

# Attentive Normalization for Conditional Image generation

Yi Wang<sup>1</sup>, Ying-Cong Chen<sup>1</sup>, Xiangyu Zhang<sup>2</sup>, Jian Sun<sup>2</sup>, Jiaya Jia<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>MEGVII Technology



# What is Conditional Image generation?

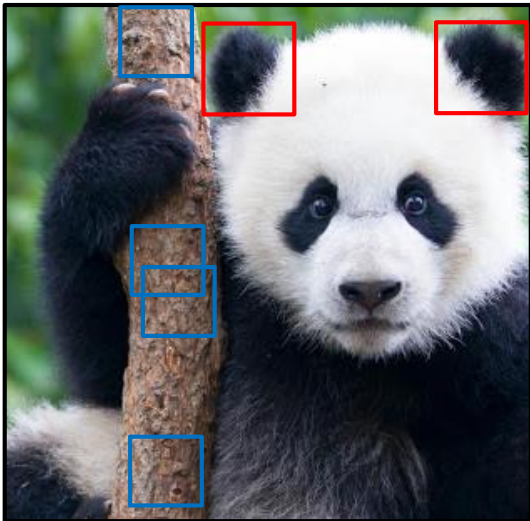
**panda**  
Input Class Label

Generation  

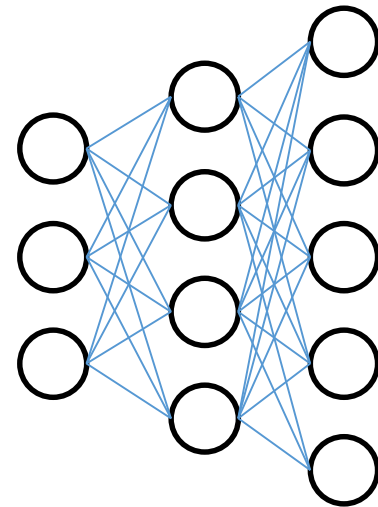



Generative model

# Long-range dependency in image generation

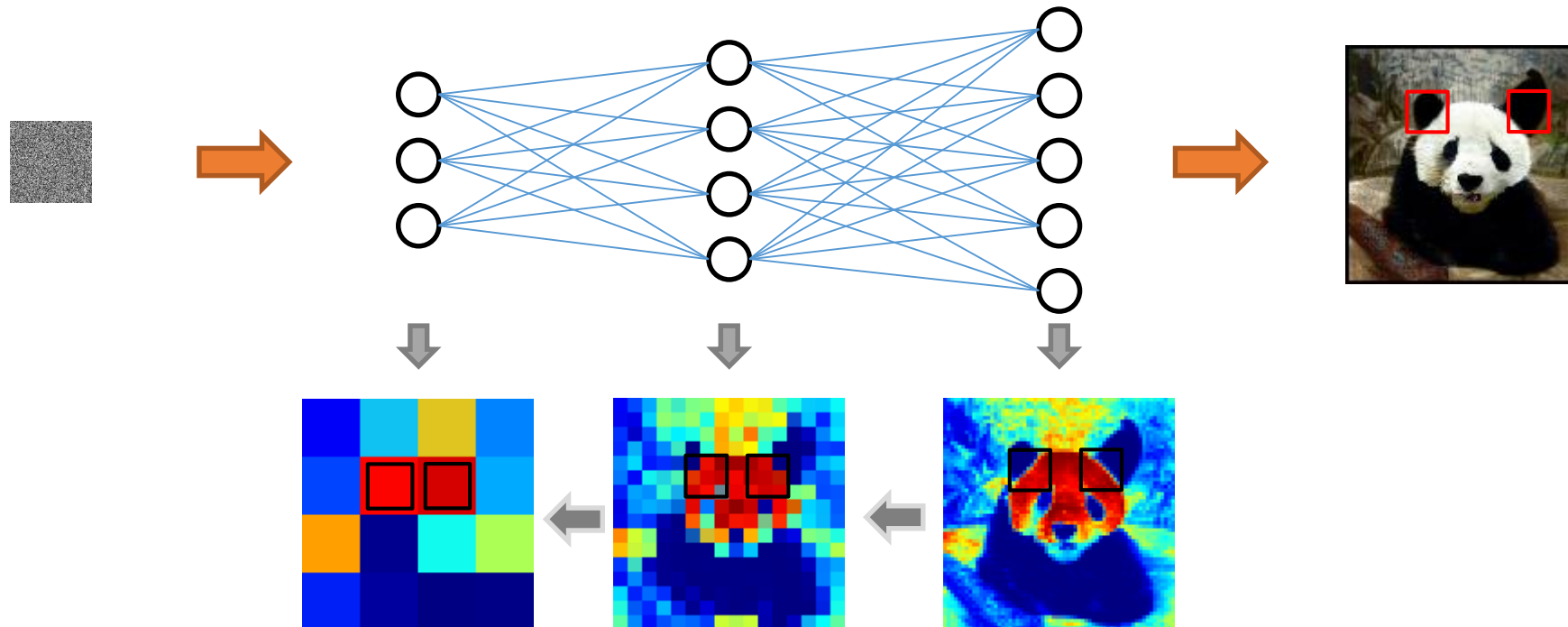


- **Standard convolutional neural network:**
  - **Modelling image contents in a hierarchical manner.**

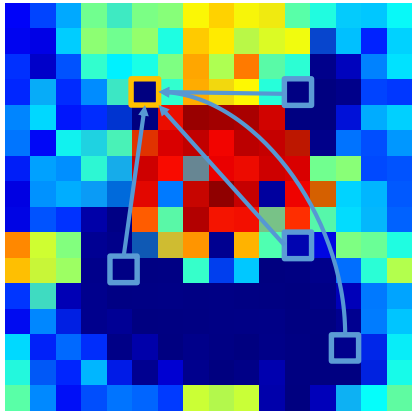


# Long-range dependency in image generation

- **Standard convolutional neural network:**
  - **Long-range dependency is conduct in a Markov chain.**



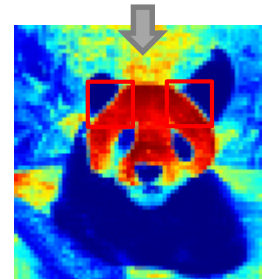
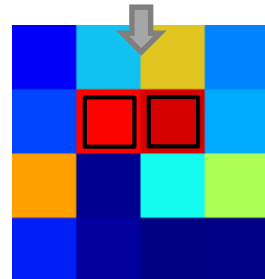
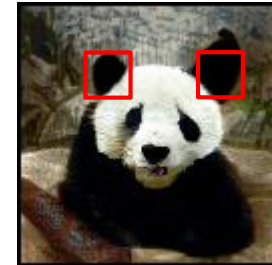
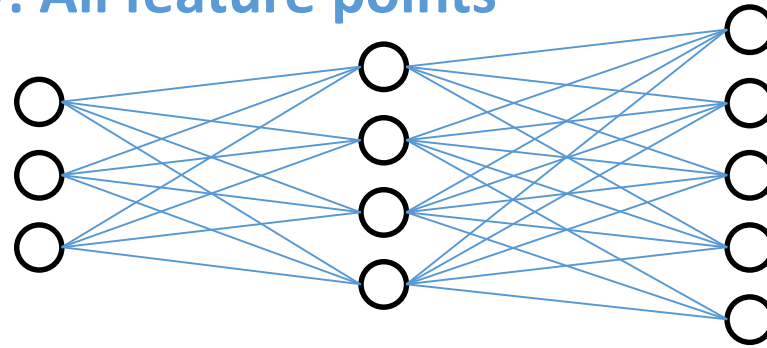
# Prior work: self-attention



**Self attention: reconstructing each feature point using the weighted sum of all feature points<sup>[1]</sup>.**

**Query: every feature point**

**Key: All feature points**



[1] Zhang, Han, et al. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318*, 2018.

# Our method: Attentive Normalization

## Core idea:

We **normalize** the input feature maps **spatially** according to the **semantic layouts** predicted from them.

Regional normalization

Semantic layout learning

## Empirical observations to backup our method:

- A feature map can be viewed as a composition of multiple semantic entities<sup>[3,4]</sup>.
- The deep layers in a neural network capture high-level semantics of the input images<sup>[5]</sup>.

[3] Greff, Klaus, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. In *NeurIPS*, 2017.

[4] Sabour, Sara, Nicholas Frosst, and Geoffrey E. Hinton. Dynamic routing between capsules. In *NeurIPS*, 2017.

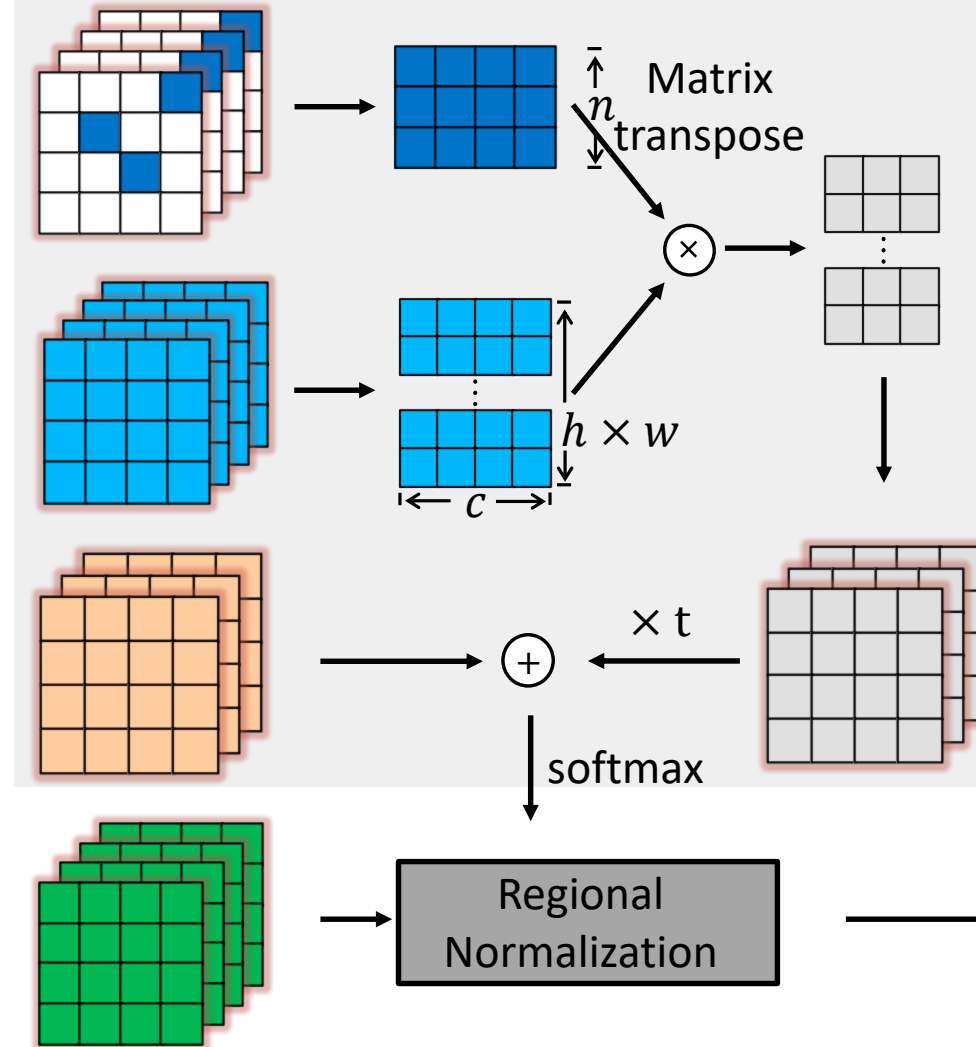
[5] Le, Quoc V. Building high-level features using large scale unsupervised learning. In *ICASSP*, 2013.



# Our method: Attentive Normalization

Semantic layout learning  
+  
Regional Normalization

Input feature maps  $X$



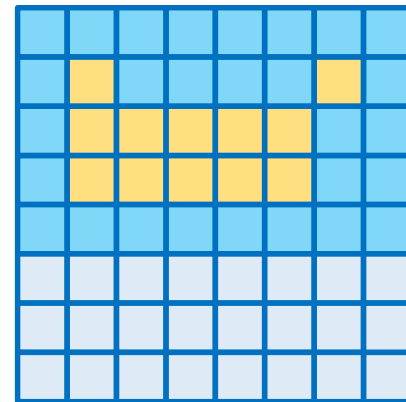
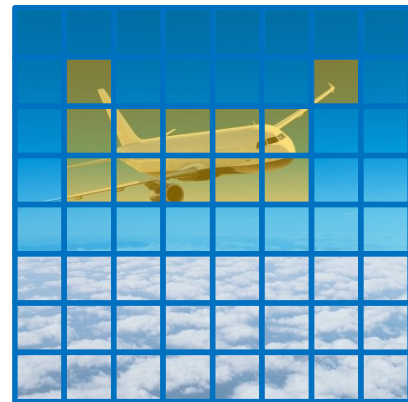
- Forward flow
- Convolution
- $\otimes$  Matrix multiplication
- Semantic layout learning






# Semantic layout learning

An image is composed of  $n$  semantic entities.

Each feature point of the image, it is determined by at least one entity.



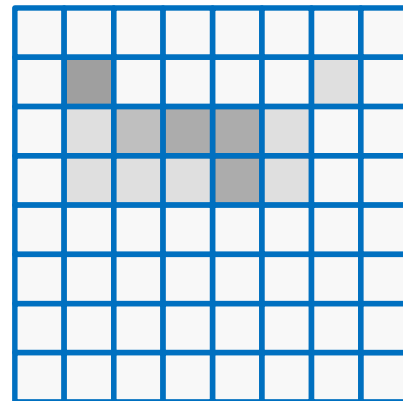
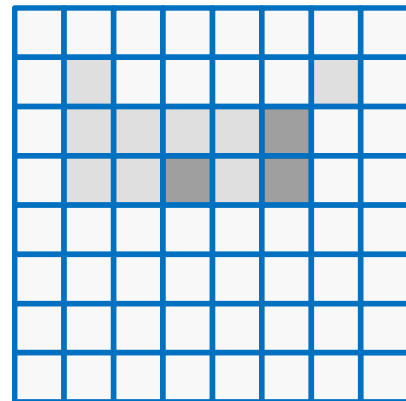
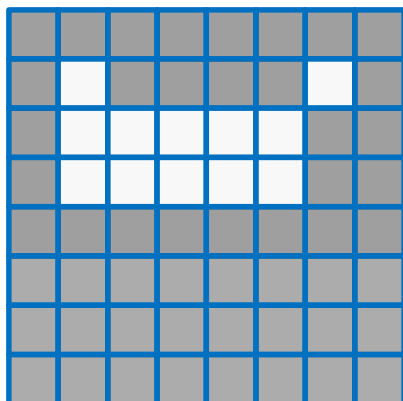
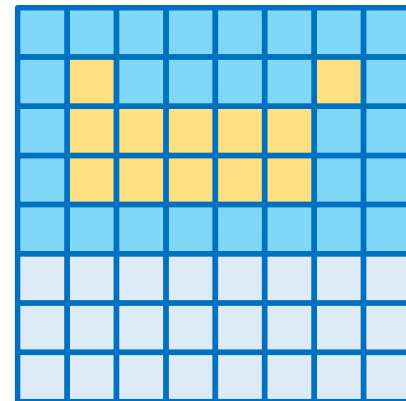
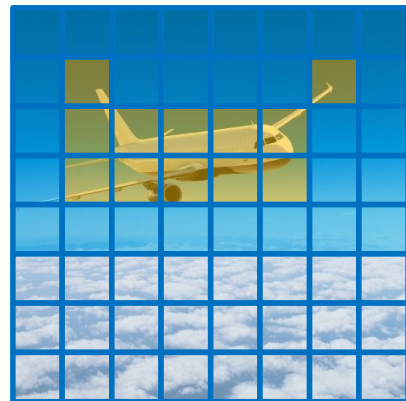
-  Sky
-  Plane
-  Cloud





# Semantic layout learning

How to group feature points of an image according to their correlation to the semantic entities.



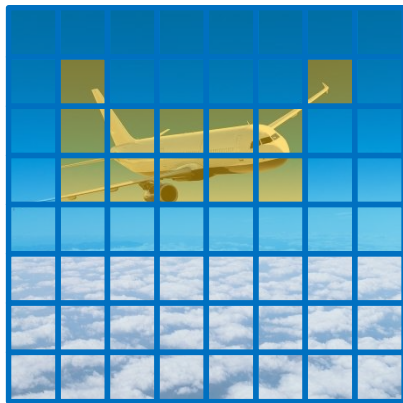
- Sky
- Plane
- Cloud

- Learned entities**
- Background
  - Round thing
  - Square thing
  - Flat thing



# Semantic layout learning

How to get these semantic entities?



Learned entities

Background

Round thing

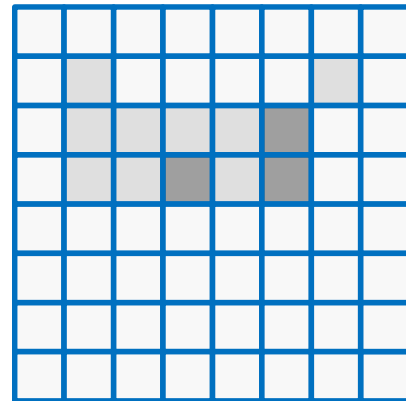
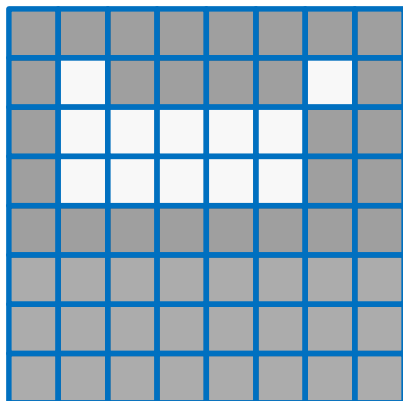
...

Square thing

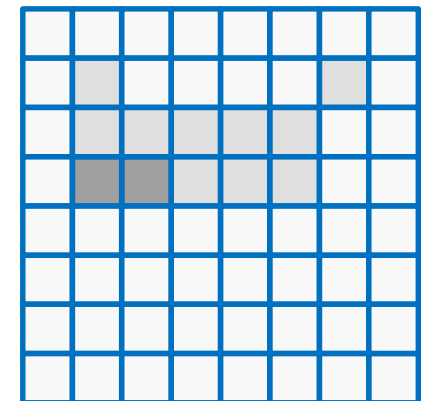
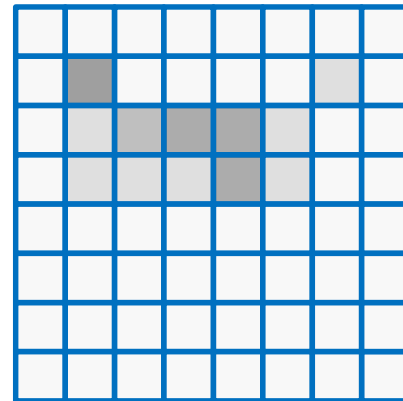
Flat thing

## Implementations

- A convolutional layer with  $n$  filters ( $n$  denotes a predefined class number)



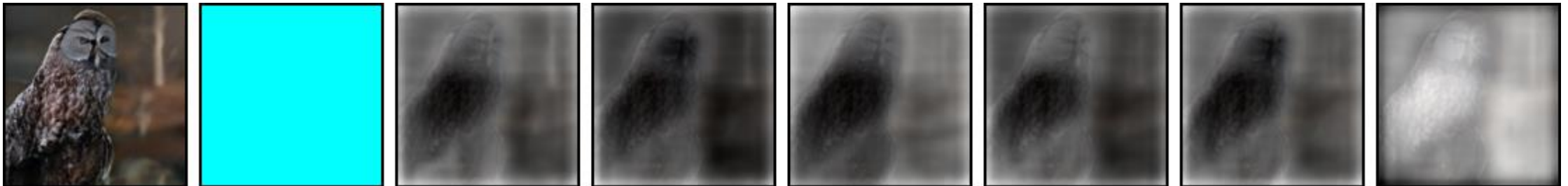
...



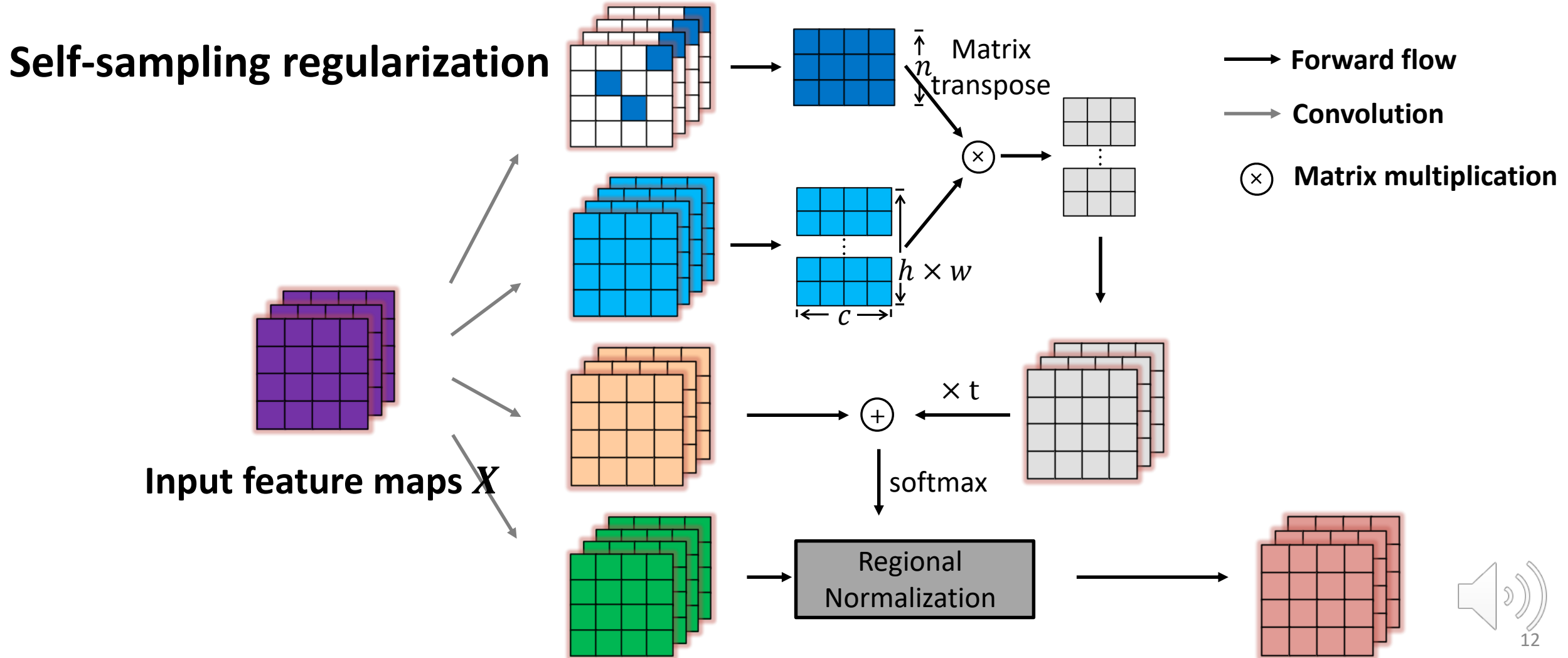
# Challenges in optimizing SSL

## Trivial solutions for learning entities directly

- It tends to group all feature points with a single semantic entities.
- No protocols are set to ban useless semantic entities.

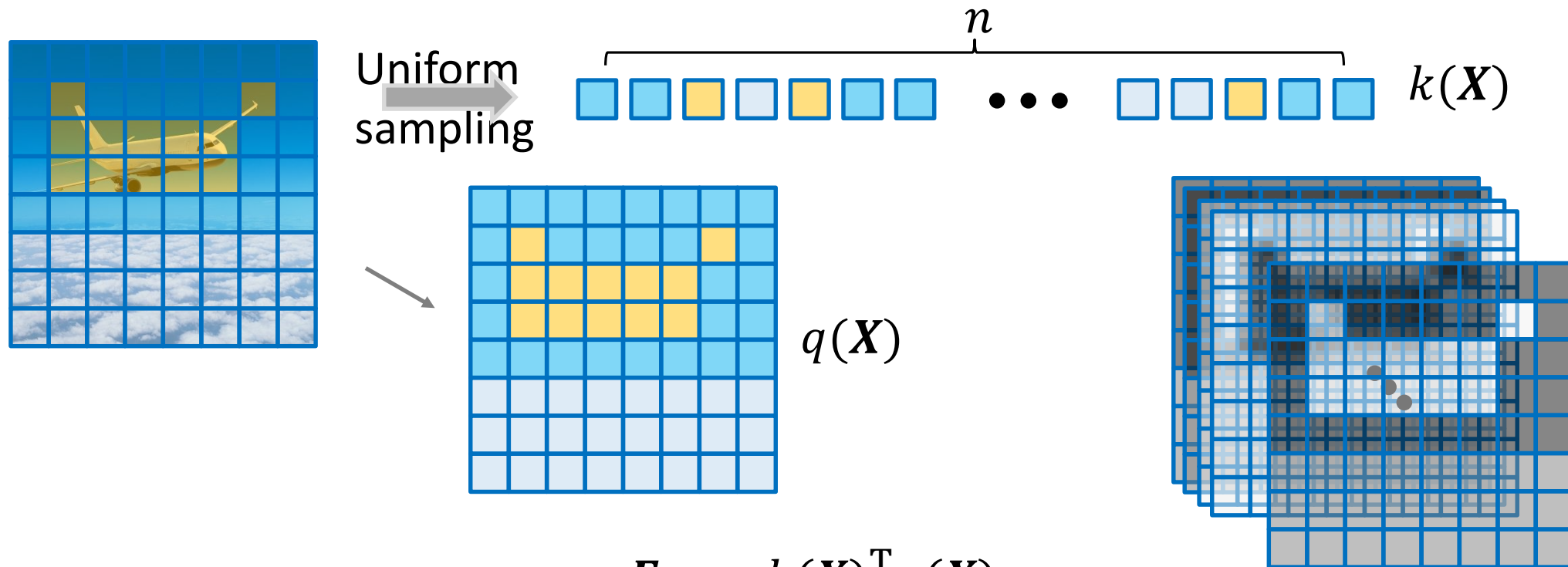


# Our method: Attentive Normalization



# Self-Sampling Regularization (SSR)

Regularizing semantics learning with a self-sampling branch

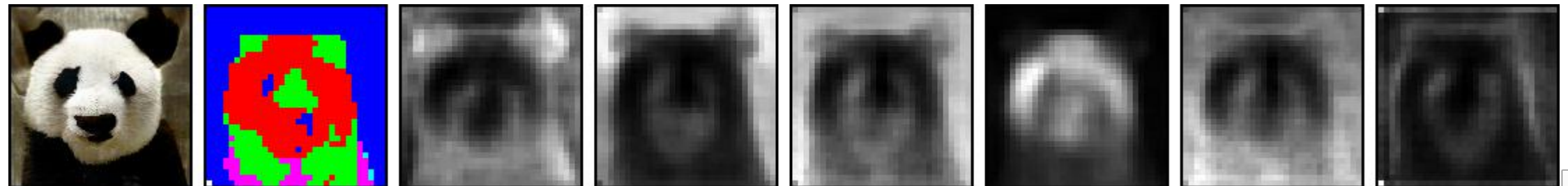
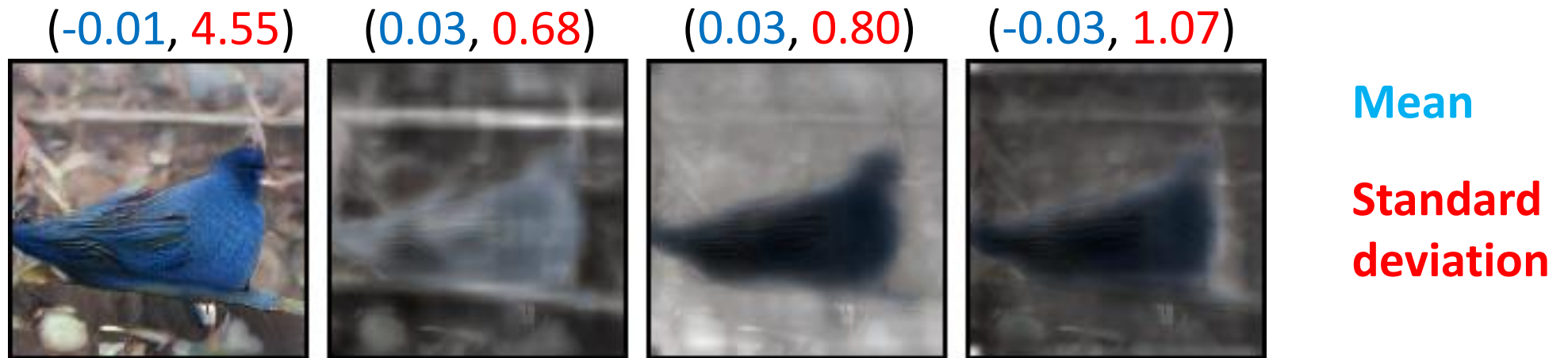


$$F_{i,j} = k(X)_i^T q(X)_j$$

where  $F \in \mathcal{R}^{h \times w \times n}$ ,  $q(X)$  are also translated feature maps.  $i$  and  $j$  denote pixel location. We set  $\#\{i\} = n$  and  $\#\{j\} = h \times w$ .

# Analysis: learned semantic layouts

The predicted semantic layout indicates regions with high inner coherence in semantics.



Generated

The highlighted regions indicated by the learned semantic layouts

# Analysis: complexity analysis

## Complexity Analysis

The computational complexity of Attentive Normalization:  $O(nMHW C)$

The computational complexity of Self-Attention:  $O(N(H^2W^2C + HWC^2))$

Module (ms)	128 x 128	256 x 256	512 x 512	1024 x 1024
AN (n=16)	<b>0.73</b>	<b>2.24</b>	<b>9.46</b>	<b>37.68</b>
Self-attention <sup>[1]</sup>	5.21	79.42	-	-

All fed tensors are with the same batch size 1 and channel number 32. Resolutions are different.

**'-' stands for evaluation time unmeasurable due to out-of-memory in GPU.**

Running environment: Pytorch 1.1.0, 4 CPUs, 1 TiTAN 2080 GPU, 32GB Memory.

[1] Zhang, Han, et al. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318*, 2018.



# Quantitative Results

## Class-conditional image generation On ImageNet (128x128):

	Itr x 1K↓	FID↓	Intra FID↓	IS↑
<b>AC-GAN</b> <sup>[6]</sup>	/	/	260.0	28.5
<b>SN-GAN</b> <sup>[7]</sup>	1000	27.62	92.4	36.80
<b>SN-GAN*</b> <sup>[1]</sup>	1000	22.96	/	42.87
<b>SA-GAN</b> <sup>[1]</sup>	1000	18.65	83.7	<b>52.52</b>
<b>Ours</b>	<b>880</b>	<b>17.84</b>	<b>83.4</b>	46.57

[1] Zhang, Han, et al. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318*, 2018.

[6] Odena, Augustus, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *ICML*, 2017.

[7] Miyato, Takeru, and Masanori Koyama. cGANs with projection discriminator. *arXiv preprint arXiv:1802.05637*. 2018.





# Qualitative Results

## Class-conditional image generation



Drilling platform (540)



Agaric (992)



Schooner (780)

## Image inpainting



## Categorical interpolation



Blenheim spaniel



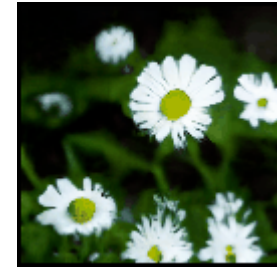
indigo hunting



coffee



Panda



Flower



# Thanks

