

# Generative Modeling of Infinite Occluded Objects for Compositional Scene Representation

Jinyang Yuan, Bin Li, Xiangyang Xue

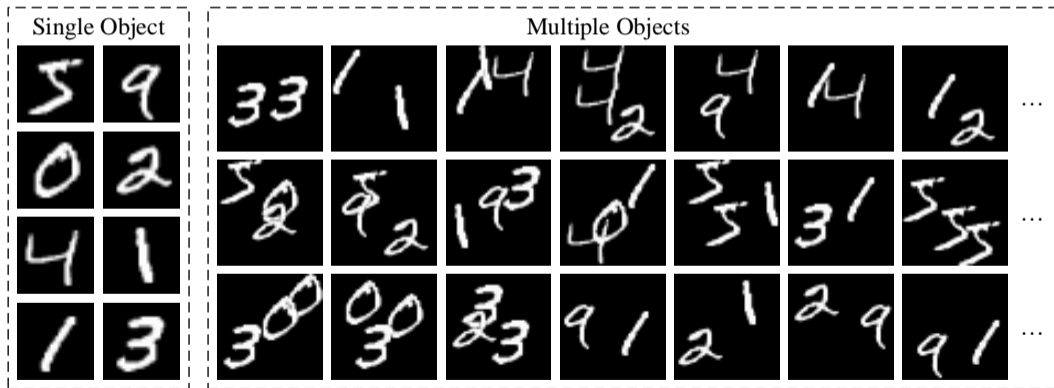
Fudan University

*{yuanjinyang, libin, xyxue}@fudan.edu.cn*

Jun 12, 2019

# Compositional Scene Representation

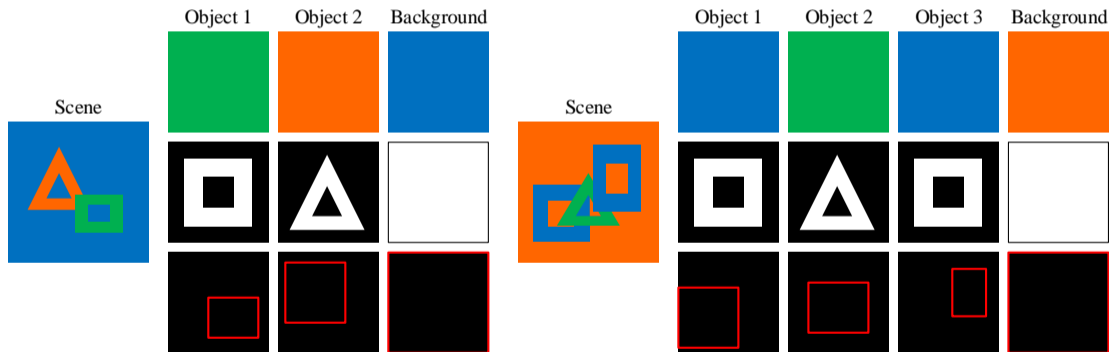
- **Scenes** are composed of **objects** and **background**
- The **combinations** of objects and background are **diverse**
- A **single representation** for the entire scene is **relatively complex**



# Compositional Scene Representation

Compositional scene representation is desirable

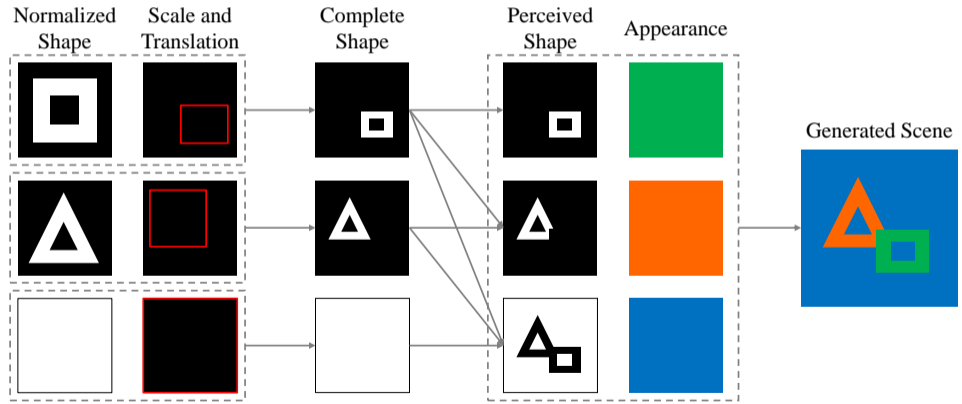
- Lower representation **complexity**
- Higher **generalizability** to novel scenes



# Generative Modeling of Infinite Occluded Objects

Two major difficulties

- The **number of objects** is unknown
- The perceived objects may be incomplete due to **occlusions**



# Generative Modeling of Infinite Occluded Objects

Background:  $k = 0$ ,      Objects:  $k \geq 1$

Latent Representation  $\mathbf{s}_{\cdot k} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \text{diag}(\tilde{\boldsymbol{\sigma}}^2)), \quad k \geq 0$

Presence (number of objects)  $\nu_k \sim \text{Beta}(\alpha, 1), \quad z_k^{\text{ind}} \sim \text{Ber}(\prod_{k'=1}^k \nu_{k'}), \quad k \geq 1$

Complete Shape  $z_{n,k}^{\text{dep}} \sim \text{Ber}(f_{\text{stn}}(\underbrace{f_{\text{shp}}(\mathbf{s}_{\cdot k}^{\text{shp}})}_{\text{normalized shape}}, \underbrace{\mathbf{s}_{\cdot k}^{\text{stn}}}_n), \quad k \geq 1$   
normalized shape    scale and translation

Perceived Shape (occlusions)  $\rho_{n,k} = \begin{cases} z_k^{\text{ind}} z_{n,k}^{\text{dep}} \prod_{k'=1}^{k-1} (1 - z_{k'}^{\text{ind}} z_{n,k'}^{\text{dep}}), & k \geq 1 \\ 1 - \sum_{k'=1}^{\infty} \rho_{n,k'}, & k = 0 \end{cases}$

Appearance  $\mathbf{a}_{n,k} = \begin{cases} f_{\text{apc}}^{\text{obj}}(\mathbf{s}_{\cdot k}^{\text{apc}}), & k \geq 1 \\ f_{\text{apc}}^{\text{back}}(\mathbf{s}_{\cdot k}^{\text{apc}}), & k = 0 \end{cases}$

Generated Scene  $\mathbf{x}_n \sim \sum_{k=0}^{\infty} \rho_{n,k} \mathcal{N}(\mathbf{a}_{n,k}, \hat{\sigma}^2 \mathbf{I})$

- Parameters are inferred by long short-term memories (**LSTMs**)
- Each object and background are updated **sequentially** and **iteratively**
- The LSTMs imitate the procedure of **coordinate ascent**

$$q(\mathbf{h}|\mathbf{x}) = q(\mathbf{s}_{.0}^{\text{apc}}) \prod_{k=1}^K \left( q(\mathbf{s}_{.k}^{\text{stn}}) q(\mathbf{s}_{.k}^{\text{shp}} | \mathbf{s}_{.k}^{\text{stn}}) q(\mathbf{s}_{.k}^{\text{apc}} | \mathbf{s}_{.k}^{\text{stn}}) q(\nu_k | \mathbf{s}_{.k}^{\text{stn}}) q(z_k^{\text{ind}} | \mathbf{s}_{.k}^{\text{stn}}) \prod_{n=1}^N q(z_{n,k}^{\text{dep}} | \mathbf{s}_{.k}^{\text{shp}}, \mathbf{s}_{.k}^{\text{stn}}) \right)$$

$$q(\mathbf{s}_{.k}^* | \mathbf{s}_{.k}^{\text{stn}}) = \mathcal{N}(\mathbf{s}_{.k}^*; \boldsymbol{\mu}_{.k}^*, \text{diag}(\boldsymbol{\sigma}_{.k}^{*2}))$$

$$q(\nu_k | \mathbf{s}_{.k}^{\text{stn}}) = \text{Beta}(\nu_k; \tau_{1,k}, \tau_{2,k})$$

$$q(z_k^{\text{ind}} | \mathbf{s}_{.k}^{\text{stn}}) = \text{Ber}(z_k^{\text{ind}}; \zeta_k)$$



$$q(z_{n,k}^{\text{dep}} | \mathbf{s}_{.k}^{\text{shp}}, \mathbf{s}_{.k}^{\text{stn}}) = \text{Ber}(z_{n,k}^{\text{dep}}; \xi_{n,k})$$



**Table:** Comparison of segregation and counting performance *with* existence of occlusion.

Data set	N-EM [Greff et al., 2017]			AIR [Eslami et al., 2016]			Proposed		
	AMI	MSE	OCA	AMI	MSE	OCA	AMI	MSE	OCA
Gray-S	77.3%	10e-3	56.2%	85.4%	6.5e-3	80.9%	<b>94.6%</b>	<b>2.9e-3</b>	<b>90.5%</b>
Gray-M	30.5%	22e-3	13.5%	62.8%	9.0e-3	66.0%	<b>71.1%</b>	<b>7.5e-3</b>	<b>77.6%</b>
RGB1-S	81.8%	5.6e-3	74.2%	95.3%	2.4e-3	88.8%	<b>98.3%</b>	<b>1.1e-3</b>	<b>95.1%</b>
RGB1-M	57.0%	9.4e-3	16.3%	78.2%	3.5e-3	67.9%	<b>82.0%</b>	<b>3.1e-3</b>	<b>74.8%</b>
RGB2-S	66.2%	9.0e-3	60.8%	85.7%	3.7e-3	84.4%	<b>92.3%</b>	<b>2.2e-3</b>	<b>86.3%</b>
RGB2-M	34.9%	13e-3	12.5%	64.1%	4.8e-3	69.8%	<b>67.9%</b>	<b>4.7e-3</b>	<b>71.0%</b>
RGB3-S	29.6%	21e-3	7.44%	91.3%	3.9e-3	90.3%	<b>97.4%</b>	<b>1.4e-3</b>	<b>92.5%</b>
RGB3-M	15.4%	22e-3	2.30%	67.5%	5.4e-3	60.5%	<b>77.9%</b>	<b>3.8e-3</b>	<b>68.6%</b>
RGB4-S	24.7%	20e-3	10.3%	86.7%	4.0e-3	78.3%	<b>90.7%</b>	<b>2.5e-3</b>	<b>83.3%</b>
RGB4-M	3.82%	32e-3	2.35%	56.9%	6.3e-3	58.2%	<b>67.9%</b>	<b>4.6e-3</b>	<b>77.3%</b>



-  Eslami, S., Heess, N., Weber, T., Tassa, Y., Szepesvari, D., Kavukcuoglu, K., and Hinton, G. E. (2016).  
Attend, infer, repeat: Fast scene understanding with generative models.  
*In Advances in Neural Information Processing Systems (NeurIPS)*, pages 3225–3233.
-  Greff, K., van Steenkiste, S., and Schmidhuber, J. (2017).  
Neural expectation maximization.  
*In Advances in Neural Information Processing Systems (NeurIPS)*, pages 6691–6701.

Thank You!