

# **Positional Encoding as Spatial Inductive Bias in GANs**

CVPR 2021

2021. 7. 8  
Dajin Han

# **Index**

- **Introduction**
- **Preliminaries**
- **Motivation**
- **Methods**
- **Experiments**
- **QA**

# Introduction

Key Contents

# Key Contents

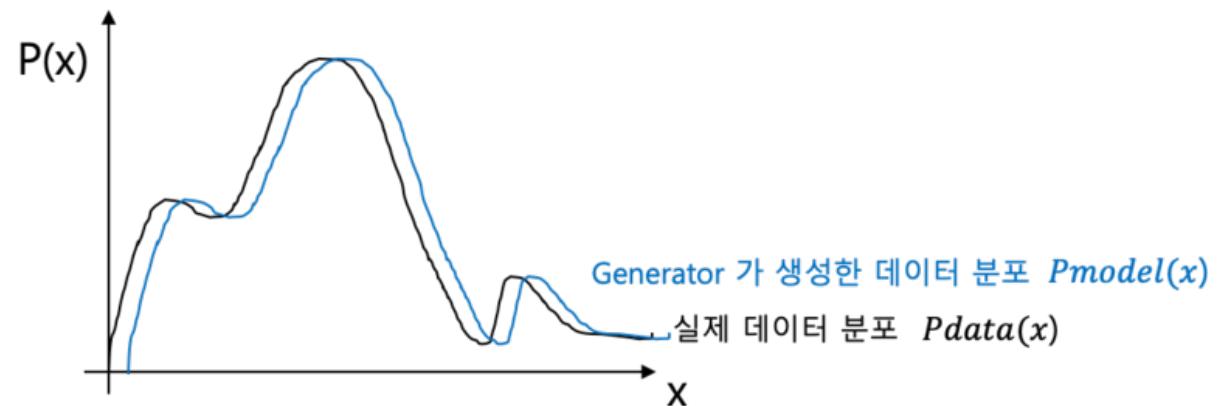
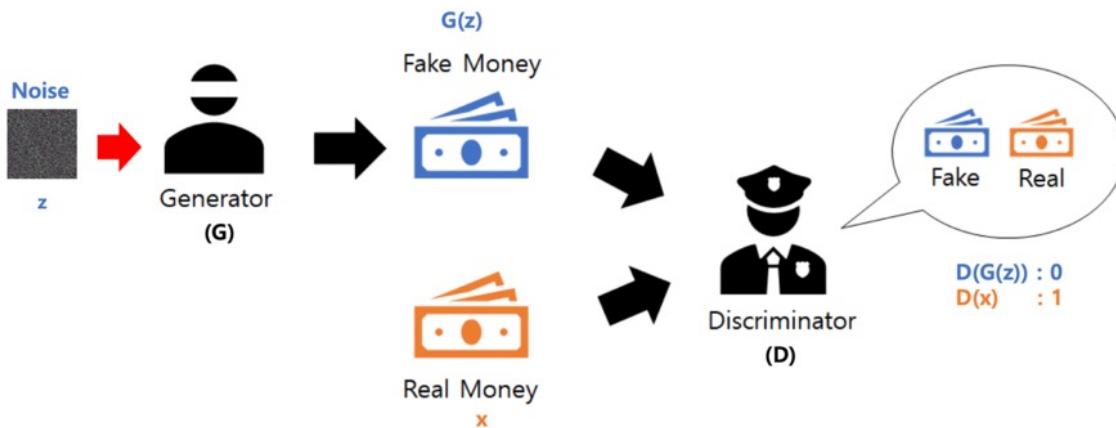
- Zero padding works as **implicit positional encoding**
  - Unintended design
  - Pros : Global structure information without center region
  - Cons : **Unbalanced spatial info**, Limits the diversity of a generative model
- Alternative **explicit positional encodings**
  - Cartesian Grid
  - Sinusoidal Positional Encoding
- New multi-scale training strategy
  - MS-PIE : **Scaling encoding**, not image

# Preliminaries

GAN, Positional Encoding, Padding effect

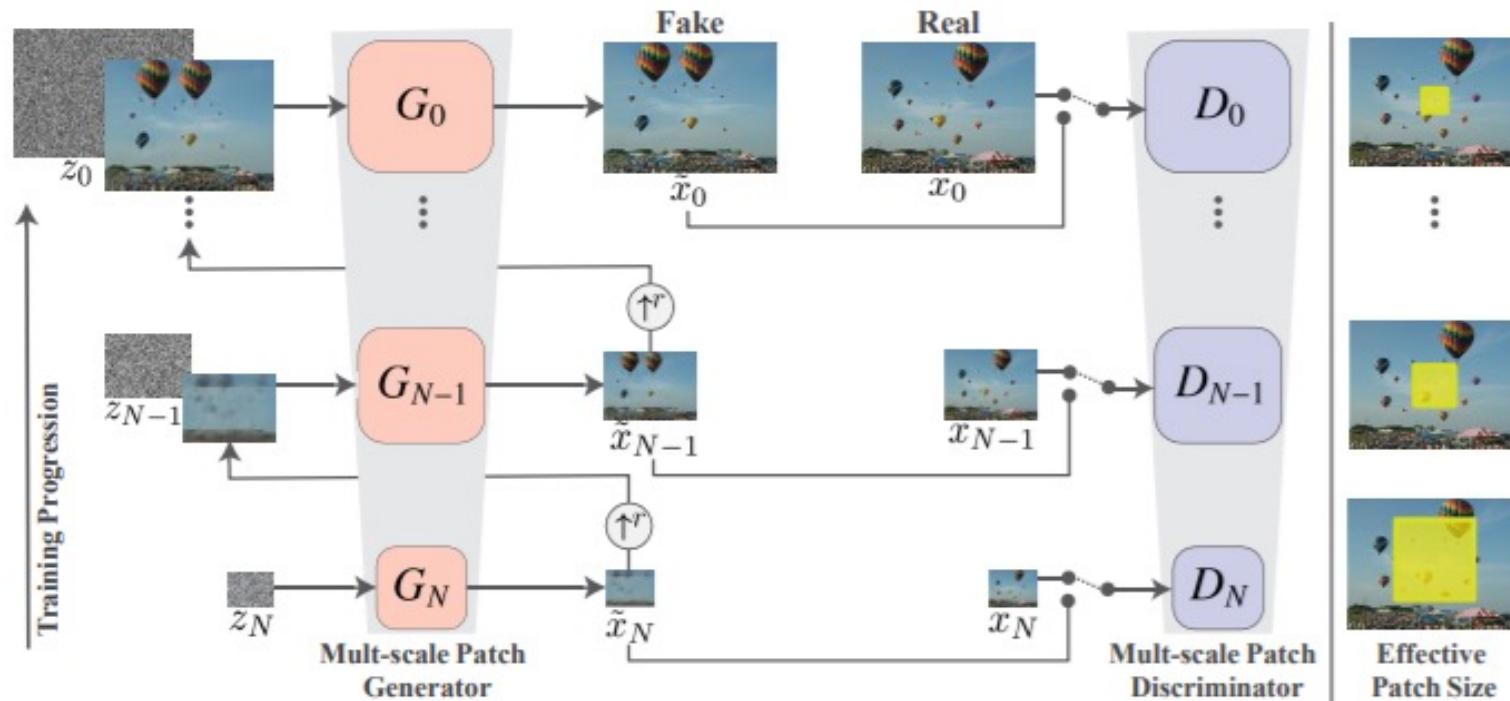
# GAN

- Generative Adversarial Network
- Minmax problem
- Training data distribution of training dataset



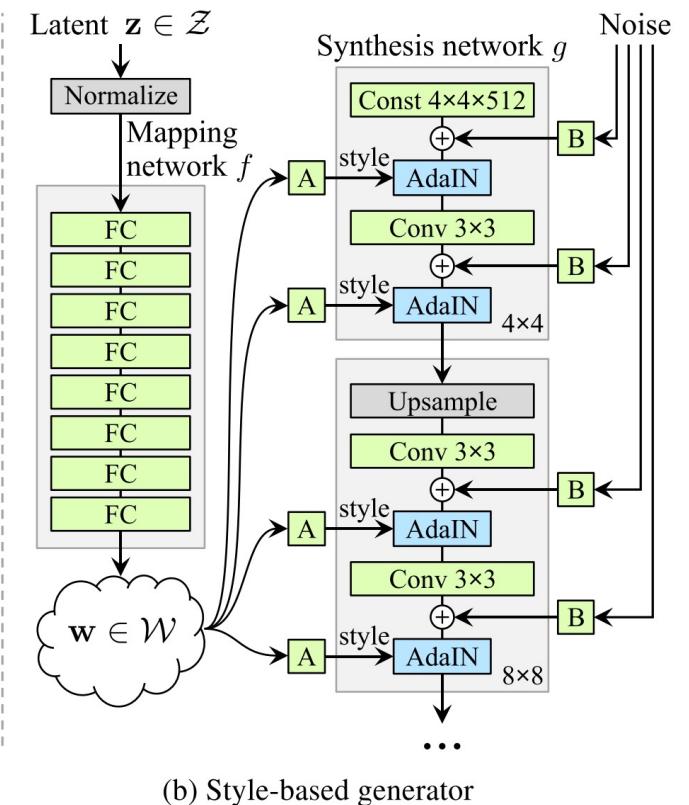
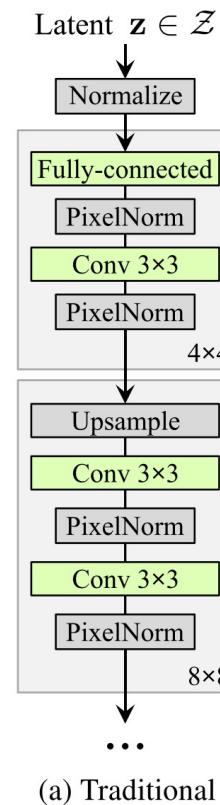
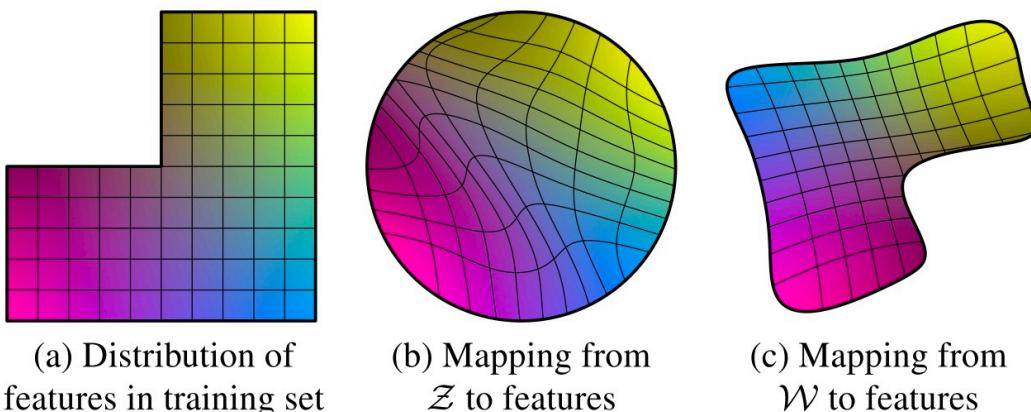
# SinGAN

- Learning a Generative Model from a Single Natural Image
- Pyramid of fully convolutional GANs
- Coarse to Fine



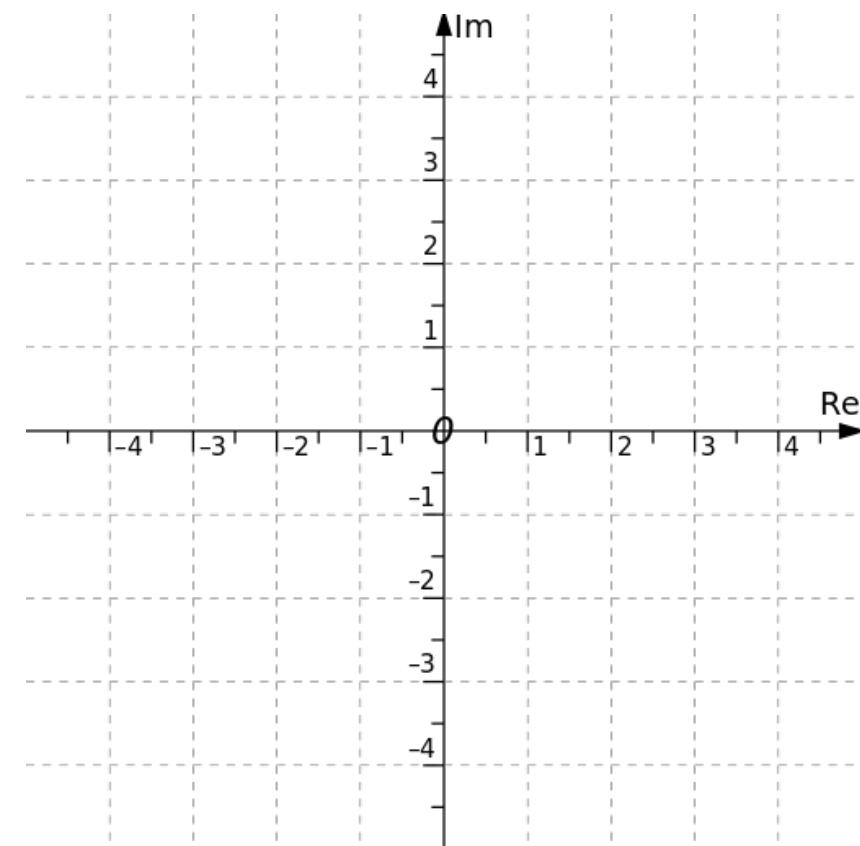
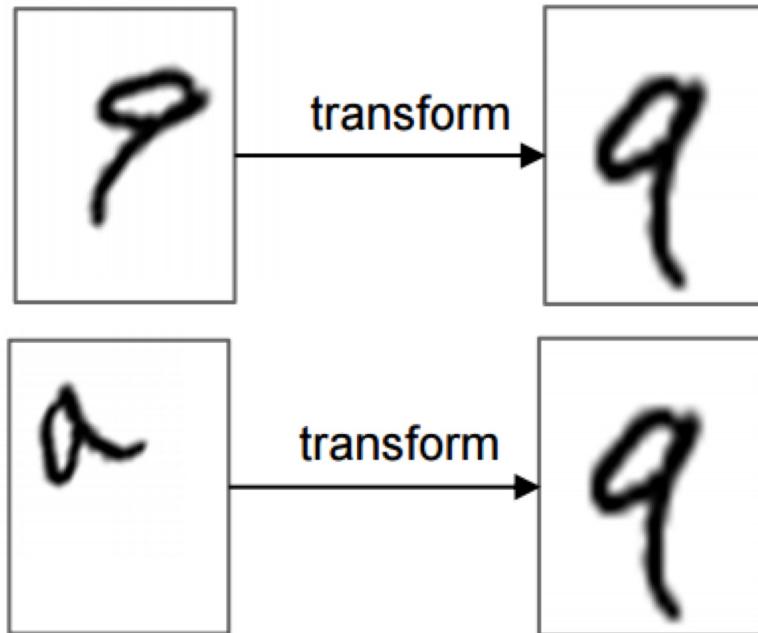
# StyleGAN

- A Style-Based Generator Architecture for GANs
- Intermediate latent vector
- More complex data distribution



# Cartesian Grid

- Spatial Transformer Network
- Spatially-Invariant

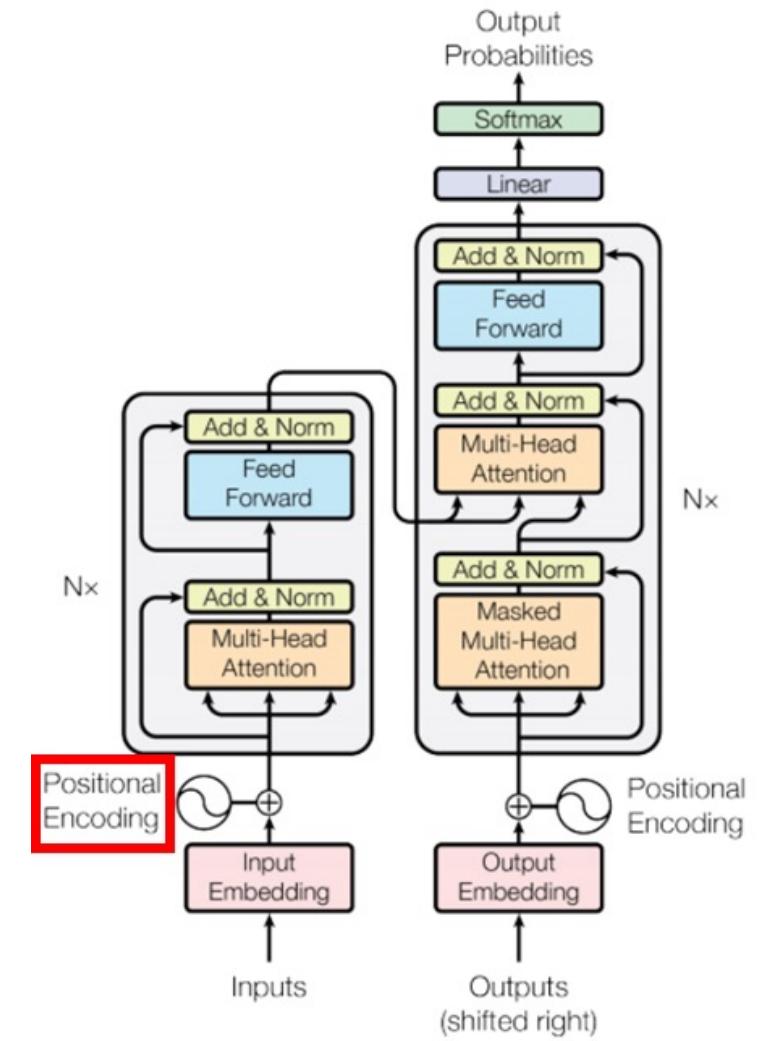
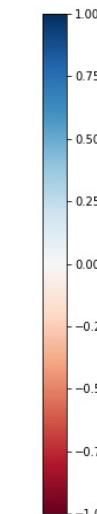
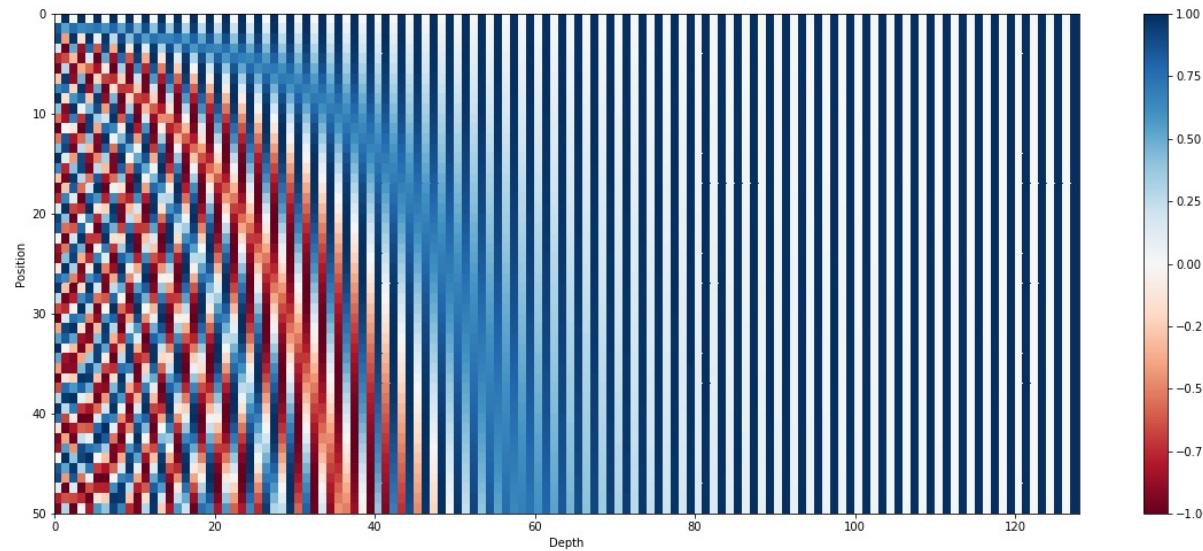


# Sinusoidal Positional Encoding

- NLP – Transformer
- Sequence of the input

$$PE_{(pos, 2k)} = \sin \left( pos / 10000^{2k/d} \right)$$

$$PE_{(pos, 2k+1)} = \cos \left( pos / 10000^{2k/d} \right)$$



# Motivation

Effect of zero padding, Statistical Property

# GAN Results

- Small receptive field, but good performance
- Weak center region
- Spatial bias spread from the **border to center** -> zero padding?



# GAN without Padding

- Fail to capture **global structure**
- Zero padding works as **positional encoding**



*Original Image*



**(a) SinGAN**



**(b) SinGAN w/o padding**

# Statistical property w/o padding

- $\mathbb{E}(y_i)$  : Distributional property of each location  $i$  in the convolutional feature map
- $R(y_i, y_j)$  : Relationship between two spatial locations

# Statistical property - $\mathbb{E}(y_i)$

- $\mathbb{E}(y_i)$  : Distributional property of each location  $i$  in the convolutional feature map
- Keep spatially identical expectation value

$$\mathbb{E}(\vec{y}_i^{(1)}) = \sum_k w_k^{(1)} \int_{-\infty}^{+\infty} x_k p(x_k) dx_k + b^{(1)}$$

$$= \sum_k w_k^{(1)} \mathbb{E}(x_k) + b^{(1)} = b^{(1)},$$

$$X_{\vec{i}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

$x_k$  : input feature

$p(\cdot)$  : probability density function

$w_k, b$  : parameters in the convolutional layer

$g(\cdot)$  : LeakyReLU

$$\mathbb{E}(\vec{y}_i^{(2)}) = \sum_k w_k^{(2)} \int_{-\infty}^{+\infty} g(y_k^{(1)}) p(y_k^{(1)}) dy_k^{(1)} + b^{(2)}$$

$$= \sum_k w_k^{(2)} \cdot (\gamma \mathbb{C}_1 + \mathbb{C}_2) + b^{(2)},$$

# Statistical property - $R(y_i, y_j)$

- $R(y_i, y_j)$  : Relationship between two spatial locations
- Irrelevant to the absolute position vector

$$\begin{aligned} R(\vec{y_i}, \vec{y_j}) &= \mathbb{E}(\vec{y_i} \vec{y_j}) \\ &= \sum_{x_l \in X_{\vec{i}} \cap X_{\vec{j}}} w_{k_l} w_{t_l} \mathbb{E}(x_l^2) \\ &= R(\vec{i} - \vec{j}). \end{aligned}$$

$$X_{\vec{i}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

$x_k$  : input feature

$w_k, b$  : parameters in the convolutional layer

# Statistical property w/o padding

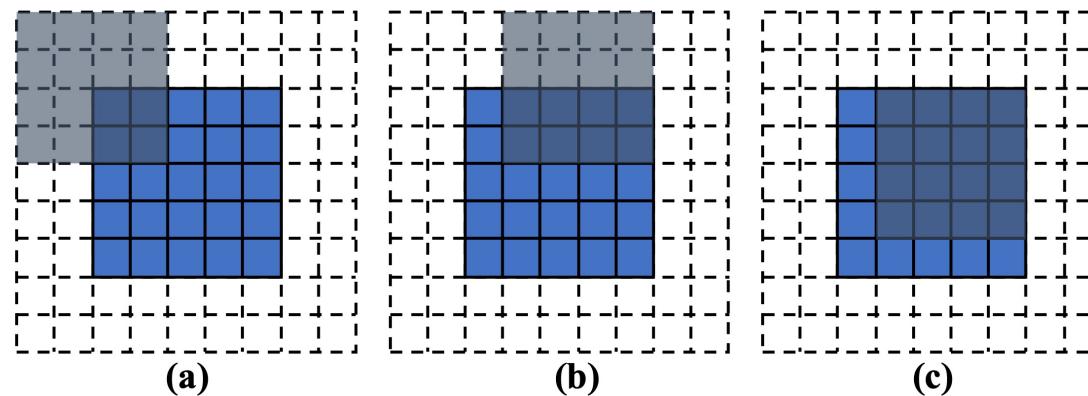
- $\mathbb{E}(y_i)$  : Distributional property of each location  $i$  in the convolutional feature map
- $R(y_i, y_j)$  : Relationship between two spatial locations
- **Weak stationary stochastic process (fail to capture global structure)**
- **Learn only small patch**

# Statistical property with padding

- Unique encoding on corners

$$\mathbb{E}(y_{\vec{i}}) = \sum_k w_k (\gamma \mathbb{C}_1 + \mathbb{C}_2) \mathbb{1}(x_k \notin Pad) + b,$$

$$R(y_{\vec{i}}, y_{\vec{j}}) = \sum_{x_l \in X_{\vec{i}} \cap X_{\vec{j}}} w_{k_l} w_{t_l} \mathbb{E}(x_l^2) \mathbb{1}(x_l \notin Pad),$$



[ ] zero padding    ■ feature map    ■ convolution kernel

# Method

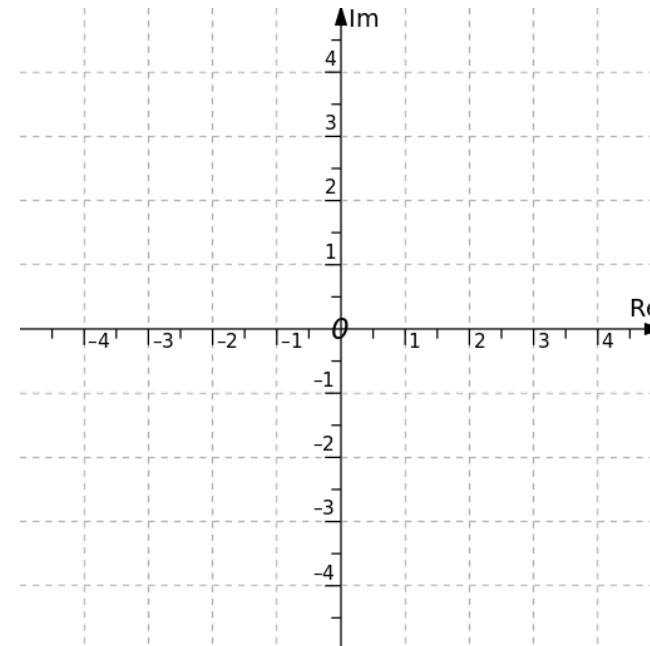
Learnable Constant Input, Cartesian Spatial Grid, Sinusoidal Positional Encoding, MS-PIE

# Learnable Constant Input

- StyleGAN
- “Learnable” constant as the input of generator
- Unclear and lacks explicit prior on image space

# Cartesian Spatial Grid

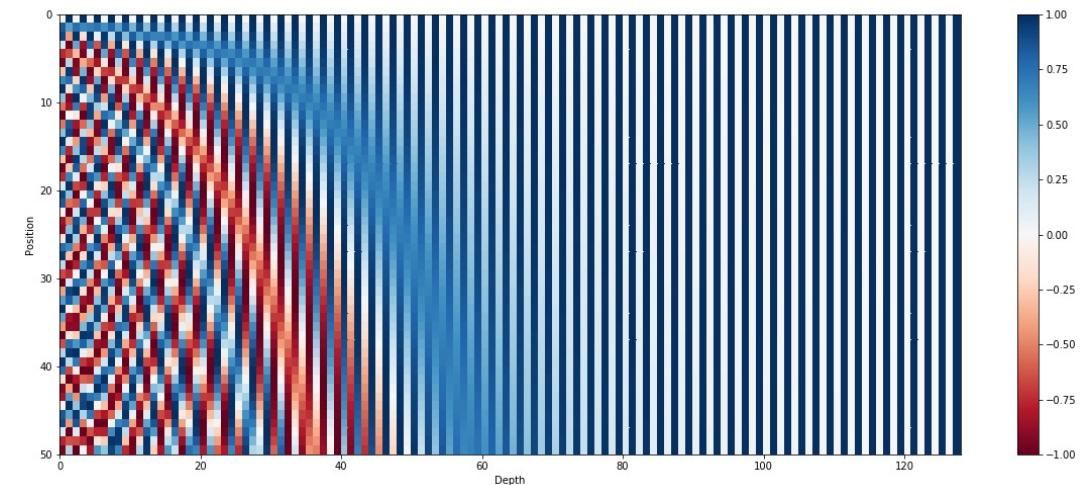
- Spatial Transformer Network (transformation)
- $[-1, -1]$  for top left,  $[1, 1]$  for bottom right
- Larger input scale, Closer distances between adjacent points
- **Robust to align global structure**



# Sinusoidal Positional Encoding

- NLP – Transformer
- Combination of sin and cos function
- **No spatial anchor**

$\underbrace{[\sin(\omega_0 i), \cos(\omega_0 i), \dots]}_{height\ dimension}, \underbrace{[\sin(\omega_0 j), \cos(\omega_0 j), \dots]}_{width\ dimension},$



# Positional Encodings

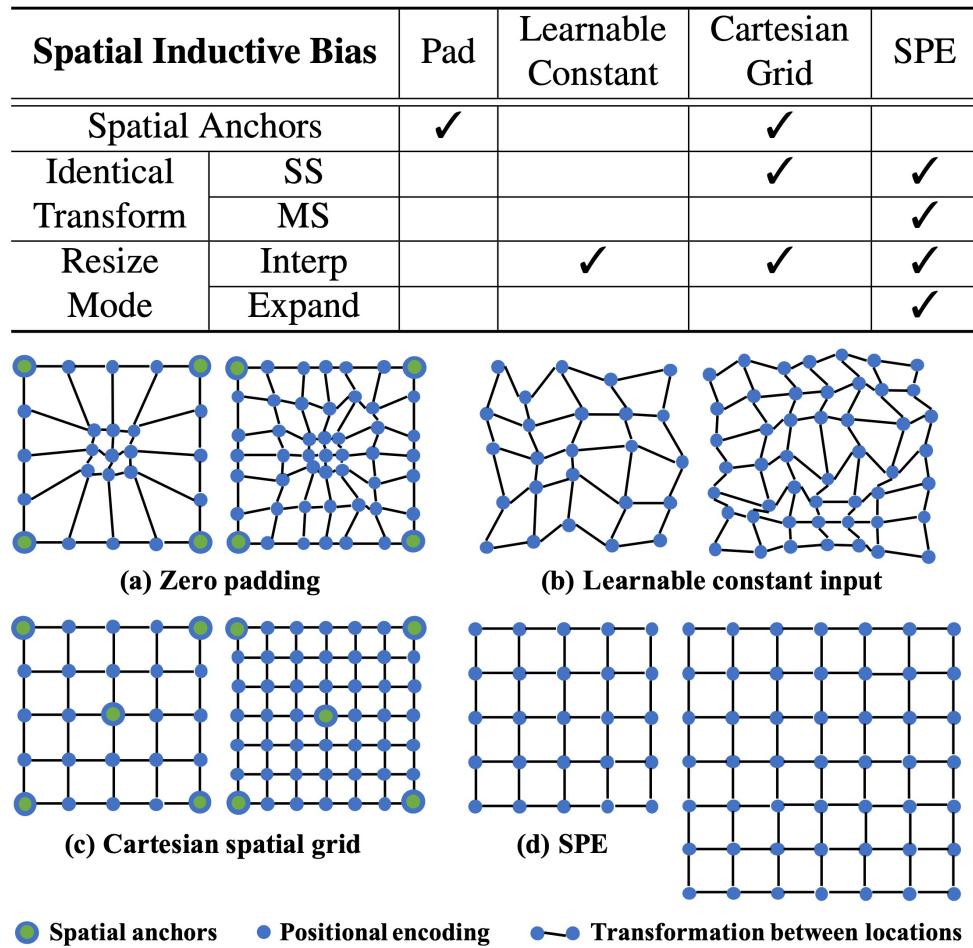


Table 1: Summary of spatial inductive bias defined by different positional encodings. ‘Identical Transform’ with ‘SS’ means the transformation between locations is spatially identical at a single scale. The ‘MS’ row shows whether the transformation is scale-invariant. ‘Interp’ is the traditional interpolation while ‘Expand’ denotes the expansion in SPE.

Figure 4: Illustration for the 2D spatial space defined by different positional encodings. For each encoding, we present two spatial spaces at  $5 \times 5$  and  $7 \times 7$  scale. (d) shows how SPE naturally expands its space with consistent transformation.

# MS-PIE

- Multi-Scale training with Positional Encoding
- Two methods
  - with zero padding
  - w/o zero padding

# MS-PIE – with zero padding

- Standard Model with zero padding
  - Capture global structure well
  - **Weak representation on central region**
- SPE (Spatial Positional Encoding)
  - Much **precise spatial inductive bias over dynamic input space**

# MS-PIE – w/o zero padding

- Standard Model without zero padding
  - Fail to capture global structure
  - Weak representation on central region
- CSG (Cartesian Spatial Grid)
  - Fixed anchors align global structure

# Experiments

Effect of zero padding, StyleGAN

# Experiments – zero padding

Training configuration	SWD ( $\times 10^3$ ) ↓			
	128	64	32	Avg.
(a) PGGAN	<b>3.162</b>	<b>4.285</b>	<b>5.000</b>	<b>4.149</b>
(b) + Remove padding	11.169	6.945	7.488	8.534
(c) + SPE, w/o padding	4.555	6.164	6.365	5.694

Training configuration	SWD ( $\times 10^3$ ) ↓			
	128	64	32	Avg.
(a) DCGAN w/ conv	<b>11.82</b>	17.30	28.10	19.07
(b) + No padding	22.21	27.26	44.24	31.24
(c) + SPE, w/o padding	14.75	<b>16.52</b>	<b>22.53</b>	<b>17.93</b>

# Experiments - StyleGAN

Table 2: Results based on  $256^2$  StyleGAN2 with a channel multiplier of two [22] in FFHQ dataset [21, 24]. The Fréchet inception distance (FID) are reported on two best snapshots before the discriminator has been shown with 10M and 20M images, respectively. The precision and recall are reported on the training snapshot with best FID.  $\uparrow$  indicates higher is better, and  $\downarrow$  indicates lower is better.

Training configuration	FID@256 $\downarrow$		Precision (%) $\uparrow$	Recall (%) $\uparrow$
	10M	20M		
(a) StyleGAN2-C2-256	6.32	<b>5.62</b>	<b>76.85</b>	50.41
(b) Deconv → Up-Conv	<b>5.79</b>	5.68	76.78	50.46
(c) + Remove padding	6.45	6.13	73.71	51.73
(d) + Cartesian grid	6.31	6.07	73.01	<b>52.90</b>
(e) + SPE	6.40	5.86	72.94	<b>52.87</b>

# Experiments – MS-PIE

Table 5: Main results for our MS-PIE with  $256^2$  StyleGAN2 in the FFHQ dataset. The precision and recall are calculated at the same scale as the Fréchet inception distance (FID). ‘C2’ indicates the channel multiplier in the generator is two.

Training configuration		Resize	FID@512↓ 20M 25M		Precision (%)↑	Recall (%)↑	FID@256↓ 20M 25M		Precision (%)↑	Recall (%)↑
(a) StyleGAN2-C2-256			—	—	—	—	5.62	5.56	75.92	51.24
(b) StyleGAN2-C2-512		—	3.47	3.41	75.88	54.61	5.00	4.91	75.65	54.58
MS-PIE w/ padding	(c) Learnable constant input	Interp	3.46	3.35	73.84	55.77	4.82	4.50	72.75	55.42
	(d) Cartesian spatial grid	Interp	3.59	3.50	73.28	56.16	<b>4.74</b>	4.71	73.34	55.07
	(e) SPE-interp	Interp	3.41	3.15	<b>74.13</b>	56.88	4.99	4.73	73.28	<b>56.93</b>
	(f) SPE-expand	Expand	<b>3.31</b>	<b>2.93</b>	73.51	<b>57.32</b>	4.79	<b>4.27</b>	<b>73.48</b>	55.69
	(g) SPE-expand-C1	Expand	3.65	3.40	73.05	56.45	5.54	4.83	73.59	54.29
MS-PIE w/o padding	(h) Learnable constant input	Interp	4.99	4.01	72.81	54.35	5.67	5.11	72.37	55.52
	(i) Cartesian spatial grid	Interp	<b>3.96</b>	<b>3.76</b>	<b>73.26</b>	<b>54.71</b>	<b>5.30</b>	<b>5.09</b>	70.74	56.06
	(j) SPE-interp	Interp	4.80	4.23	73.11	54.63	5.89	5.38	71.21	<b>56.21</b>
	(k) SPE-expand	Expand	4.46	4.17	73.05	51.07	6.08	5.59	<b>72.65</b>	49.74



Original Image

(a) SinGAN

(b) SinGAN w/ Cartesian grid

(c) SinGAN w/ SPE

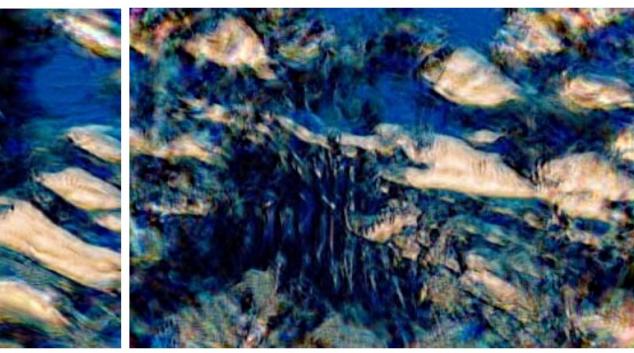
# Results



*Original Image*



**(c) SinGAN w/ padding**



**(d) SinGAN w/o padding**



**(e) SinGAN w/ Cartesian Grid**



**(f) SinGAN w/ SPE**

# QA