

Rotating Features for Object Discovery

Sindy Löwe, Phillip Lippe, Francesco Locatello, Max Welling

NeurIPS 2023 Oral

PRESENTER: MINGHAO LIU

2023/12/3

Outline

1 / **Background**

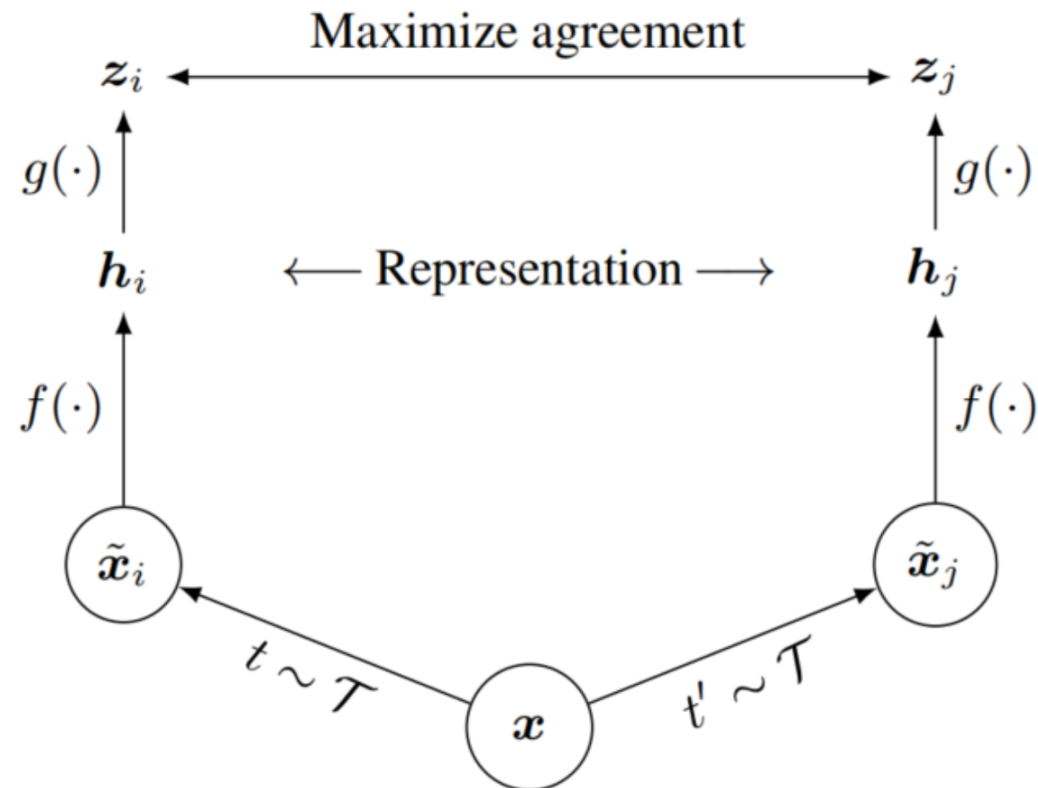
2 / **Method**

3 / **Experiments**

Background

A Simple Framework for Contrastive Learning of Visual Representations

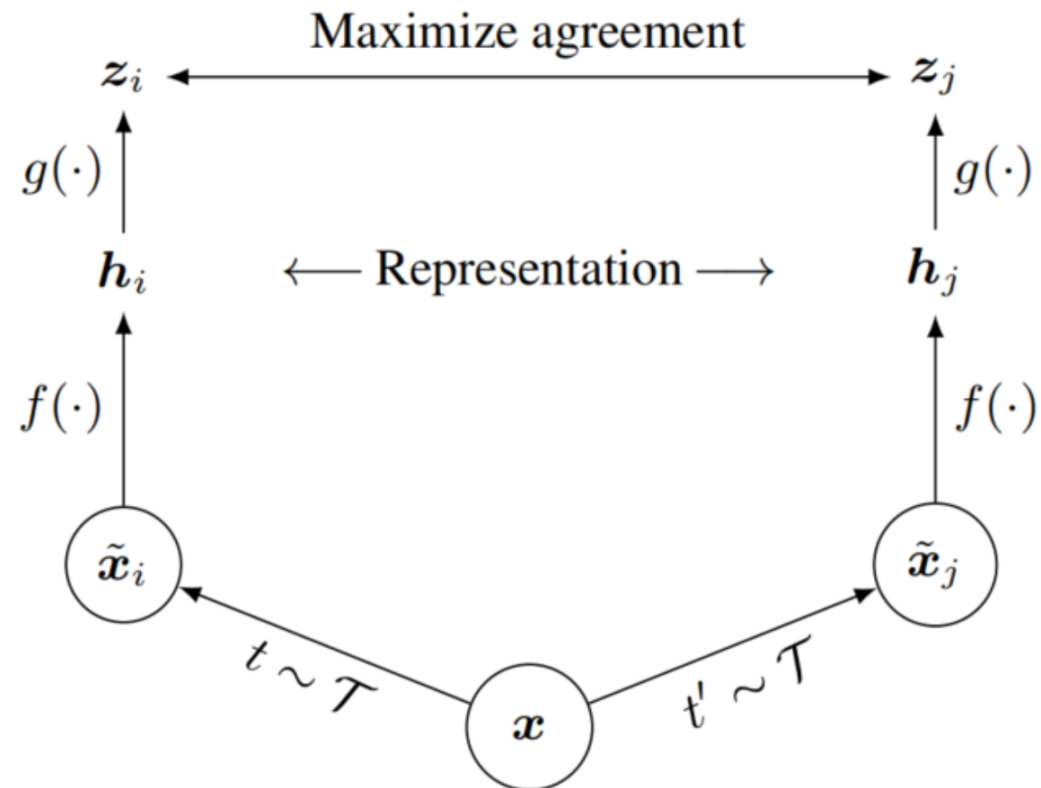
■ SimCLR (ICML 2020) Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹



Background

■ SimCLR

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}, \quad (1)$$



Background

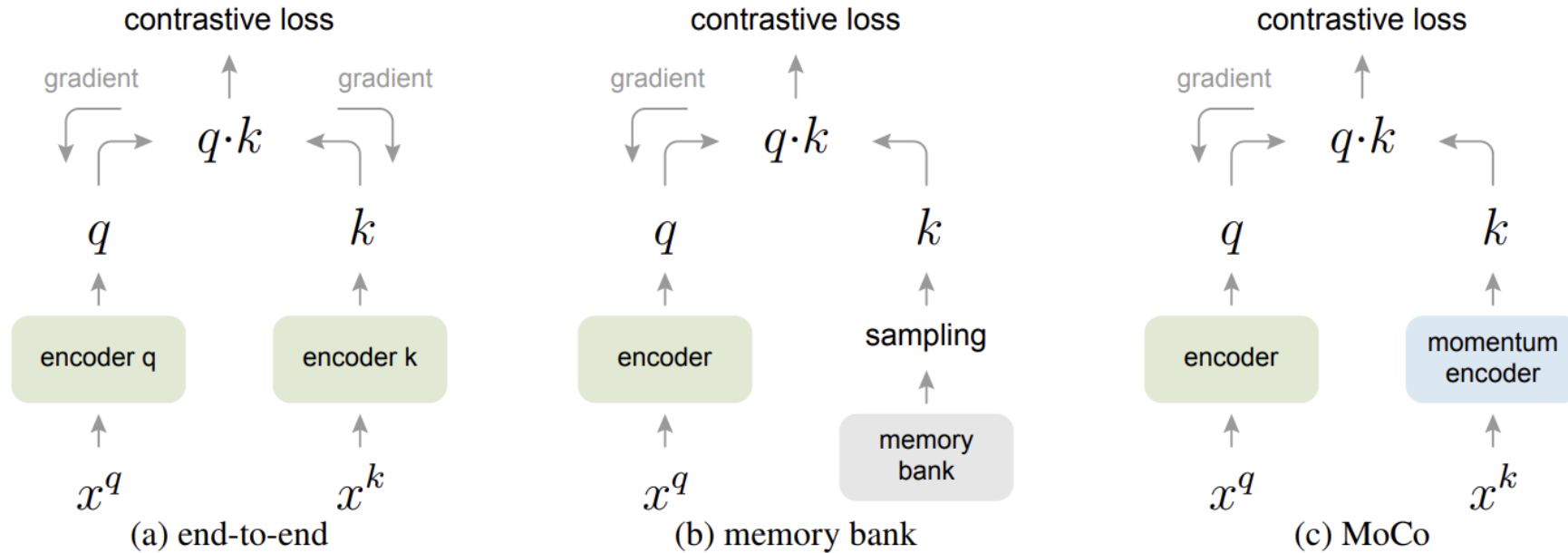
Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

Facebook AI Research (FAIR)

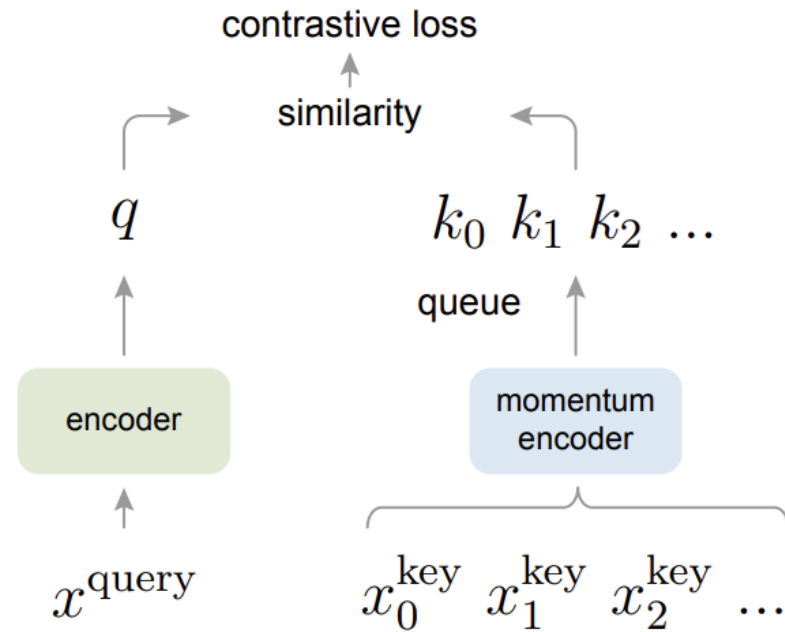
MoCo (CVPR 2020)

$$L_q = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=0}^K \exp(q \cdot k_i / \tau)}$$



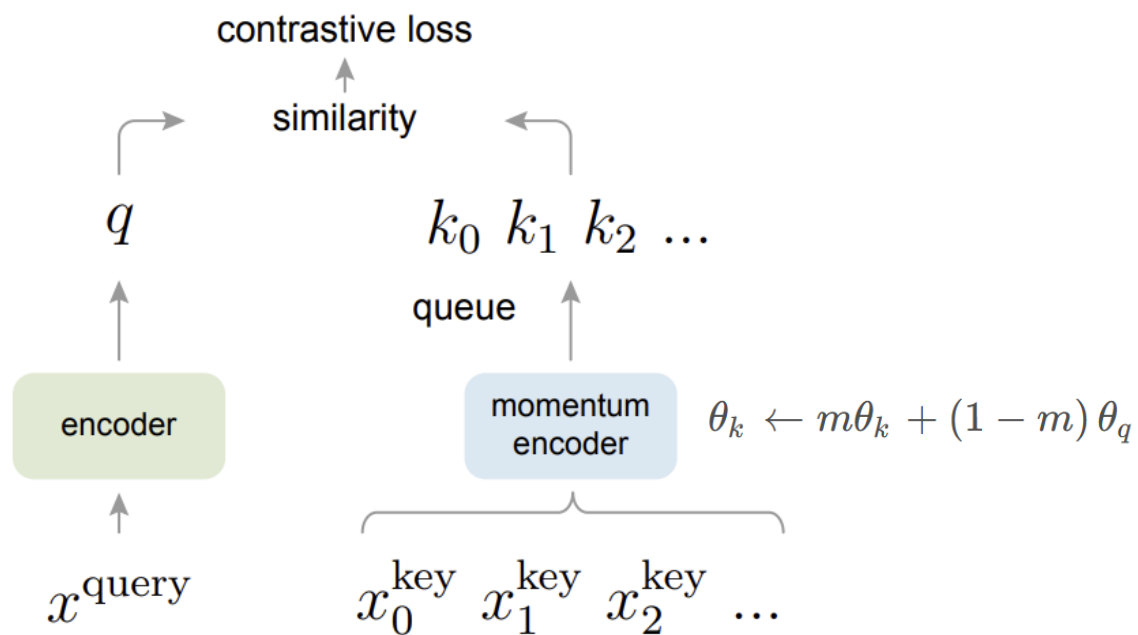
Background

■ MoCo



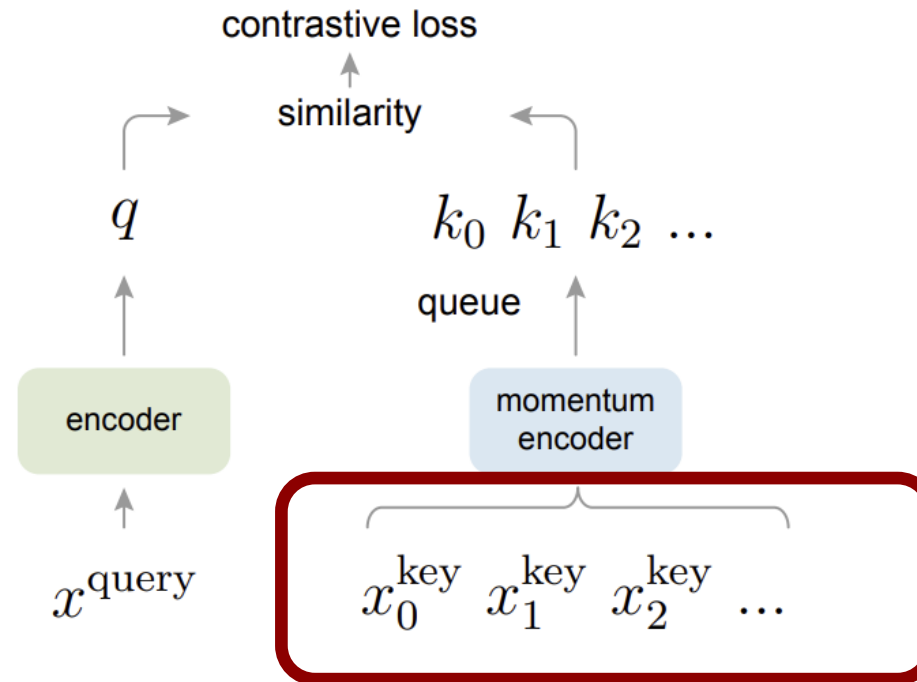
Background

■ MoCo



Background

■ MoCo



Background

■ MoCo v3

Algorithm 1 MoCo v3: PyTorch-like Pseudocode

```
# f_q: encoder: backbone + proj mlp + pred mlp
# f_k: momentum encoder: backbone + proj mlp
# m: momentum coefficient
# tau: temperature

for x in loader: # load a minibatch x with N samples
    x1, x2 = aug(x), aug(x) # augmentation
    q1, q2 = f_q(x1), f_q(x2) # queries: [N, C] each
    k1, k2 = f_k(x1), f_k(x2) # keys: [N, C] each

    loss = ctr(q1, k2) + ctr(q2, k1) # symmetrized
    loss.backward()

    update(f_q) # optimizer update: f_q
    f_k = m*f_k + (1-m)*f_q # momentum update: f_k

# contrastive loss
def ctr(q, k):
    logits = mm(q, k.t()) # [N, N] pairs
    labels = range(N) # positives are in diagonal
    loss = CrossEntropyLoss(logits/tau, labels)
    return 2 * tau * loss
```

Background

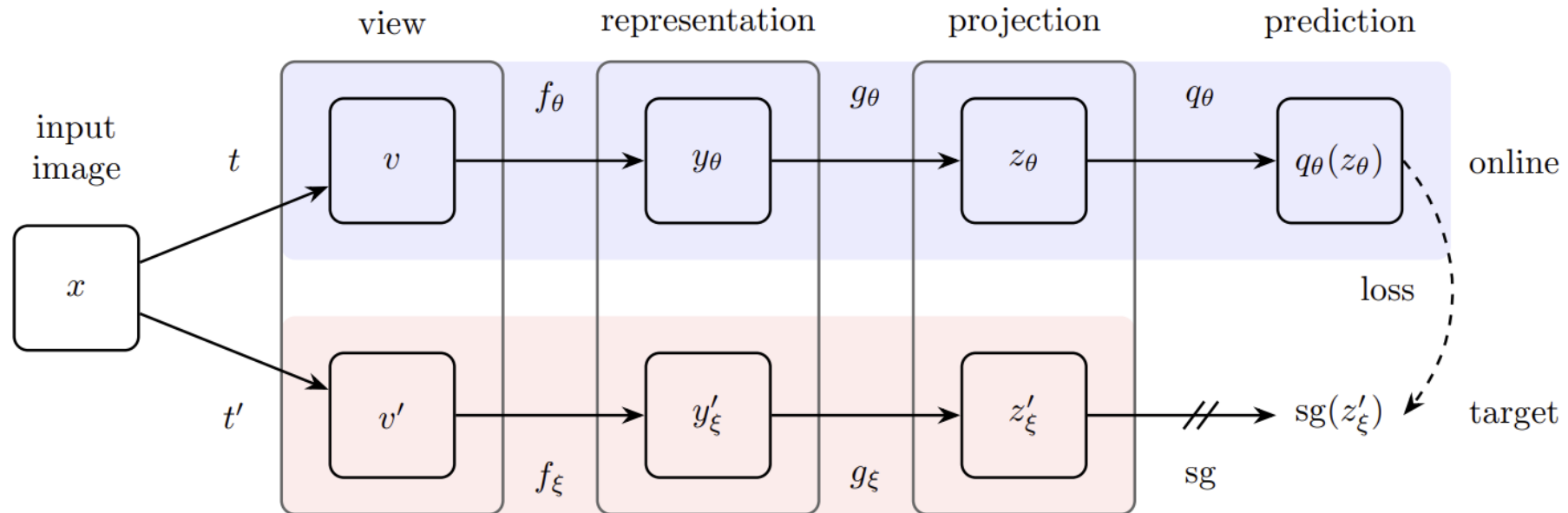
■ BYOL (NeurIPS 2020)

Bootstrap Your Own Latent A New Approach to Self-Supervised Learning

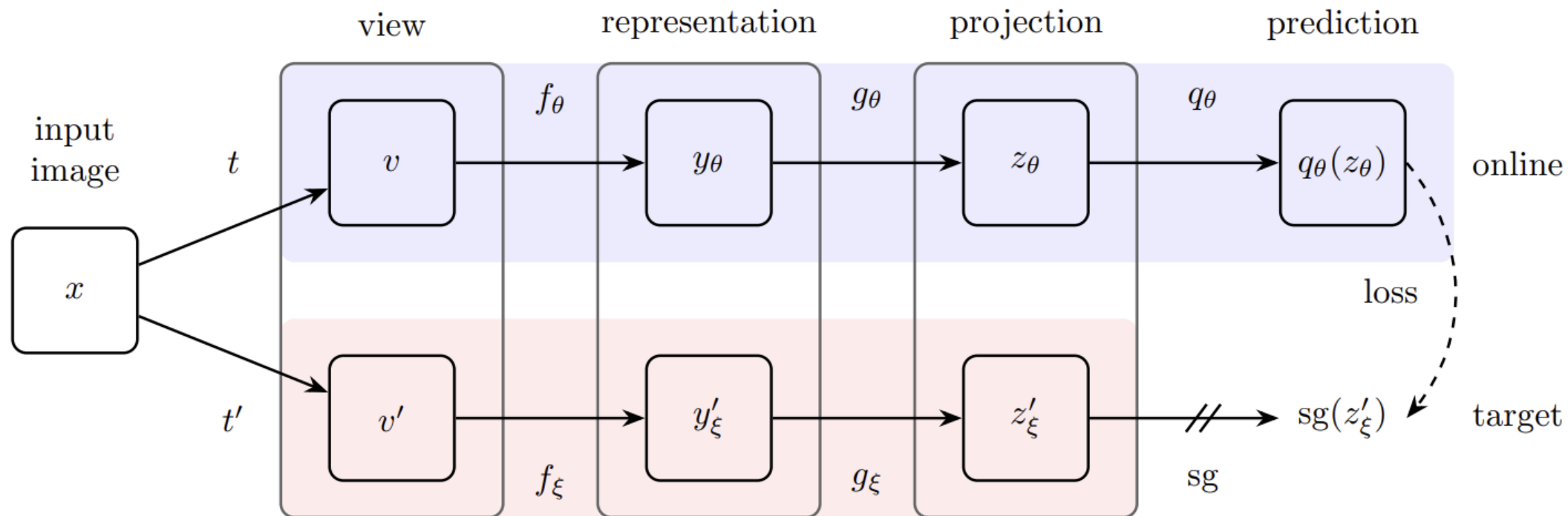
Jean-Bastien Grill^{*1}, Florian Strub^{*1}, Florent Alché^{*1}, Corentin Tallec^{*1}, Pierre H. Richemond^{*1,2}
 Elena Buchatskaya¹, Carl Doersch¹, Bernardo Avila Pires¹, Zhaohan Daniel Guo¹
 Mohammad Gheshlaghi Azar¹, Bilal Piot¹, Koray Kavukcuoglu¹, Rémi Munos¹, Michal Valko¹

¹DeepMind

²Imperial College



BYOL



BYOL

Algorithm 1: BYOL: Bootstrap Your Own Latent

Inputs :

\mathcal{D} , \mathcal{T} , and \mathcal{T}' set of images and distributions of transformations
 θ , f_θ , g_θ , and q_θ initial online parameters, encoder, projector, and predictor
 ξ , f_ξ , g_ξ initial target parameters, target encoder, and target projector
 optimizer optimizer, updates online parameters using the loss gradient
 K and N total number of optimization steps and batch size
 $\{\tau_k\}_{k=1}^K$ and $\{\eta_k\}_{k=1}^K$ target network update schedule and learning rate schedule

```

1 for  $k = 1$  to  $K$  do
2    $\mathcal{B} \leftarrow \{x_i \sim \mathcal{D}\}_{i=1}^N$       // sample a batch of  $N$  images
3   for  $x_i \in \mathcal{B}$  do
4      $t \sim \mathcal{T}$  and  $t' \sim \mathcal{T}'$       // sample image transformations
5      $z_1 \leftarrow g_\theta(f_\theta(t(x_i)))$  and  $z_2 \leftarrow g_\theta(f_\theta(t'(x_i)))$       // compute projections
6      $z'_1 \leftarrow g_\xi(f_\xi(t'(x_i)))$  and  $z'_2 \leftarrow g_\xi(f_\xi(t(x_i)))$       // compute target projections
7      $l_i \leftarrow -2 \cdot \left( \frac{\langle q_\theta(z_1), z'_1 \rangle}{\|q_\theta(z_1)\|_2 \cdot \|z'_1\|_2} + \frac{\langle q_\theta(z_2), z'_2 \rangle}{\|q_\theta(z_2)\|_2 \cdot \|z'_2\|_2} \right)$       // compute the loss for  $x_i$ 
8   end
9    $\delta\theta \leftarrow \frac{1}{N} \sum_{i=1}^N \partial_\theta l_i$       // compute the total loss gradient w.r.t.  $\theta$ 
10   $\theta \leftarrow \text{optimizer}(\theta, \delta\theta, \eta_k)$       // update online parameters
11   $\xi \leftarrow \tau_k \xi + (1 - \tau_k) \theta$       // update target parameters
12 end

```

Output : encoder f_θ

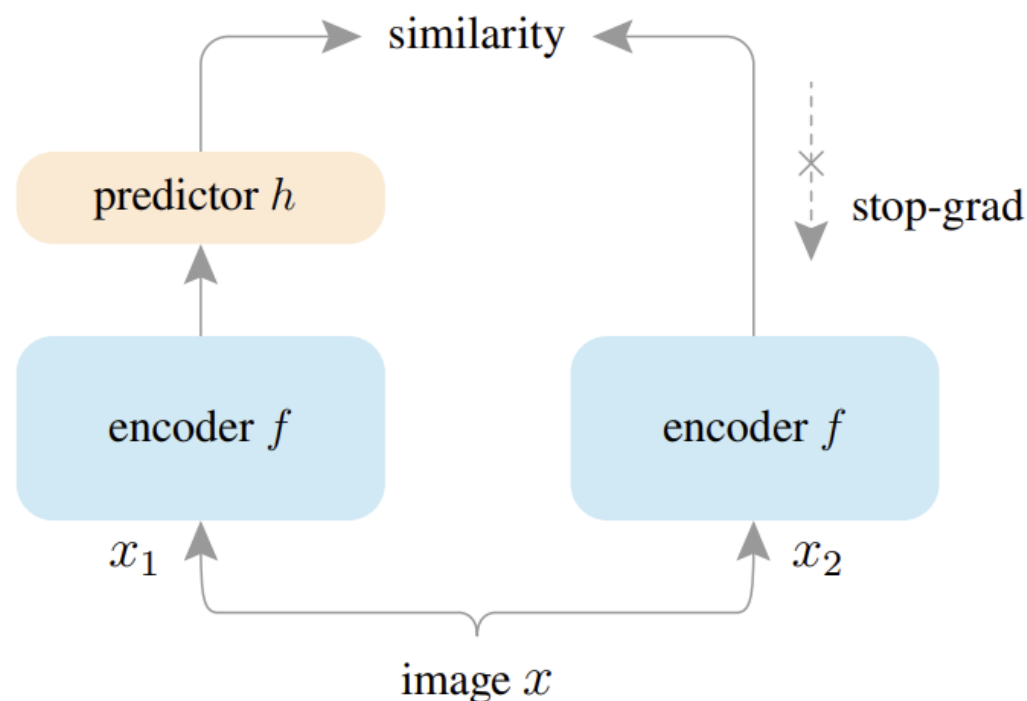
Background

Exploring Simple Siamese Representation Learning

Xinlei Chen Kaiming He

Facebook AI Research (FAIR)

■ SimSiam



Algorithm 1 SimSiam Pseudocode, PyTorch-like

```
# f: backbone + projection mlp
# h: prediction mlp

for x in loader: # load a minibatch x with n samples
    x1, x2 = aug(x), aug(x) # random augmentation
    z1, z2 = f(x1), f(x2) # projections, n-by-d
    p1, p2 = h(z1), h(z2) # predictions, n-by-d

    L = D(p1, z2)/2 + D(p2, z1)/2 # loss

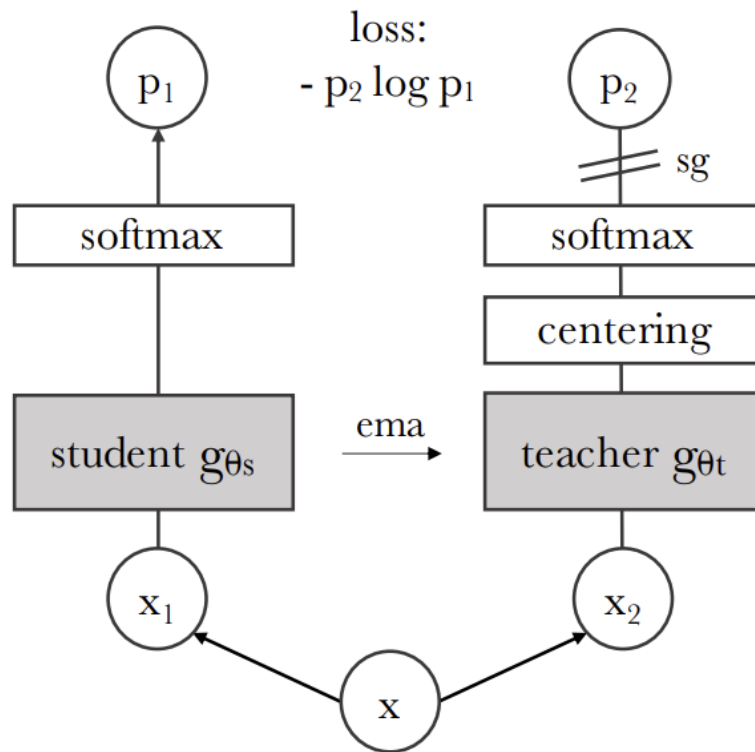
    L.backward() # back-propagate
    update(f, h) # SGD update

def D(p, z): # negative cosine similarity
    z = z.detach() # stop gradient

    p = normalize(p, dim=1) # l2-normalize
    z = normalize(z, dim=1) # l2-normalize
    return -(p*z).sum(dim=1).mean()
```

Background

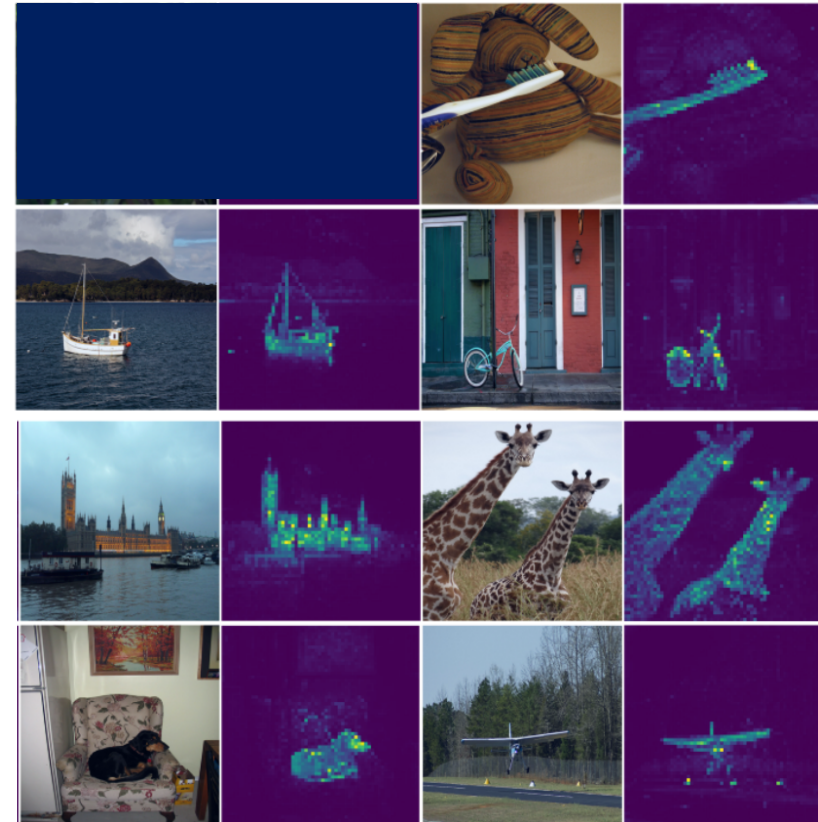
■ DINO(CVPR 2021)



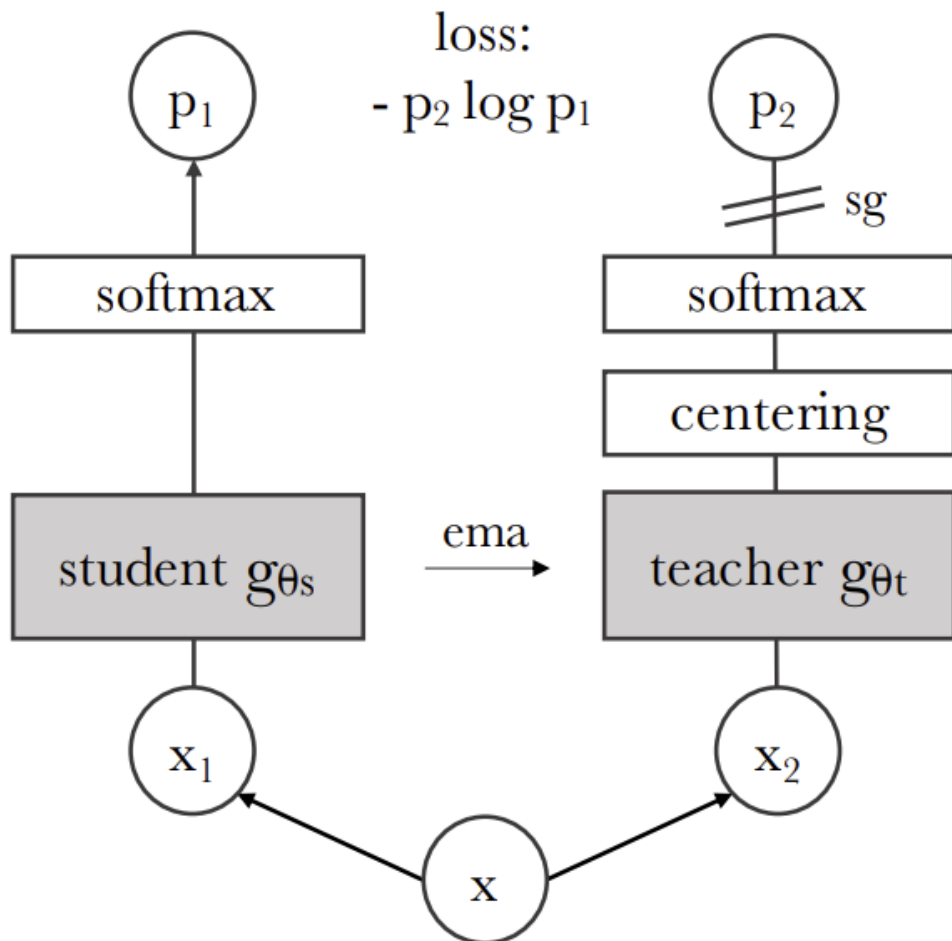
Emerging Properties in Self-Supervised Vision Transformers

Mathilde Caron^{1,2} Hugo Touvron^{1,3} Ishan Misra¹ Hervé Jegou¹
 Julien Mairal² Piotr Bojanowski¹ Armand Joulin¹

¹ Facebook AI Research ² Inria* ³ Sorbonne University



DINO



Algorithm 1 DINO PyTorch pseudocode w/o multi-crop.

```

# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# l, m: network and center momentum rates
gt.params = gs.params
for x in loader: # load a minibatch x with n samples
    x1, x2 = augment(x), augment(x) # random views

    s1, s2 = gs(x1), gs(x2) # student output n-by-K
    t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

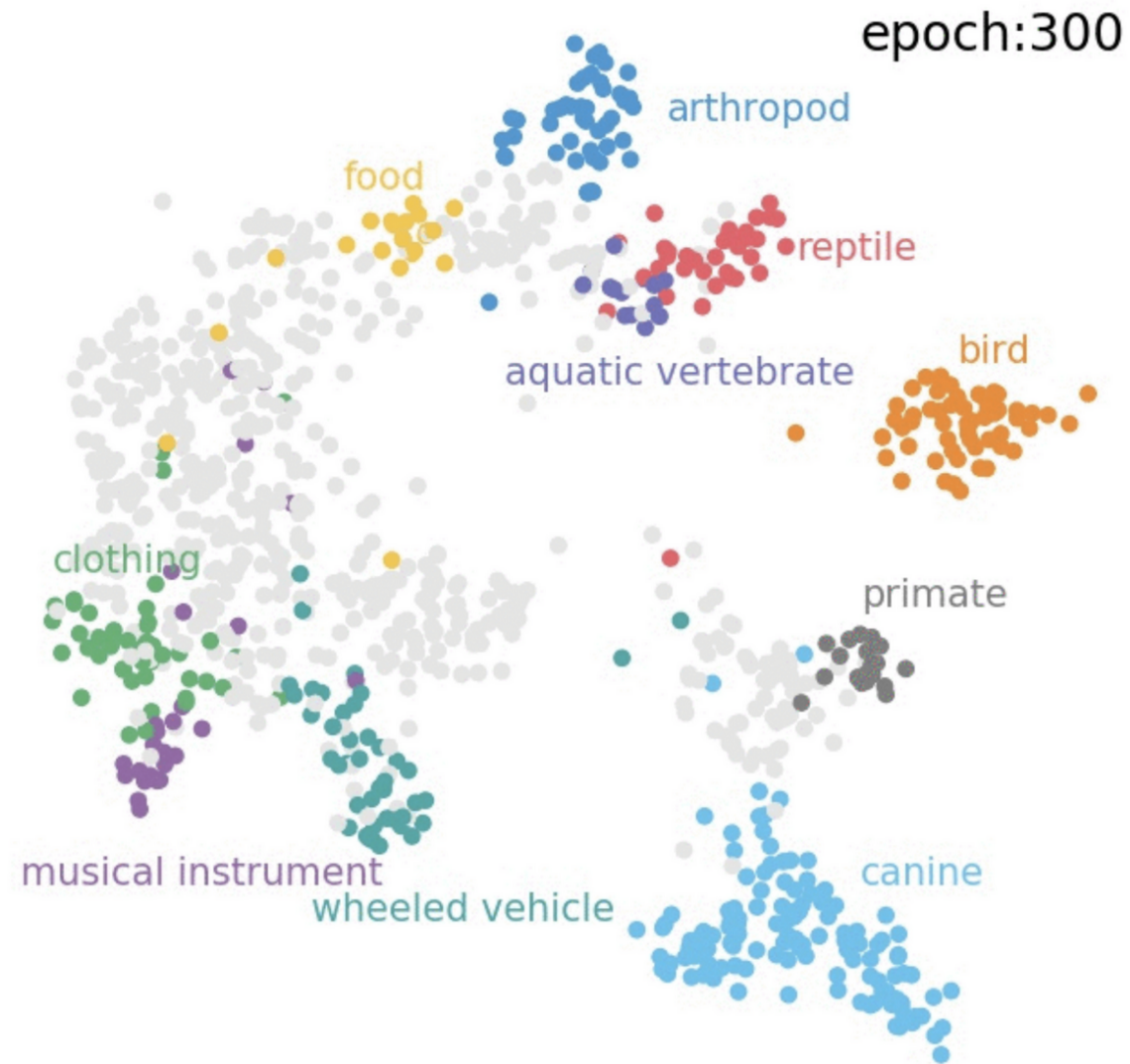
    loss = H(t1, s2)/2 + H(t2, s1)/2
    loss.backward() # back-propagate

    # student, teacher and center updates
    update(gs) # SGD
    gt.params = l*gt.params + (1-l)*gs.params
    C = m*C + (1-m)*cat([t1, t2]).mean(dim=0)

def H(t, s):
    t = t.detach() # stop gradient
    s = softmax(s / tps, dim=1)
    t = softmax((t - C) / tpt, dim=1) # center + sharpen
    return - (t * log(s)).sum(dim=1).mean()

```

DINO



Background

■ Object Discovery



Background

- **Object-centric representations**
- **Decompose scenes in terms of abstract building blocks**

Unsupervised Conditional Slot Attention for Object Centric Learning

Avinash Kori[†] Francesco Locatello^{*} Francesca Toni[†] Ben Glocker[†]

[†] Department of Computing, Imperial College London
a.kori21@imperial.ac.uk

Background

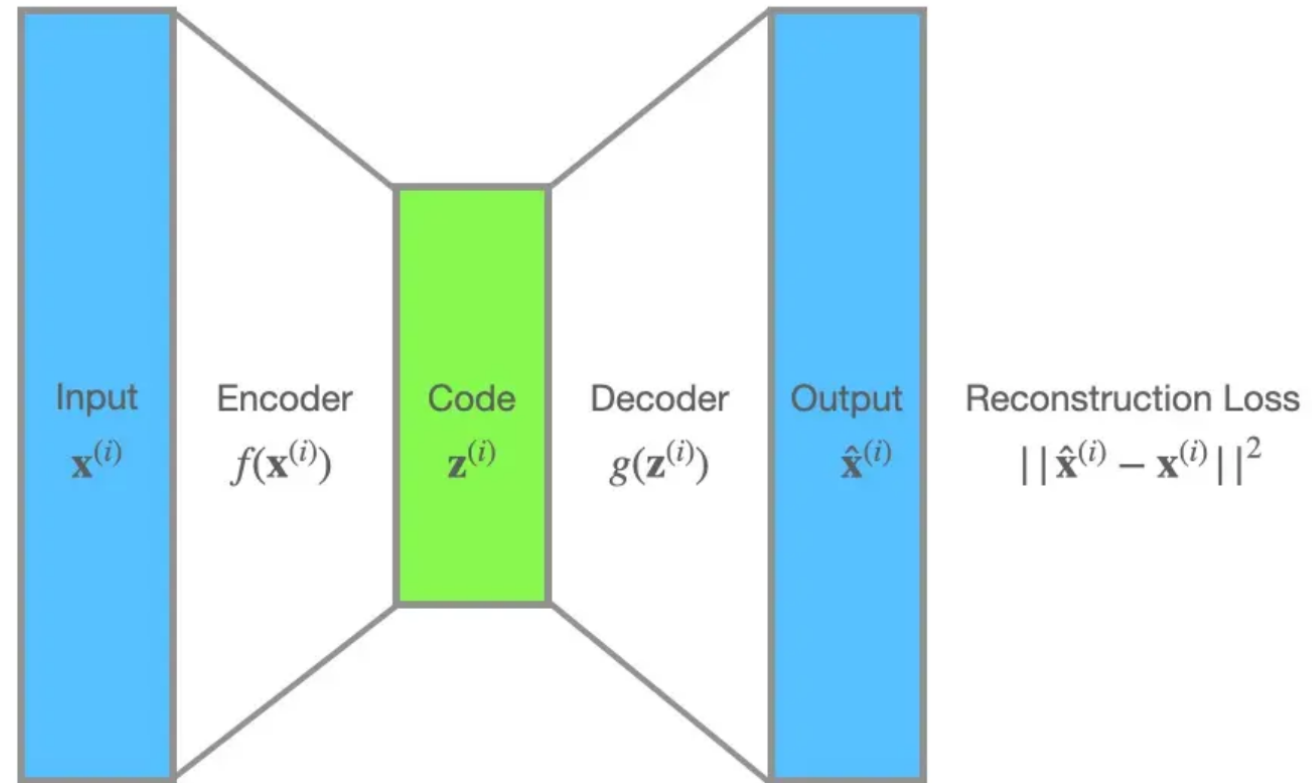
- **MONet**
- **Slot-Based Methods**
- **Multi Objects**
- **Components VAE**
- **Recurrent attention network**

MONet: Unsupervised Scene Decomposition and Representation

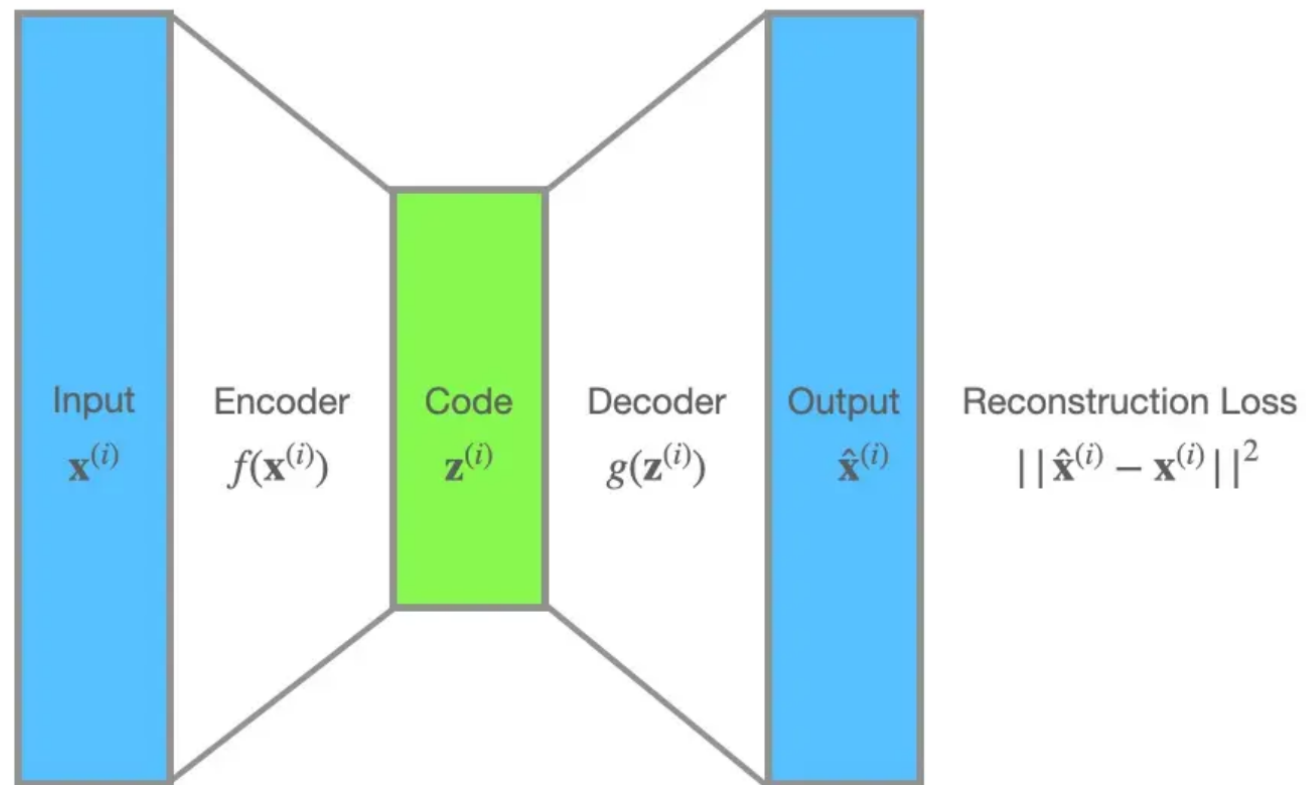
Christopher P. Burgess, Loic Matthey, Nicholas Watters,
Rishabh Kabra, Irina Higgins, Matt Botvinick, Alexander Lerchner

DeepMind
London, United Kingdom

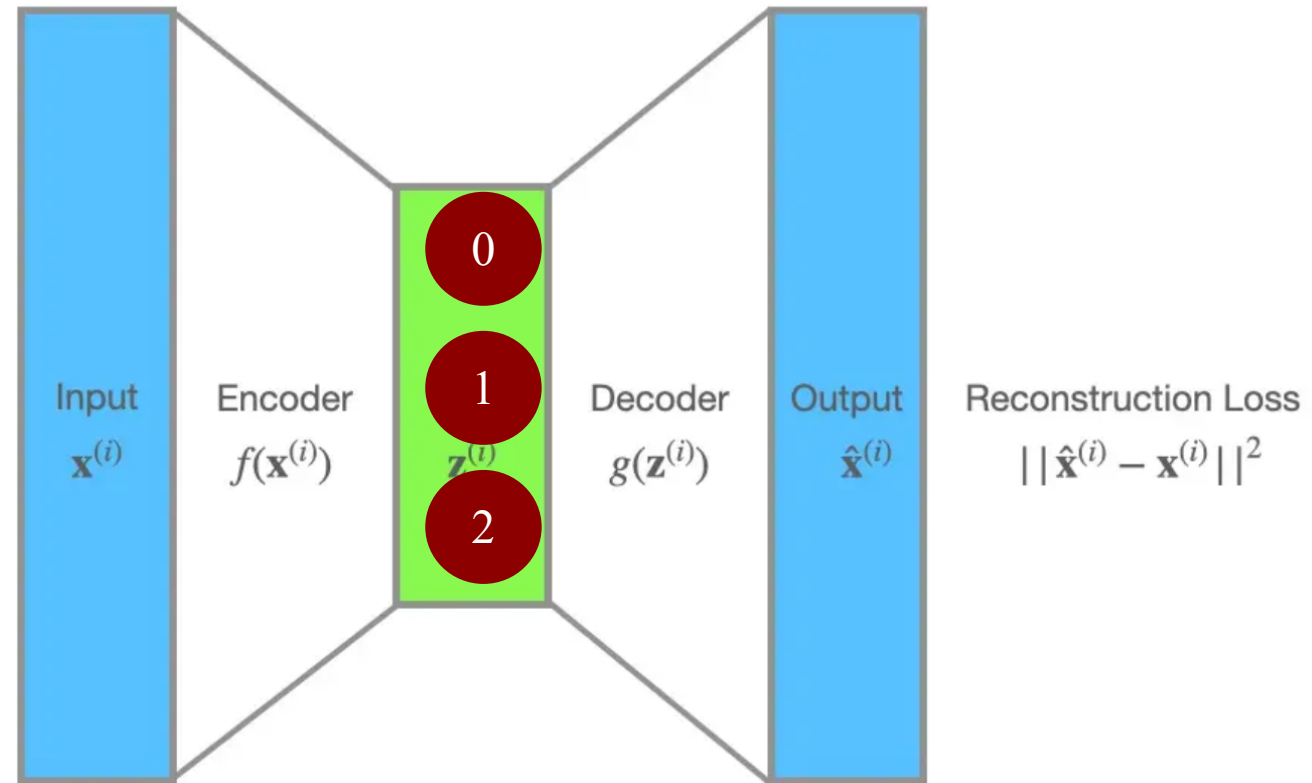
VAE



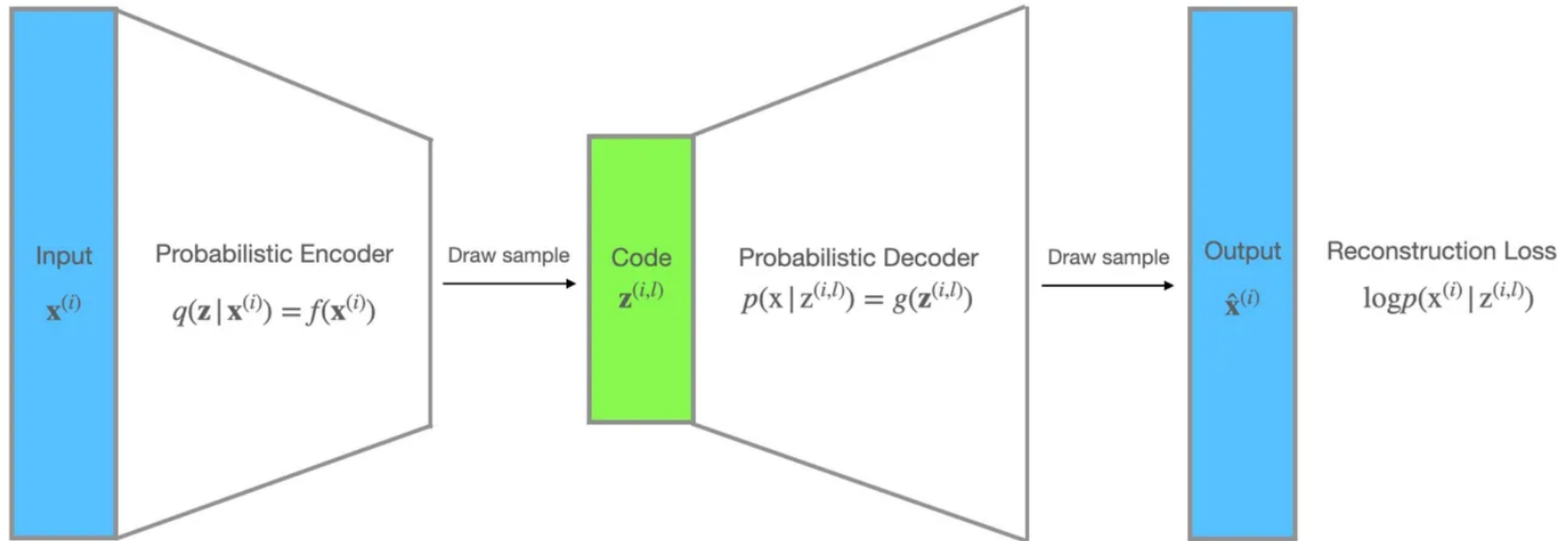
VAE



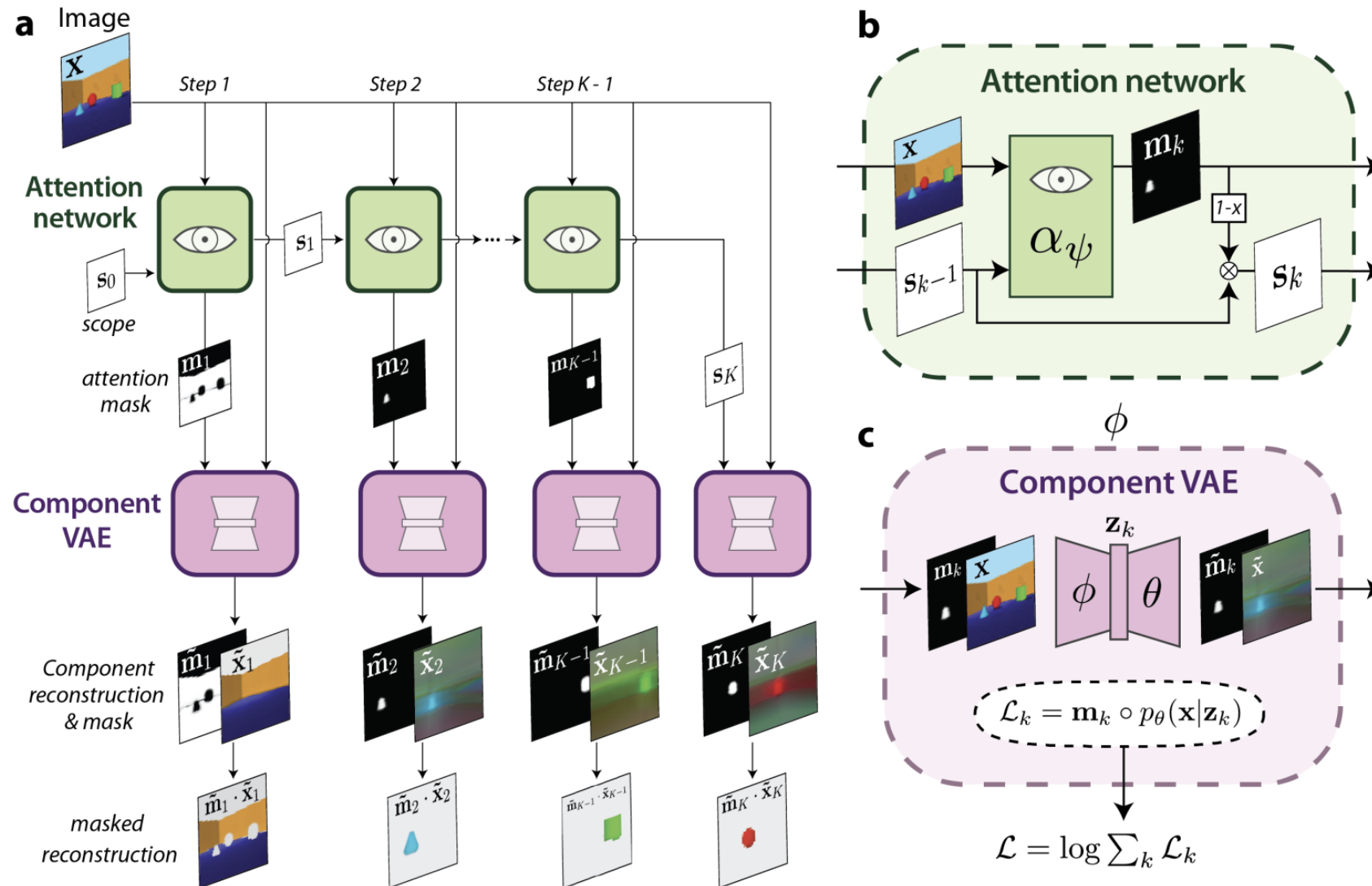
VAE



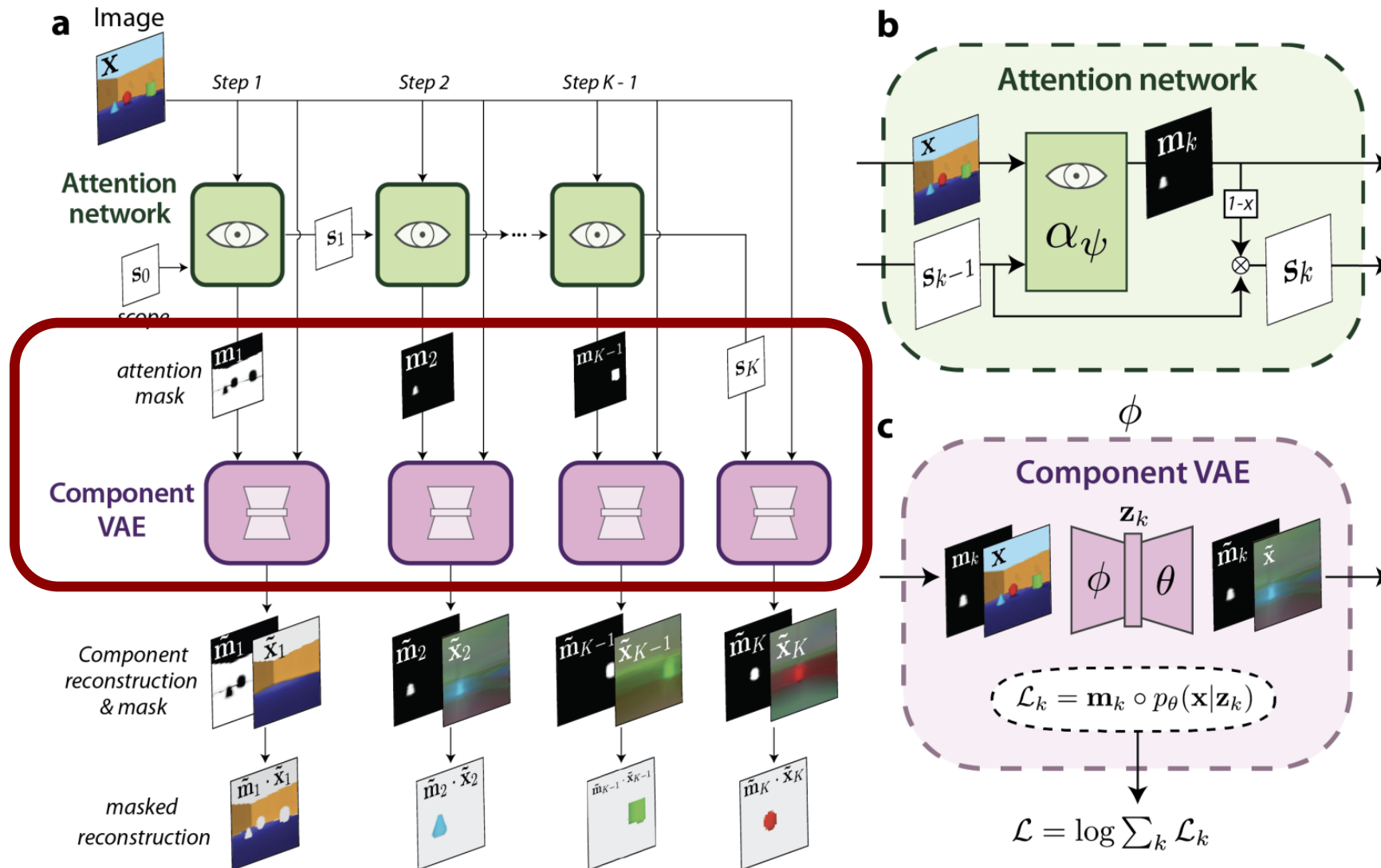
VAE



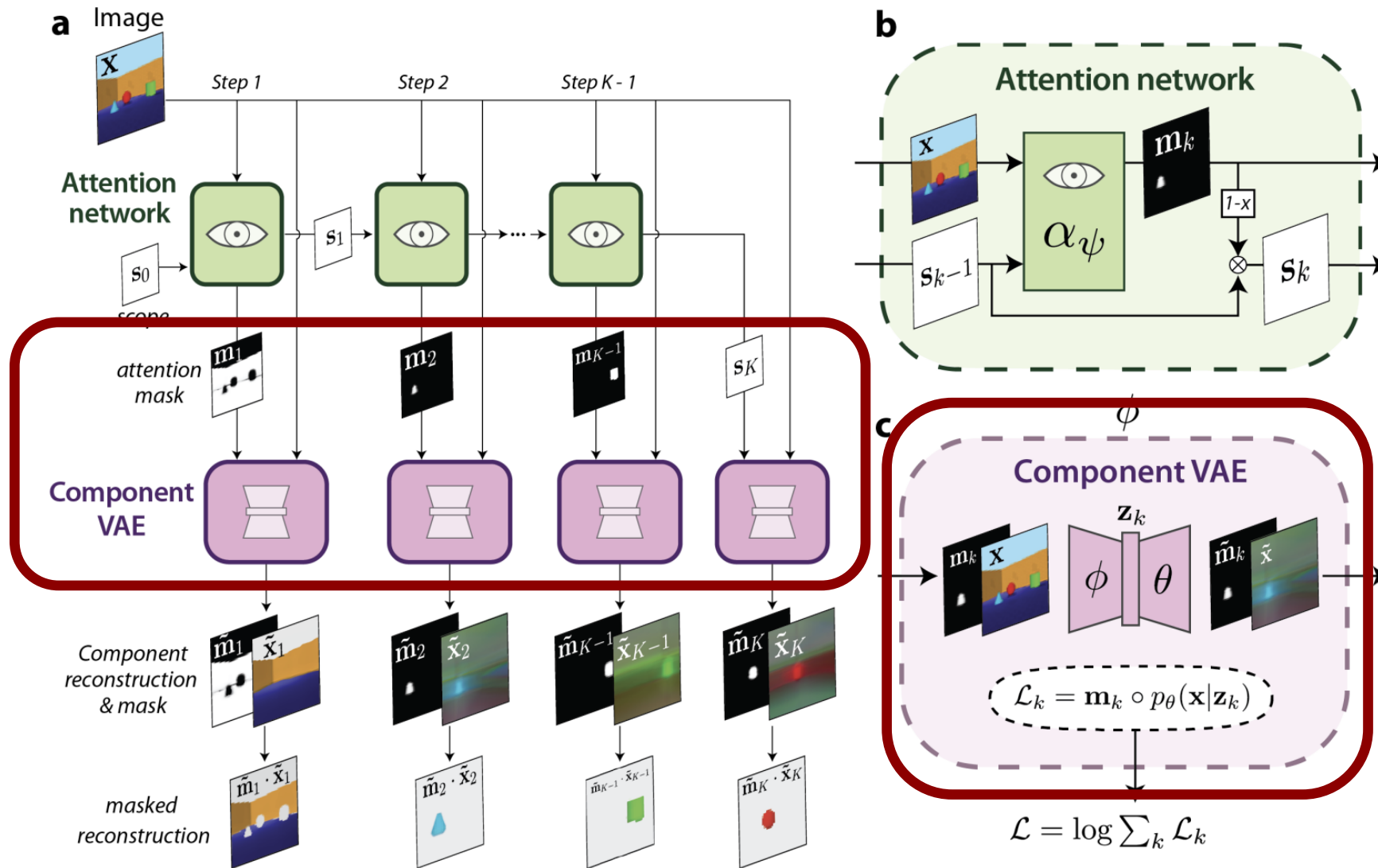
MONet



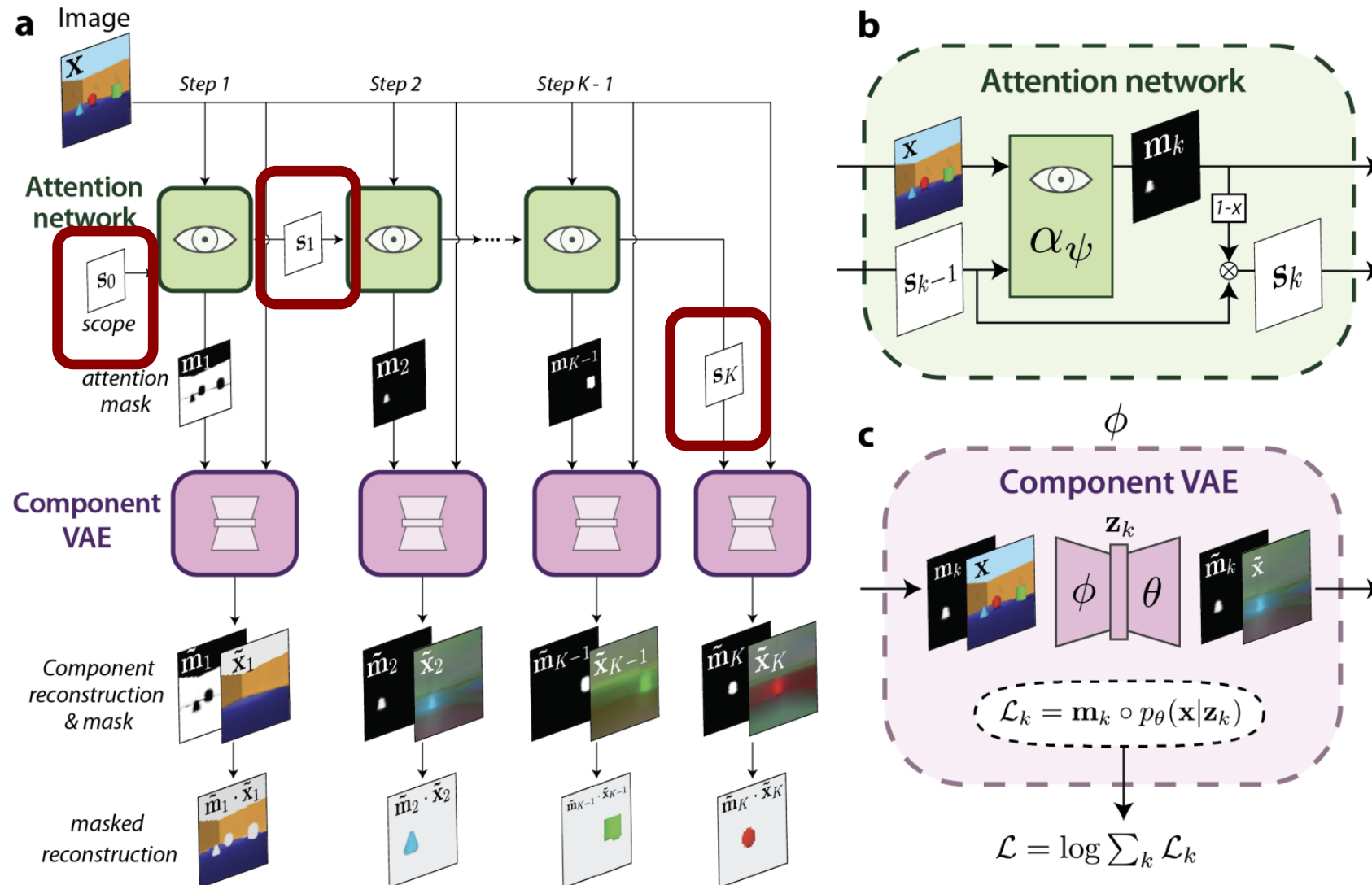
MONet



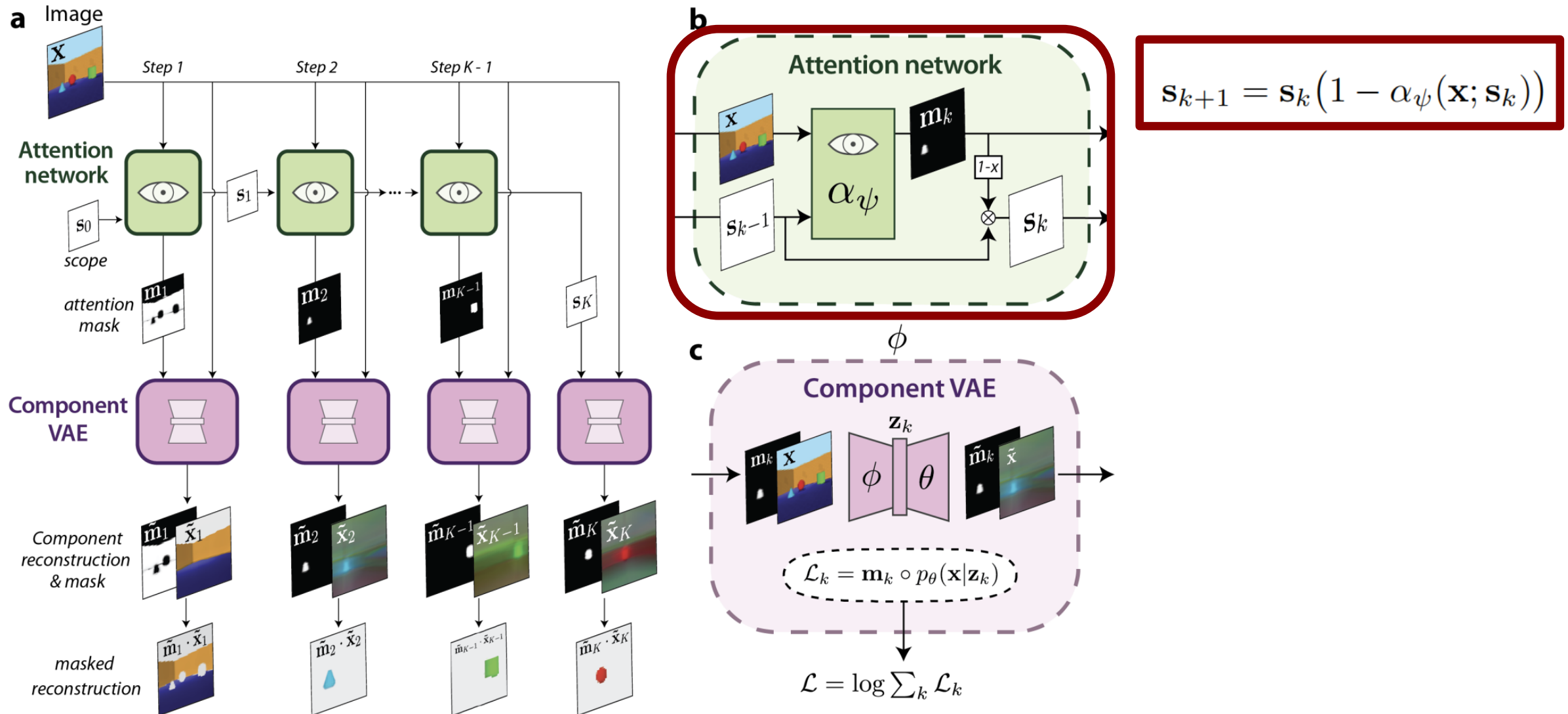
MONet



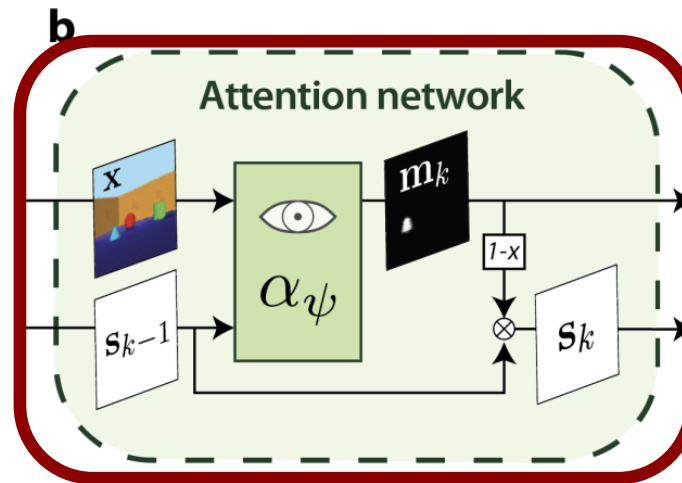
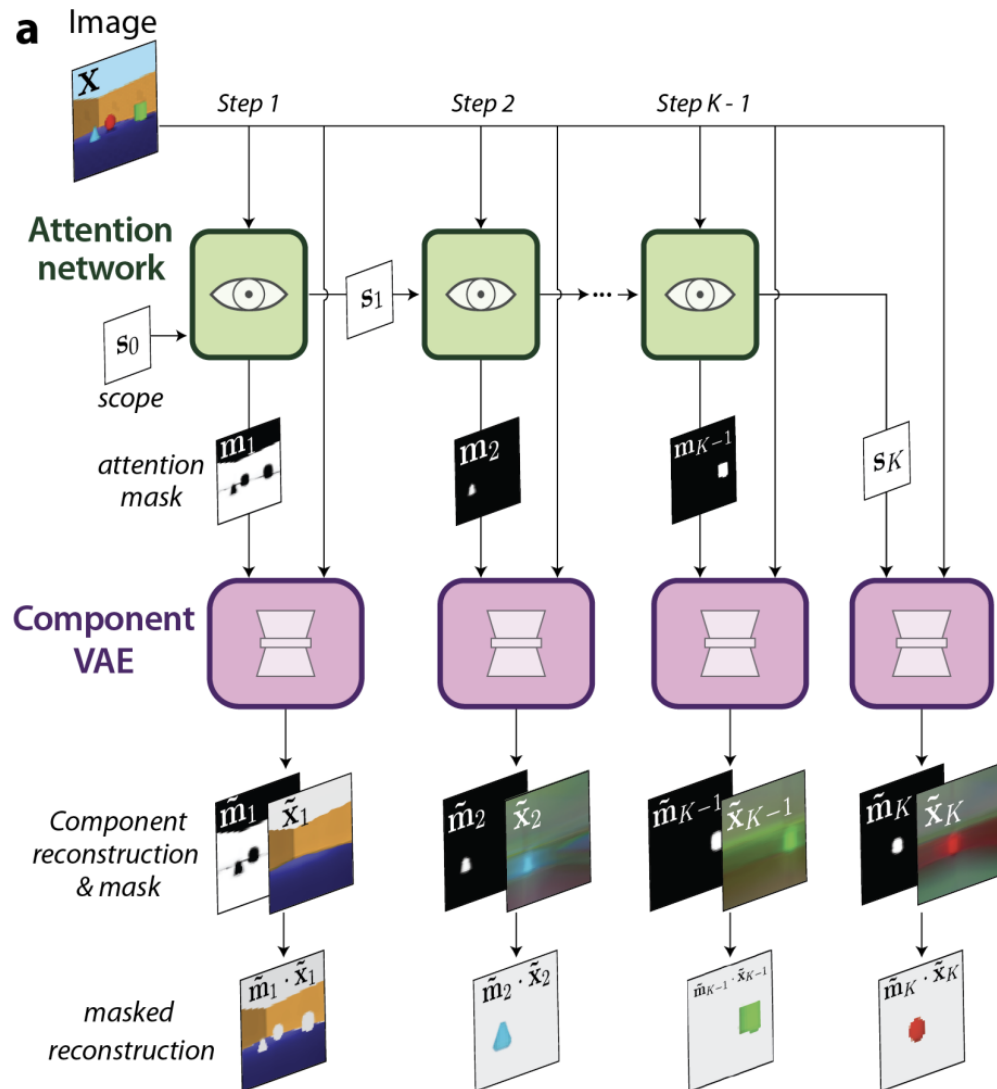
MONet



MONet

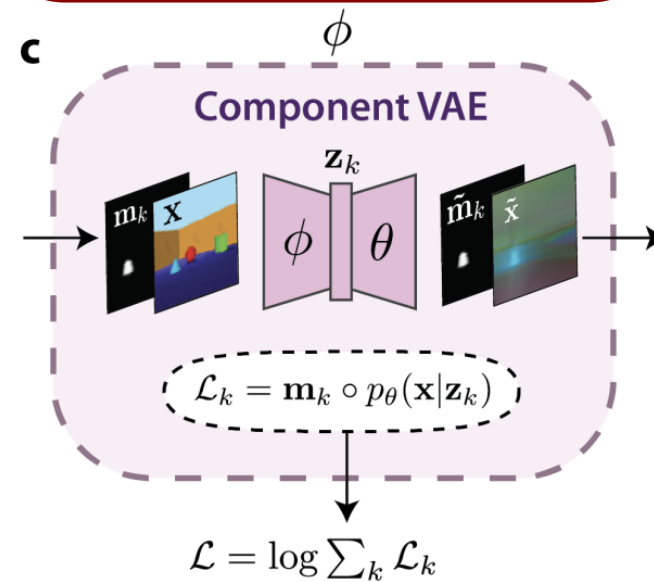


MONet

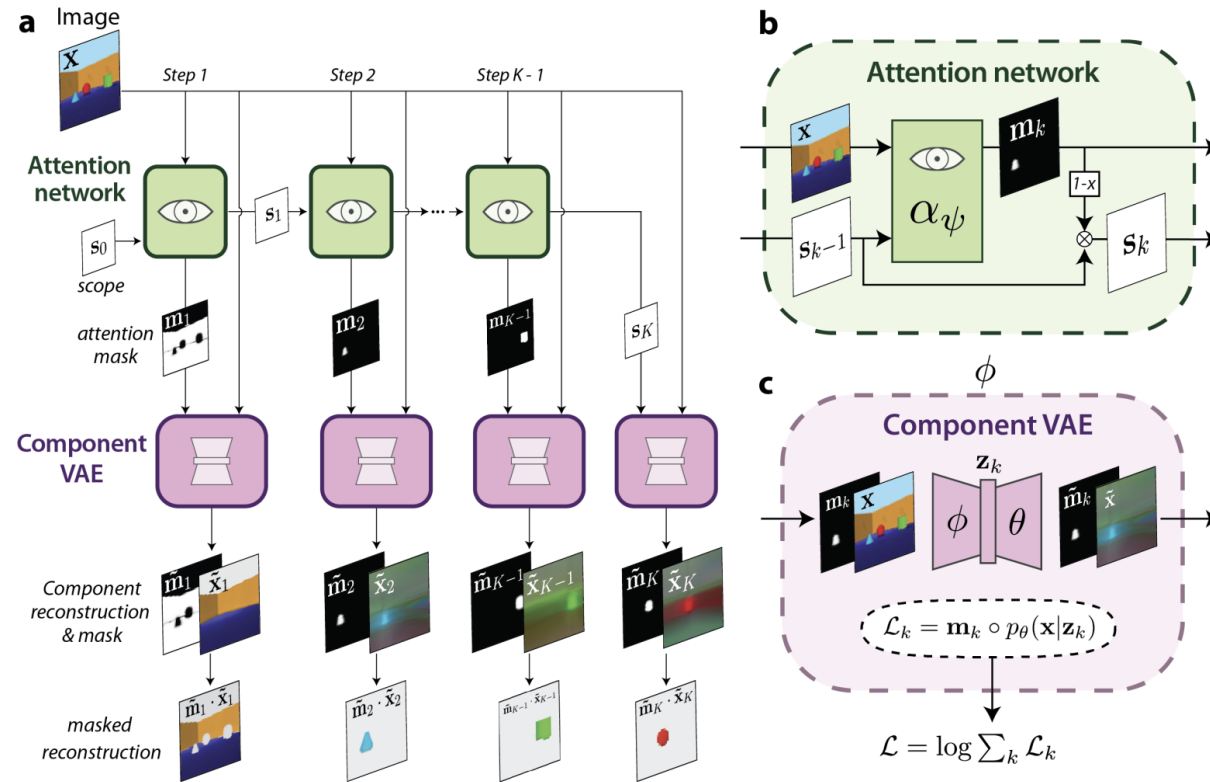


$$s_{k+1} = s_k (1 - \alpha_\psi(\mathbf{x}; s_k))$$

$$m_k = s_{k-1} \alpha_\psi(\mathbf{x}; s_{k-1})$$

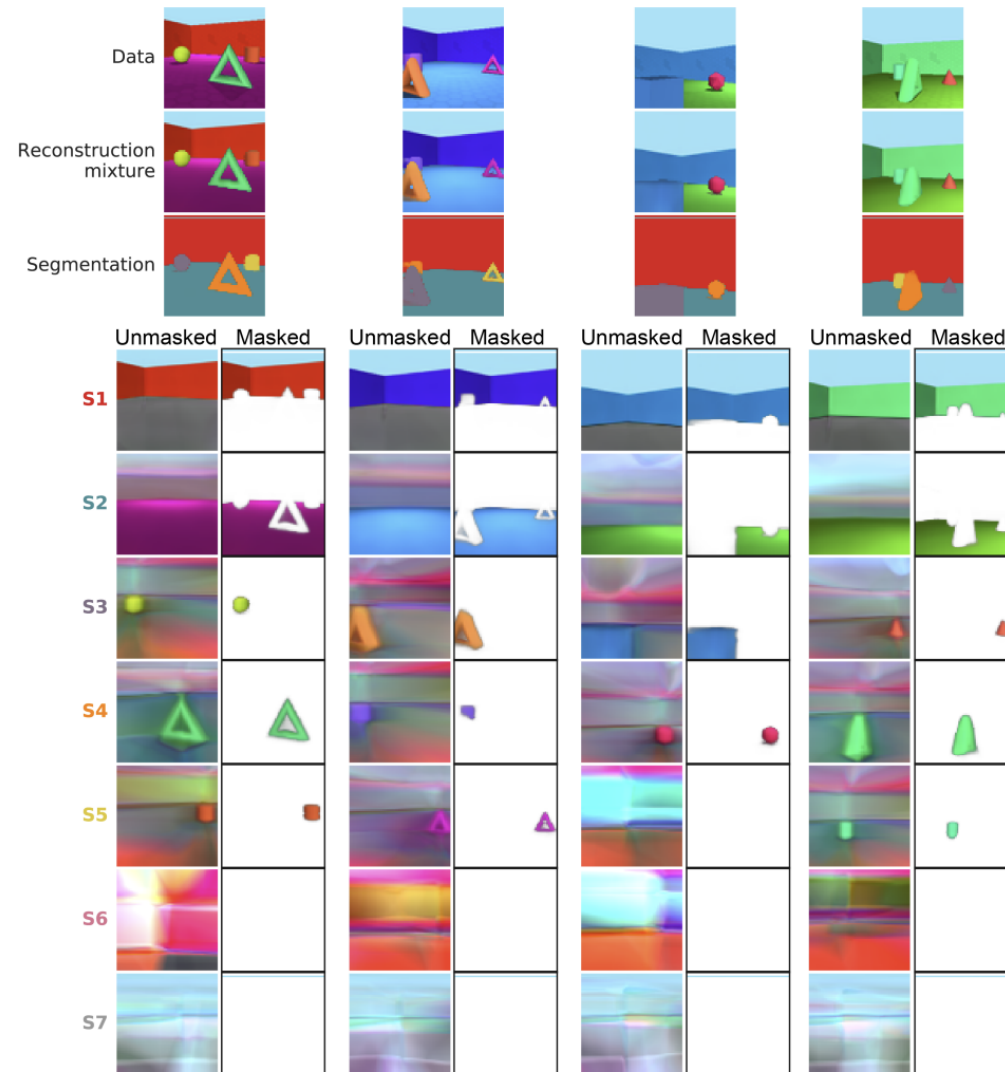


MONet



$$\mathcal{L}(\phi; \theta; \psi; \mathbf{x}) = -\log \sum_{k=1}^K \mathbf{m}_k p_\theta(\mathbf{x}|\mathbf{z}_k) + \beta D_{KL} \left(\prod_{k=1}^K q_\phi(\mathbf{z}_k|\mathbf{x}, \mathbf{m}_k) \parallel p(\mathbf{z}) \right) + \gamma D_{KL} (q_\psi(\mathbf{c}|\mathbf{x}) \parallel p_\theta(\mathbf{c}|\{\mathbf{z}_k\}))$$

MONet



Background

■ SlotCon (NeurIPS 2022)

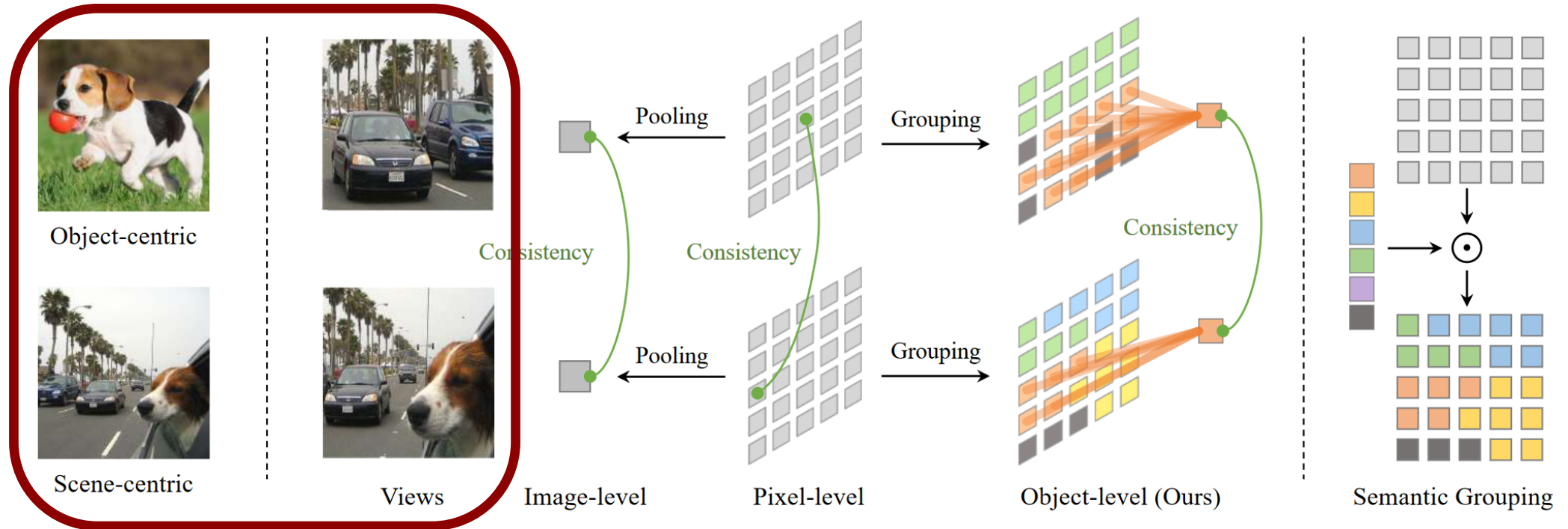
Self-Supervised Visual Representation Learning with Semantic Grouping

Xin Wen¹ Bingchen Zhao^{2,3} Anlin Zheng^{1,4} Xiangyu Zhang⁴ Xiaojuan Qi¹

¹University of Hong Kong ²University of Edinburgh ³LunarAI ⁴MEGVII Technology

{wenxin, xjqj}@eee.hku.hk zhaobc.gmail.com
{zhenganlin, zhangxiangyu}@megvii.com

SlotCon



SlotCon



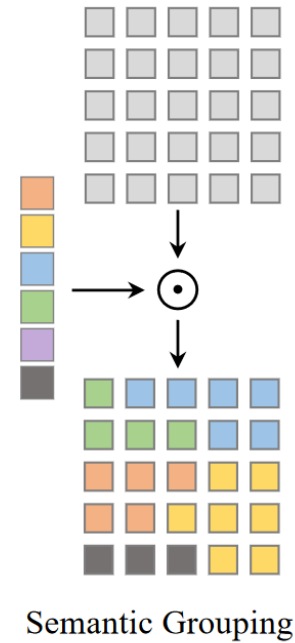
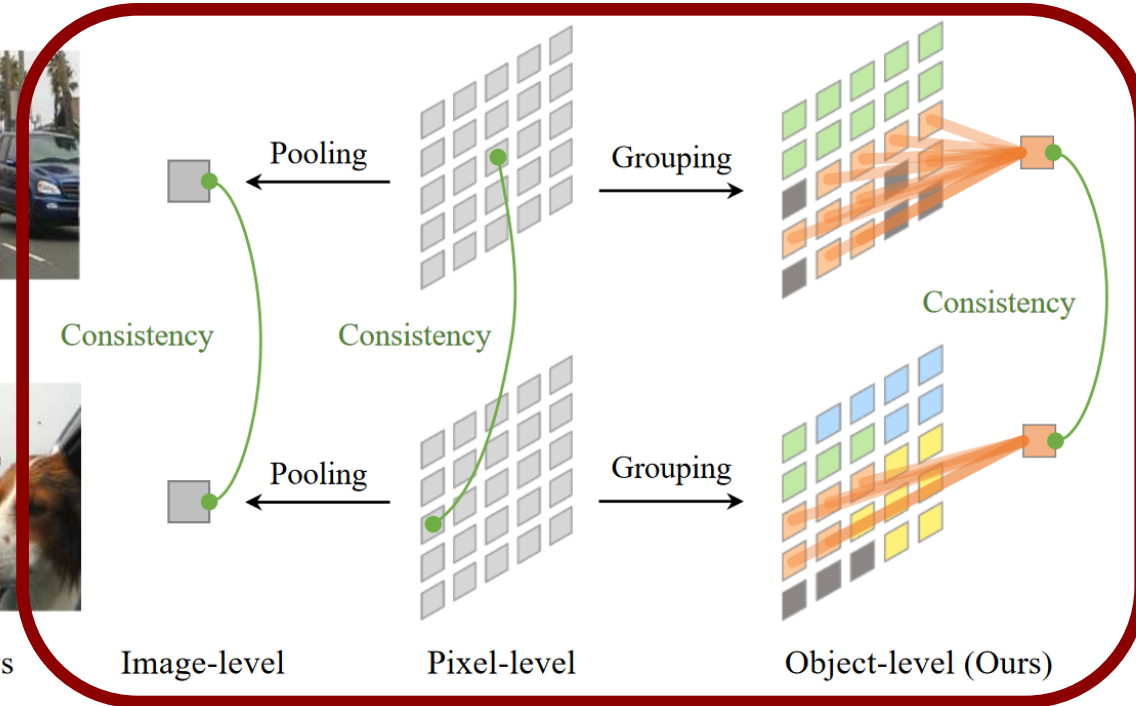
Object-centric



Scene-centric



Views



SlotCon



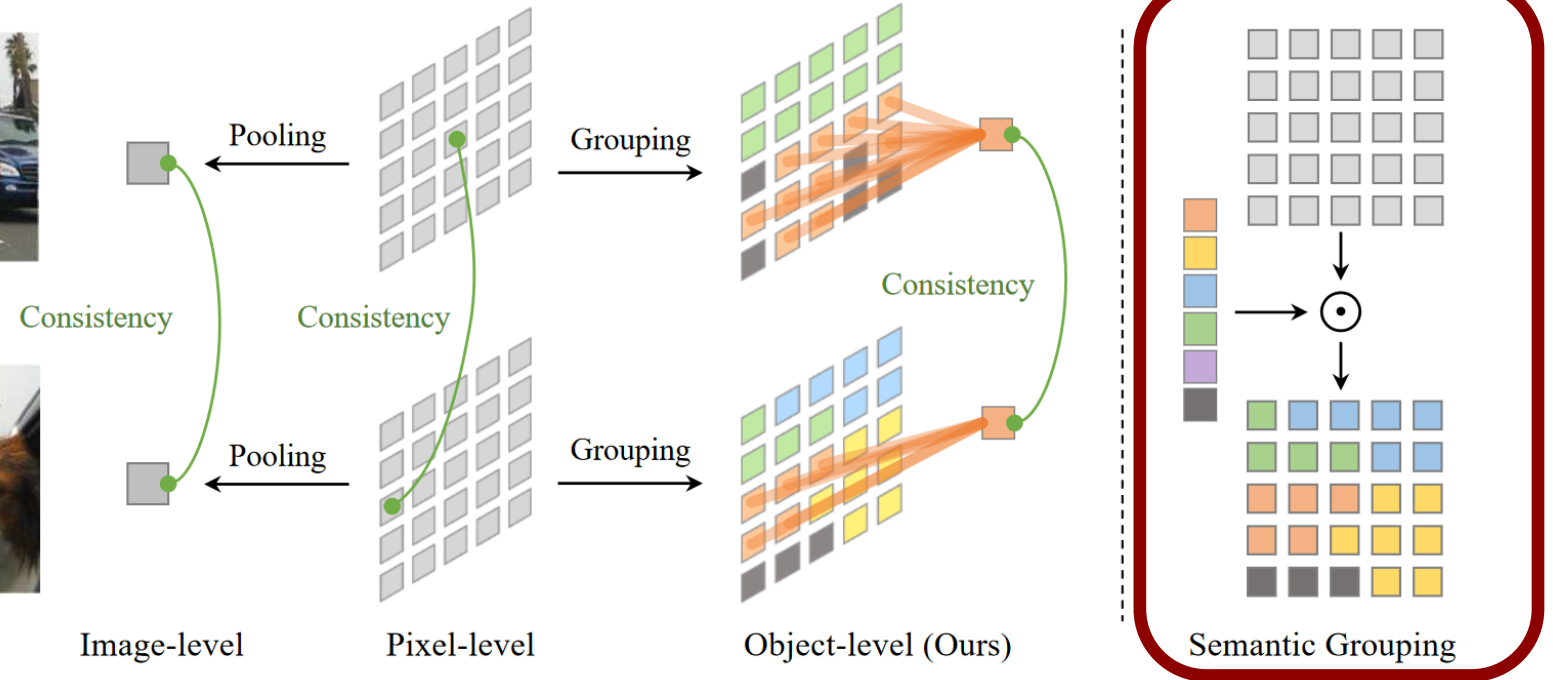
Object-centric



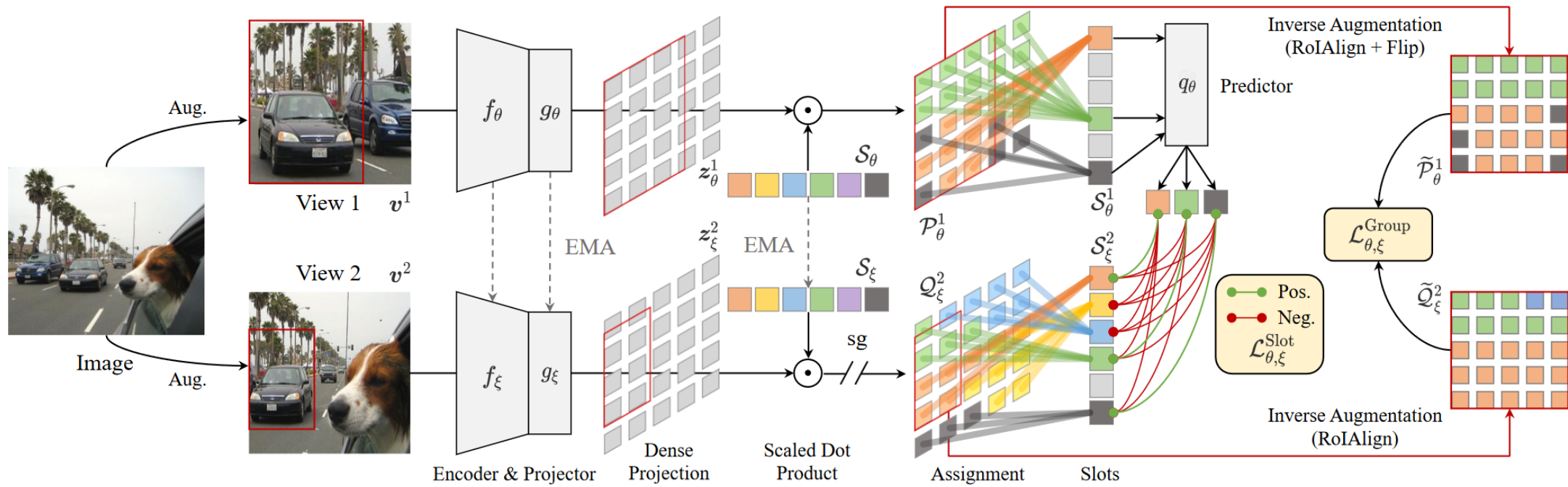
Scene-centric



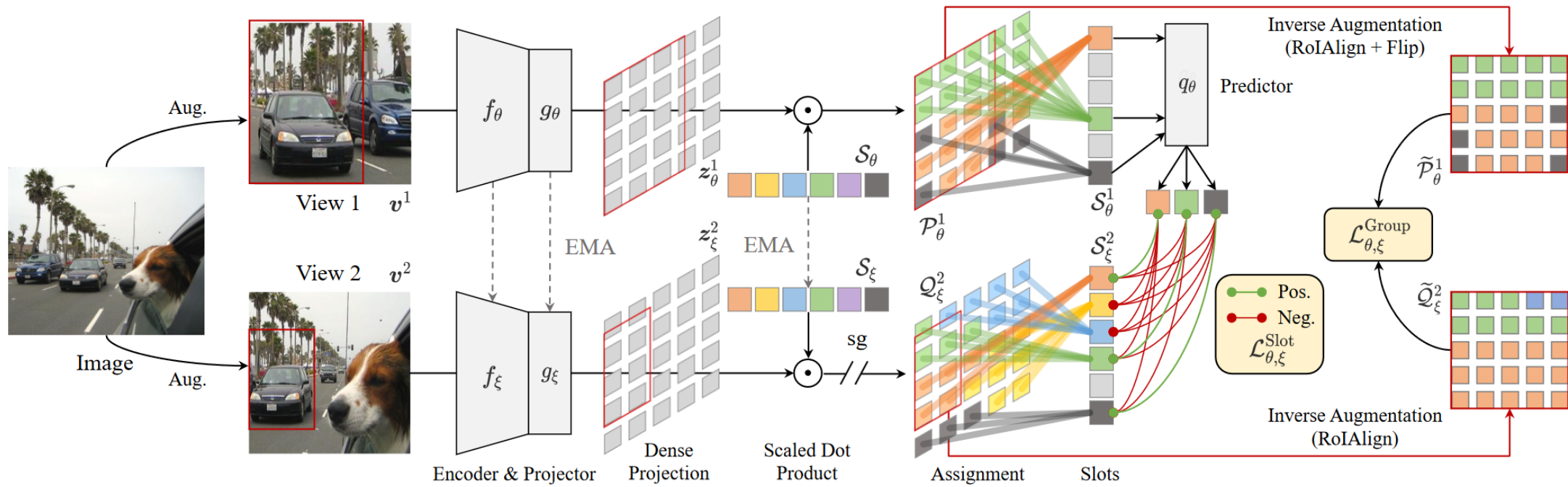
Views



SlotCon



SlotCon



Background

- CAE

- Complex-Valued Autoencoders

- Simple Design

Complex-Valued Autoencoders for Object Discovery

Sindy Löwe

UvA-Bosch Delta Lab, University of Amsterdam

loewe.sindy@gmail.com

Phillip Lippe

QUVA Lab, University of Amsterdam

p.lippe@uva.nl

Maja Rudolph

Bosch Center for AI

maja.rudolph@us.bosch.com

Max Welling

UvA-Bosch Delta Lab, University of Amsterdam

m.welling@uva.nl

Background

■ CAE

■ Complex-Valued Autoencoders

■ Simple Design

Complex-Valued Autoencoders for Object Discovery

Sindy Löwe

UvA-Bosch Delta Lab, University of Amsterdam

loewe.sindy@gmail.com

Phillip Lippe

QUVA Lab, University of Amsterdam

p.lippe@uva.nl

Maja Rudolph

Bosch Center for AI

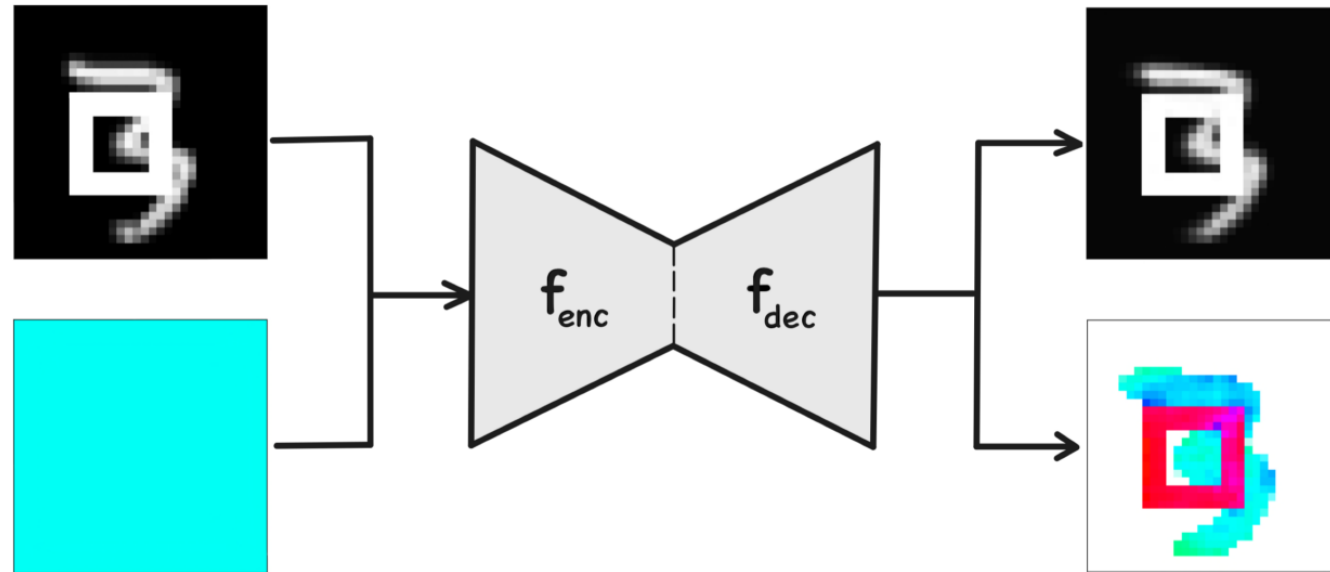
maja.rudolph@us.bosch.com

Max Welling

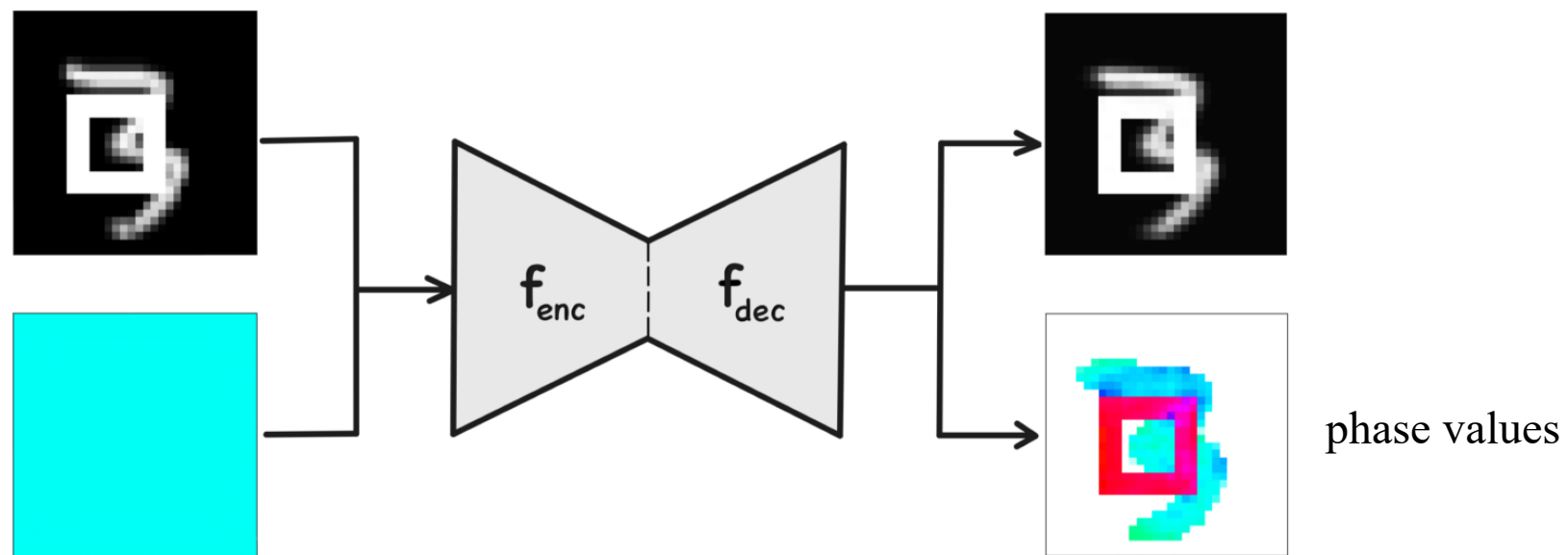
UvA-Bosch Delta Lab, University of Amsterdam

m.welling@uva.nl

CAE



CAE



CAE

■ Synchronization

■ Additive operations $\psi = f_{\mathbf{w}}(\mathbf{z}) = f_{\mathbf{w}}(\text{Re}(\mathbf{z})) + f_{\mathbf{w}}(\text{Im}(\mathbf{z})) \cdot i \in \mathbb{C}^{d_{\text{out}}}$

■ Desynchronization $m_{\psi} = |\psi| + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}} \quad \varphi_{\psi} = \arg(\psi) + \mathbf{b}_{\varphi} \in \mathbb{R}^{d_{\text{out}}}$

■ More control over the precise phase shifts:

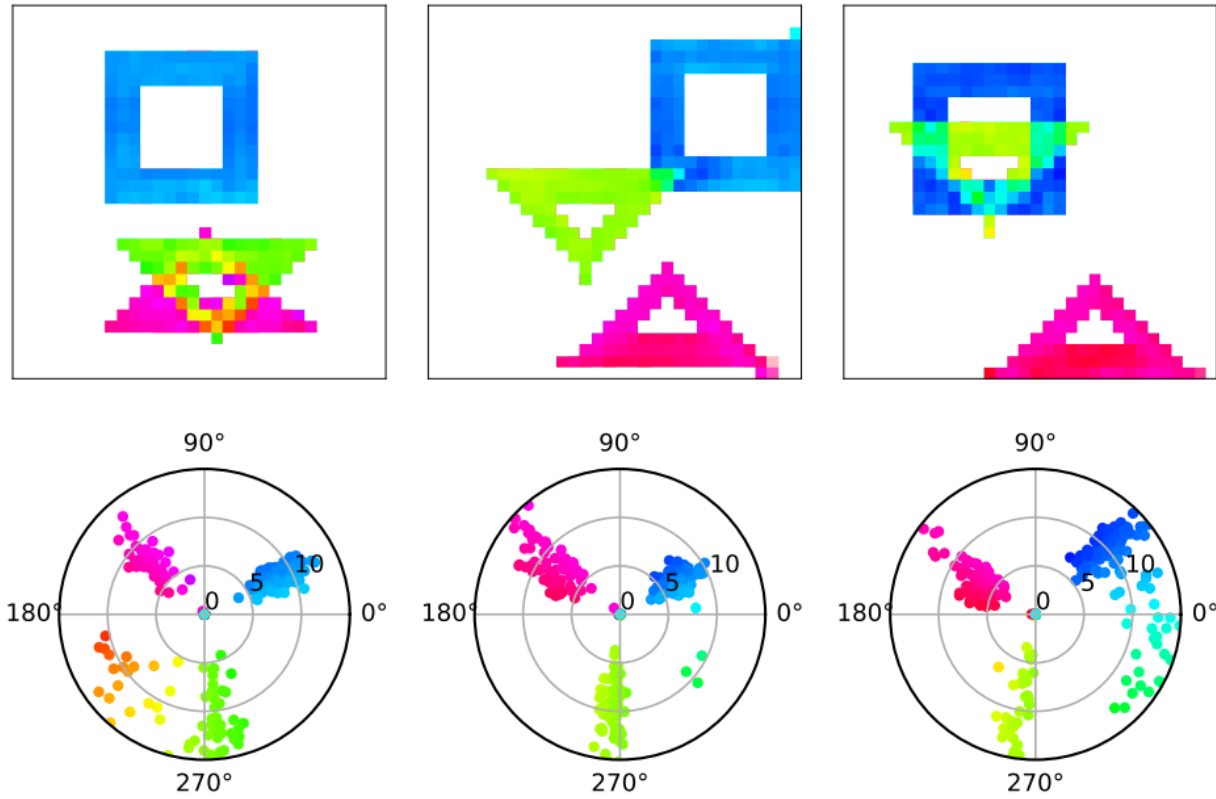
■ Gating $\chi = f_{\mathbf{w}}(|\mathbf{z}|) + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}}$
 $m_{\mathbf{z}} = 0.5 \cdot m_{\psi} + 0.5 \cdot \chi \in \mathbb{R}^{d_{\text{out}}}$

CAE

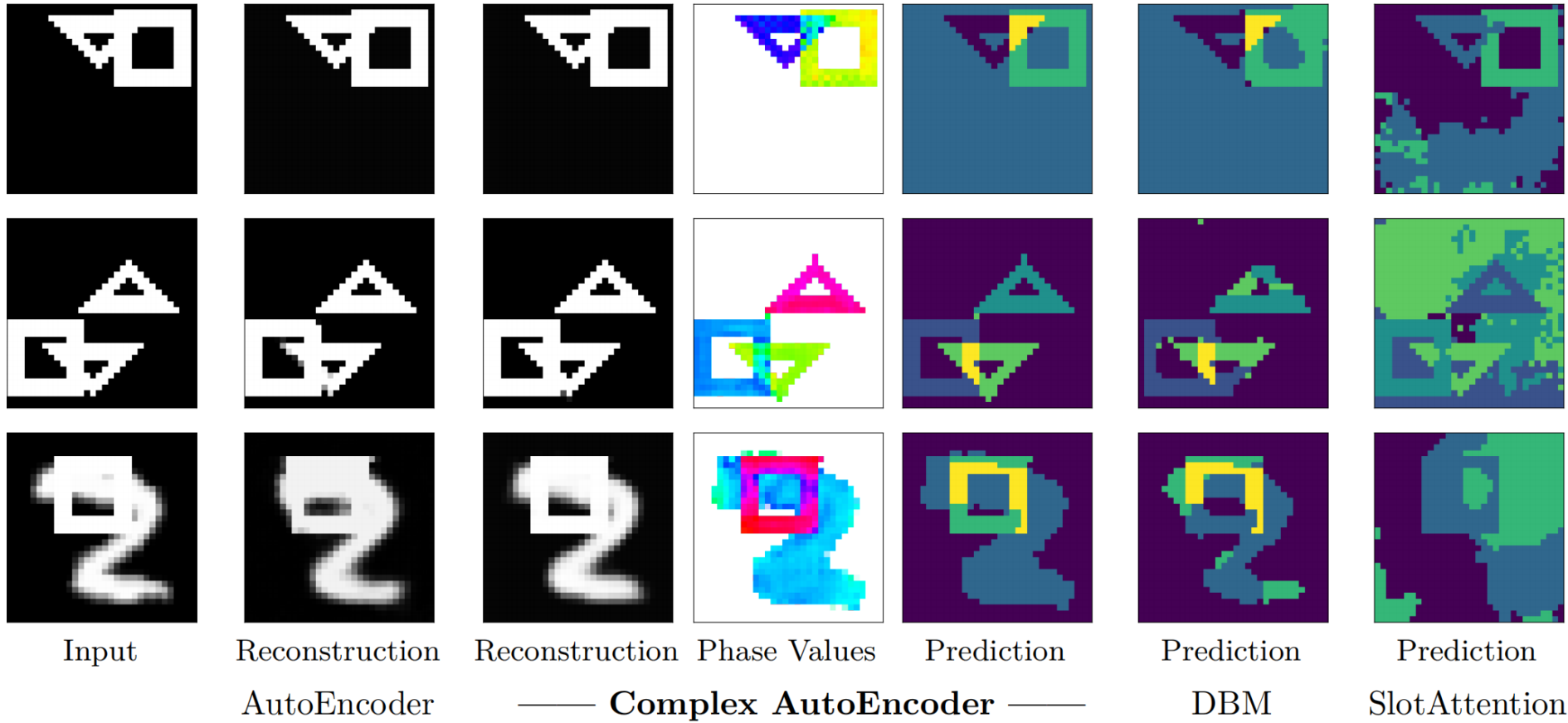
■ Complex-Valued Activation Function

$$\mathbf{z}' = \text{ReLU}(\text{BatchNorm}(\mathbf{m}_z)) \circ e^{i\varphi_\psi} \in \mathbb{C}^{d_{\text{out}}}$$

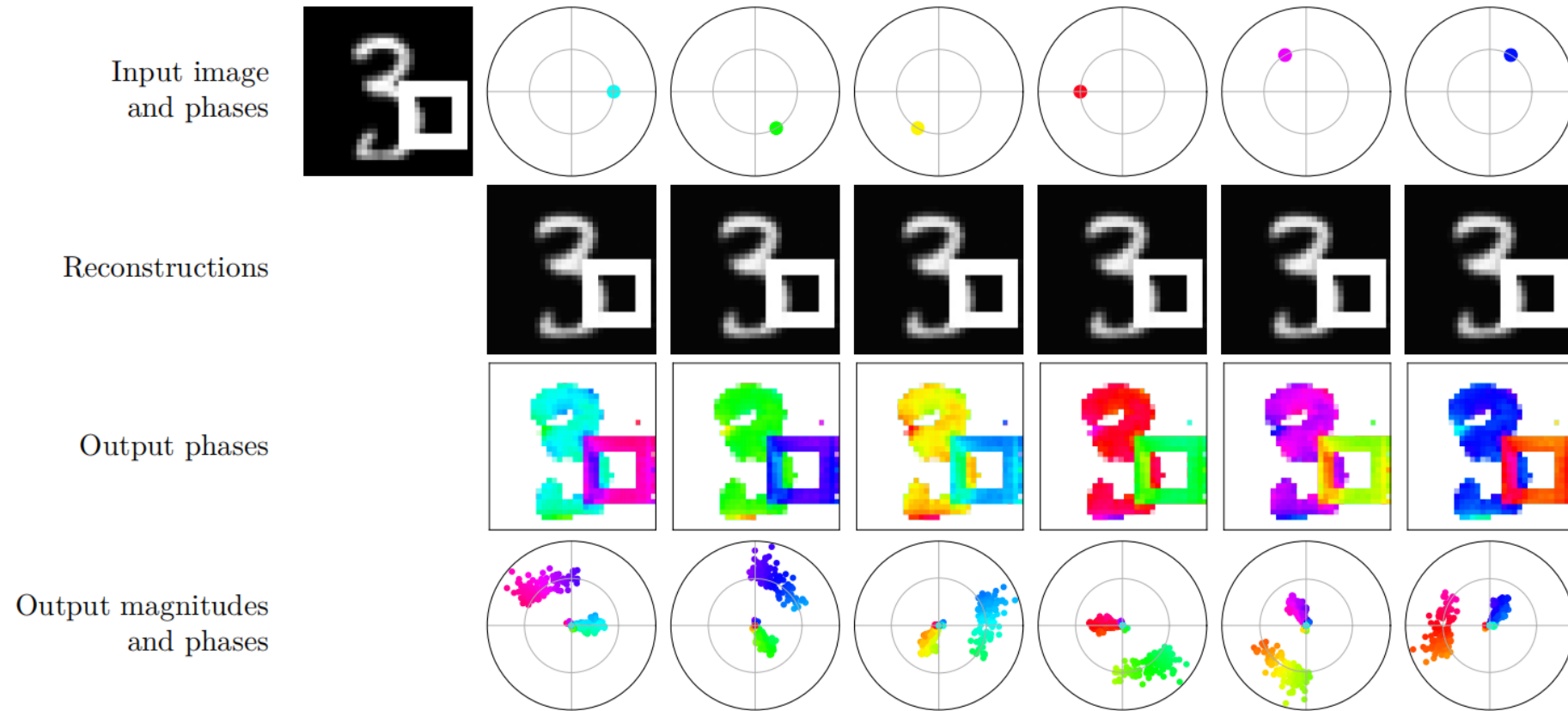
CAE



CAE



CAE

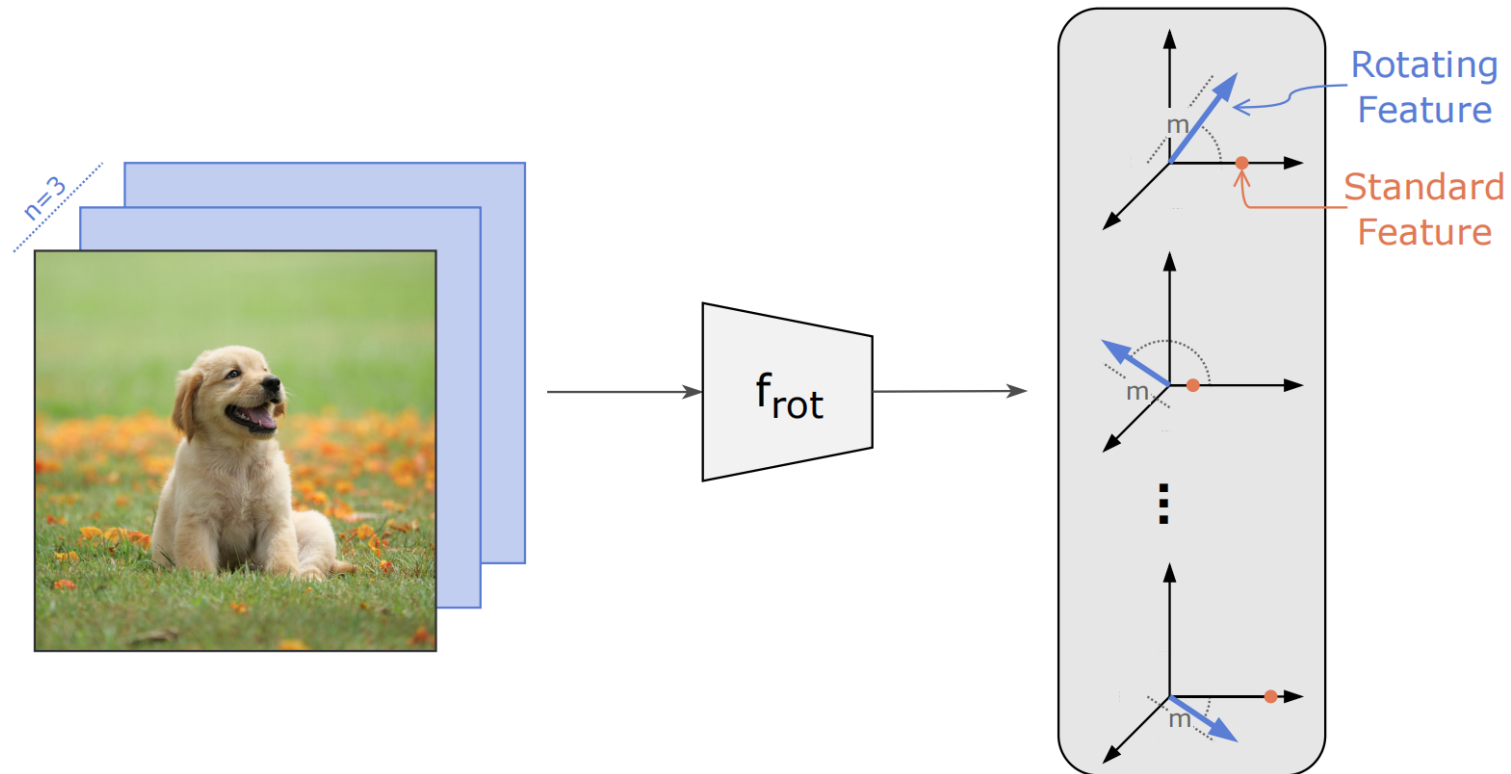


Outline

- Background
- Method
- Experiments
- Conclusion

Method

■ Rotating Features



Method

■ Rotating Features

$$\mathbf{z}_{\text{rotating}} \in \mathbb{R}^{n \times d} \quad \|\mathbf{z}_{\text{rotating}}\|_2 \in \mathbb{R}^d$$

Method

■ Rotating Features

$$\mathbf{z}_{\text{in}} \in \mathbb{R}^{n \times d_{\text{in}}} \quad \mathbf{w} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}} \quad \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

$$\boldsymbol{\psi} = f_{\mathbf{w}}(\mathbf{z}_{\text{in}}) + \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

Method

■ Rotating Features

$$\mathbf{z}_{\text{in}} \in \mathbb{R}^{n \times d_{\text{in}}} \quad \mathbf{w} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}} \quad \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

$$\boldsymbol{\psi} = f_{\mathbf{w}}(\mathbf{z}_{\text{in}}) + \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

$$\boldsymbol{\chi} = f_{\mathbf{w}}(\|\mathbf{z}_{\text{in}}\|_2) \in \mathbb{R}^{d_{\text{out}}}$$

$$\mathbf{m}_{\text{bind}} = 0.5 \cdot \|\boldsymbol{\psi}\|_2 + 0.5 \cdot \boldsymbol{\chi} \in \mathbb{R}^{d_{\text{out}}}$$

Method

■ Rotating Features

$$\mathbf{z}_{\text{in}} \in \mathbb{R}^{n \times d_{\text{in}}} \quad \mathbf{w} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}} \quad \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

$$\boldsymbol{\psi} = f_{\mathbf{w}}(\mathbf{z}_{\text{in}}) + \mathbf{b} \in \mathbb{R}^{n \times d_{\text{out}}}$$

$$\boldsymbol{\chi} = f_{\mathbf{w}}(\|\mathbf{z}_{\text{in}}\|_2) \in \mathbb{R}^{d_{\text{out}}}$$

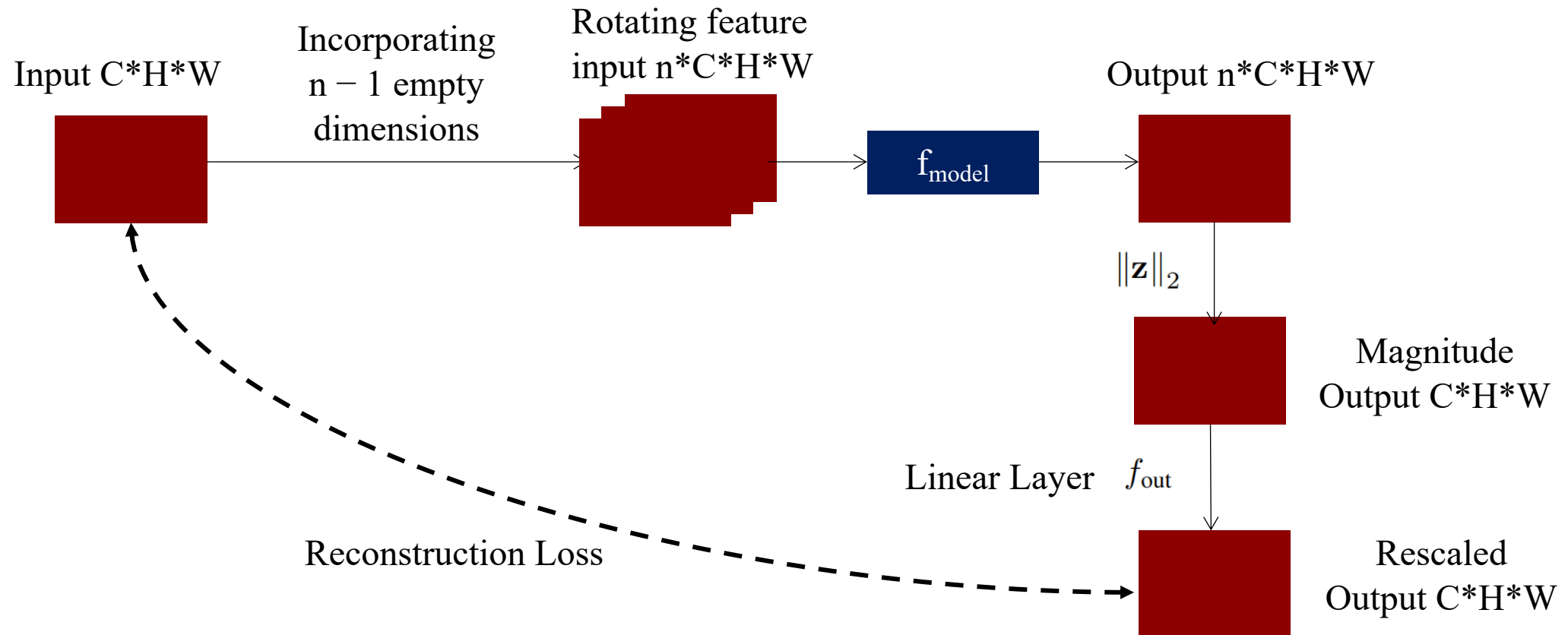
$$\mathbf{m}_{\text{bind}} = 0.5 \cdot \|\boldsymbol{\psi}\|_2 + 0.5 \cdot \boldsymbol{\chi} \in \mathbb{R}^{d_{\text{out}}}$$

$$\mathbf{m}_{\text{out}} = \text{ReLU}(\text{BatchNorm}(\mathbf{m}_{\text{bind}})) \in \mathbb{R}^{d_{\text{out}}}$$

$$\mathbf{z}_{\text{out}} = \frac{\boldsymbol{\psi}}{\|\boldsymbol{\psi}\|_2} \cdot \mathbf{m}_{\text{out}} \in \mathbb{R}^{n \times d_{\text{out}}}$$

Method

■ Training Process



Method

- **Evaluating Object Separation in Rotating Features**

Method

■ Evaluating Object Separation in Rotating Features

$$\mathbf{z}_{\text{norm}} = \frac{\mathbf{z}}{\|\mathbf{z}\|_2} \in \mathbb{R}^{n \times c \times h \times w}$$

Method

■ Evaluating Object Separation in Rotating Features

$$\mathbf{w}_{\text{eval}}^{i,j,l} = \begin{cases} 1 & \text{if } \|\mathbf{z}\|_2^{i,j,l} > t \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{z}_{\text{eval}} = \frac{\sum_{i=1}^c \mathbf{w}_{\text{eval}}^i \circ \mathbf{z}_{\text{norm}}^i}{\sum_{i=1}^c \mathbf{w}_{\text{eval}}^i + \varepsilon} \in \mathbb{R}^{n \times h \times w}$$

Outline

- Background
- Method
- **Experiments**
- Conclusion

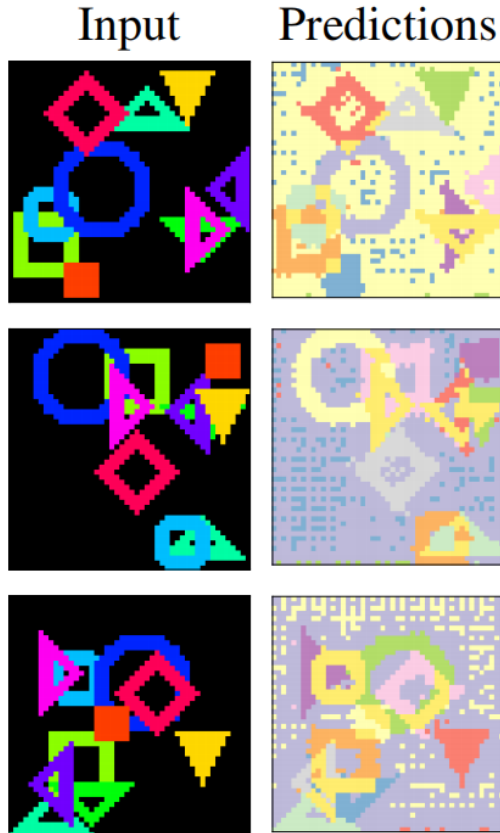
Experiments

■ Rotating Features Applied to Real-World Images



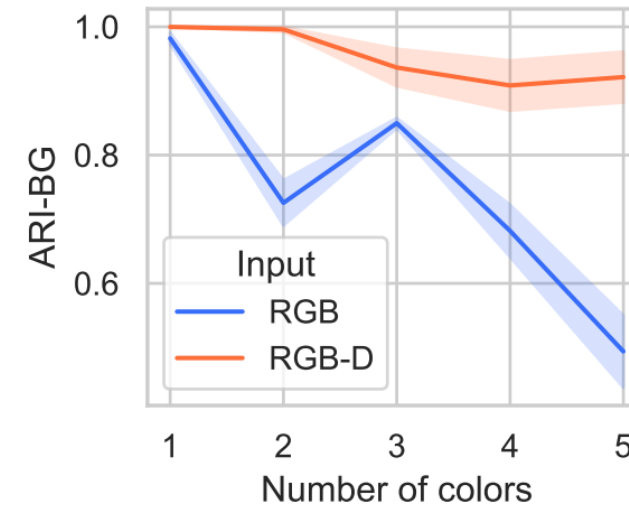
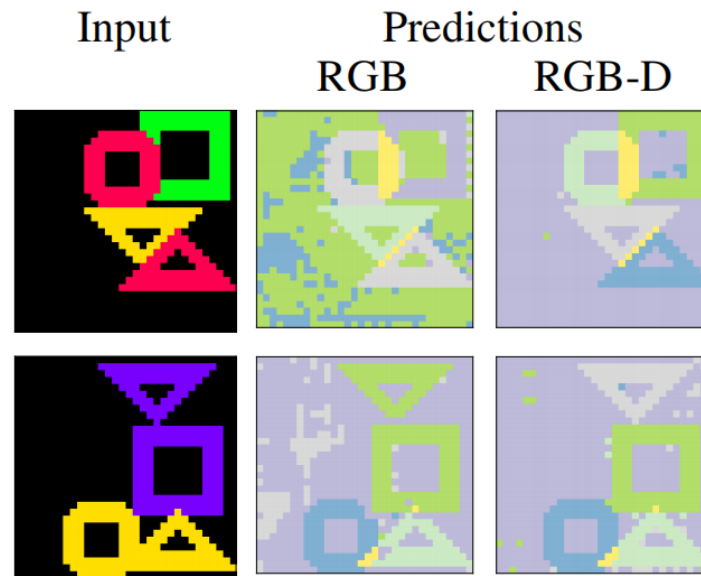
Experiments

■ Rotating Features Can Represent More Objects



Experiments

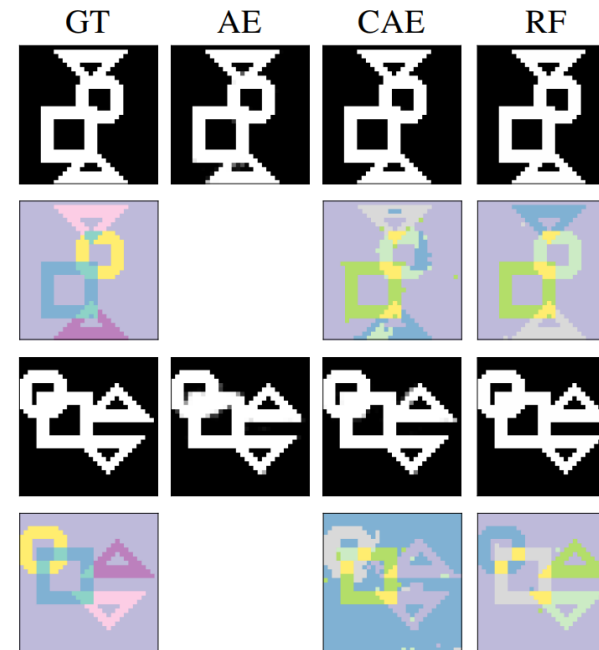
■ Rotating Features Are Applicable to Multi-Channel Images



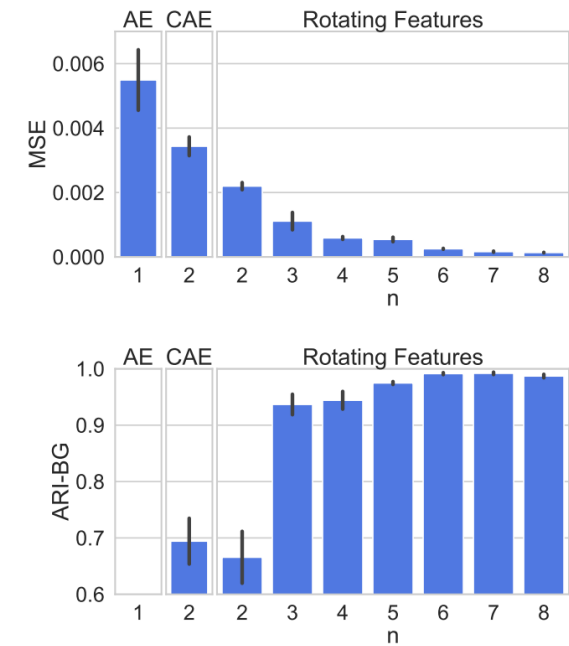
Experiments

■ Object Discovery

Model	$MBO_i \uparrow$	$MBO_c \uparrow$
Block Masks	0.247	0.259
Slot Attention	0.222 ± 0.008	0.237 ± 0.008
SLATE	0.310 ± 0.004	0.324 ± 0.004
Rotating Features –DINO	0.282 ± 0.006	0.320 ± 0.006
DINO k -means	0.363	0.405
DINO CAE	0.329 ± 0.009	0.374 ± 0.010
DINOSAUR Transformer	0.440 ± 0.008	0.512 ± 0.008
DINOSAUR MLP	0.395 ± 0.000	0.409 ± 0.000
Rotating Features	0.407 ± 0.001	0.460 ± 0.001



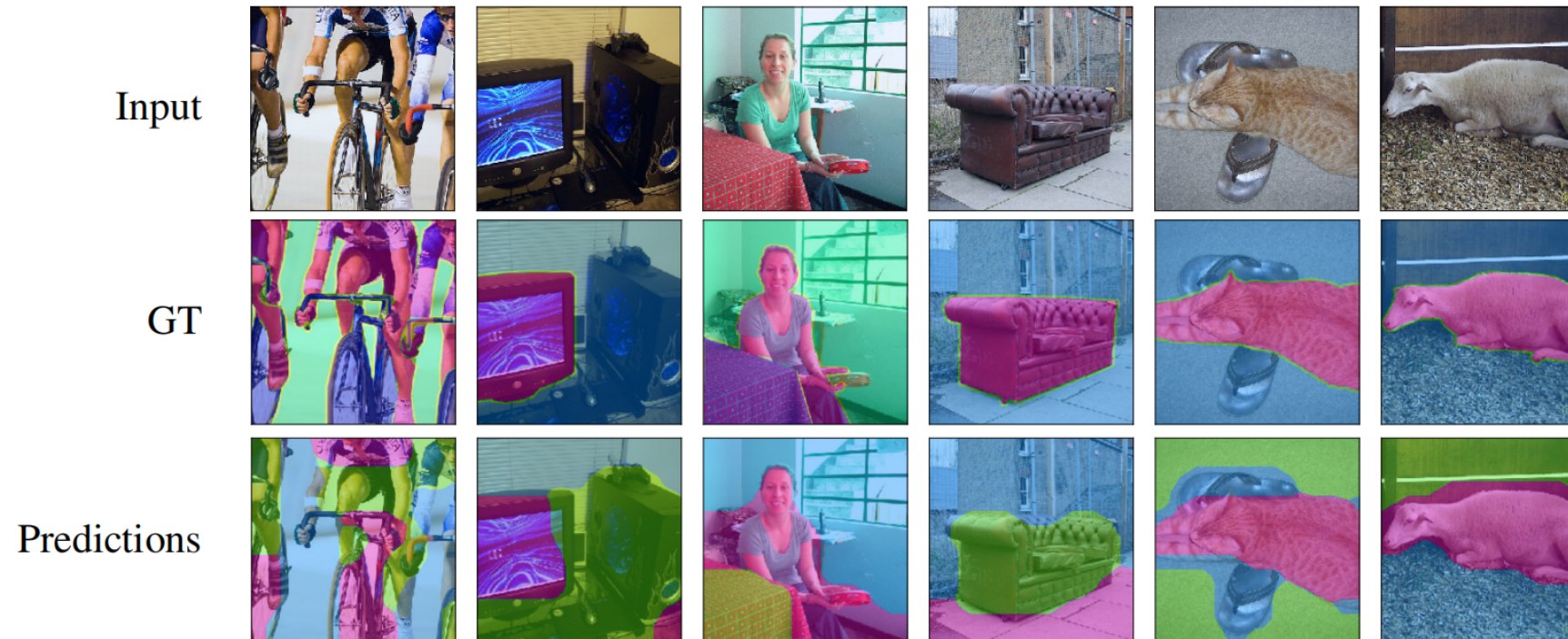
(a)



(b)

Experiments

■ Rotating Features Are Applicable to Real-World Images



Thanks!