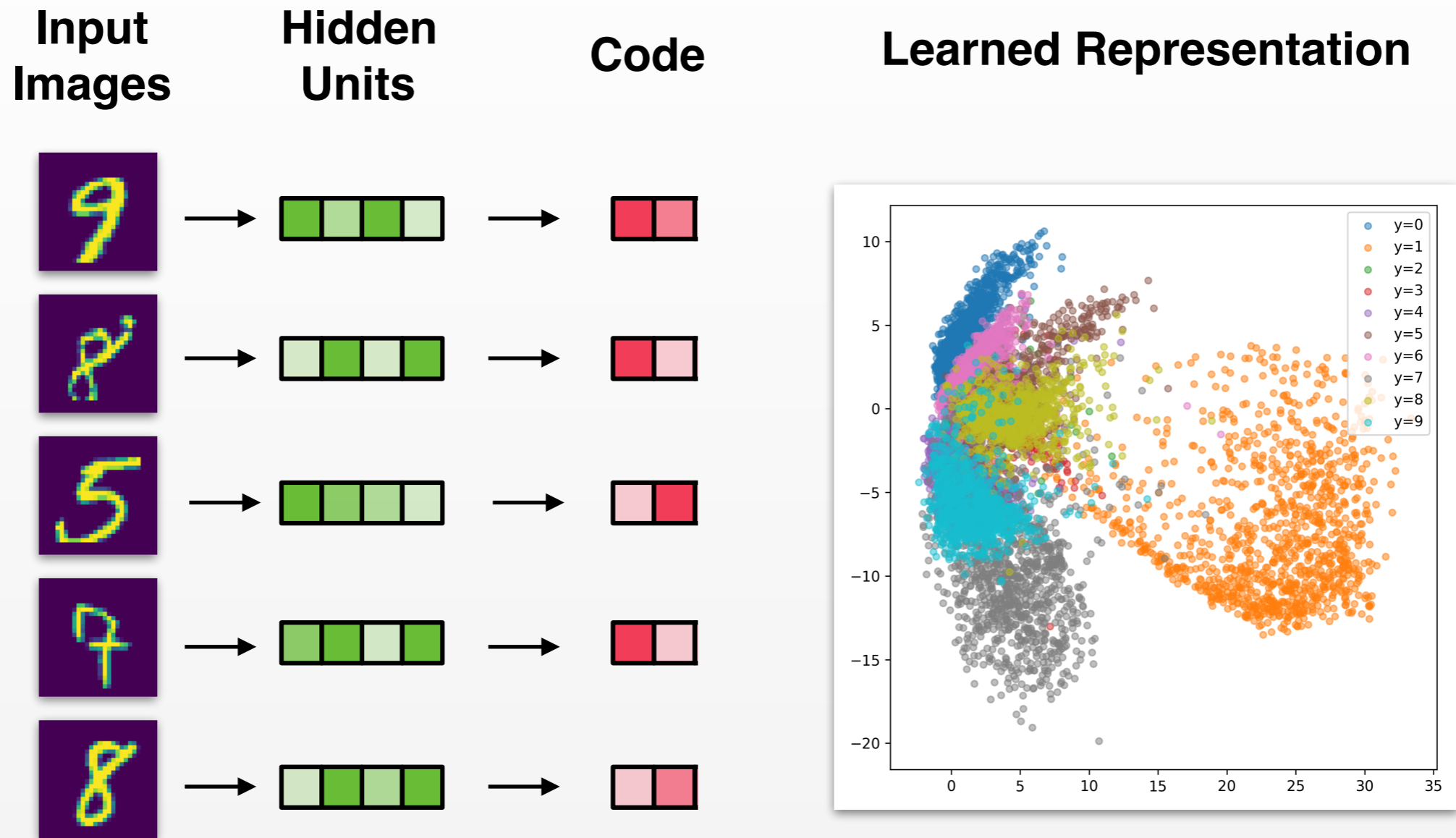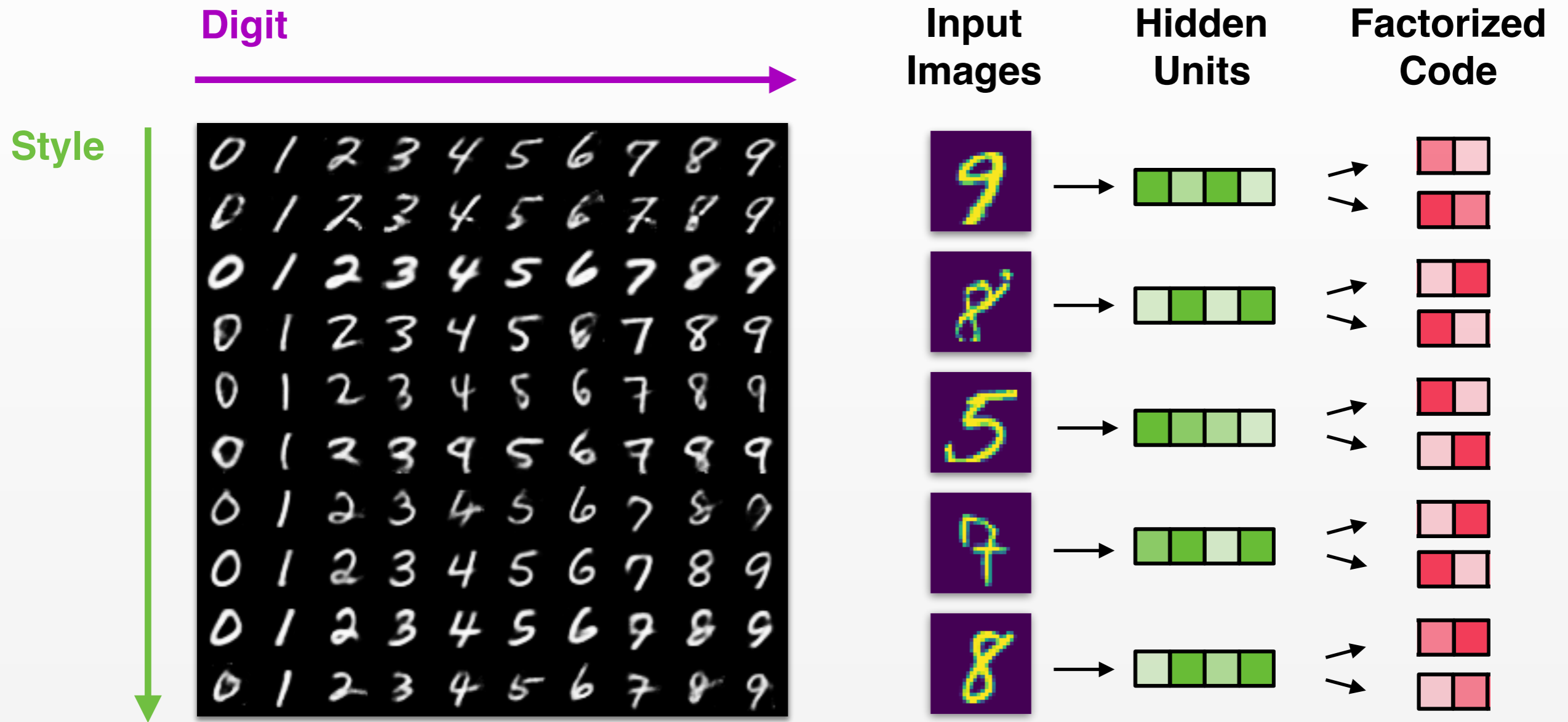
# Autoencoders

**Input Images**    **Hidden Units**    **Code**    **Learned Representation**

# Autoencoders



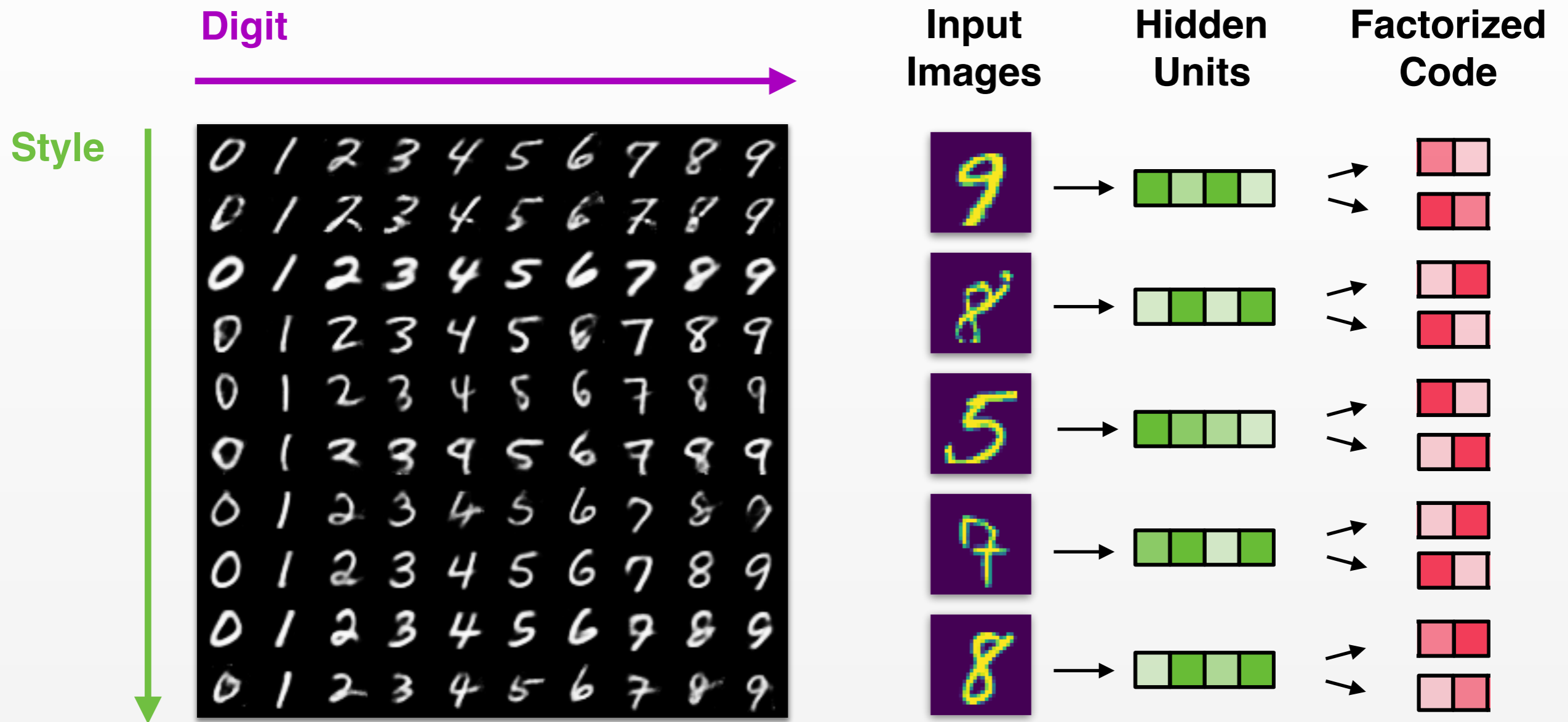**Input Images**  **Hidden Units**  **Code**  **Learned Representation**

**Entangled Representation:** Individual dimensions in code encode some unknown combination of features in the data.
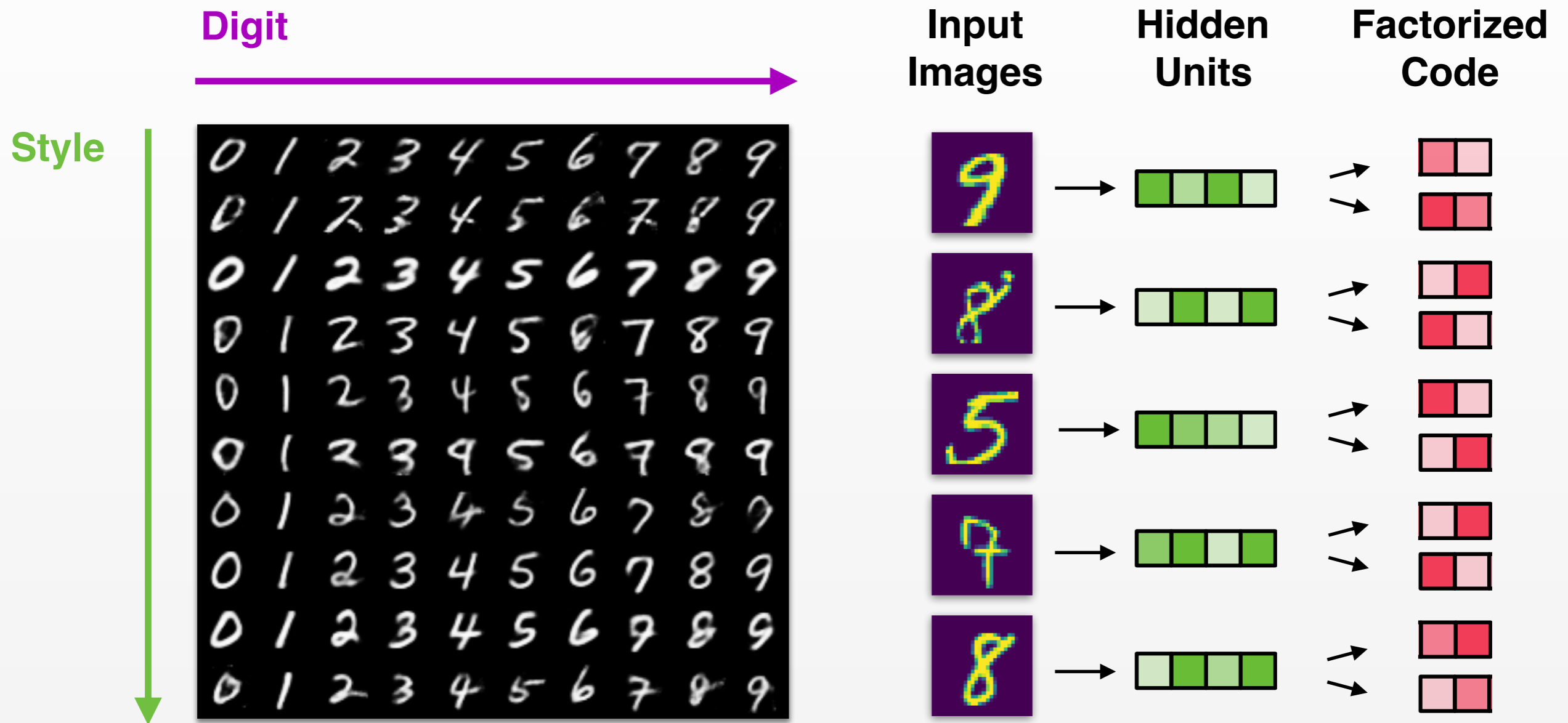
# Disentangled Representations

**Digit**

**Style**

# Disentangled Representations

# Disentangled Representations



**Goal:** Learn features that correspond to *distinct* factors of variation

# Disentangled Representations



**Goal:** Learn features that correspond to *distinct* factors of variation

**One Notion of Utility:** Statistical independence

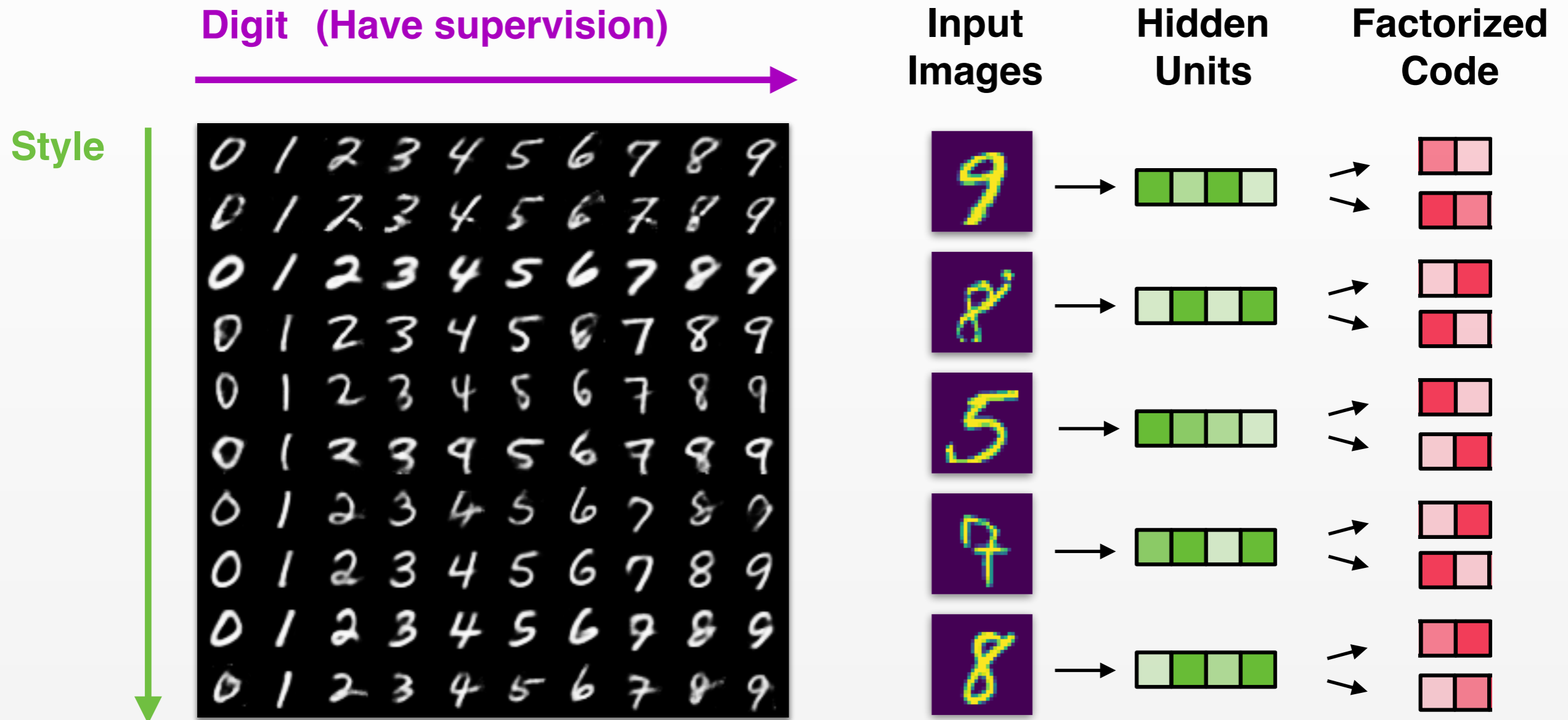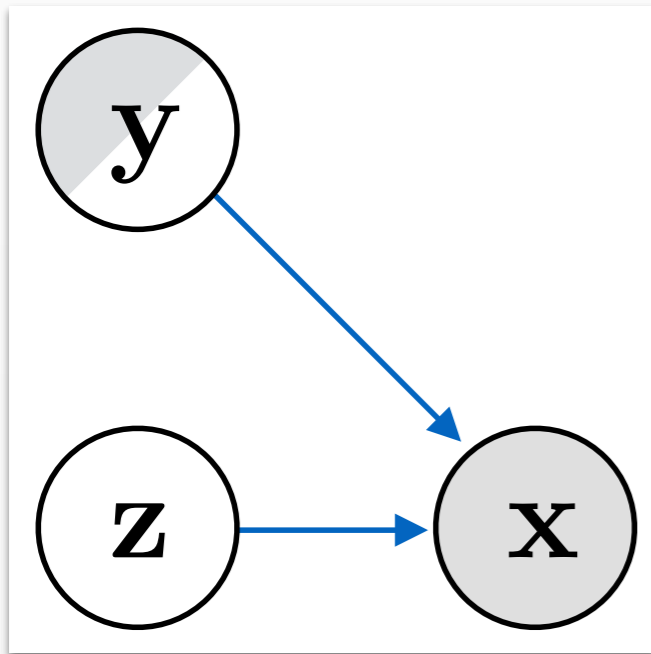# Disentangled Representations



**Goal:** Learn features that correspond to *distinct* factors of variation

**One Notion of Utility:** Statistical independence

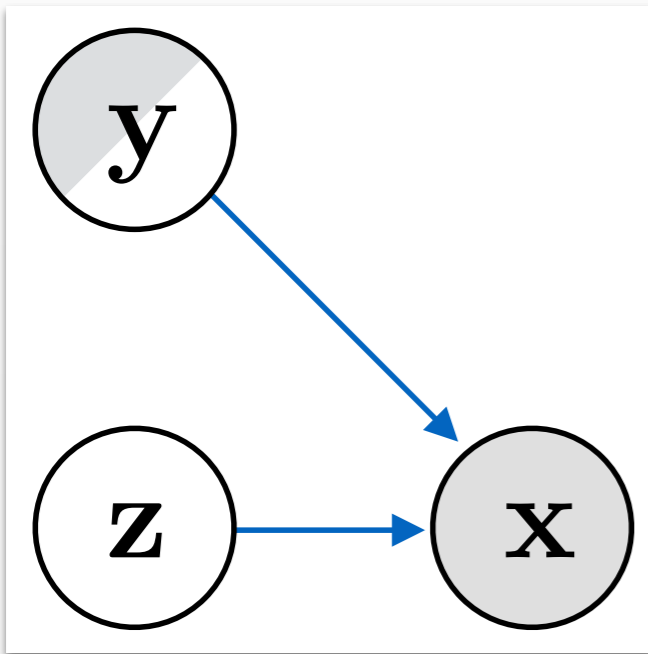# Semi-supervised Learning

Generative Model (Decoder)

$$p_\theta(x \mid y, z)$$



Assume independence between digit **y** and style **z**



(a)

Figure 2: (a) Vis label varied. analogies for the label git.

$\lambda$   $\varepsilon$

$q$   $y$   $p$

To train deep gen the variation point $x^i$ to the EL

$\lambda$   $\varepsilon$

128

$q$   $z$   $p$

[Kingma, Mohamed, Jimenez-Rezende, Welling, *Semi-supervised Learn...* NIPS 2014]

# Semi-supervised Learning



**Generative Model (Decoder)**

$$p_\theta(x \mid y, z)$$



Assume independence
between digit **y** and style **z**

**Inference Model (Encoder)**

$$q_\phi(y, z \mid x)$$



Infer **y** from pixels **x**,
and **z** from **y** and **x**

(a)

Figure 2: (a) Vis... label varied. ... analogies for the... label... git.

To train deep gen... the variation... point $x^i$ to the EL...

[Kingma, Mohamed, Jimenez-Rezende, Welling, *Semi-supervised Learn... NIPS 2014]

# Semi-supervised Lear...

**Generative Model (Decoder)**

$$p_\theta(x\,|\,y,z)$$



Assume independence
between digit **y** and style **z**

**Inference Model (Encoder)**

$$q_\phi(y,z|x)$$



Infer **y** from pixels **x**,
and **z** from **y** and **x**



Separate interpretable **y**
from "nuisance" variables **z**

[Kingma, Mohamed, Jimenez-Rezende, Welling, *Semi-supervised Learn... ee...* NIPS 2014

# Semi-supervised Learning

**Generative Model (Decoder)**

$$p_\theta(x \mid y, z)$$



Assume independence
between digit **y** and style **z**

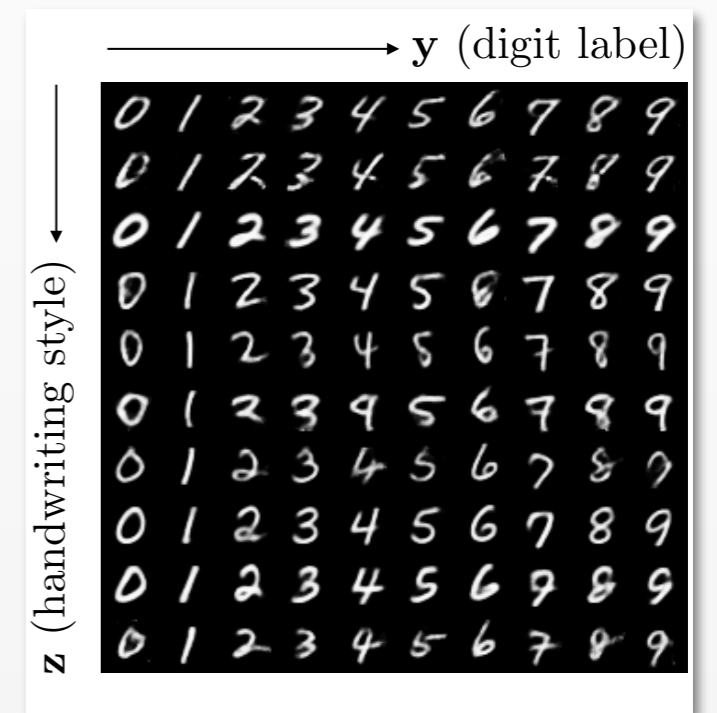**Inference Model (Encoder)**

$$q_\phi(y, z \mid x)$$



Infer **y** from pixels **x**,
and **z** from **y** and **x**



Separate interpretable **y**
from "nuisance" variables **z**

**Hypothesis:** Assuming a statistical independence
under the prior induces disentangled representations.

[Kingma, Mohamed, Jimenez-Rezende, Welling, *Semi-supervised Learning*, NIPS 2014]

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch



https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch

**PROB TORCH**

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch

PROB TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch



https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch



https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch

PROB
TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch

PROB TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch



PROB
TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch



PROB TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```
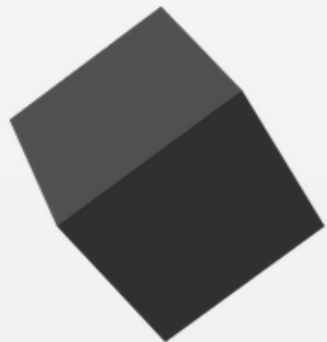
## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Probabilistic Torch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```
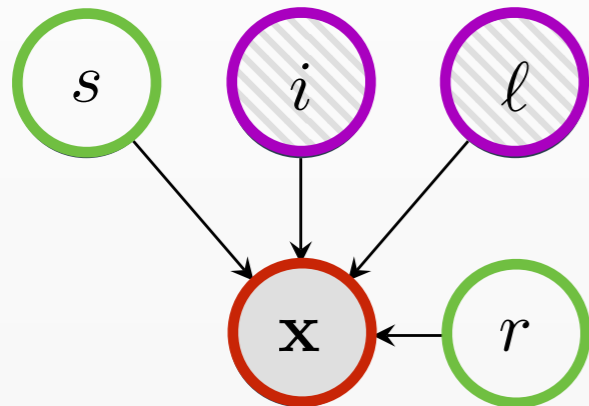
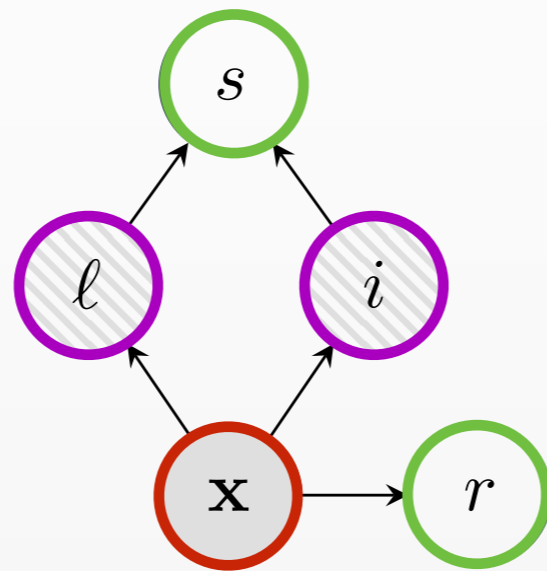## Probabilistic Torch



PROB
TORCH

https://github.com/probtorch/probtorch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Edward



## Probabilistic Torch

# Deep Probabilistic Programs

## Inference Model (Encoder)

```python
class Encoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Generative Model (Decoder)

```python
class Decoder(torch.nn.Module):
  def __init__(self, x_sz, h_sz, y_sz, z_sz):
    # intializes 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
```

## Edward

## Probabilistic Torch

PROB TORCH

## Pyro

# Example: Yale B Faces

**Generative Model**



**Inference Model**



**Semi-supervised:**
**Identity (i), Lighting (l)**

**Unsupervised:**
**Shading (s), Reflectance (r)**

**Data:**
**Pixels (x)**

Original    Reconstruction

Same Lighting, Different Identity



Same Identity, Different Lighting



[Siddharth, Paige, van de Meent, Desmaison, Wood, Goodman, Kohli, Torr, NIPS 2017]

# Example: Multiple MNIST Digits



Inference Model

Generative Model

Semi-supervised

Data

Unsupervised

Multiple MNIST digits

[Siddharth, Paige, van de Meent, Desmaison, Wood, Goodman, Kohli, Torr, NIPS 2017]

# Example: Multiple MNIST Digits



Encoder

Decoder

— **Semi-supervised**

— **Data**

**Unsupervised**

Data

Decomposition

Accuracy

| $\frac{M}{M+N}$ | Count Error (%) | |
|---|---|---|
| | w/o MNIST | w/ MNIST |
| 0.1 | 85.45 ($\pm$ 5.77) | 76.33 ($\pm$ 8.91) |
| 0.5 | 93.27 ($\pm$ 2.15) | 80.27 ($\pm$ 5.45) |
| 1.0 | 99.81 ($\pm$ 1.81) | 84.79 ($\pm$ 5.11) |

[Siddharth, Paige, van de Meent, Desmaison, Wood, Goodman, Kohli, Torr, NIPS 2017]

# Deep Probabilistic Models

## Example: Modeling Laboratory Mice with Deep State Space Models



[Johnson, Duvenaud, Wiltschko, Adams, Datta, NIPS 2016]

# Deep Probabilistic Models

Example: Modeling Laboratory Mice with Deep State Space Models



[Johnson, Duvenaud, Wiltschko, Adams, Datta, NIPS 2016]

## Data Science ❤️ Probabilistic Modeling + Deep Learning?

• Structured priors model problem domain

• Bayesian inference for uncertainty estimates

• Neural likelihood models for data such as text and images

• Neural inference models predict values for latent variables

# Today: Unsupervised Le...

**Generative Model (Decoder)**

$$p_\theta(x \mid y, z)$$



Assume independence
between digit **y** and style **z**

**Inference Model (Encoder)**

$$q_\phi(y, z \mid x)$$



Infer **y** from pixels **x**,
and **z** from **y** and **x**

Separate interpretable **y**
from "nuisance" variables **z**

**Hypothesis:** Assuming a statistical independence
under the prior induces disentangled representations.

# Today: Unsupervised Le

Generative Model (Decoder)

$$p_\theta(x \mid y, z)$$

Inference Model (Encoder)

$$q_\phi(y, z \mid x)$$



Assume independence
between digit **y** and style **z**

Infer **y** from pixels **x**,
and **z** from **y** and **x**

Separate interpretable **y**
from "nuisance" variables **z**

**Problem:** Unsupervised learning (with same model)
does *not* disentangle digit from style

# Variational Autoencoders

(a.k.a. Deep Latent-Variable Models)

# Autoencoders

| Input Images | Hidden Units | Code | Hidden Units | Reconstructed Images |
|:---:|:---:|:---:|:---:|:---:|
| $x_n$ | $h_n$ | $z_n$ | $\hat{h}_n$ | $\hat{x}_n$ |
| 784 (28 x 28) | 256 | 2-50 | 256 | 784 (28 x 28) |

# Autoencoders



**Input Images**  **Hidden Units**  **Code**  **Hidden Units**  **Reconstructed Images**

Encoder $z(x; \phi)$

Decoder $\hat{x}(z; \theta)$

# Autoencoders



Input Images | Hidden Units | Code | Hidden Units | Reconstructed Images

Objective: Maximize Reconstruction Quality

$$\max_{\phi,\theta} \frac{1}{N} \sum_{n=1}^{N} \sum_{p=1}^{P} x_{n,p} \log \hat{x}_{n,p}(z(x_n;\phi);\theta)$$

# Variational Autoencoders

# Variational Autoencoders



Encoder (Inference Network)

$$q_\phi(z \mid x)$$

# Variational Autoencoders



| Input Images | Hidden Units | Mean Std | Code (Random) | Hidden Units | Likelihood Mean |
|---|---|---|---|---|---|

Encoder (Inference Network)

$$q_\phi(\boldsymbol{z} \mid \boldsymbol{x})$$

Decoder (Likelihood)

$$p_\theta(\boldsymbol{x} \mid \boldsymbol{z})$$

# Variational Autoencoders



| Input Images | Hidden Units | Mean Std | Code (Random) | Hidden Units | Likelihood Mean |
|---|---|---|---|---|---|

Encoder (Inference Network)

$q_\phi(\boldsymbol{z} \mid \boldsymbol{x})$

Prior

$p(\boldsymbol{z})$

Decoder (Likelihood)

$p_\theta(\boldsymbol{x} \mid \boldsymbol{z})$

# Variational Inference

$$\mathscr{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{z} \mid \boldsymbol{x})) \Big]$$

# Variational Inference

Objective: Maximize Evidence Lower Bound

$$\mathscr{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[\log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \,|\, \boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{z} \,|\, \boldsymbol{x}))\Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Bigg[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[\log \frac{p_\theta(\boldsymbol{x})p_\theta(\boldsymbol{z} \,|\, \boldsymbol{x})}{q_\phi(\boldsymbol{z} \,|\, \boldsymbol{x})}\bigg]\Bigg]$$

# Variational Inference

Objective: Maximize Evidence Lower Bound

$$\mathscr{L}(\theta, \phi) := \mathbb{E}_{p^{\mathrm{data}}(\boldsymbol{x})}\Big[ \log p_\theta(\boldsymbol{x}) - \mathrm{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{z} \mid \boldsymbol{x}))\Big]$$

$$= \mathbb{E}_{p^{\mathrm{data}}(\boldsymbol{x})}\Bigg[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[ \log \frac{p_\theta(\boldsymbol{x}) p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \bigg]\Bigg]$$

$$= \mathbb{E}_{p^{\mathrm{data}}(\boldsymbol{x})}\Bigg[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \bigg]\Bigg] \quad \text{(computable)}$$

# Variational Inference

$$\mathscr{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{z} \mid \boldsymbol{x})) \Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[ \log \frac{p_\theta(\boldsymbol{x}) p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \right] \right]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \right] \right] \quad \text{(computable)}$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[ \log p_\theta(\boldsymbol{x} \mid \boldsymbol{z}) + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \right] \right]$$
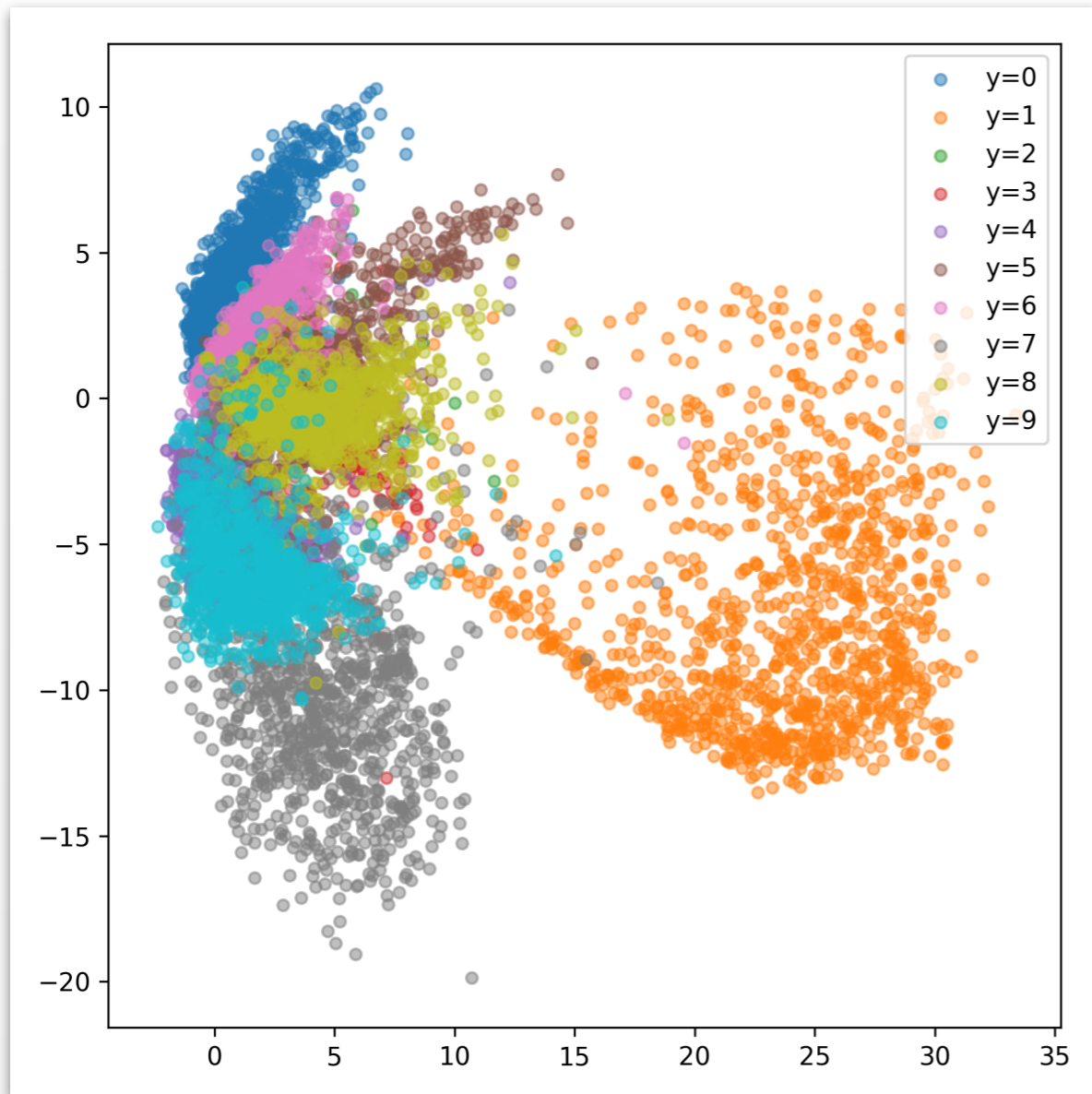
# Variational Inference

Objective: Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{z} \mid \boldsymbol{x})) \Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\Big[ \log \frac{p_\theta(\boldsymbol{x}) p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \Big] \Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\Big[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \Big] \Big] \quad \text{(computable)}$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\Big[ \log p_\theta(\boldsymbol{x} \mid \boldsymbol{z}) + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \Big] \Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z})] - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p(\boldsymbol{z})) \Big]$$

# Variational Inference

Objective: Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[\log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \| p_\theta(\boldsymbol{z} \mid \boldsymbol{x}))\Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Bigg[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[\log \frac{p_\theta(\boldsymbol{x})p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\bigg]\Bigg]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Bigg[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\bigg]\Bigg] \quad \text{(computable)}$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Bigg[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\bigg[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z}) + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\bigg]\Bigg]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Bigg[\underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z})]}_{\text{Reconstruction Error}} - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \| p(\boldsymbol{z}))\Bigg]$$

# Variational Inference

Objective: Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) := \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[\log p_\theta(\boldsymbol{x}) - \text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \| p_\theta(\boldsymbol{z} \mid \boldsymbol{x}))\Big]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x})p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\right]\right]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\right]\right] \quad \text{(computable)}$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\left[\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z}) + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}\right]\right]$$

$$= \mathbb{E}_{p^{\text{data}}(\boldsymbol{x})}\Big[\underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z})]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \| p(\boldsymbol{z}))}_{\text{KL Regularization}}\Big]$$

# Regular vs Variational Autoencoders

Autoencoder (2-dim)

VAE (2-dim)



Arbitrary scale (no "typical" values)

Well-defined scale  (-4 < z < 4)

KL regularization constrains values of latent codes

# Regular vs Variational Autoencoders

Autoencoder (50-dim TSNE)

VAE (50-dim TSNE)

# Regular vs Variational Autoencoders

### Autoencoder (50-dim TSNE)



### VAE (50-dim TSNE)



Representations are still entangled

# Unsupervised vs Semi-Supervised



**Unsupervised, Entangled**

**Semi-supervised, Disentangled**

# Unsupervised vs Semi-Supervised

**Unsupervised, Entangled**



**Semi-supervised, Disentangled**



Latent code **z** represents both style and digit

# Unsupervised vs Semi-Supervised

**Unsupervised, Entangled**

**Semi-supervised, Disentangled**



Latent code *z* represents
both style and digit

Style variable *z* is conditionally
independent from digit *y* (*)

# Unsupervised vs Semi-Supervised

**Unsupervised, Entangled**



**Semi-supervised, Disentangled**



Latent code **z** represents both style and digit

Style variable **z** is conditionally independent from digit **y** (*)

(*but not without supervision)

# Unsupervised vs Semi-Supervised

**Unsupervised, Entangled**



**Semi-supervised, Disentangled**



In both models, prior on $z$ assumes uncorrelated dimensions

$$p(\boldsymbol{z}) = \prod_d p(\boldsymbol{z}_d)$$

# Learning Statistically Independent Factors

# Deep Latent-Variable Models

### Generative Model

### Inference Model

$$p(\boldsymbol{z})$$

$$p_\theta(\boldsymbol{x})$$

$$q_\phi(\boldsymbol{z})$$

$$q(\boldsymbol{x})$$

Prior
(*known*)

Data Distribution
(*learned*)

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

# Deep Latent-Variable Models

## Generative Model

$p(z)$

$p_\theta(x)$

Prior
(*known*)

Data Distribution
(*learned*)

$$x_n \sim p^{\text{data}}(x)$$

Data Sampled from
Unknown Distribution

## Inference Model

$q_\phi(z)$

$q(x)$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

$$q(x) := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}(x)$$
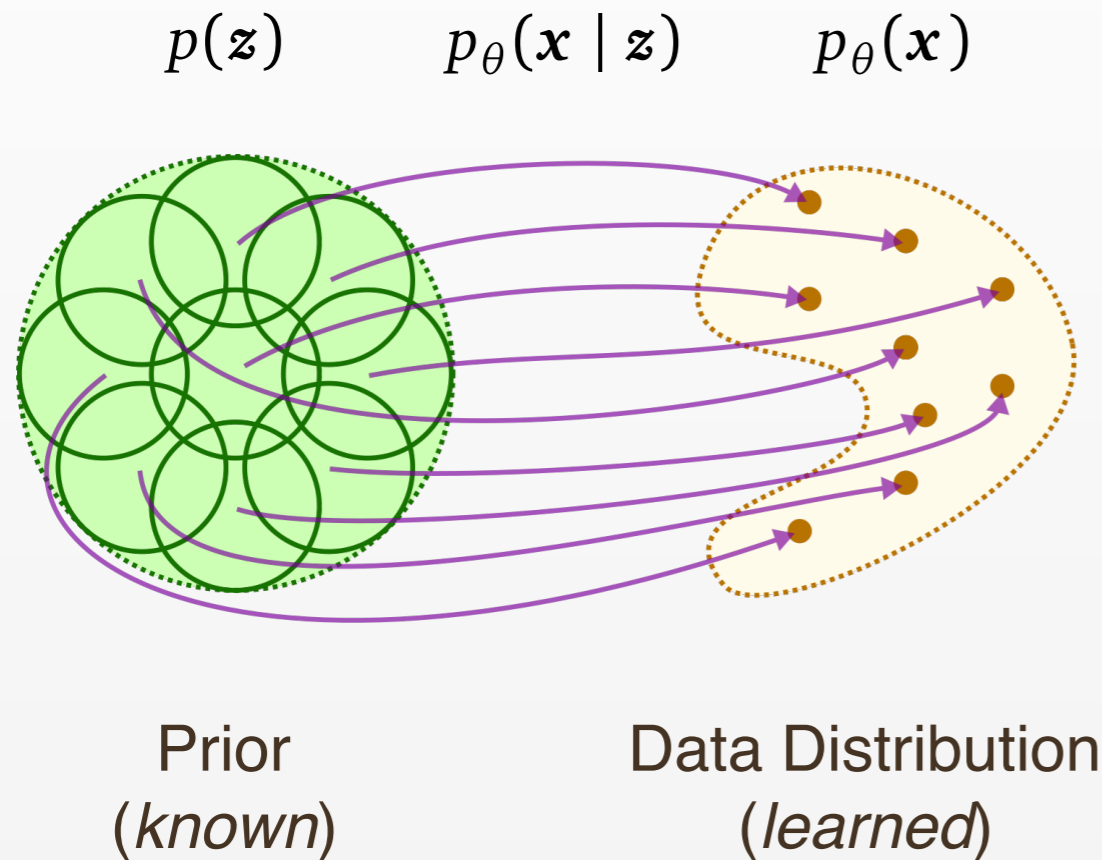
Approximation of
Data Distribution

# Deep Latent-Variable Models

## Generative Model

$p(\boldsymbol{z})$

$p_\theta(\boldsymbol{x})$

Prior
(*known*)

Data Distribution
(*learned*)

$$\boldsymbol{x}_n \sim p^{\mathrm{data}}(\boldsymbol{x})$$

Data Sampled from
Unknown Distribution

## Inference Model

$q_\phi(\boldsymbol{z})$

$q(\boldsymbol{x})$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

$$q(\boldsymbol{x}) := \frac{1}{N} \sum_{n=1}^{N} \delta_{\boldsymbol{x}_n}(\boldsymbol{x})$$

Approximation of
Data Distribution

# Deep Latent-Variable Models



## Generative Model

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

Prior
(*known*)

Data Distribution
(*learned*)

$$x_n \sim p^{\text{data}}(x)$$

Data Sampled from
Unknown Distribution

## Inference Model

$q_\phi(z)$      $q(x)$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

$$q(x) := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}(x)$$

Approximation of
Data Distribution

# Deep Latent-Variable Models

## Generative Model

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



Prior
(*known*)

Data Distribution
(*learned*)

## Inference Model

$q_\phi(z)$     $q_\phi(z \mid x)$     $q(x)$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

$$x_n \sim p^{\text{data}}(x)$$

Data Sampled from
Unknown Distribution
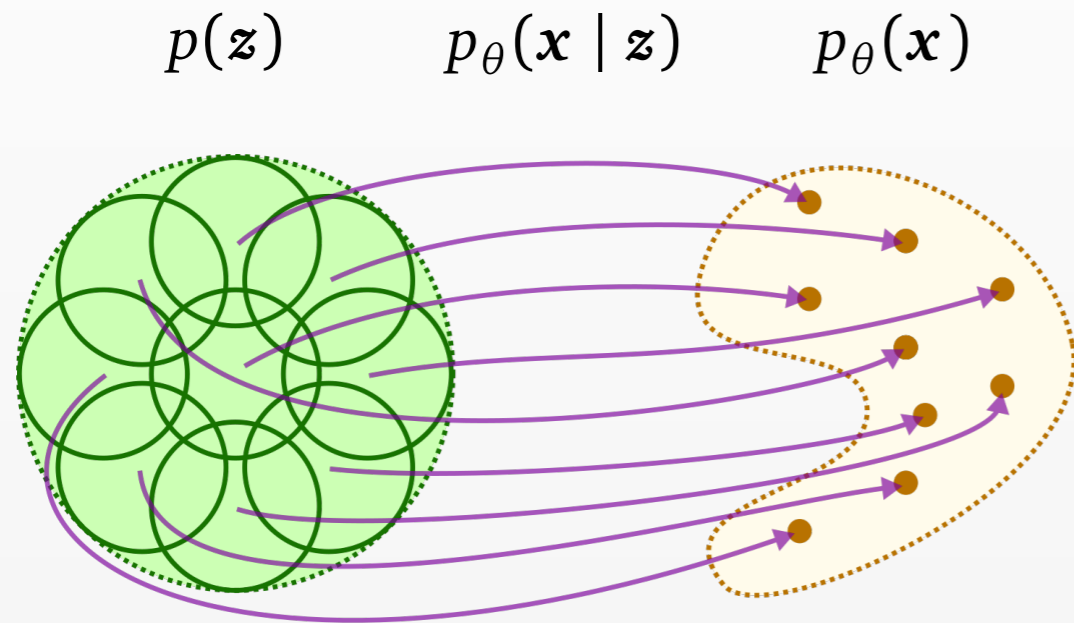
$$q(x) := \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}(x)$$

Approximation of
Data Distribution

# Deep Latent-Variable Models

## Generative Model

$p(\boldsymbol{z})$   $p_\theta(\boldsymbol{x} \mid \boldsymbol{z})$   $p_\theta(\boldsymbol{x})$



Prior
(*known*)

Data Distribution
(*learned*)

## Inference Model

$q_\phi(\boldsymbol{z})$   $q_\phi(\boldsymbol{z} \mid \boldsymbol{x})$   $q(\boldsymbol{x})$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

# Deep Latent-Variable Models

## Generative Model

$$p(z) \qquad p_\theta(x \mid z) \qquad p_\theta(x)$$



Prior
(*known*)

Data Distribution
(*learned*)

## Inference Model

$$q_\phi(z) \qquad q_\phi(z \mid x) \qquad q(x)$$

Inference Marginal
(*learned*)

Empirical Distribution
(*known*)

Constraints: Models are Equivalent when

$$p_\theta(x) = q(x) \qquad q_\phi(z) = p(z) \qquad p_\theta(z \mid x) = p_\phi(z \mid x) \; \forall \, x$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(\boldsymbol{x})} \left[ \underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \left[\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z})\right]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p(\boldsymbol{z}))}_{\text{KL Regularization}} \right]$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(x)}\left[ \underbrace{\mathbb{E}_{q_\phi(z|x)}\left[\log p_\theta(x \mid z)\right]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(z \mid x) \parallel p(z))}_{\text{KL Regularization}} \right]$$

**Alternate View:** KL Divergence Between Two Models

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(\boldsymbol{x})} \left[ \underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_\theta(\boldsymbol{x} \mid \boldsymbol{z}) \right]}_{\text{Reconstruction Error}} - \underbrace{\mathrm{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \parallel p(\boldsymbol{z}))}_{\text{KL Regularization}} \right]$$

**Alternate View:** KL Divergence Between Two Models

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(\boldsymbol{x}, \boldsymbol{z})} \left[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \right]$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathscr{L}(\theta, \phi) = \mathbb{E}_{q(x)}\left[\underbrace{\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x \mid z)]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(z \mid x) \| p(z))}_{\text{KL Regularization}}\right]$$

**Alternate View:** KL Divergence Between Two Models

$$\mathscr{L}(\theta, \phi) = \mathbb{E}_{q_\phi(x,z)}\left[\log \frac{p_\theta(x, z)}{q_\phi(z \mid x)}\right] = \mathbb{E}_{q_\phi(x,z)}\left[\log \frac{p_\theta(x, z)}{q_\phi(z \mid x)} + \log \frac{q(x)}{q(x)}\right]$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(x)}\left[ \underbrace{\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x \mid z)]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(z \mid x) \| p(z))}_{\text{KL Regularization}} \right]$$

**Alternate View:** KL Divergence Between Two Models

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x, z)}{q_\phi(z \mid x)} \right] = \mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x, z)}{q_\phi(z \mid x)} + \log \frac{q(x)}{q(x)} \right]$$

$$= \mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x, z)}{q_\phi(x, z)} \right] + \mathbb{E}_{q(x)}[\log q(x)]$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(x)}\left[ \underbrace{\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x \mid z)]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(z \mid x) \parallel p(z))}_{\text{KL Regularization}} \right]$$

**Alternate View:** KL Divergence Between Two Models

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x,z)}{q_\phi(z \mid x)} \right] = \mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x,z)}{q_\phi(z \mid x)} + \log \frac{q(x)}{q(x)} \right]$$

$$= \underbrace{\mathbb{E}_{q_\phi(x,z)}\left[ \log \frac{p_\theta(x,z)}{q_\phi(x,z)} \right]}_{-\text{KL}\left(q_\phi(x,z) \parallel p_\theta(x,z)\right)} + \mathbb{E}_{q(x)}[\log q(x)]$$

# Implicit Tradeoffs in the Variational Objective

**Classical View:** Maximize Evidence Lower Bound

$$\mathscr{L}(\theta, \phi) = \mathbb{E}_{q(\boldsymbol{x})} \left[ \underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})} [\log p_\theta(\boldsymbol{x} \mid \boldsymbol{z})]}_{\text{Reconstruction Error}} - \underbrace{\text{KL}(q_\phi(\boldsymbol{z} \mid \boldsymbol{x}) \parallel p(\boldsymbol{z}))}_{\text{KL Regularization}} \right]$$

**Alternate View:** KL Divergence Between Two Models

$$\mathscr{L}(\theta, \phi) = \mathbb{E}_{q_\phi(\boldsymbol{x},\boldsymbol{z})} \left[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} \right] = \mathbb{E}_{q_\phi(\boldsymbol{x},\boldsymbol{z})} \left[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})} + \log \frac{q(\boldsymbol{x})}{q(\boldsymbol{x})} \right]$$

$$= \underbrace{\mathbb{E}_{q_\phi(\boldsymbol{x},\boldsymbol{z})} \left[ \log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{x}, \boldsymbol{z})} \right]}_{-\text{KL}(q_\phi(\boldsymbol{x},\boldsymbol{z}) \parallel p_\theta(\boldsymbol{x},\boldsymbol{z}))} + \underbrace{\mathbb{E}_{q(\boldsymbol{x})} [\log q(\boldsymbol{x})]}_{H_q(\boldsymbol{x}) = \log N}$$

# Implicit Tradeoffs in the Variational Objective

$$\mathscr{L}(\theta, \phi) := -\mathrm{KL}\big(q_\phi(\boldsymbol{z}, \boldsymbol{x}) \,\big\|\, p_\theta(\boldsymbol{x}, \boldsymbol{z})\big) = \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\right]$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Implicit Tradeoffs in the Variational Objective

$$\mathscr{L}(\theta, \phi) := -\mathrm{KL}\big(q_\phi(\boldsymbol{z}, \boldsymbol{x}) \,\big\|\, p_\theta(\boldsymbol{x}, \boldsymbol{z})\big) = \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{p_\theta(\boldsymbol{x})p(\boldsymbol{z})} + \log \frac{q_\phi(\boldsymbol{z})q(\boldsymbol{x})}{q_\phi(\boldsymbol{z}, \boldsymbol{x})} + \log \frac{p_\theta(\boldsymbol{x})}{q(\boldsymbol{x})} + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z})}\right]$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Implicit Tradeoffs in the Variational Objective

$$\mathcal{L}(\theta, \phi) := -\text{KL}\big(q_\phi(\boldsymbol{z}, \boldsymbol{x}) \,\big|\big|\, p_\theta(\boldsymbol{x}, \boldsymbol{z})\big) = \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\left[\log \frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{p_\theta(\boldsymbol{x})p(\boldsymbol{z})} + \log \frac{q_\phi(\boldsymbol{z})q(\boldsymbol{x})}{q_\phi(\boldsymbol{z}, \boldsymbol{x})} + \log \frac{p_\theta(\boldsymbol{x})}{q(\boldsymbol{x})} + \log \frac{p(\boldsymbol{z})}{q_\phi(\boldsymbol{z})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}, \boldsymbol{x})}\bigg[\underbrace{\log \frac{p_\theta(\boldsymbol{z} \,|\, \boldsymbol{x})}{p(\boldsymbol{z})}}_{\text{①}} - \underbrace{\log \frac{q_\phi(\boldsymbol{z} \,|\, \boldsymbol{x})}{q_\phi(\boldsymbol{z})}}_{\text{②}}\bigg] - \underbrace{\text{KL}\big(q(\boldsymbol{x}) \,||\, p_\theta(\boldsymbol{x})\big)}_{\text{③}} - \underbrace{\text{KL}\big(q_\phi(\boldsymbol{z}) \,\big|\big|\, p(\boldsymbol{z})\big)}_{\text{④}}$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Implicit Tradeoffs in the Variational Objective

$$\mathscr{L}(\theta, \phi) := -\mathrm{KL}\big(q_\phi(z,x) \,\big\|\, p_\theta(x,z)\big) = \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x,z)}{q_\phi(z,x)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x,z)}{p_\theta(x)p(z)} + \log \frac{q_\phi(z)q(x)}{q_\phi(z,x)} + \log \frac{p_\theta(x)}{q(x)} + \log \frac{p(z)}{q_\phi(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z\mid x)}{p(z)}}_{①} - \underbrace{\log \frac{q_\phi(z\mid x)}{q_\phi(z)}}_{②}\right] - \underbrace{\mathrm{KL}\big(q(x)\,\|\,p_\theta(x)\big)}_{③} - \underbrace{\mathrm{KL}\big(q_\phi(z)\,\big\|\,p(z)\big)}_{④}$$

$$p_\theta(x) = q(x)$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Implicit Tradeoffs in the Variational Objective

$$\mathcal{L}(\theta, \phi) := -\mathrm{KL}\big(q_\phi(z, x) \,\big\|\, p_\theta(x, z)\big) = \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x, z)}{q_\phi(z, x)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x, z)}{p_\theta(x)p(z)} + \log \frac{q_\phi(z)q(x)}{q_\phi(z, x)} + \log \frac{p_\theta(x)}{q(x)} + \log \frac{p(z)}{q_\phi(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z \mid x)}{p(z)}}_{\textcircled{1}} - \underbrace{\log \frac{q_\phi(z \mid x)}{q_\phi(z)}}_{\textcircled{2}}\right] - \underbrace{\mathrm{KL}\big(q(x) \,\|\, p_\theta(x)\big)}_{\textcircled{3}} - \underbrace{\mathrm{KL}\big(q_\phi(z) \,\big\|\, p(z)\big)}_{\textcircled{4}}$$

$$p_\theta(x) = q(x) \qquad q_\phi(z) = p(z)$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Implicit Tradeoffs in the Variational Objective

$$\mathcal{L}(\theta, \phi) := -\mathrm{KL}\big(q_\phi(z, x) \,\big\|\, p_\theta(x, z)\big) = \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x, z)}{q_\phi(z, x)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\log \frac{p_\theta(x, z)}{p_\theta(x)p(z)} + \log \frac{q_\phi(z)q(x)}{q_\phi(z, x)} + \log \frac{p_\theta(x)}{q(x)} + \log \frac{p(z)}{q_\phi(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z \mid x)}{p(z)}}_{①} - \underbrace{\log \frac{q_\phi(z \mid x)}{q_\phi(z)}}_{②}\right] - \underbrace{\mathrm{KL}\big(q(x) \,\|\, p_\theta(x)\big)}_{③} - \underbrace{\mathrm{KL}\big(q_\phi(z) \,\big\|\, p(z)\big)}_{④}$$

$$p_\theta(z \mid x) = p_\phi(z \mid x) \; \forall \, x \qquad\qquad p_\theta(x) = q(x) \qquad\qquad q_\phi(z) = p(z)$$

[Hoffman, Johnson, NIPS AABI Workshop 2016], [Chen, Li, Grosse, Duvenaud, Arxiv 2018]

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Relaxing Constraints in the Objective

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z,x)} \left[ \underbrace{\log \frac{p_\theta(z \mid x)}{p(z)}}_{\text{①}} - \underbrace{\log \frac{q_\phi(z \mid x)}{q_\phi(z)}}_{\text{②}} \right] - \underbrace{\text{KL}(q(x) \,\|\, p_\theta(x))}_{\text{③}} - \underbrace{\text{KL}(q_\phi(z) \,\|\, p(z))}_{\text{④}}$$

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]
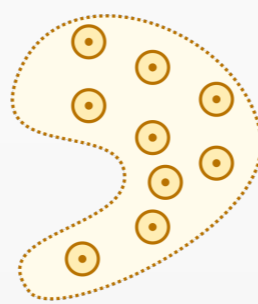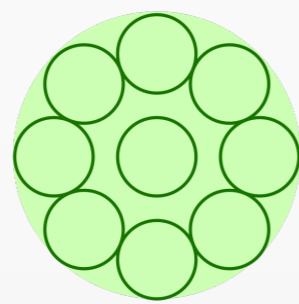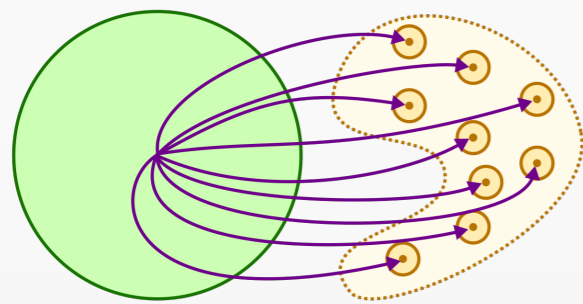
# Relaxing Constraints in the Objective

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z \mid x)}{p(z)}}_{①} - \underbrace{\log \frac{q_\phi(z \mid x)}{q_\phi(z)}}_{②}\right] - \underbrace{\mathrm{KL}(q(x) \| p_\theta(x))}_{③} - \underbrace{\mathrm{KL}(q_\phi(z) \| p(z))}_{④}$$



② + ③ + ④

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Relaxing Constraints in the Objective

$$\mathscr{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z,x)}\left[ \log \frac{p_\theta(z \mid x)}{p(z)} - \log \frac{q_\phi(z \mid x)}{q_\phi(z)} \right] - \mathrm{KL}\big(q(x) \,\|\, p_\theta(x)\big) - \mathrm{KL}\big(q_\phi(z) \,\|\, p(z)\big)$$

①       ②       ③       ④



② + ③ + ④       ① + ③ + ④

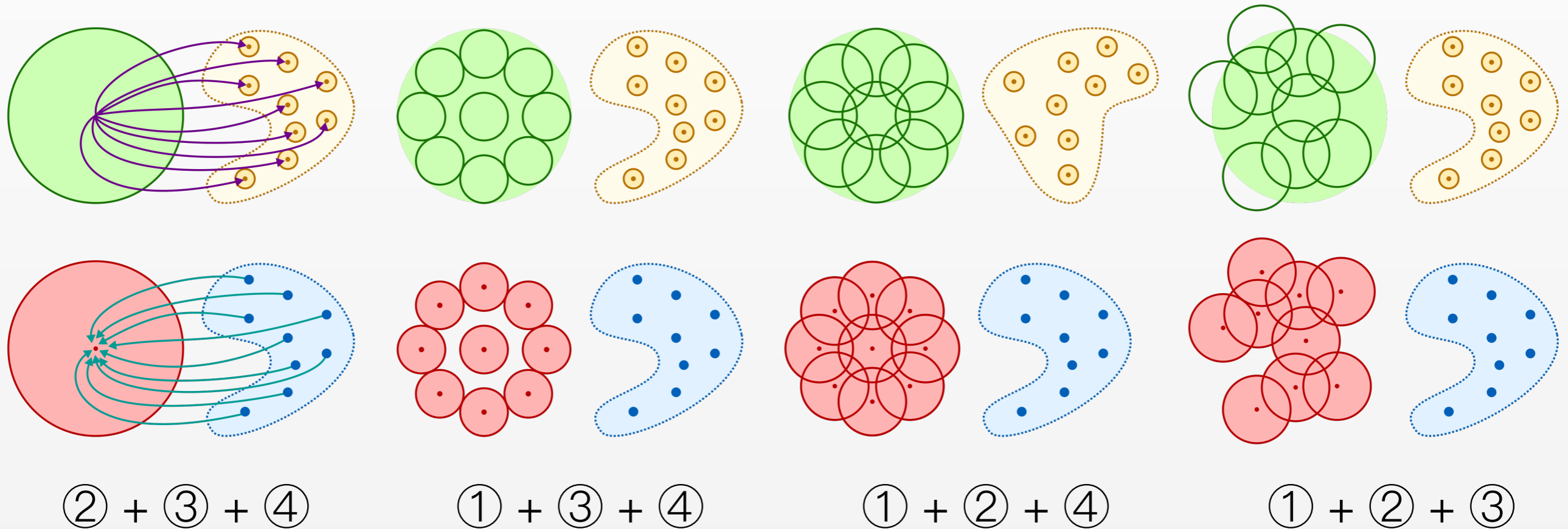[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Relaxing Constraints in the Objective

$$\mathscr{L}(\theta,\phi) = \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z\mid x)}{p(z)}}_{①} - \underbrace{\log \frac{q_\phi(z\mid x)}{q_\phi(z)}}_{②}\right] - \underbrace{\mathrm{KL}\big(q(x)\,\|\,p_\theta(x)\big)}_{③} - \underbrace{\mathrm{KL}\big(q_\phi(z)\,\|\,p(z)\big)}_{④}$$



② + ③ + ④          ① + ③ + ④          ① + ② + ④

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Relaxing Constraints in the Objective

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(\boldsymbol{z},\boldsymbol{x})}\left[ \underbrace{\log \frac{p_\theta(\boldsymbol{z} \mid \boldsymbol{x})}{p(\boldsymbol{z})}}_{\text{①}} - \underbrace{\log \frac{q_\phi(\boldsymbol{z} \mid \boldsymbol{x})}{q_\phi(\boldsymbol{z})}}_{\text{②}} \right] - \underbrace{\mathrm{KL}\big(q(\boldsymbol{x}) \,\|\, p_\theta(\boldsymbol{x})\big)}_{\text{③}} - \underbrace{\mathrm{KL}\big(q_\phi(\boldsymbol{z}) \,\|\, p(\boldsymbol{z})\big)}_{\text{④}}$$



② + ③ + ④          ① + ③ + ④          ① + ② + ④          ① + ② + ③

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Relaxing Constraints in the Objective

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z,x)}\left[\underbrace{\log \frac{p_\theta(z \mid x)}{p(z)}}_{①} - \underbrace{\log \frac{q_\phi(z \mid x)}{q_\phi(z)}}_{②}\right] - \underbrace{\mathrm{KL}(q(x) \| p_\theta(x))}_{③} - \underbrace{\mathrm{KL}(q_\phi(z) \| p(z))}_{④}$$



② + ③ + ④      ① + ③ + ④      ① + ② + ④      ① + ② + ③

Idea: Relax ①, ②, and ③ in favor of ④

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\text{KL}\big(q_\phi(\boldsymbol{z})\,\big|\big|\,p(\boldsymbol{z})\big) = -\mathbb{E}_{q_\phi(\boldsymbol{z})}\left[\log \frac{q_\phi(\boldsymbol{z})}{\prod_d q_\phi(\boldsymbol{z}_d)} + \log \frac{\prod_d q_\phi(\boldsymbol{z}_d)}{\prod_d p(\boldsymbol{z}_d)} + \log \frac{\prod_d p(\boldsymbol{z}_d)}{p(\boldsymbol{z})}\right]$$

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\mathrm{KL}\big(q_\phi(z)\,\big\|\,p(z)\big) = -\mathbb{E}_{q_\phi(z)}\left[\log\frac{q_\phi(z)}{\prod_d q_\phi(z_d)} + \log\frac{\prod_d q_\phi(z_d)}{\prod_d p(z_d)} + \log\frac{\prod_d p(z_d)}{p(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z)}\underbrace{\left[\log\frac{p(z)}{\prod_d p(z_d)} - \log\frac{q_\phi(z)}{\prod_d q_\phi(z_d)}\right]}_{\text{Ⓐ}} - \sum_d \underbrace{\mathrm{KL}\big(q_\phi(z_d)\,\big\|\,p(z_d)\big)}_{\text{Ⓑ}}.$$

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\mathrm{KL}\big(q_\phi(\boldsymbol{z})\,\big\|\,p(\boldsymbol{z})\big) = -\mathbb{E}_{q_\phi(\boldsymbol{z})}\left[\log\frac{q_\phi(\boldsymbol{z})}{\prod_d q_\phi(\boldsymbol{z}_d)} + \log\frac{\prod_d q_\phi(\boldsymbol{z}_d)}{\prod_d p(\boldsymbol{z}_d)} + \log\frac{\prod_d p(\boldsymbol{z}_d)}{p(\boldsymbol{z})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z})}\underbrace{\left[\log\frac{p(\boldsymbol{z})}{\prod_d p(\boldsymbol{z}_d)} - \log\frac{q_\phi(\boldsymbol{z})}{\prod_d q_\phi(\boldsymbol{z}_d)}\right]}_{\text{\textcircled{A}}} - \sum_d \underbrace{\mathrm{KL}\big(q_\phi(\boldsymbol{z}_d)\,\big\|\,p(\boldsymbol{z}_d)\big)}_{\text{\textcircled{B}}}.$$

Marginals should
be identical

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\mathrm{KL}\big(q_\phi(z)\,\big\|\,p(z)\big) = -\mathbb{E}_{q_\phi(z)}\left[\log\frac{q_\phi(z)}{\prod_d q_\phi(z_d)} + \log\frac{\prod_d q_\phi(z_d)}{\prod_d p(z_d)} + \log\frac{\prod_d p(z_d)}{p(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z)}\underbrace{\left[\log\frac{p(z)}{\prod_d p(z_d)} - \log\frac{q_\phi(z)}{\prod_d q_\phi(z_d)}\right]}_{\text{Ⓐ}} - \sum_d \underbrace{\mathrm{KL}\big(q_\phi(z_d)\,\big\|\,p(z_d)\big)}_{\text{Ⓑ}}.$$

Ⓐ Correlations between variables
should be identical

Ⓑ Marginals should
be identical

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\mathrm{KL}\big(q_\phi(\boldsymbol{z})\,\big\|\,p(\boldsymbol{z})\big) = -\mathbb{E}_{q_\phi(\boldsymbol{z})}\left[\log\frac{q_\phi(\boldsymbol{z})}{\prod_d q_\phi(\boldsymbol{z}_d)} + \log\frac{\prod_d q_\phi(\boldsymbol{z}_d)}{\prod_d p(\boldsymbol{z}_d)} + \log\frac{\prod_d p(\boldsymbol{z}_d)}{p(\boldsymbol{z})}\right]$$

$$= \ \mathbb{E}_{q_\phi(\boldsymbol{z})}\underbrace{\left[\log\frac{p(\boldsymbol{z})}{\prod_d p(\boldsymbol{z}_d)} - \log\frac{q_\phi(\boldsymbol{z})}{\prod_d q_\phi(\boldsymbol{z}_d)}\right]}_{\text{\textcircled{A}}} - \sum_d \underbrace{\mathrm{KL}\big(q_\phi(\boldsymbol{z}_d)\,\big\|\,p(\boldsymbol{z}_d)\big)}_{\text{\textcircled{B}}}.$$

Correlations between variables
should be identical

Marginals should
be identical

Disentanglement: $\quad p(\boldsymbol{z}) = \prod_d p(\boldsymbol{z}_d) \quad q_\phi(\boldsymbol{z}) = \prod_d q_\phi(\boldsymbol{z}_d)$

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Hierarchical Decomposition

$$-\mathrm{KL}\big(q_\phi(z) \,\big|\big|\, p(z)\big) = -\mathbb{E}_{q_\phi(z)}\left[\log \frac{q_\phi(z)}{\prod_d q_\phi(z_d)} + \log \frac{\prod_d q_\phi(z_d)}{\prod_d p(z_d)} + \log \frac{\prod_d p(z_d)}{p(z)}\right]$$

$$= \mathbb{E}_{q_\phi(z)}\left[\underbrace{\log \frac{p(z)}{\prod_d p(z_d)} - \log \frac{q_\phi(z)}{\prod_d q_\phi(z_d)}}_{\text{Ⓐ}}\right] - \sum_d \underbrace{\mathrm{KL}\big(q_\phi(z_d) \,\big|\big|\, p(z_d)\big)}_{\text{Ⓑ}}.$$

Correlations between variables should be identical      Marginals should be identical

Disentanglement: $\quad p(z) = \prod_d p(z_d) \quad q_\phi(z) = \prod_d q_\phi(z_d)$

$$\text{Ⓑ} = \mathbb{E}_{q_\phi(z_d)}\left[\underbrace{\log \frac{p(z_d)}{\prod_e p(z_{d,e})} - \log \frac{q_\phi(z_d)}{\prod_e q_\phi(z_{d,e})}}_{\text{ⓘ}}\right] - \sum_e \underbrace{\mathrm{KL}\big(q_\phi(z_{d,e}) \,\big|\big|\, p(z_{d,e})\big)}_{\text{ⓘⓘ}}$$

(Can continue decomposition for any number of levels)

[Esmaeli, Wu, Jain, Bozkurt, Siddharth, Paige, Brooks, Dy, van de Meent, Arxiv 2018]

# Generalizations of VAE objectives

Hierarchically Factorized Variational Autoencoders

$$\mathcal{L}(\theta, \phi) = \text{①} + \text{③} + \text{ⓘⓘ} + \alpha\,\text{②} + \beta\,\text{Ⓐ} + \gamma\,\text{ⓘ}$$

# Generalizations of VAE objectives

Hierarchically Factorized Variational Autoencoders

$$\mathcal{L}(\theta, \phi) = ① + ③ + ⅱ + \alpha② + \beta Ⓐ + \gamma ①$$

| Paper | Objective |
|---|---|
| Kingma and Welling [2013], Rezende et al. [2014] | $① + ② + ③ + ④$ |
| Higgins et al. [2016] | $① + ③ + \beta(② + ④)$ |
| Kumar et al. [2017] | $① + ② + ③ + \lambda④$ |
| Zhao et al. [2017] | $① + ③ + \lambda④$ |
| Gao et al. [2018] | $① + ② + ③ + ④ - \lambda②^{a}$ |
| Achille and Soatto [2018] | $① + ③ + \beta② + \gamma Ⓐ^{*}$ |
| Kim and Mnih [2018],Chen et al. [2018] | $① + ② + ③ + Ⓑ + \beta Ⓐ^{*}$ |
| HFVAE (this paper) | $① + ③ + ⅱ + \alpha② + \beta Ⓐ + \gamma ①$ |

# Generalizations of VAE objectives

Hierarchically Factorized Variational Autoencoders

$$\mathcal{L}(\theta, \phi) = ① + ③ + ⓘ\mathbf{i} + \alpha② + \beta Ⓐ + \gamma ⓘ$$

| Paper | Objective |
|---|---|
| Kingma and Welling [2013], Rezende et al. [2014] | $① + ② + ③ + ④$ |
| Higgins et al. [2016] | $① + ③ + \beta(② + ④)$ |
| Kumar et al. [2017] | $① + ② + ③ + \lambda④$ |
| Zhao et al. [2017] | $① + ③ + \lambda④$ |
| Gao et al. [2018] | $① + ② + ③ + ④ - \lambda②^{a}$ |
| Achille and Soatto [2018] | $① + ③ + \beta② + \gamma Ⓐ^{*}$ |
| Kim and Mnih [2018],Chen et al. [2018] | $① + ② + ③ + Ⓑ + \beta Ⓐ^{*}$ |
| HFVAE (this paper) | $① + ③ + ⓘ\mathbf{i} + \alpha② + \beta Ⓐ + \gamma ⓘ$ |

# Generalizations of VAE objectives

Hierarchically Factorized Variational Autoencoders

$$\mathcal{L}(\theta, \phi) = ① + ③ + ⅱ + \alpha② + \beta Ⓐ + \gamma ⓘ$$

| Paper | Objective |
|---|---|
| Kingma and Welling [2013], Rezende et al. [2014] | $① + ② + ③ + ④$ |
| Higgins et al. [2016] | $① + ③ + \beta(② + ④)$ |
| Kumar et al. [2017] | $① + ② + ③ + \lambda④$ |
| Zhao et al. [2017] | $① + ③ + \lambda④$ |
| Gao et al. [2018] | $① + ② + ③ + ④ - \lambda②^{a}$ |
| Achille and Soatto [2018] | $① + ③ + \beta② + \gamma Ⓐ^{*}$ |
| Kim and Mnih [2018],Chen et al. [2018] | $① + ② + ③ + Ⓑ + \beta Ⓐ^{*}$ |
| HFVAE (this paper) | $① + ③ + ⅱ + \alpha② + \beta Ⓐ + \gamma ⓘ$ |

# Results: CelebA



HFVAE and β-VAE are qualitatively similar
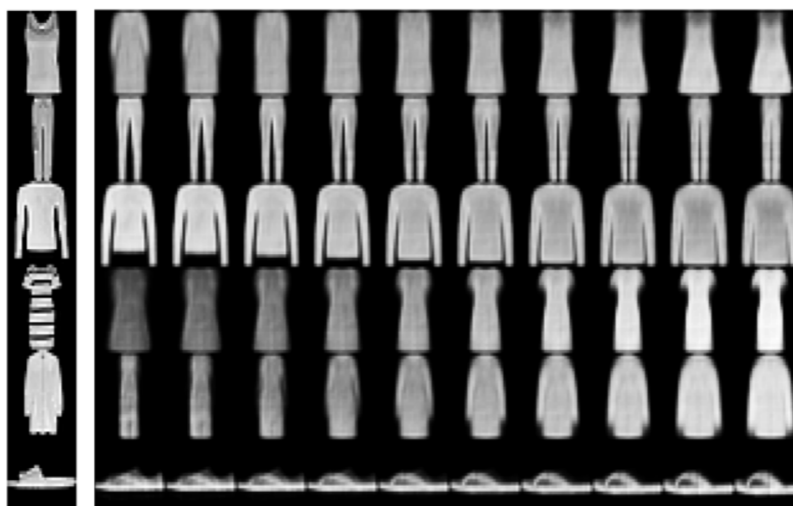
Both objectives learn (some) interpretable features

# Results: MNIST and FMNIST



MNIST HFVAE (β=12, γ=4)
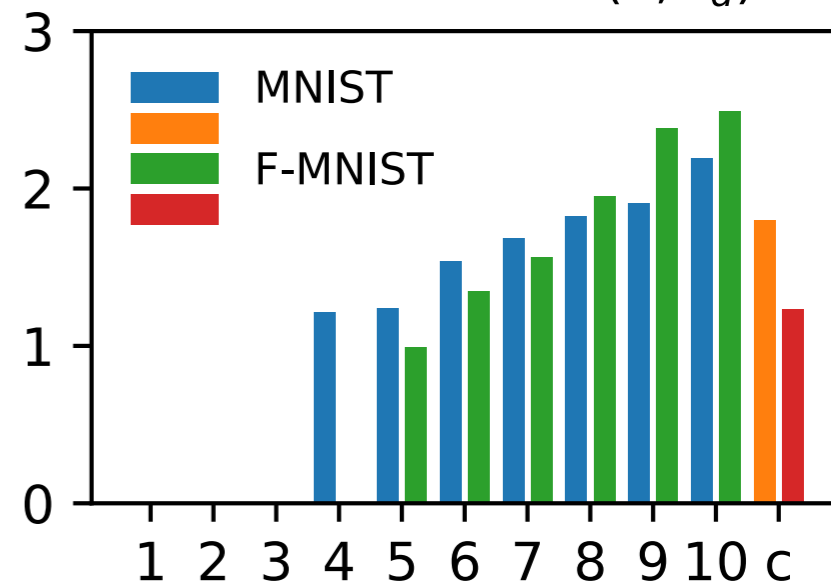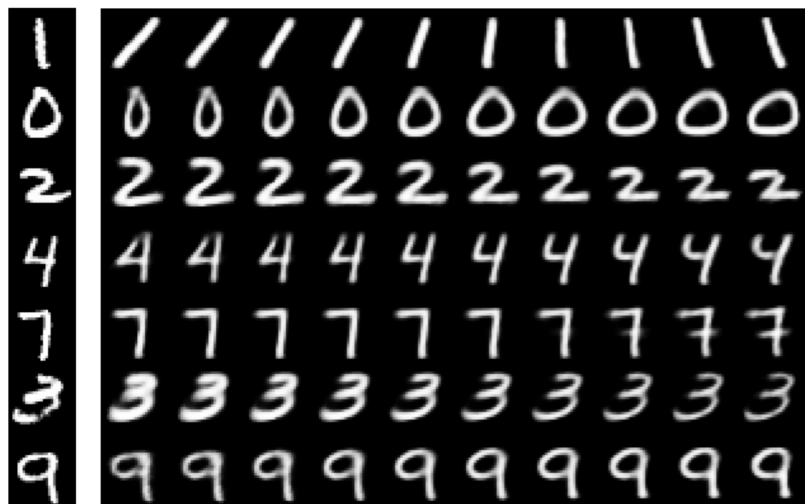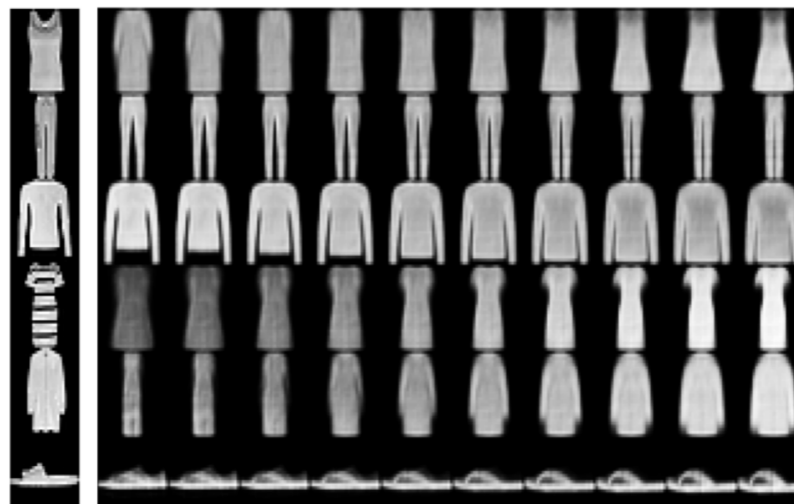
FMNIST HFVAE (β=12, γ=4)

MNIST F-MNIST $I(x; z_d)$

# Results: MNIST and FMNIST



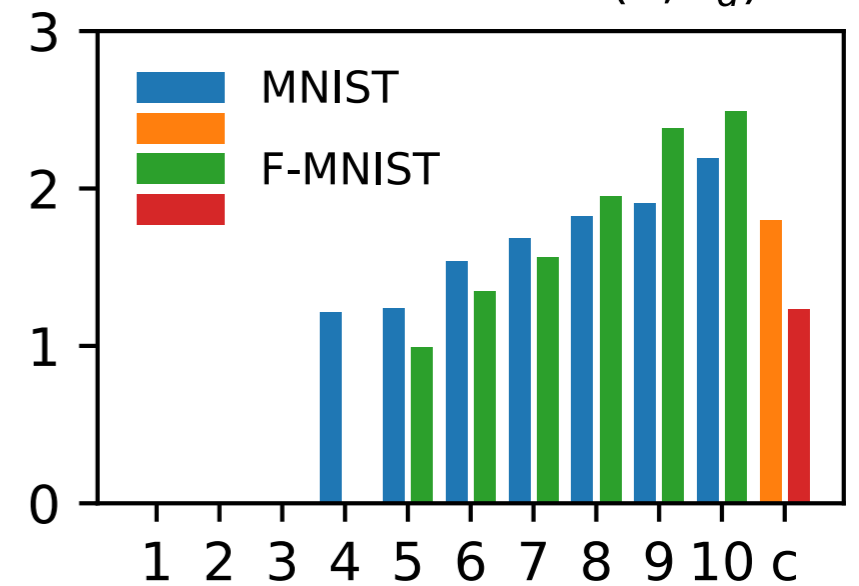MNIST HFVAE (β=12, γ=4)

FMNIST HFVAE (β=12, γ=4)

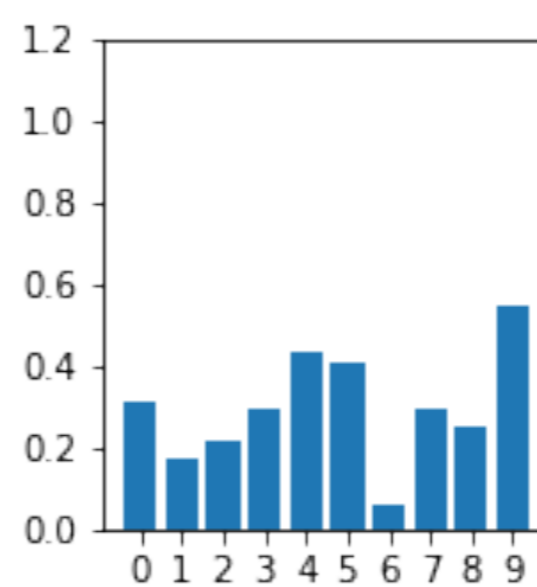MNIST F-MNIST $I(x; z_d)$

Input

$\beta$-VAE ($\beta = 4$)

HFVAE ($\beta = 12, \gamma = 4$)

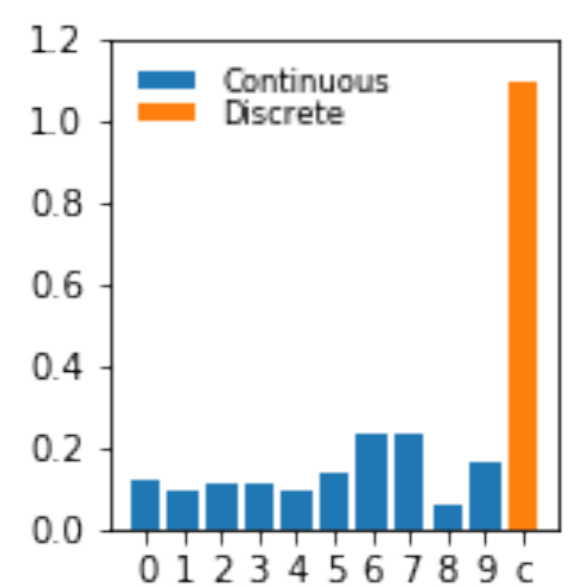$I(y; z)$ ($\beta$-VAE, $\beta = 4$)

$I(y; z)$ (HFVAE, $\beta = 12, \gamma = 4$)

# Results: Generalization

# Bonus: Learning Correlations Between Topics

## VAE (50 topics)



## HFVAE (25 + 25 topics)



| Top Words | NPMI |
|---|---|
| jesus, scripture, god, christ, bible, sin, christian, doctrine, faith, church | 0.41 |
| team, game, player, score, league, play, leafs, hockey, season, nhl | 0.37 |
| scsi, ide, controller, drive, bus, subject, lines, organization, card, problem | 0.14 |

| Top 3 Most Correlated Topics | NPMI |
|---|---|
| armenian, turk, civilian, turkish, soldier, kill, extermination, armenia, israel, jew | 0.40 |
| jesus, belief, god, christian, christ, faith, scripture, moral, truth, sin | 0.35 |
| gun, crime, people, law, defense, government, criminal, shoot, assault, fire | 0.26 |

# Deep Learning + Probabilistic Programming

Deep
Learning

Probabilistic
Programming

# Deep Learning + Probabilistic Programming

Deep
Learning

Probabilistic
Programming

Discriminative
(Data → Features)

Generative
(Variables → Data)

# Deep Learning + Probabilistic Programming

Deep
Learning

Probabilistic
Programming

Discriminative
(Data → Features)

Generative
(Variables → Data)

Large Data,
Low Uncertainty

Small Data,
High Uncertainty

# Deep Learning + Probabilistic Programming

Deep
Learning

Probabilistic
Programming

Discriminative
(Data → Features)

Generative
(Variables → Data)

Large Data,
Low Uncertainty

Small Data,
High Uncertainty

Stochastic Gradient
Descent

Many Inference
Methods

# Deep Learning + Probabilistic Programming

Deep
Learning

Probabilistic
Programming

| Discriminative (Data → Features) | | Generative (Variables → Data) |
|---|---|---|
| Large Data, Low Uncertainty | | Small Data, High Uncertainty |
| Stochastic Gradient Descent | | Many Inference Methods |

Representation
Learning

Model-based
Reasoning

# Deep Learning + Probabilistic Programming

| Deep Learning | Integrated Approaches | Probabilistic Programming |
|---|---|---|
| Discriminative (Data → Features) | | Generative (Variables → Data) |
| Large Data, Low Uncertainty | | Small Data, High Uncertainty |
| Stochastic Gradient Descent | | Many Inference Methods |
| Representation Learning | | Model-based Reasoning |

# Deep Learning + Probabilistic Programming

| Deep Learning | Integrated Approaches | Probabilistic Programming |
|---|---|---|
| Discriminative (Data → Features) | | Generative (Variables → Data) |
| Large Data, Low Uncertainty | Decision Making | Small Data, High Uncertainty |
| Stochastic Gradient Descent | | Many Inference Methods |
| Representation Learning | | Model-based Reasoning |

# Deep Learning + Probabilistic Programming

| Deep Learning | Integrated Approaches | Probabilistic Programming |
|---|---|---|
| Discriminative (Data → Features) | Autoencoders | Generative (Variables → Data) |
| Large Data, Low Uncertainty | Decision Making | Small Data, High Uncertainty |
| Stochastic Gradient Descent | | Many Inference Methods |
| Representation Learning | | Model-based Reasoning |

# Deep Learning + Probabilistic Programming

| Deep Learning | Integrated Approaches | Probabilistic Programming |
|---|---|---|
| Discriminative (Data → Features) | Autoencoders | Generative (Variables → Data) |
| Large Data, Low Uncertainty | Decision Making | Small Data, High Uncertainty |
| Stochastic Gradient Descent | Variational Inference | Many Inference Methods |

Representation Learning

Model-based Reasoning

# Thank You!

Northeastern University          Oxford          ATI

Babak Esmaeli    Hao Wu    Sarthak Jain    Alican Bozkurt    Siddharth N.    Brooks Paige

## Papers

*Structured Disentangled Representations*
B. Esmaeili, H. Wu, S. Jain, A. Bozkurt, N. Siddharth, B. Paige, D. H. Brooks, J. Dy, J.-W. van de Meent
ArXiv [https://arxiv.org/abs/1804.02086]

*Learning disentangled representations with semi-supervised deep generative models*
N. Siddharth, B. Paige, J.-W. van de Meent, A. Desmaison, F. Wood, N.D. Goodman, P. Kohli, P.H.S. Torr
NIPS 2017 [https://bit.ly/probtorch-nips-2017]

## Code

https://github.com/probtorch/probtorch