

# Structure-Aware Manipulation of Images and Videos

Sagie Benaim

School of Computer Science, Tel Aviv University



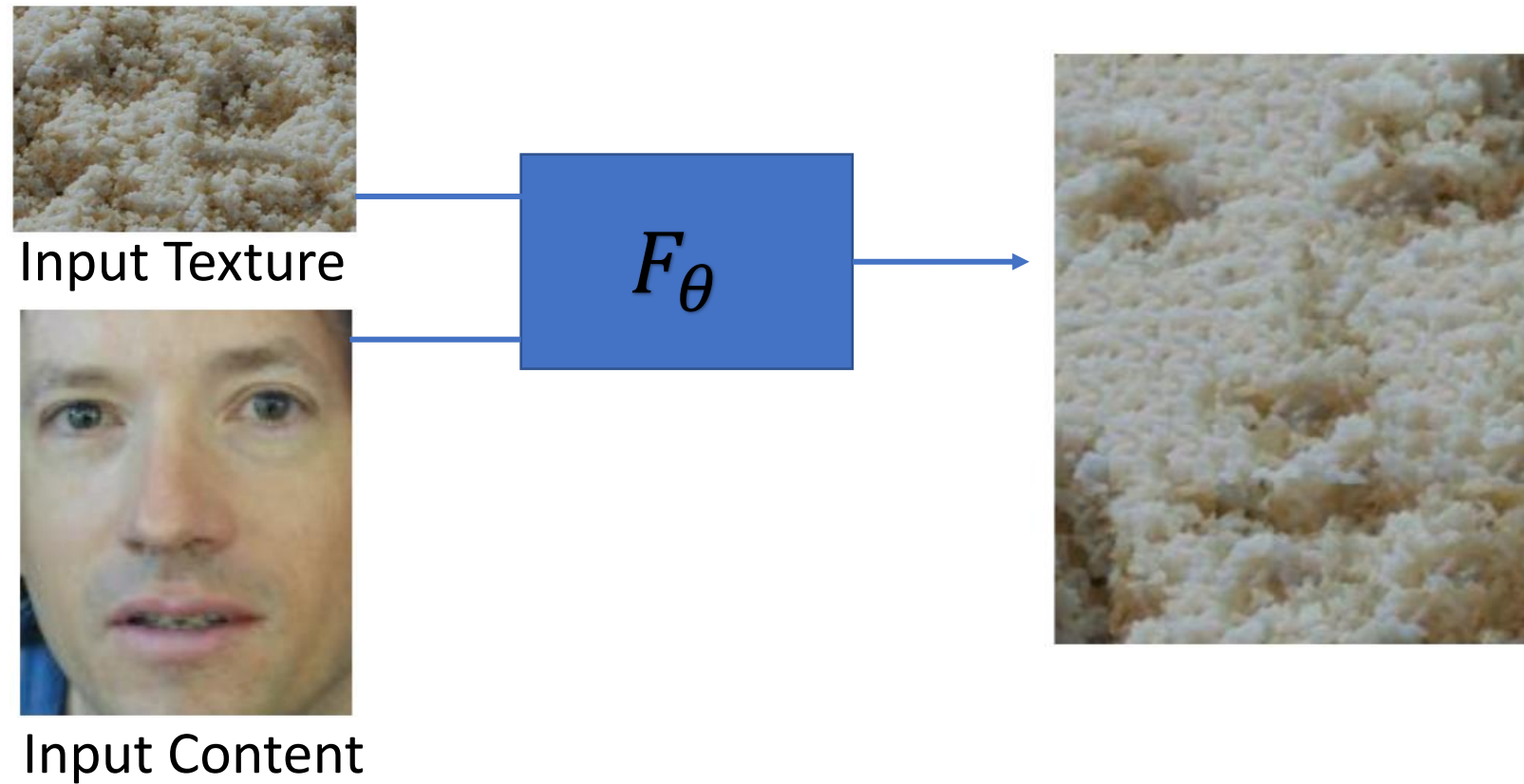
# What is a natural image?

Intelligent machines must **understand** perceived content

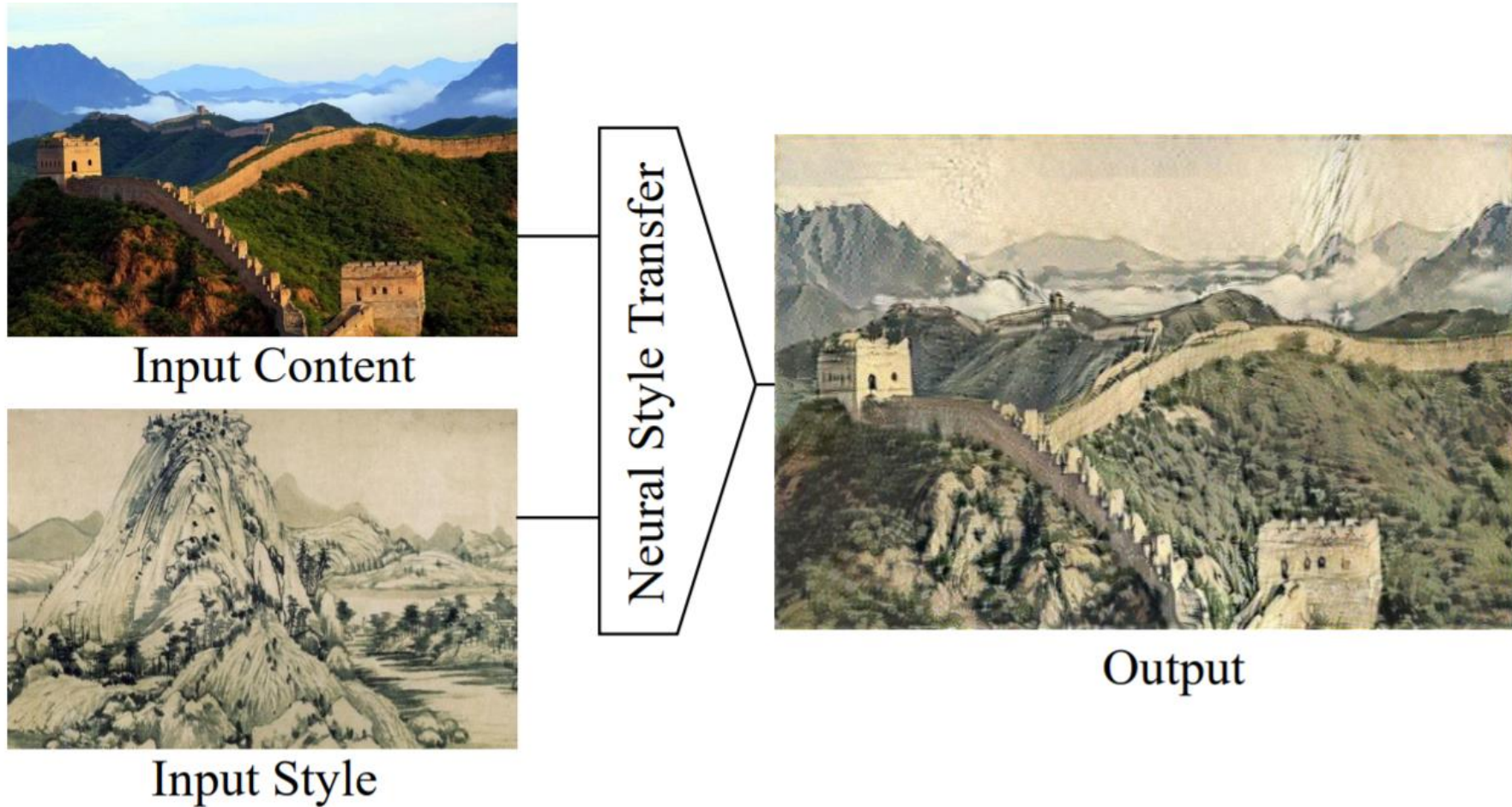


**Understanding by creating/manipulating:**  
“What I cannot create, I do not understand”  
(Richard Feynman)

# Manipulating Texture



# Manipulating Style



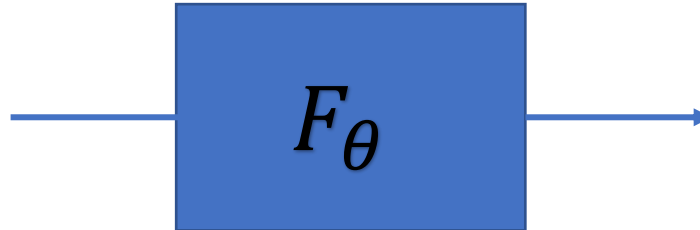
# Manipulating Structure



Target

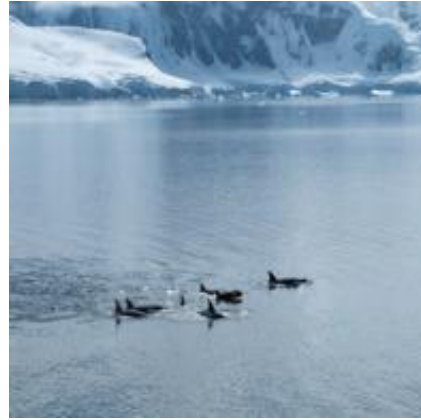


Source Structure





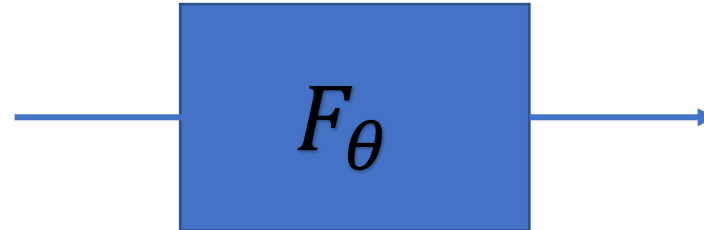
# Manipulating Structure



Target



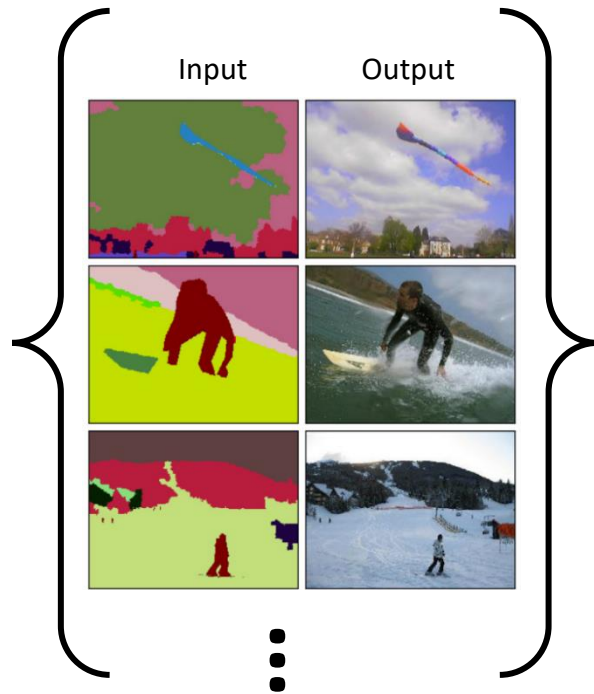
Source Structure



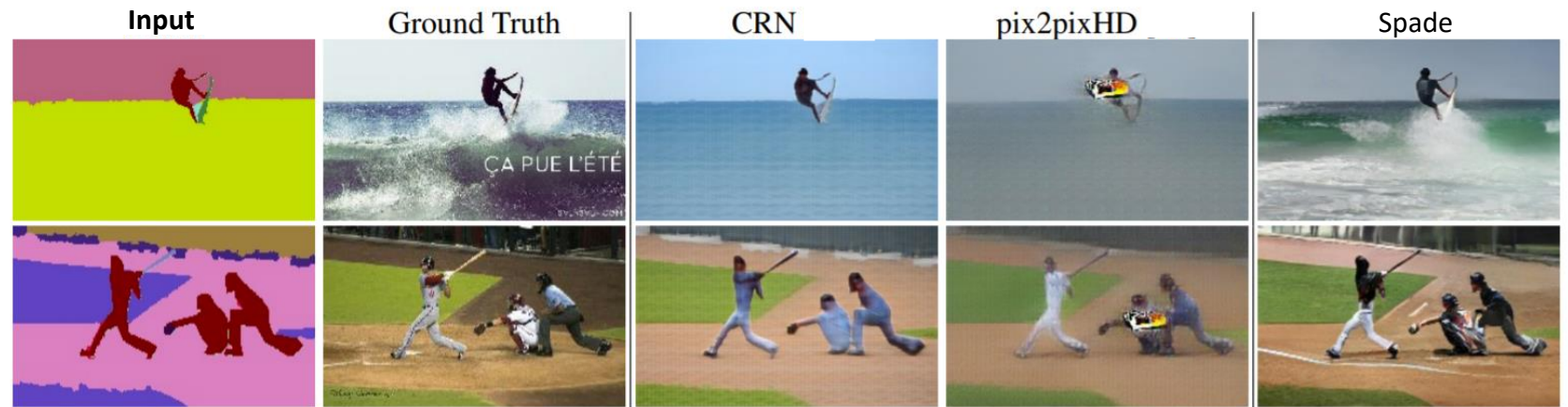
# Multi-Image Approaches

# Supervised (Paired) Setting

Train



Test





# Unsupervised (Unpaired) Setting

**A**



Faces without glasses

**B**



Faces with glasses

# Control Structure of Generated Faces (Transfer Glasses)

Common



Separate

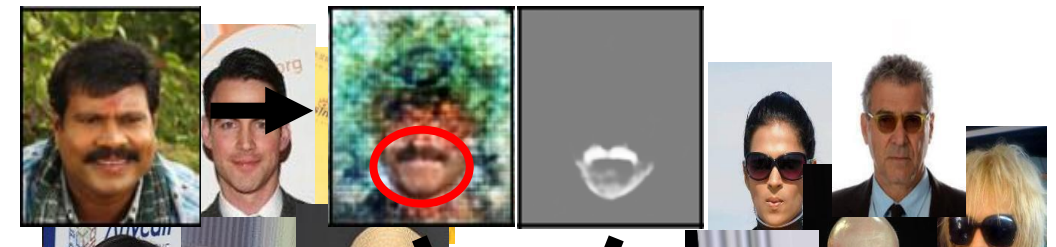
# Unsupervised Approaches

O. Press, T. Galanti, **S. Benaim**, L. Wolf.  
Emerging Disentanglement in Auto-Encoder  
Based Unsupervised Image Content Transfer.  
In **ICLR 2019**.

**S. Benaim**, M. Khaitov, T. Galanti, L. Wolf

Require a large collection of images from both domains

R. Mokady, **S. Benaim**, L. Wolf, A. Bermano.  
Mask Based Unsupervised Content Transfer.  
In **ICLR, 2020**.



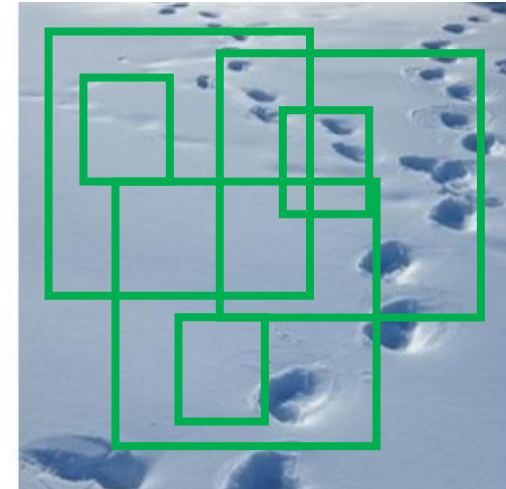
# Patch-Based Approaches



## Multi-Image Distribution



## Multi-Scale Patch Distribution



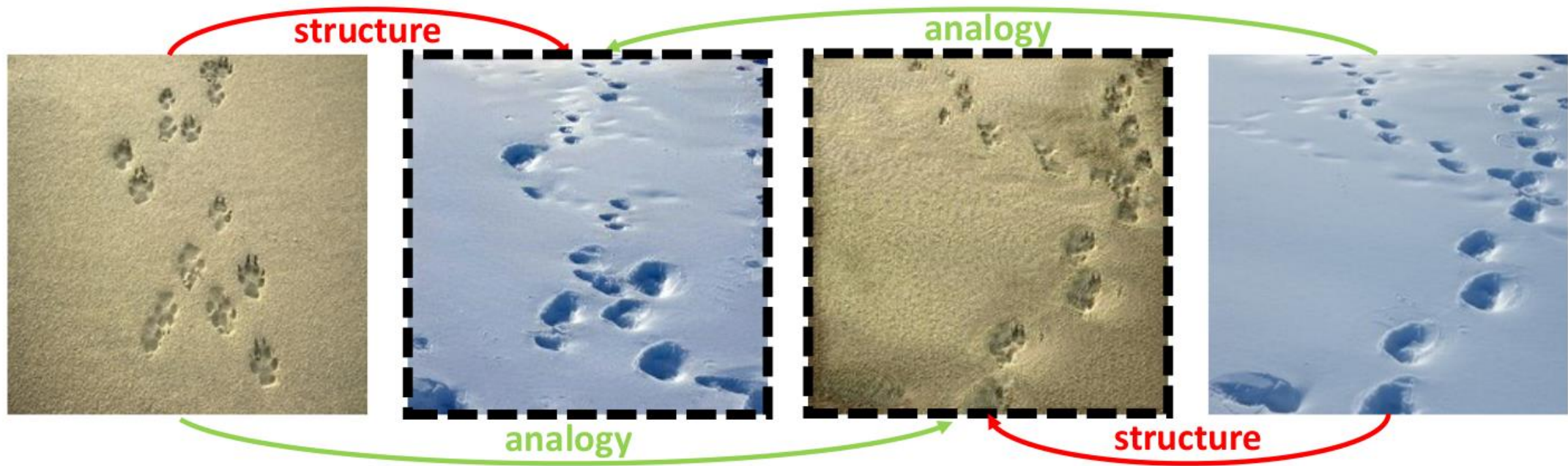
# Structural-analysis from a **Single Image Pair**

**S. Benaim\***, R. Mokady\*, A. Bermano, D Cohen-Or, L. Wolf. CGF 2020. (\*Equal contribution)





Generate an image which is **aligned** to the source image but depicts **structure** from a target image



# Structural Analogy

Target



Source

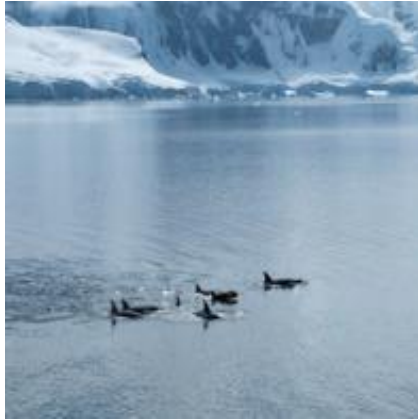


Output



# Structural Analogy

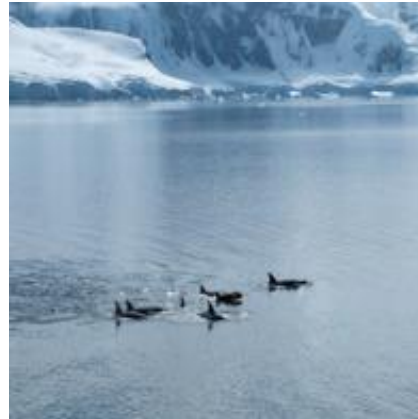
Target



Source



Output





# Structural Analogy

Target



Source



Output



# Style Transfer

# Deep Image Analogy

Style

Content

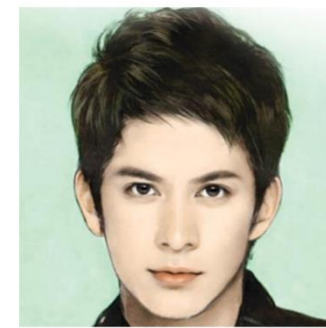
Result



Style

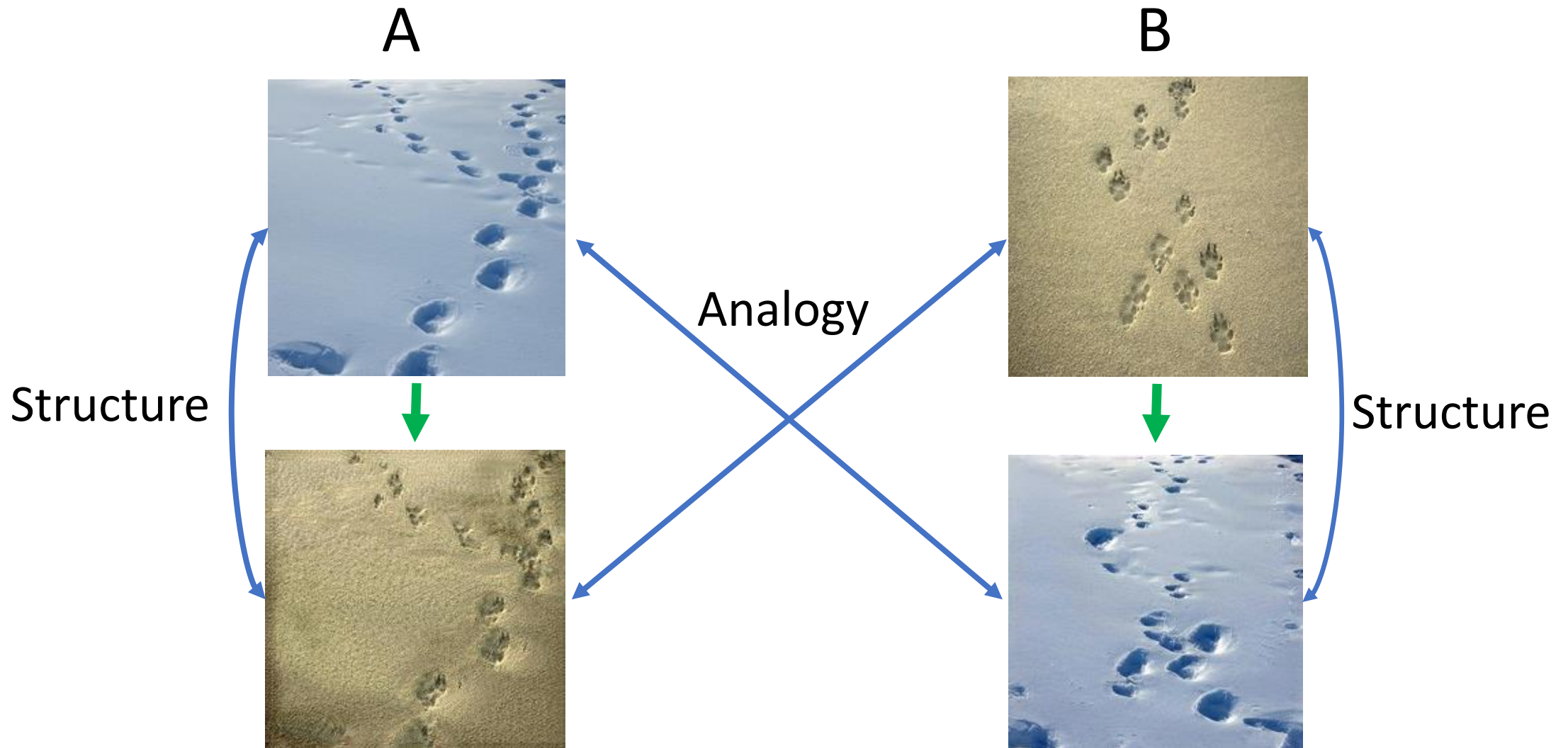
Content

Result



Cannot Change Object Shape

# Structural Analogy





# Motivation

A



B



# Motivation

A

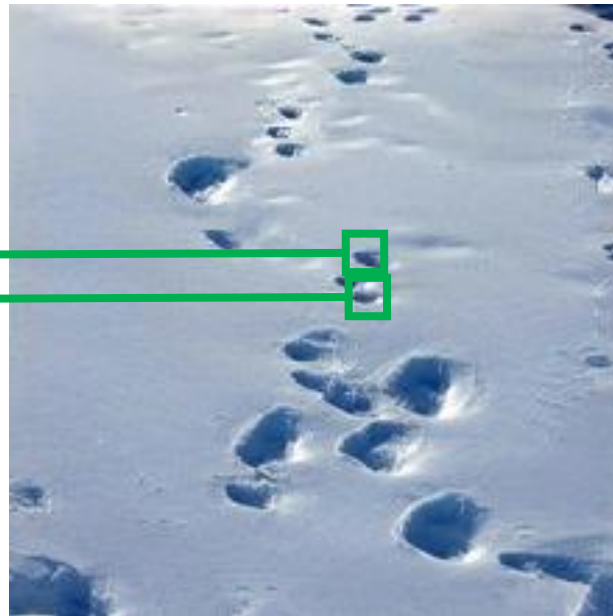
B



# Motivation

A

B



# Proposed Hierarchical Approach

Coarsest scale:

Large Patches

$\bar{a}^0$  (Unconditional)  
 $\overline{ab}^0$  (Conditional)

LEVEL = 0



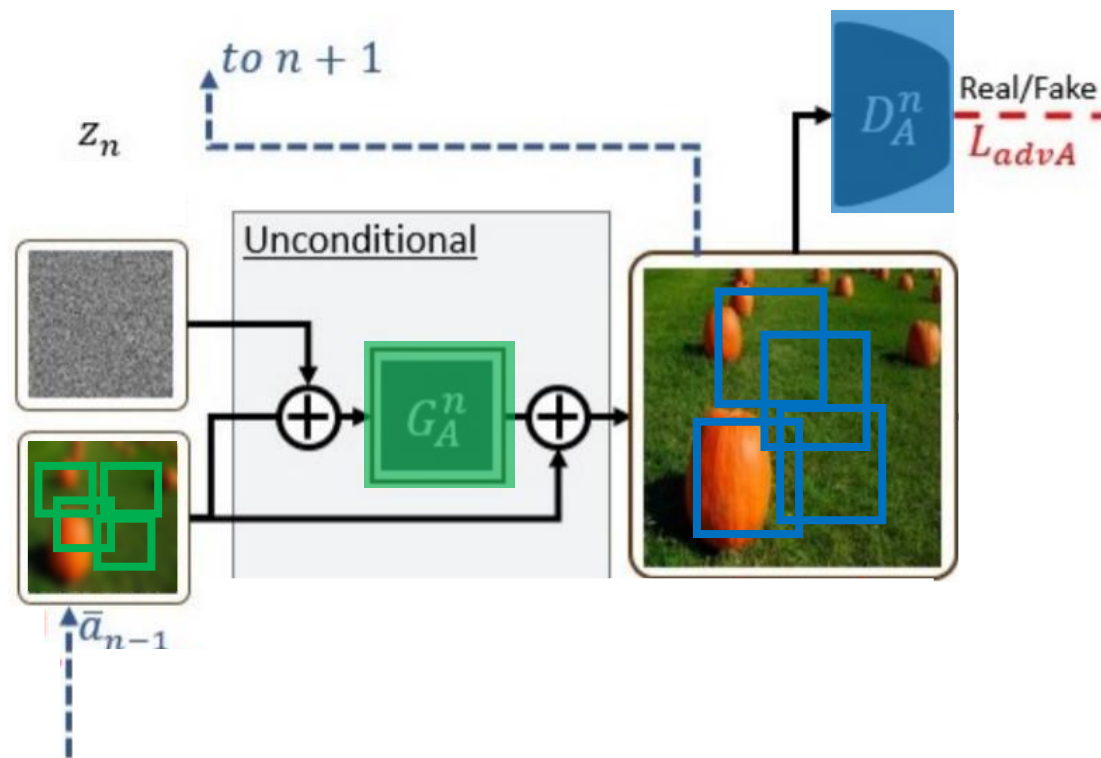
Finest scale:

Small Patches

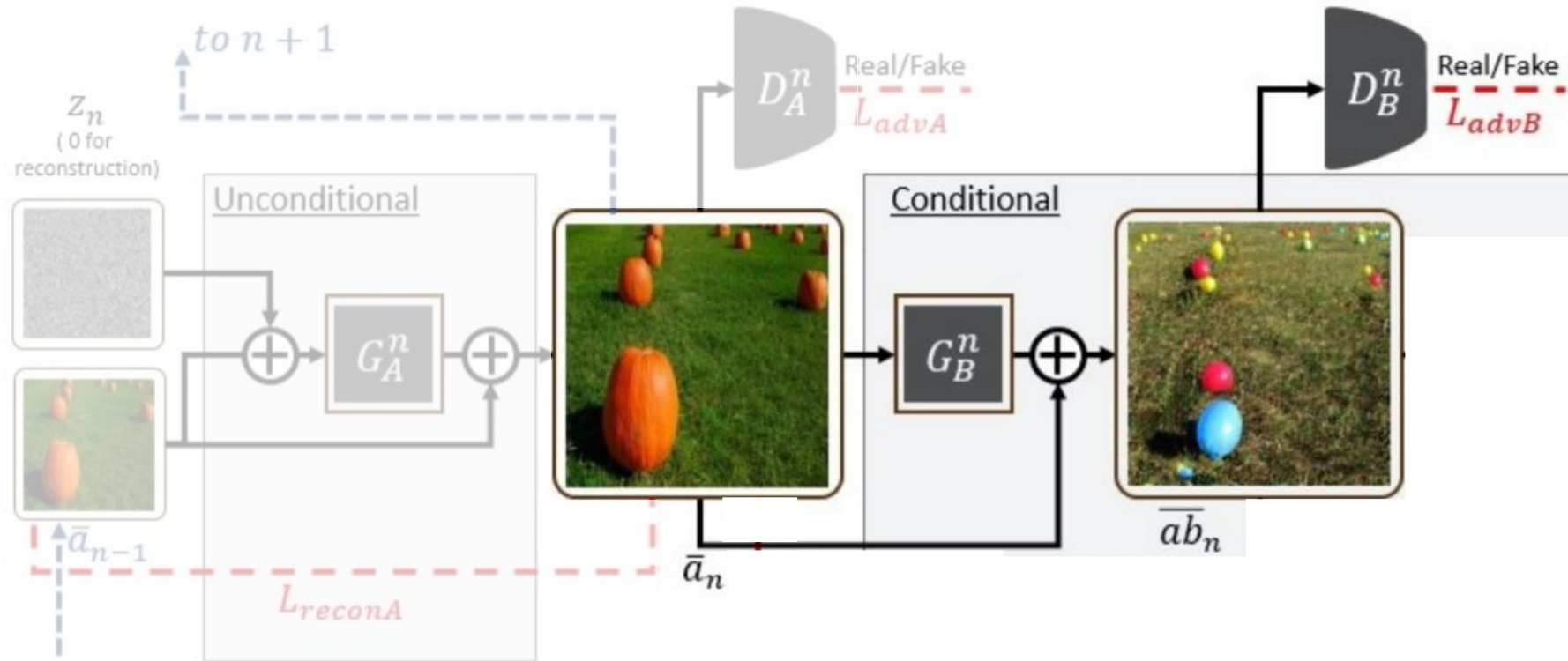
$\bar{a}^N$  (Unconditional)  
 $\overline{ab}^N$  (Conditional)

LEVEL =  $N$

# Unconditional Generation (Level $n$ )

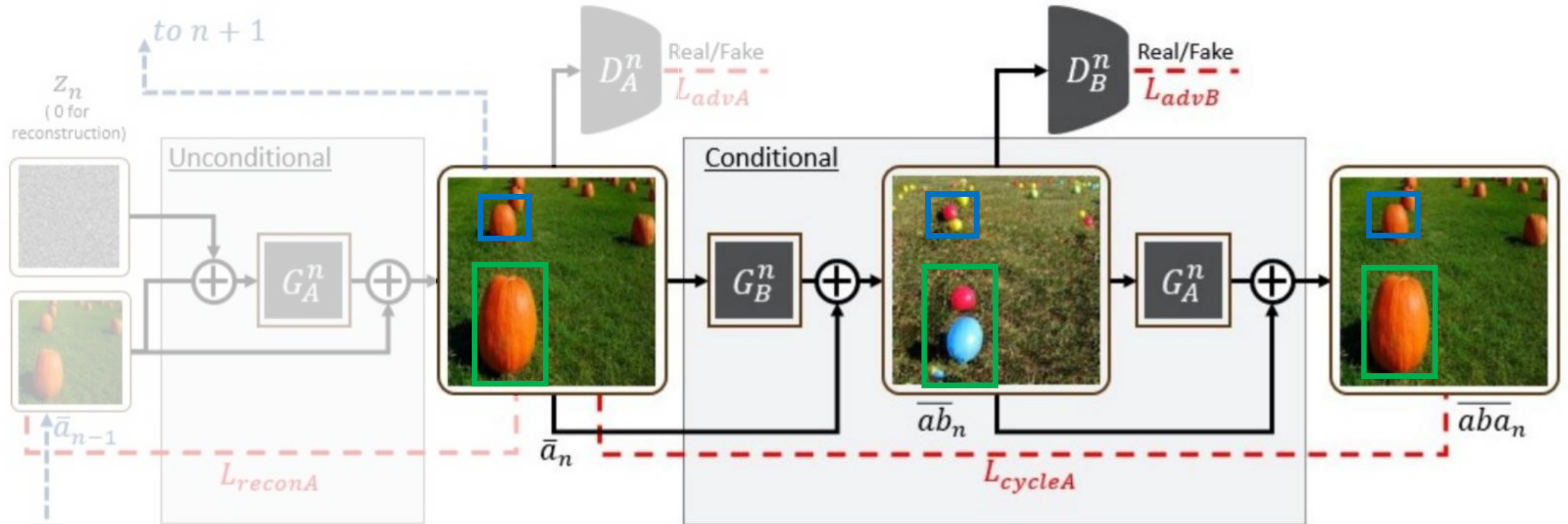


# Conditional Generation (Level n)

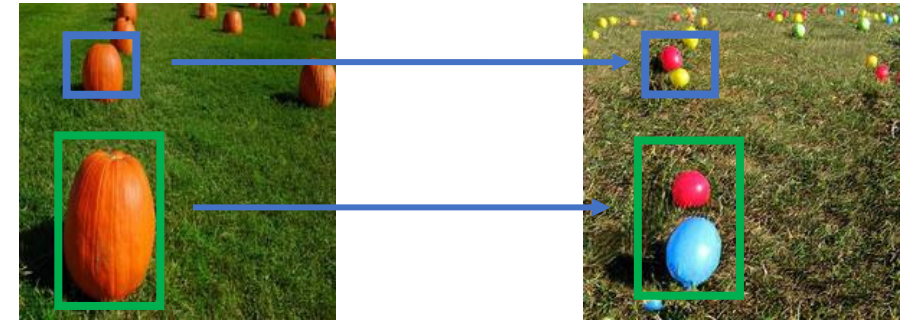
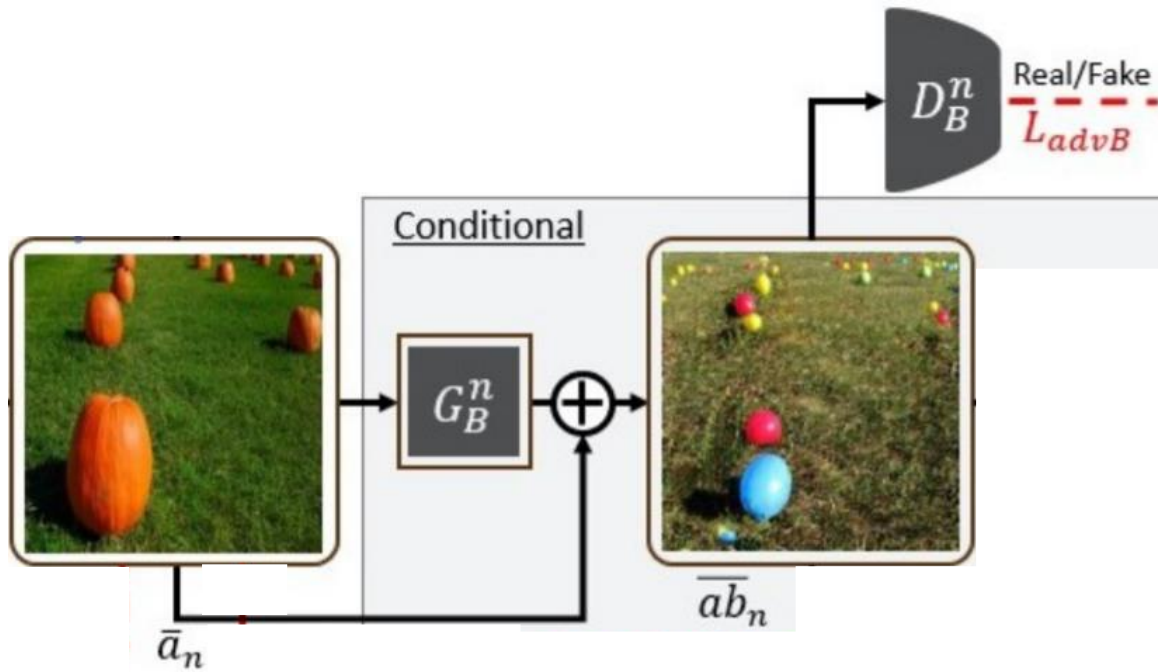




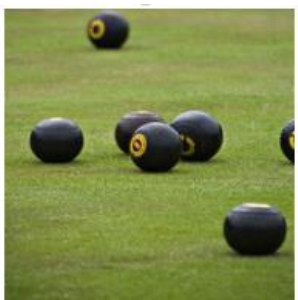
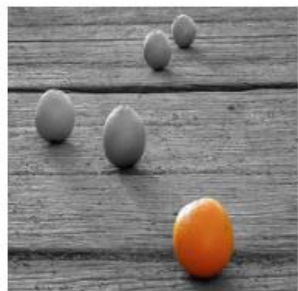
# Conditional Generation (Level n)



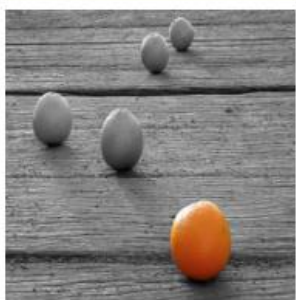
# Coarse and Mid Scales: Residual Training



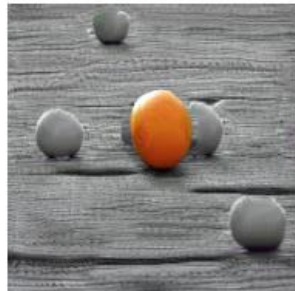
Target



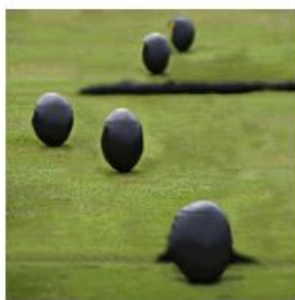
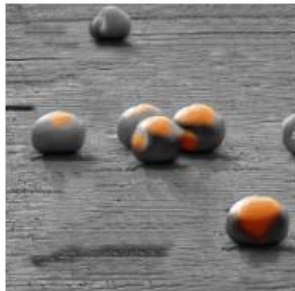
Source



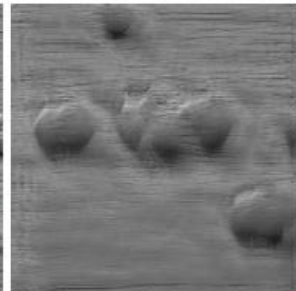
Ours



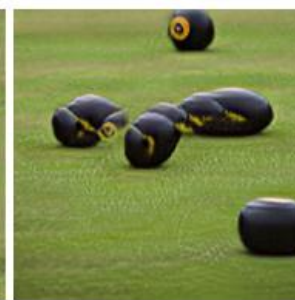
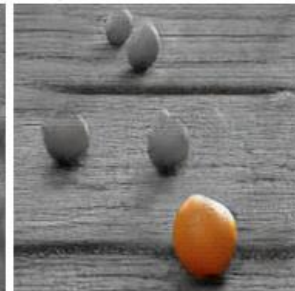
DIA



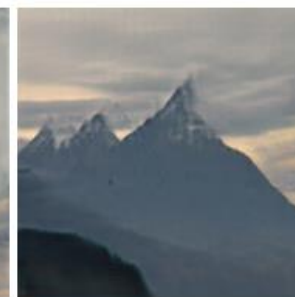
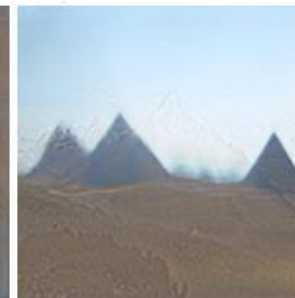
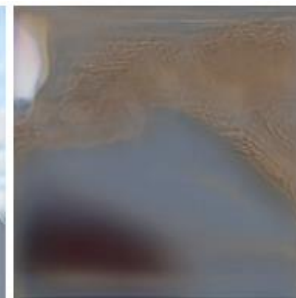
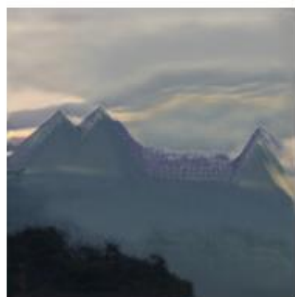
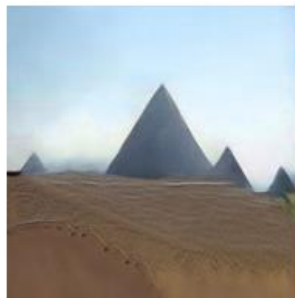
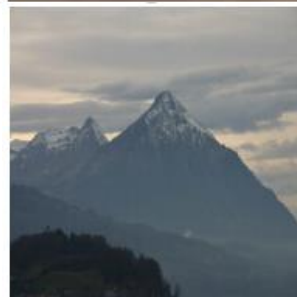
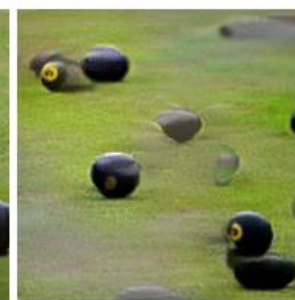
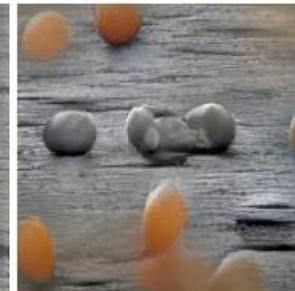
SinGAN



Cycle



Style





# Multiple Class Types

Input



Output



# Paint to Image

Input

Sketch

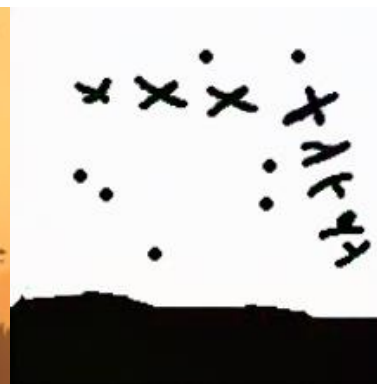
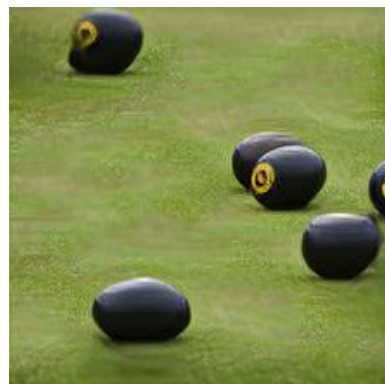
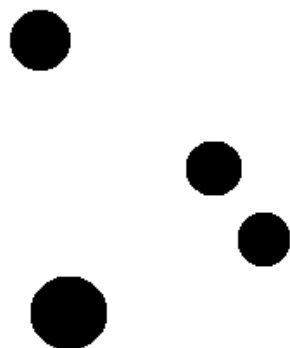
Ours



Input

Sketch

Ours



# Video Generation





Downstream Tasks?

# A Hierarchical Transformation-Discriminating Generative Model for **Few Shot Anomaly Detection**

S. Sheynin\*, **S. Benaim\***, L. Wolf. In Submission to ICCV 2021. (\*Equal contribution)

**Anomalous**



⋮



**Normal**



# A Hierarchical Transformation-Discriminating Generative Model for **Few Shot Anomaly Detection**

S. Sheynin\*, **S. Benaim\***, L. Wolf. In Submission to ICCV 2021. (\*Equal contribution)

**Anomalous**



⋮

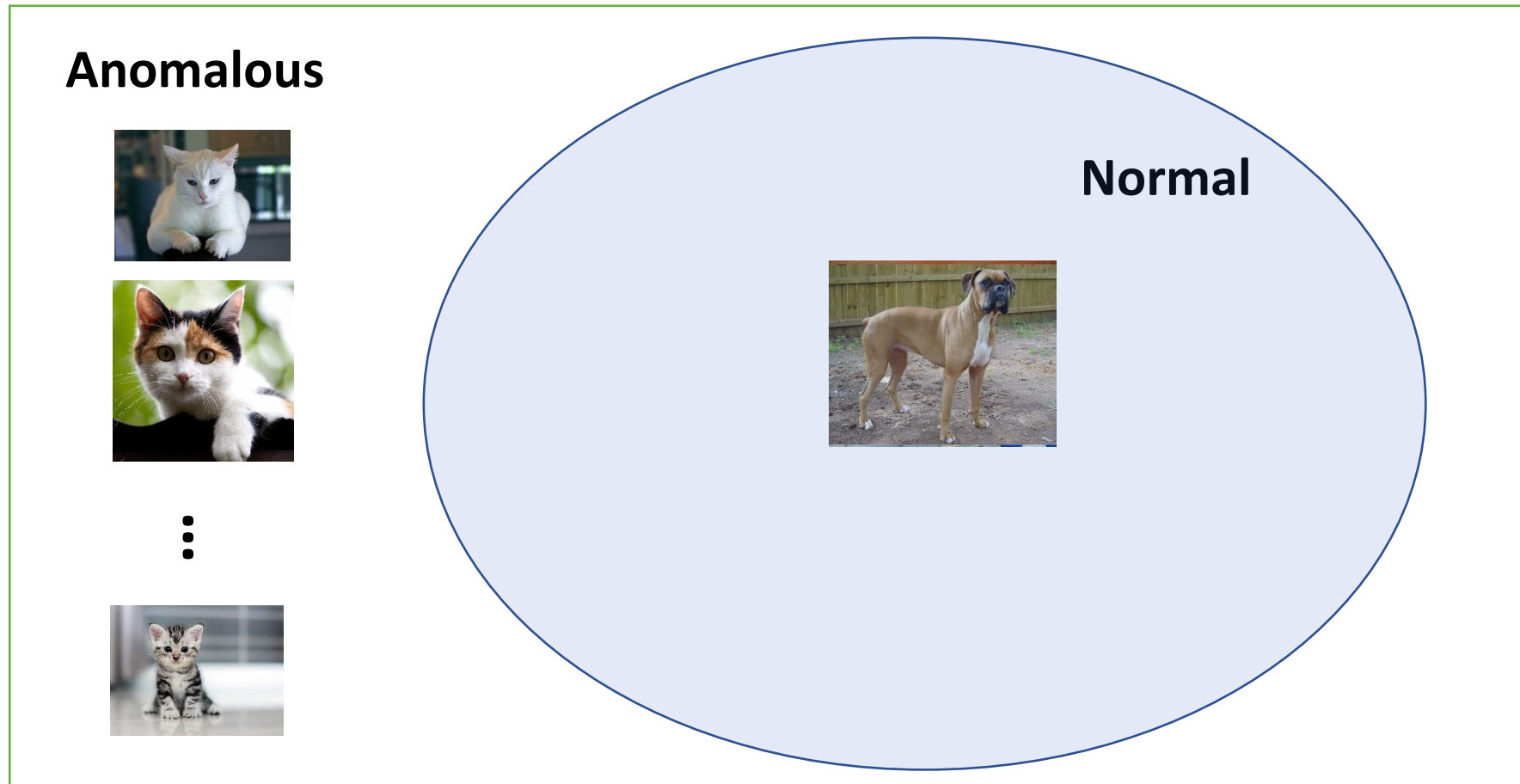


**Normal**

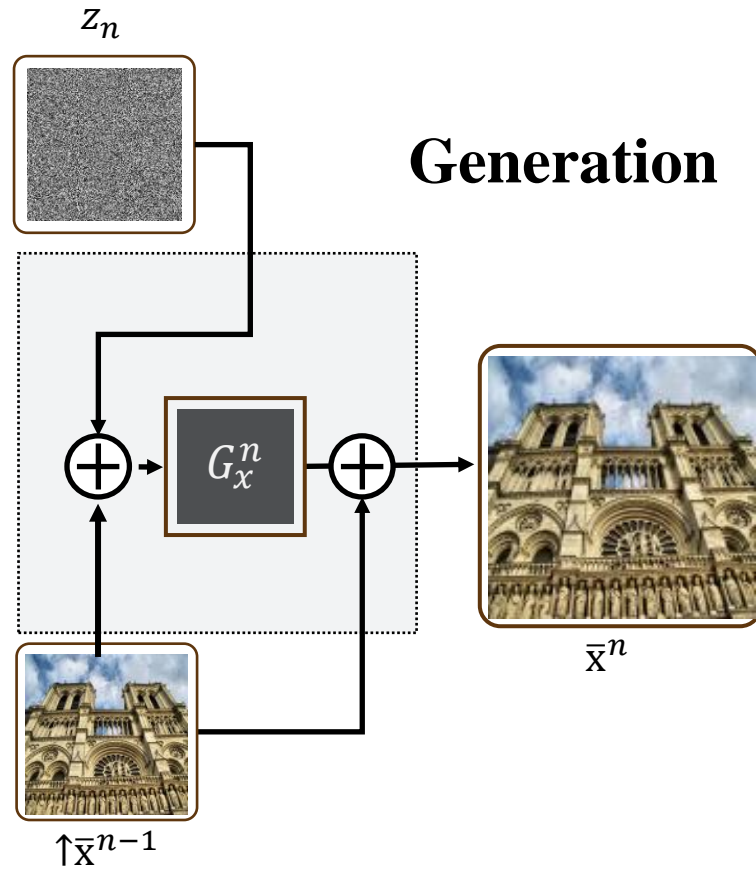


# A Hierarchical Transformation-Discriminating Generative Model for **Few Shot Anomaly Detection**

S. Sheynin\*, **S. Benaim\***, L. Wolf. In Submission to ICCV 2021. (\*Equal contribution)

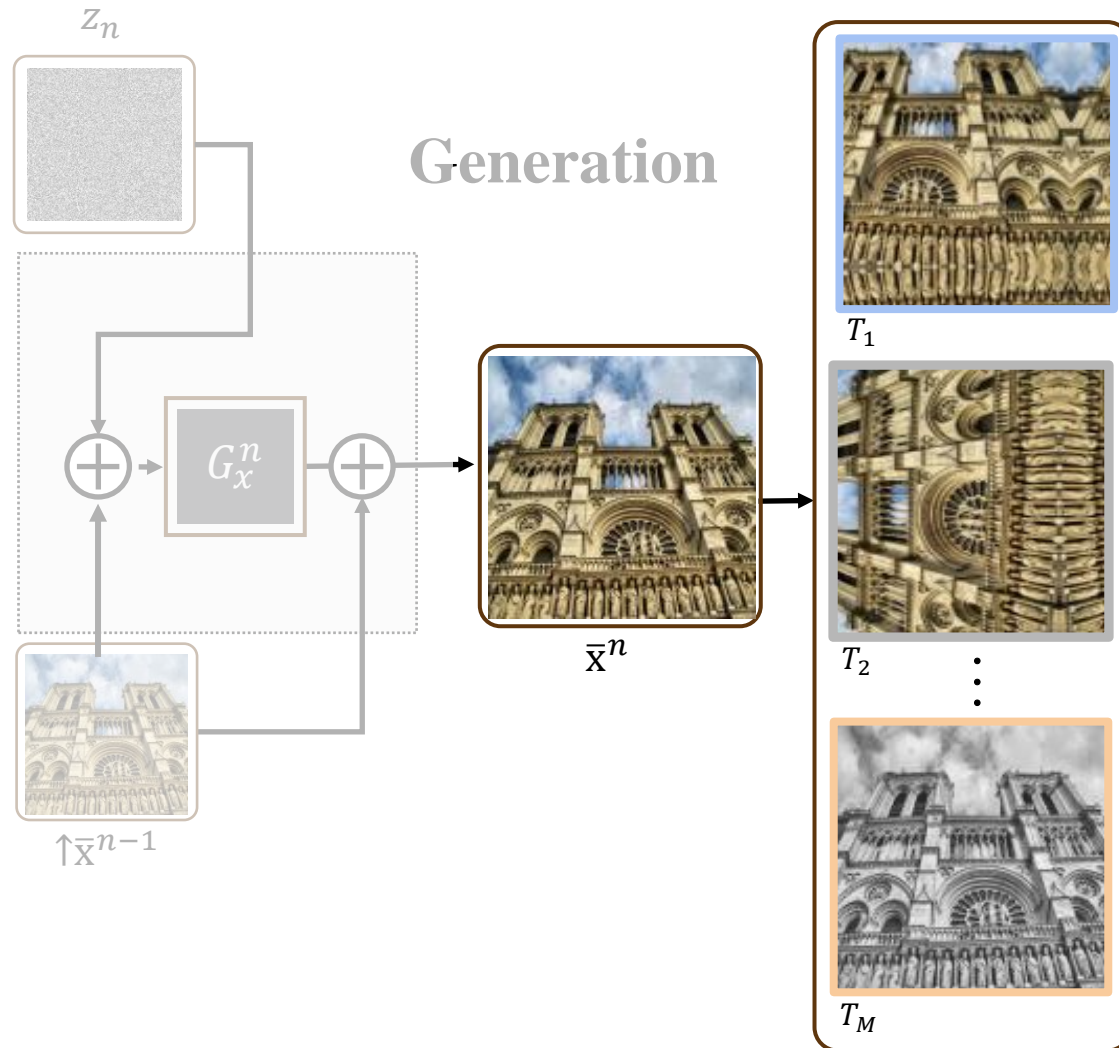


# Unconditional Generation (Level n)





# Transform Generated Sample



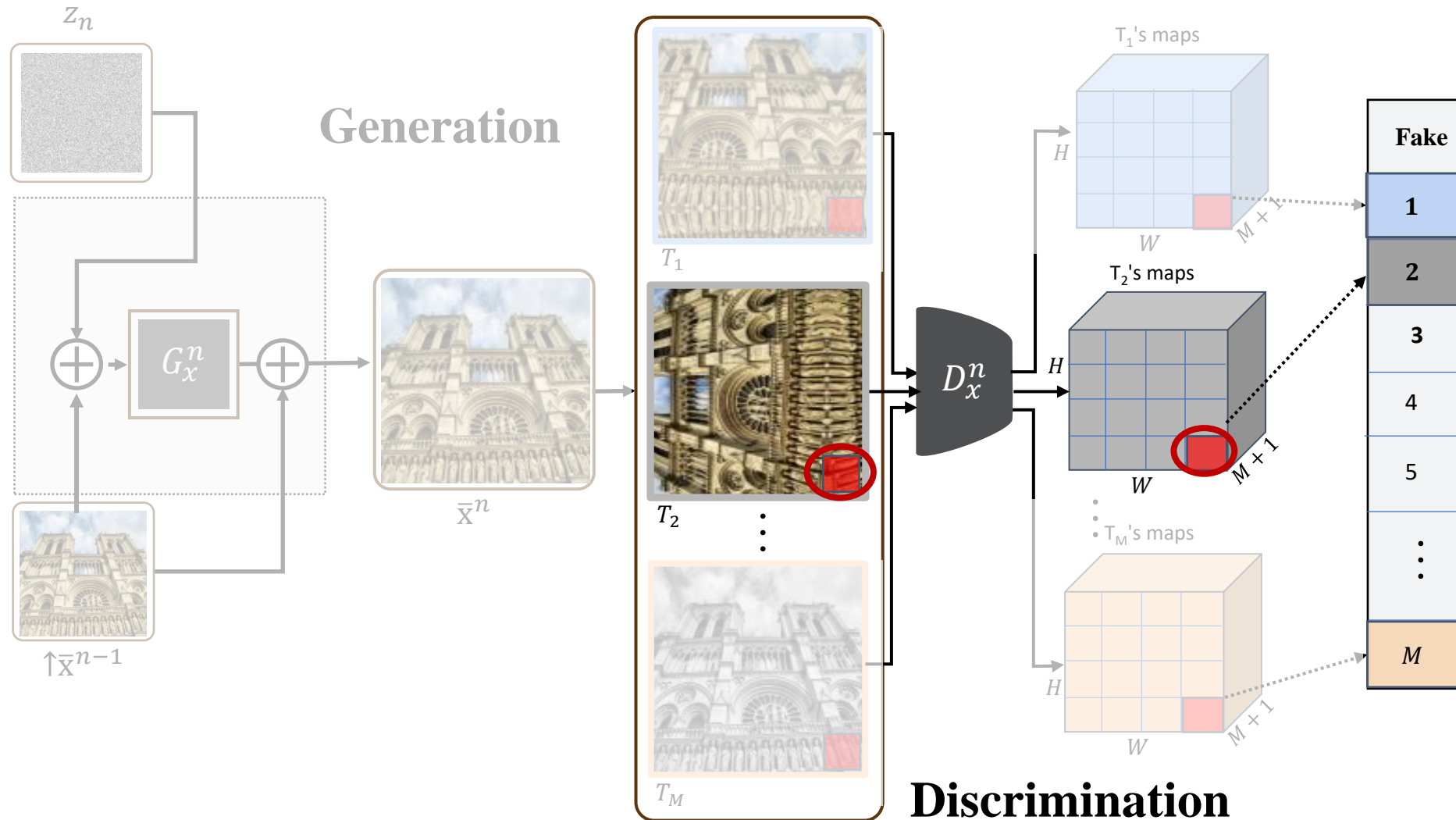
$T_1$ : Horizontal Flip, Translation (y-axis)

$T_2$ : 90° Rotation, Translation (x-axis), Translation (y-axis)

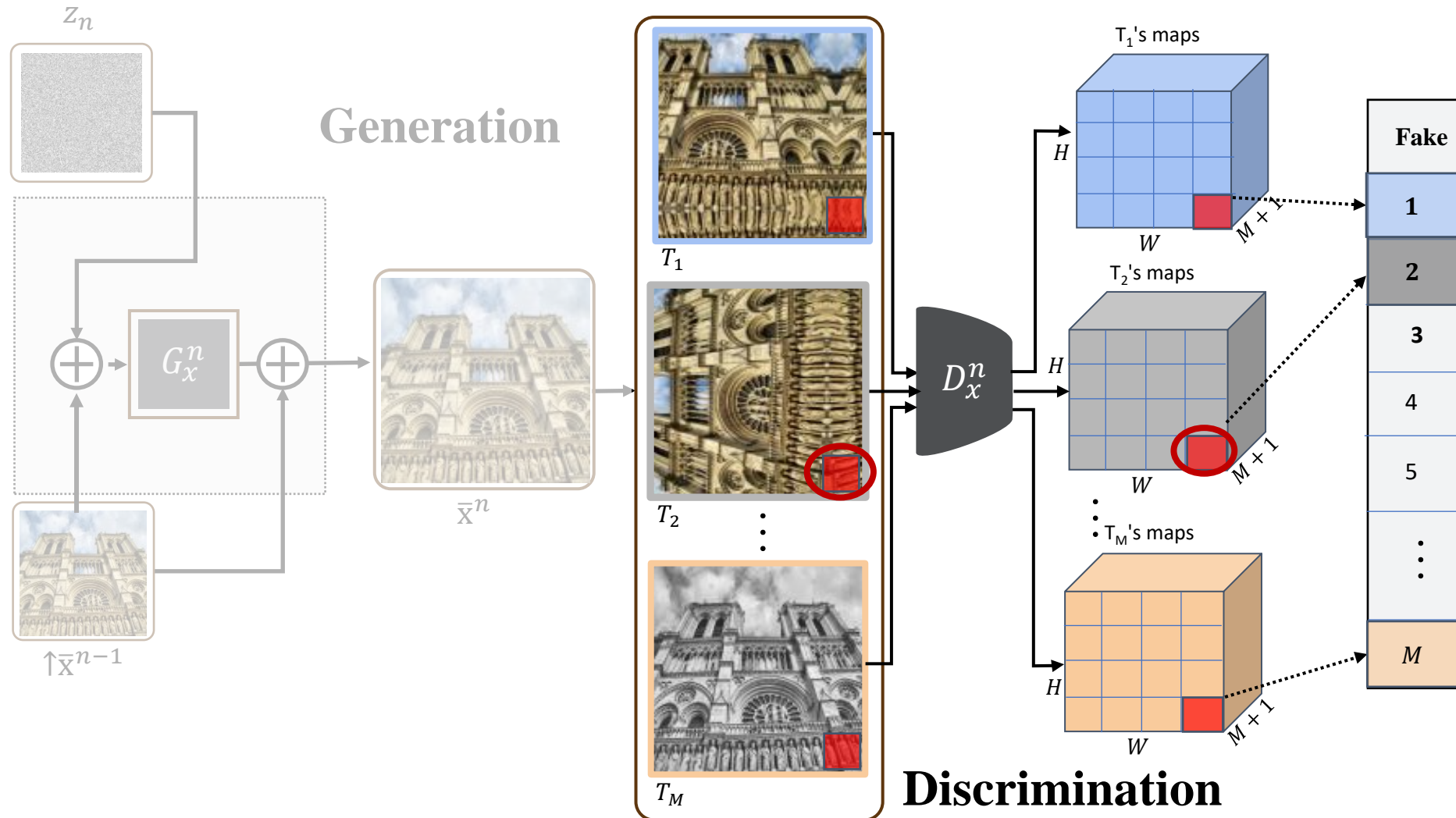
...

$T_M$ : Grayscale (y-axis)

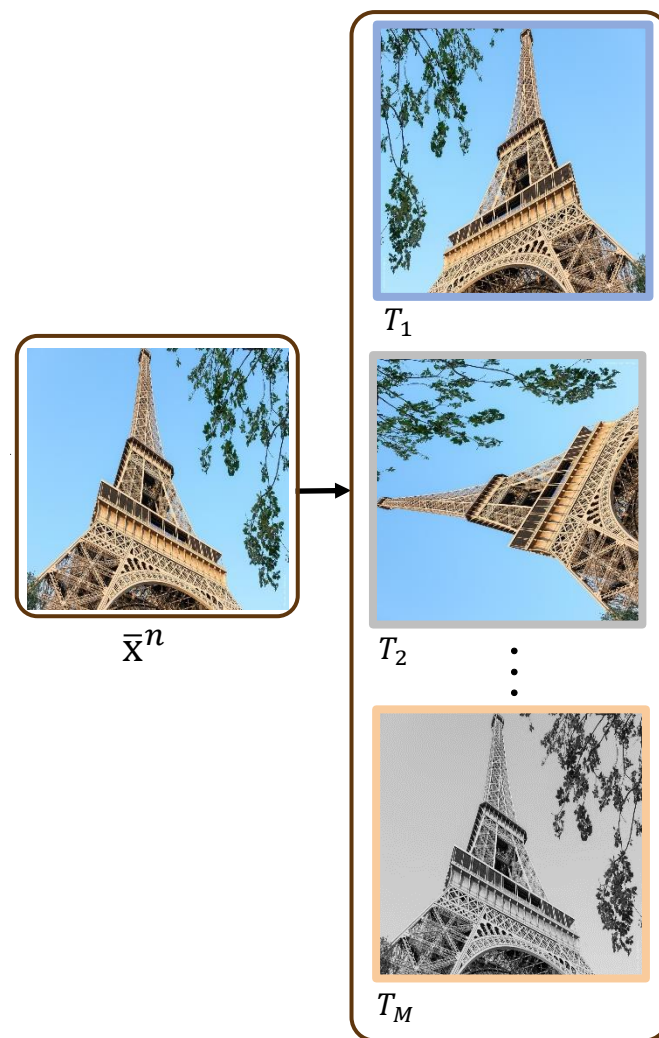
# Patch-Based Self Supervised Task



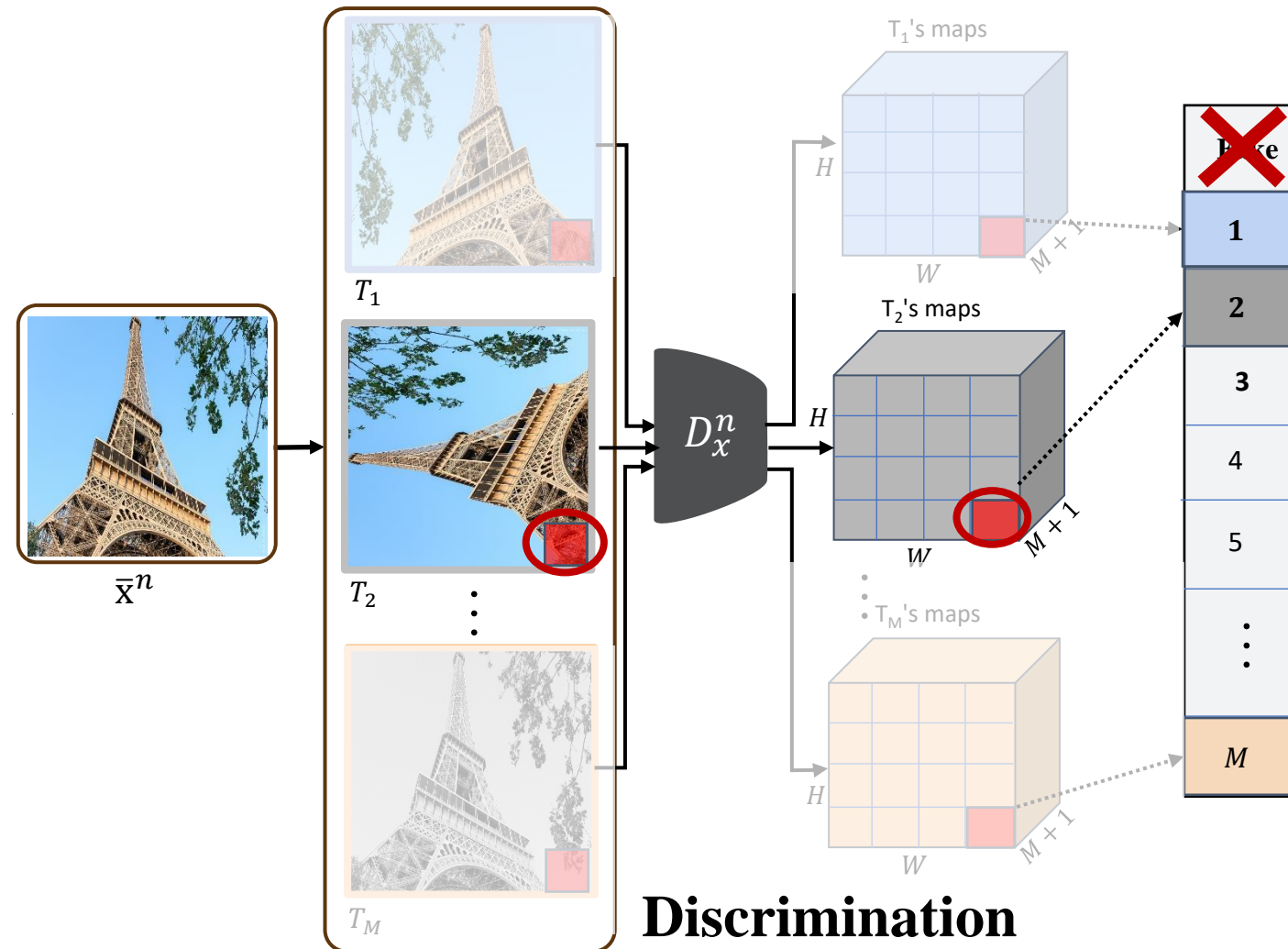
# Patch-Based Self Supervised Task



# Test Time: Anomaly Score (Scale n)

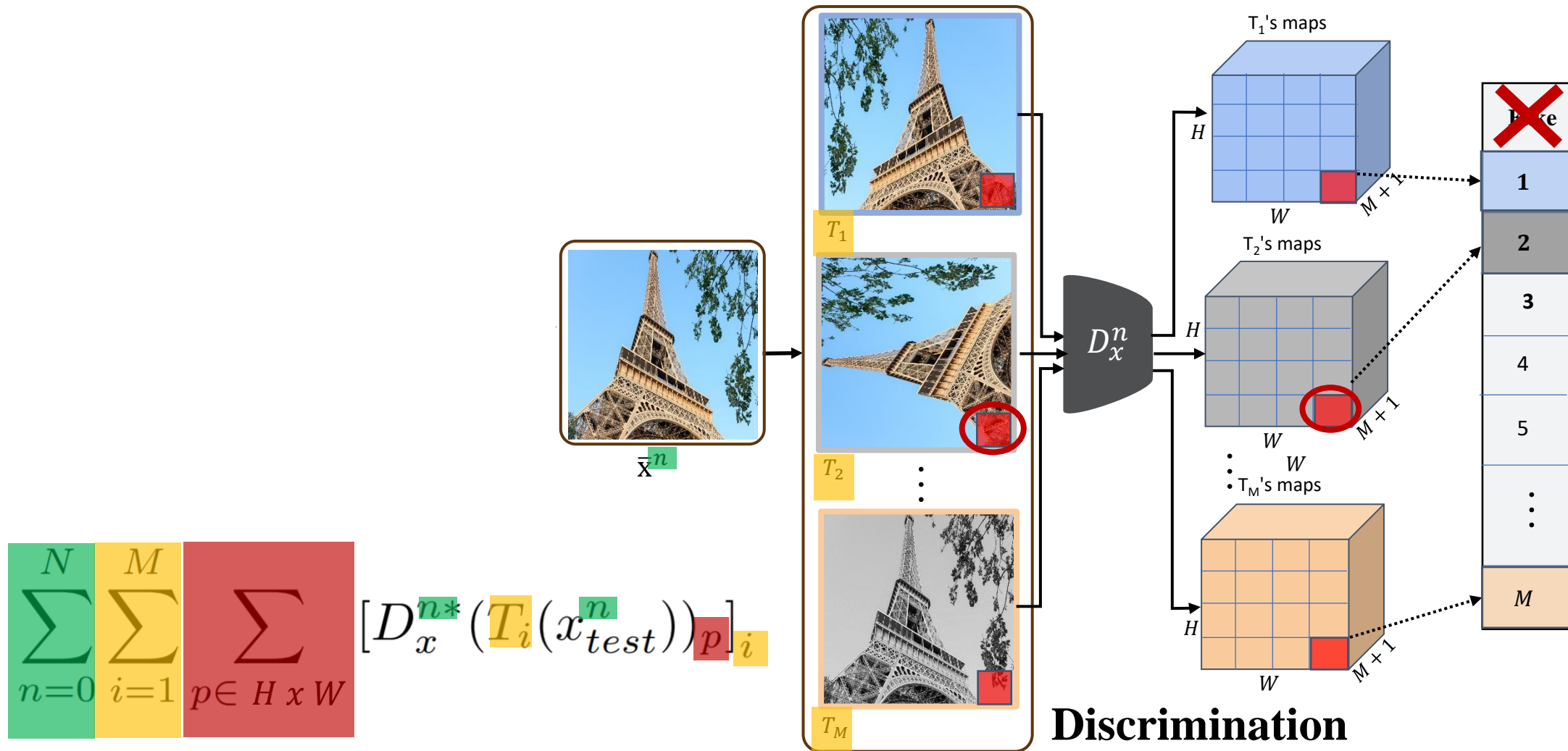


# Test Time: Anomaly Score (Scale n)

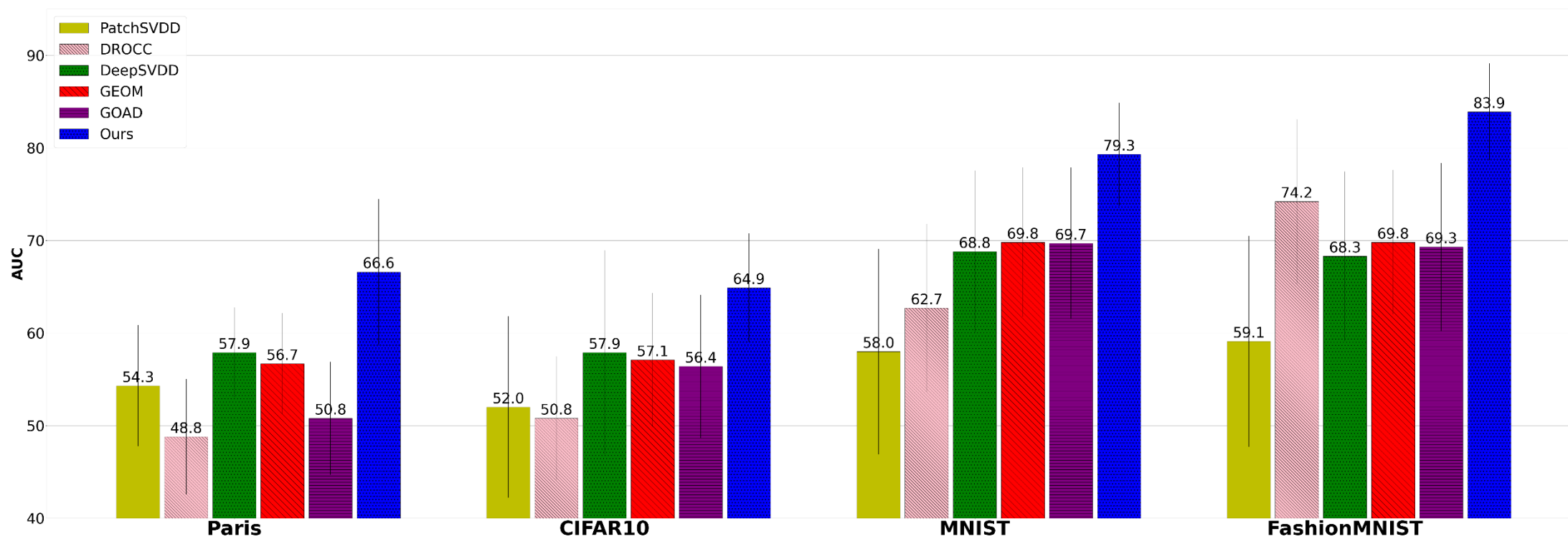




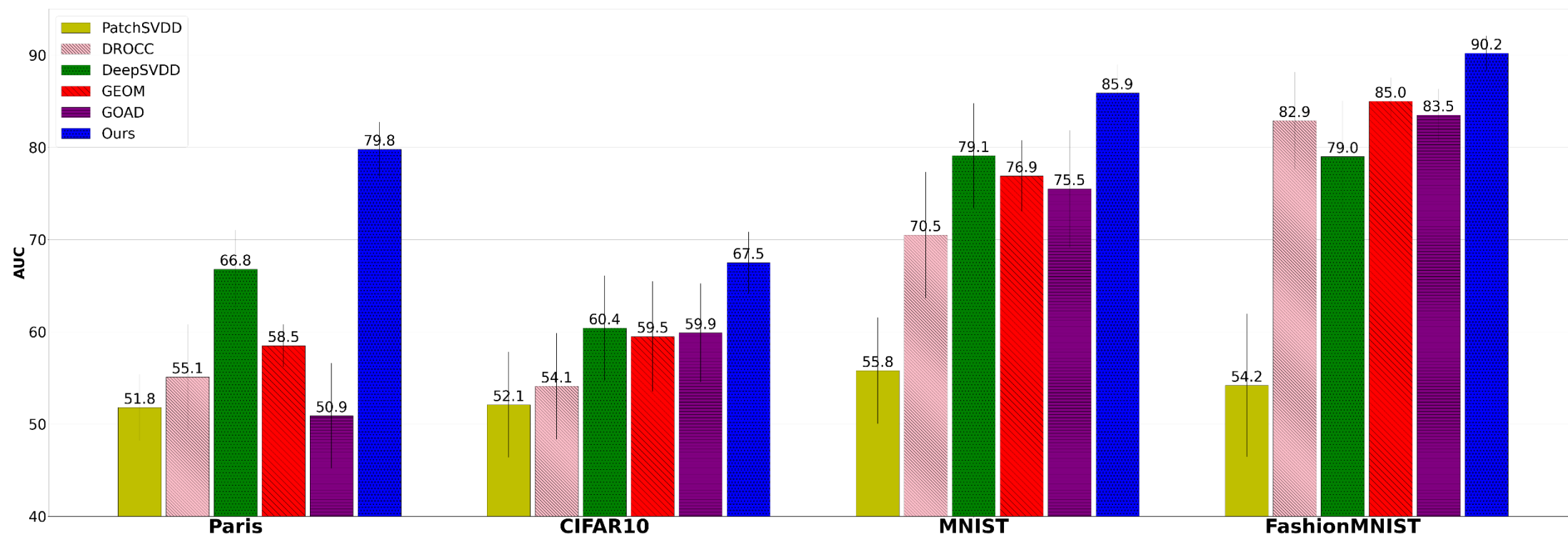
# Test Time: Anomaly Score (Scale n)



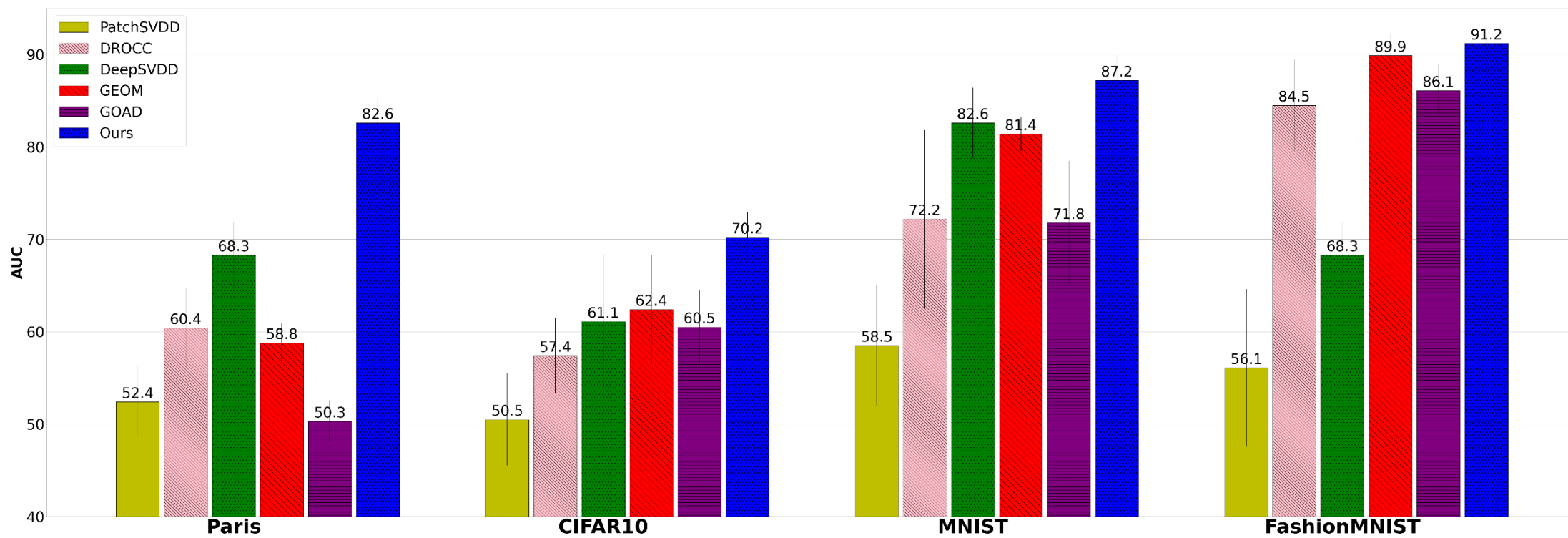
# One-Shot



# Five-Shot

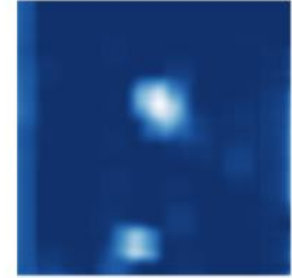
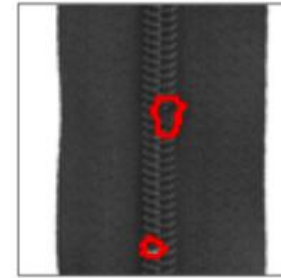
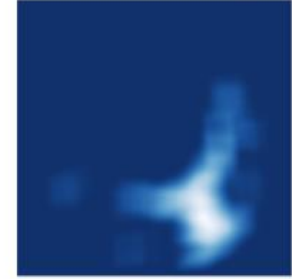
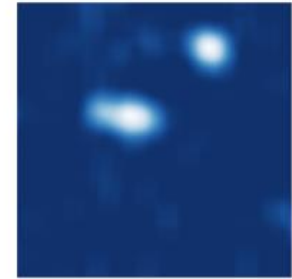
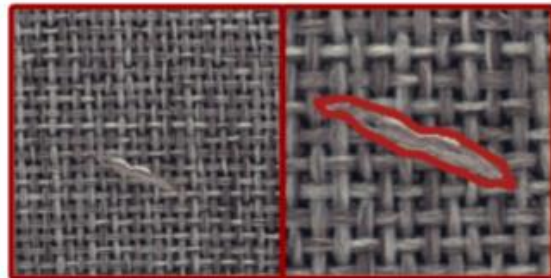
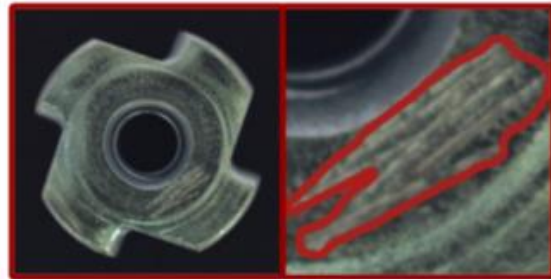


# Ten-Shot





# One Shot Defect Localization



Videos?

# Hierarchical Patch VAE-GAN: Generating Diverse Videos from a **Single Sample**

S. Gur\*, **S. Benaim\***, L. Wolf. NeurIPS 2020 (\*Equal contribution)

Real





# Hierarchical Patch VAE-GAN: Generating Diverse Videos from a Single Sample

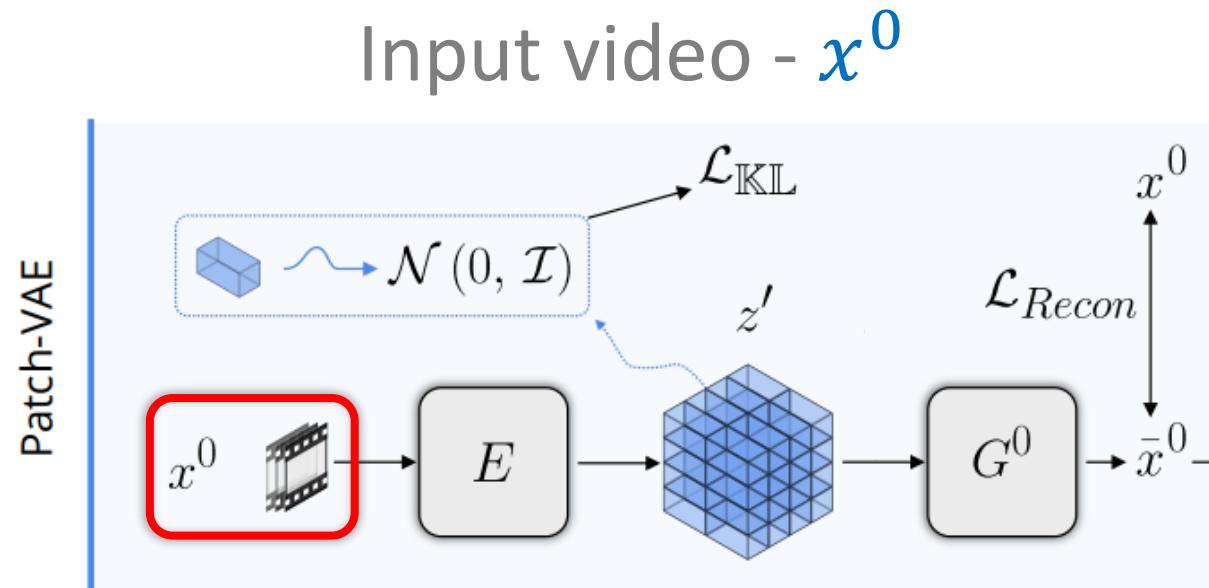
S. Gur\*, S. Benaim\*, L. Wolf. NeurIPS 2020 (\*Equal contribution)

Real

Generated Samples (13 Frames)



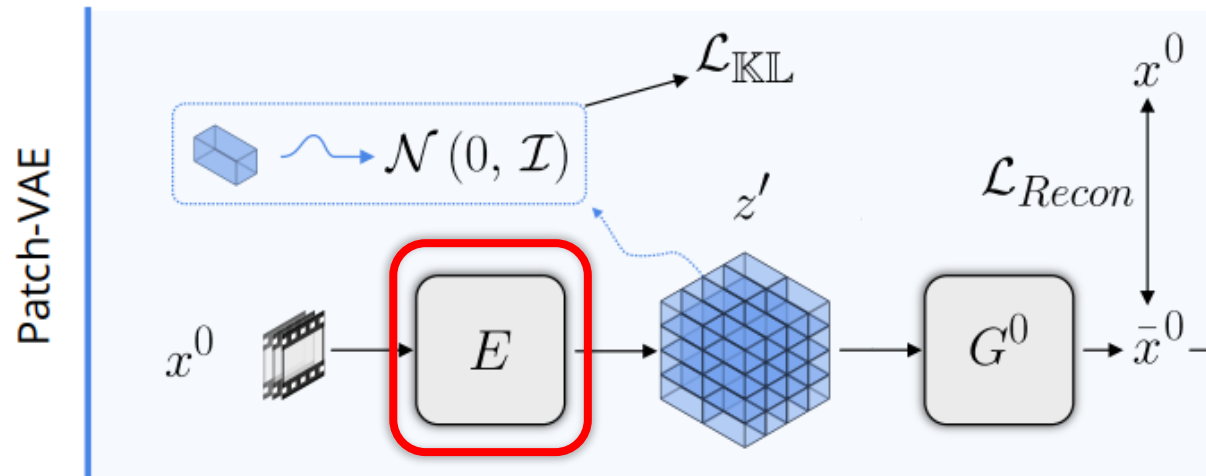
# Proposed Approach: Patch VAE



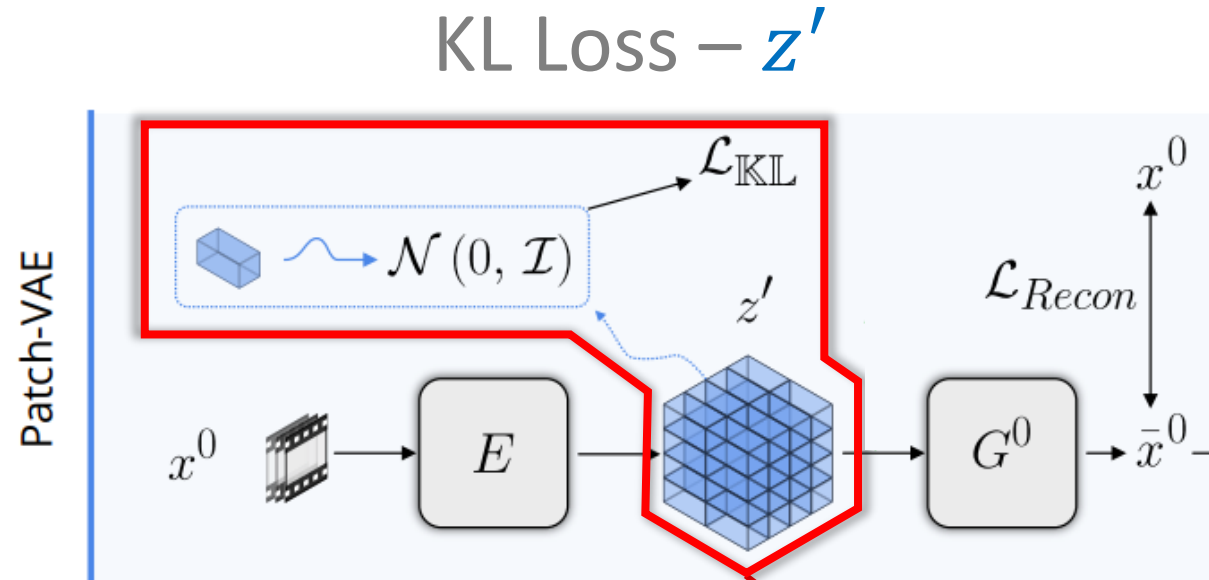


# Proposed Approach: Patch VAE

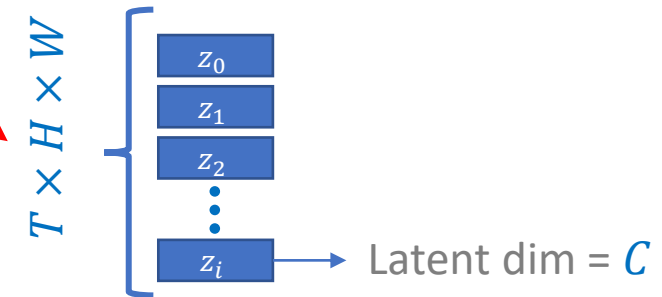
Encoder –  $E(x^0)$



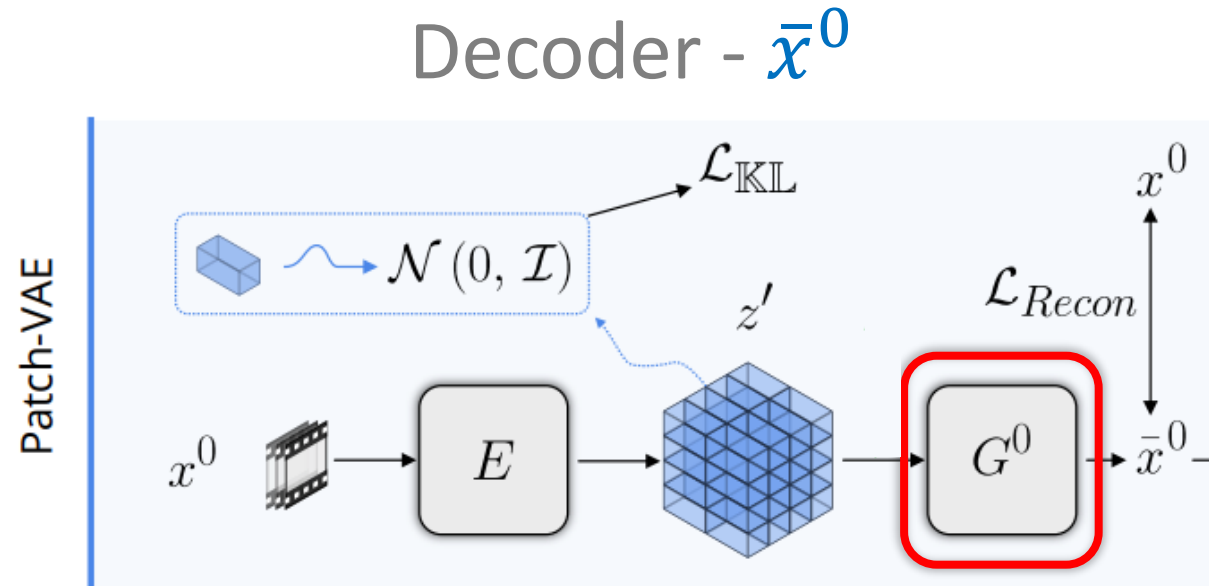
# Proposed Approach: Patch VAE



Each feature  $z_i, i = [1 \dots K], K = T \times H \times W$ ,  
in the latent space is associated with a patch  $\omega_i$

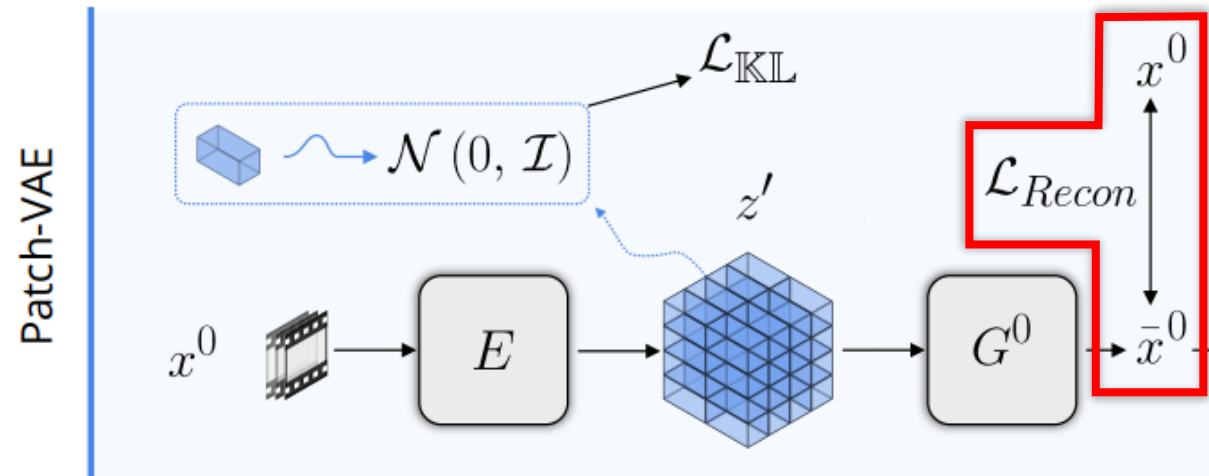


# Proposed Approach: Patch VAE



# Proposed Approach: Patch VAE

Reconstruction loss



# Proposed Approach: Hierarchical Patch VAE

Coarsest scale:  
**Low** resolution  
and frame rate

$x^0$  (Real)  
 $\bar{x}^0$  (Generated)

LEVEL = 0



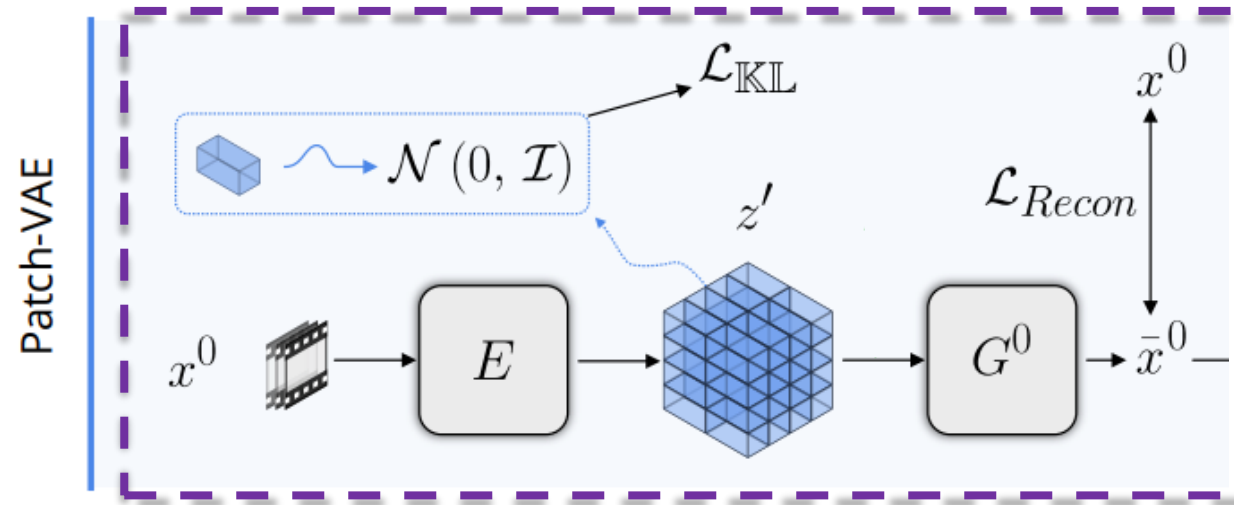
Finest scale:  
**High** resolution  
and frame rate

$x^N$  (Real)  
 $\bar{x}^N$  (Generated)

LEVEL =  $N$



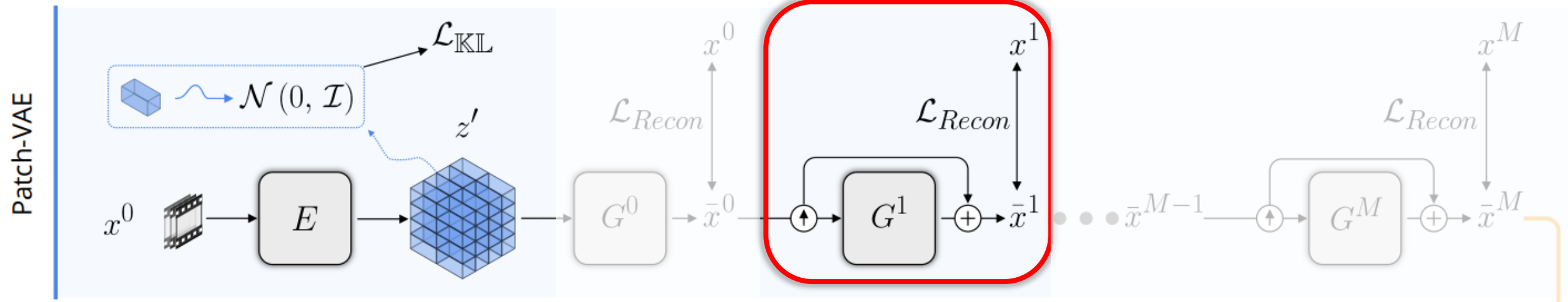
# Proposed Approach: Hierarchical Patch VAE



LEVEL = 0

# Proposed Approach: Hierarchical Patch VAE

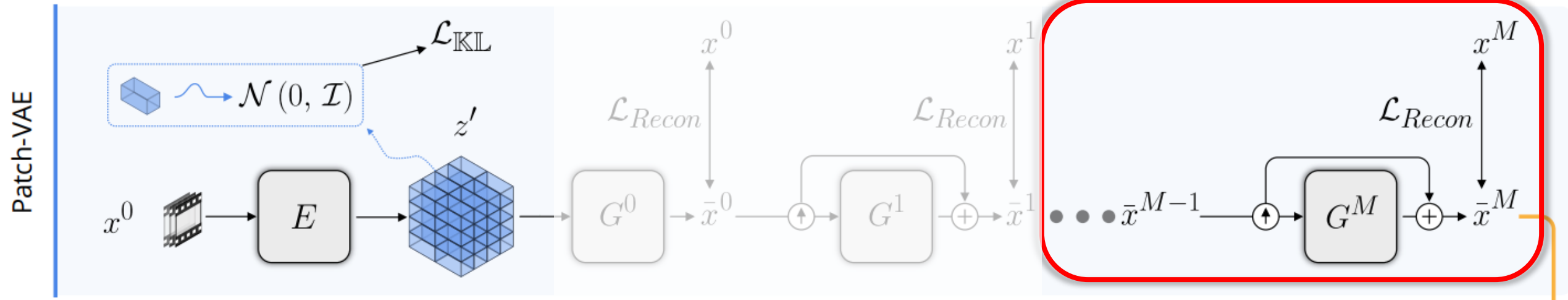
Up-sampling block -  $\bar{x}^1$



LEVEL = 1

# Proposed Approach: Hierarchical Patch VAE

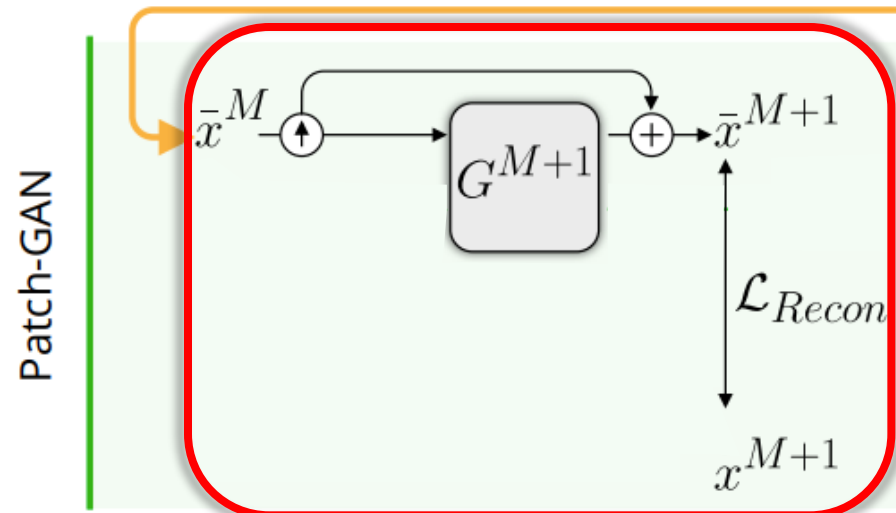
Hierarchical up-sampling up to  $\bar{x}^M$



LEVEL  $\leq M$

# Proposed Approach: Hierarchical Patch VAE GAN

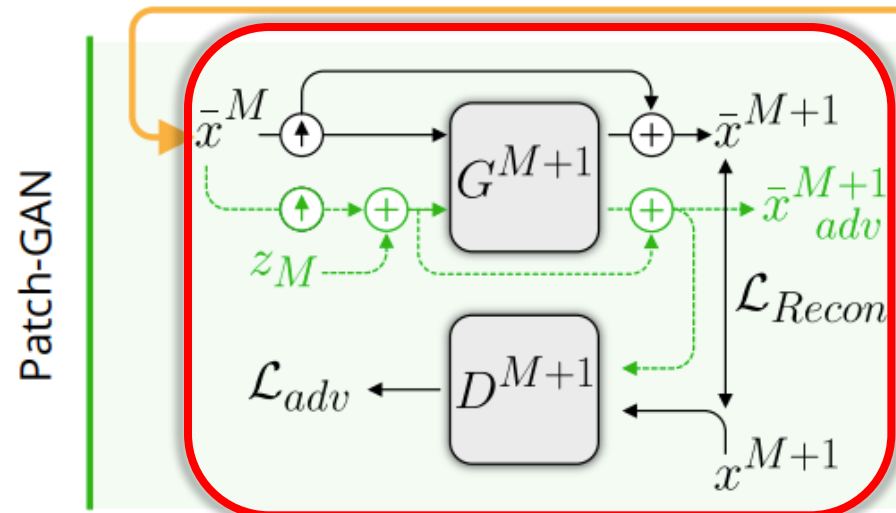
Up-sampling block  $\bar{x}^{M+1}$



LEVEL =  $M + 1$

# Proposed Approach: Hierarchical Patch VAE GAN

Adversarial training



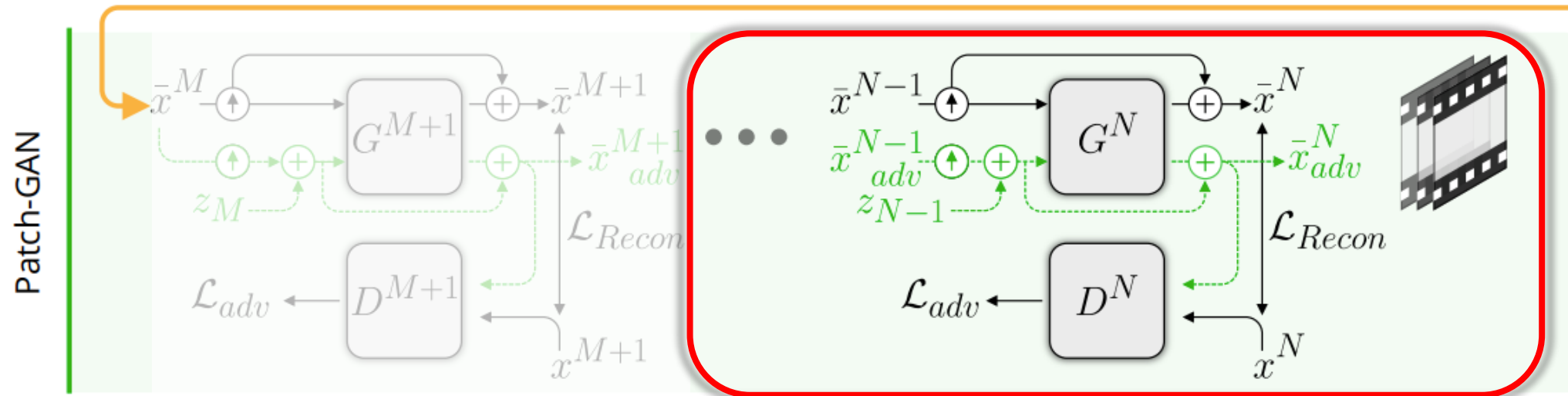
Added noise  $z_M$

LEVEL =  $M + 1$



# Proposed Approach: Hierarchical Patch VAE GAN

Hierarchical up-sampling up to final resolution  $\bar{x}^N$



$$M + 1 < \text{LEVEL} \leq N$$

# Effect of Number of patch-VAE levels

Training Video



9 Levels Total

1 p-VAE – 8 p-GAN



8 p-VAE – 1 p-GAN



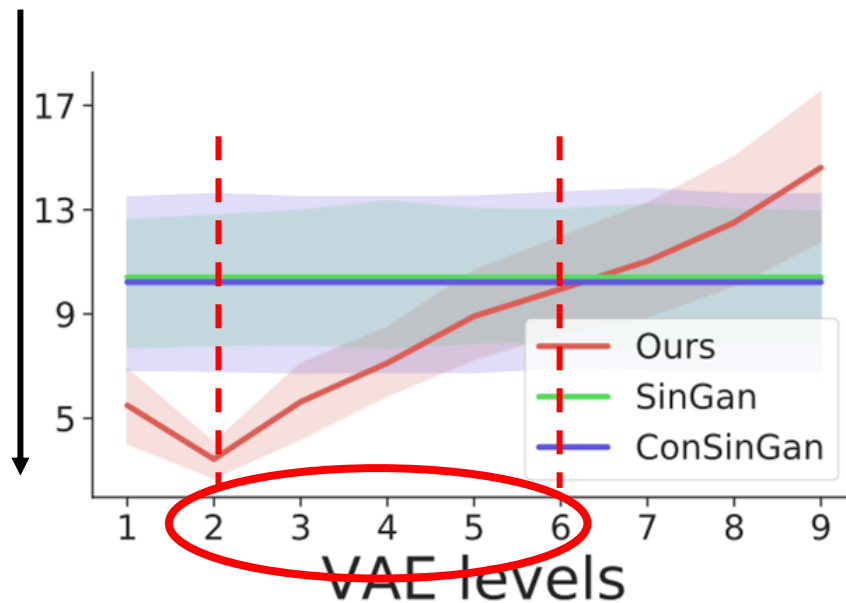
3 p-VAE – 6 p-GAN



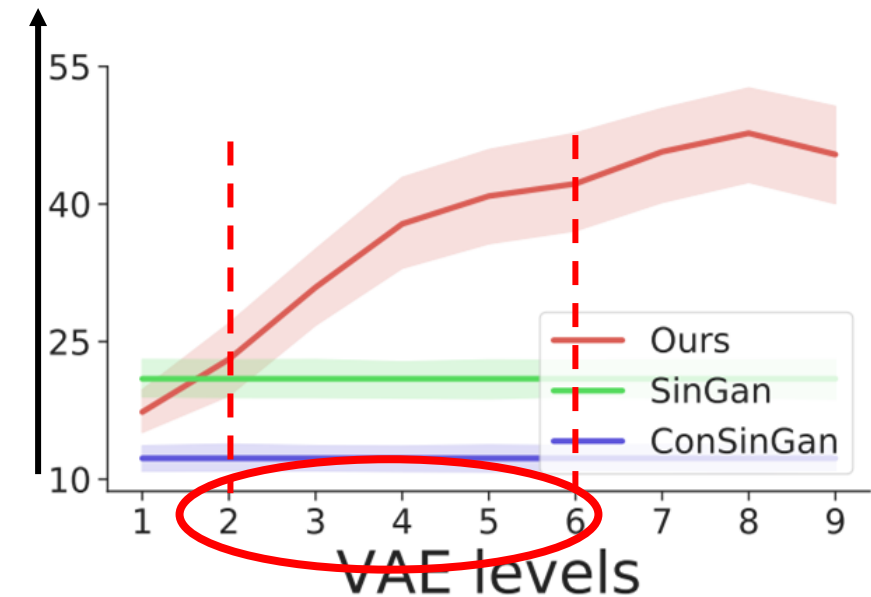
# Effect of Number of patch-VAE levels

Total of 9 layers

Quality  
(**Lower** is Better)



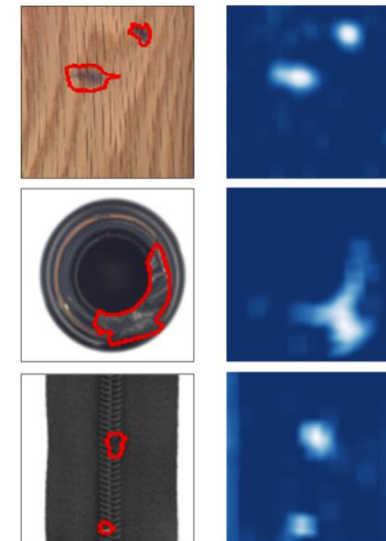
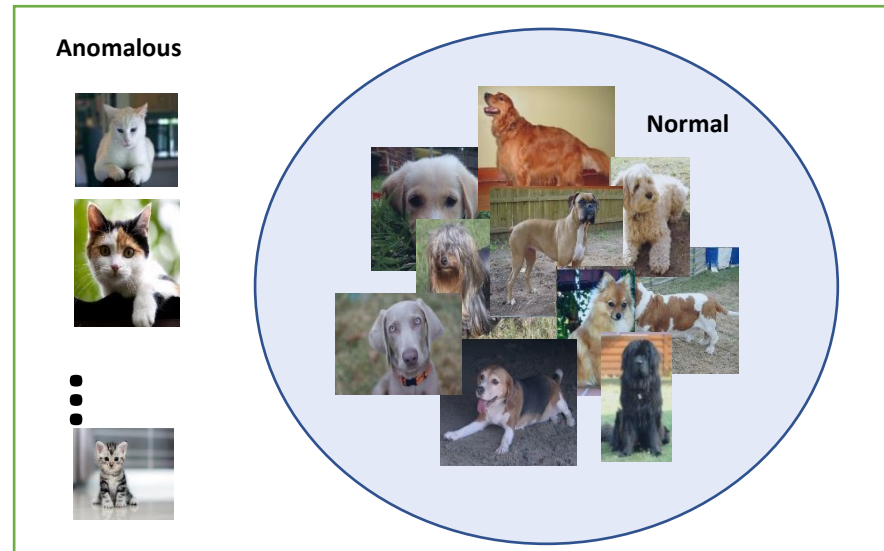
Diversity  
(**Higher** is Better)







# Manipulating ~~Structure~~ Understanding Structure



# SpeedNet: Learning the Speediness in Videos

**S. Benaim**, A. Ephrat, O. Lang, I. Mosseri, W. T. Freeman, M. Rubinstein, M. Irani, T. Dekel.  
CVPR 2020.

Slower



Normal speed



Faster



<https://speednet-cvpr20.github.io/>



# Automatically predict “speediness”

Uniform Speed Up (2x)



Adaptive speed up (2x)



---

Other Applications:

- Self-supervised action recognition
- Video retrieval

# SpeedNet

Self-supervised  
training



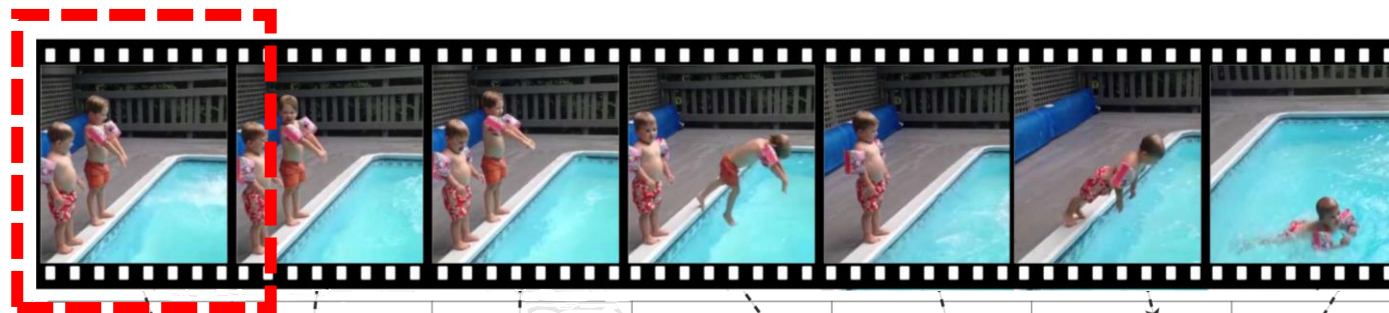
Input video



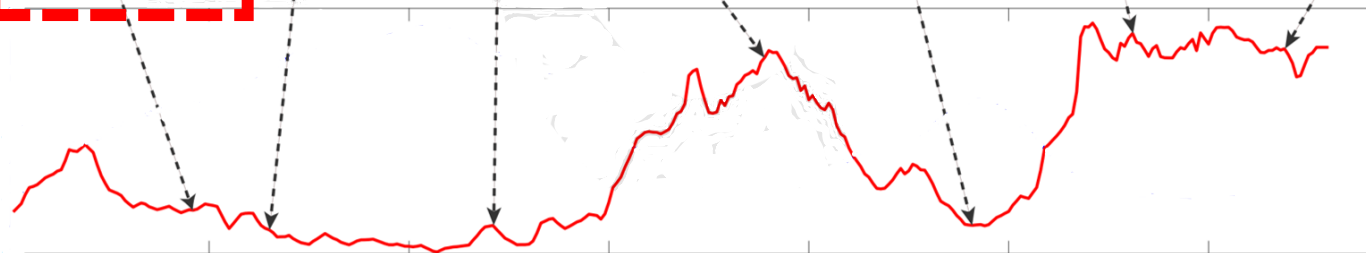
Sped Up

Inference on full  
**sped-up** video

Sped-up



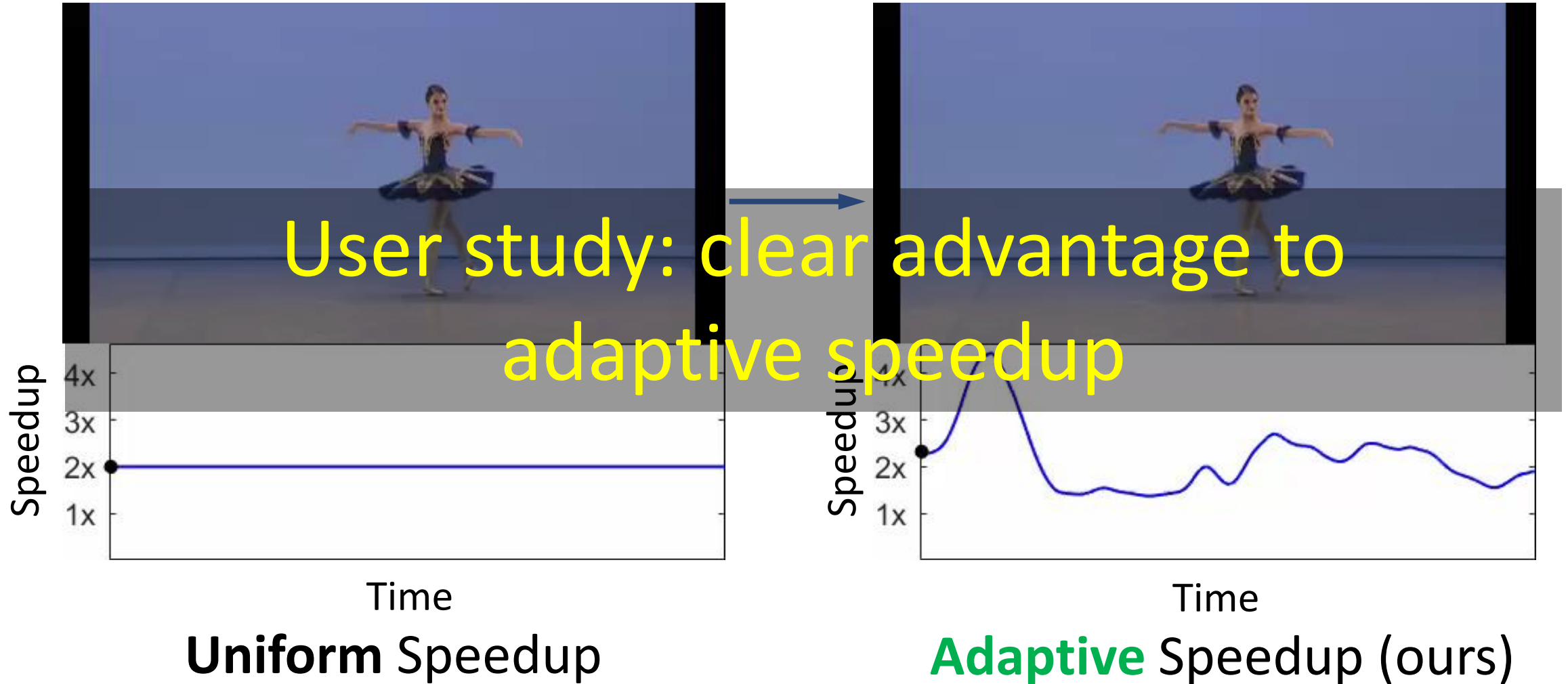
Normal speed



# Adaptive video speedup

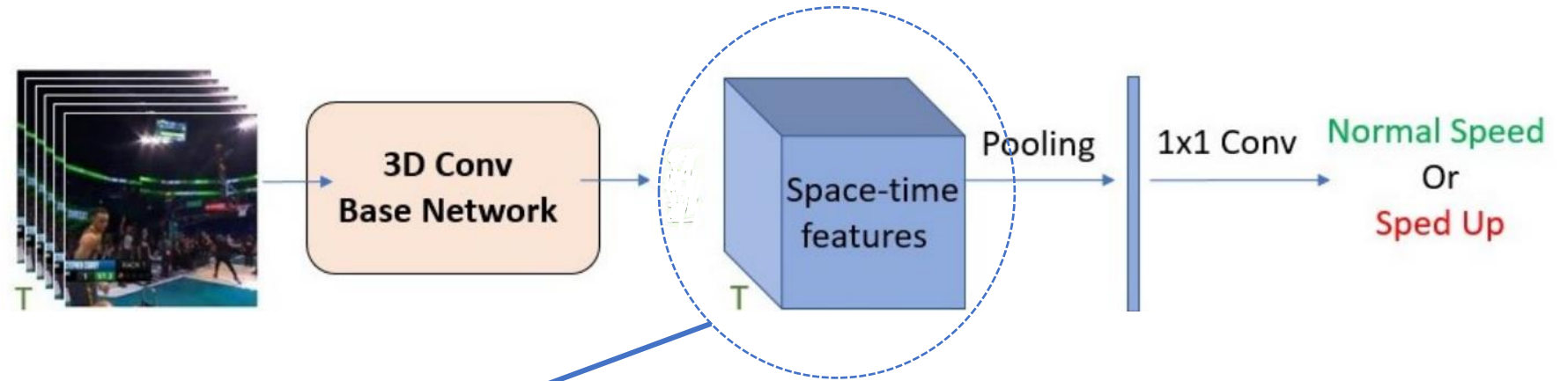
Total time =  $\frac{1}{2}$  input time

Total time =  $\frac{1}{2}$  input time



# Other self supervised tasks

Train SpeedNet



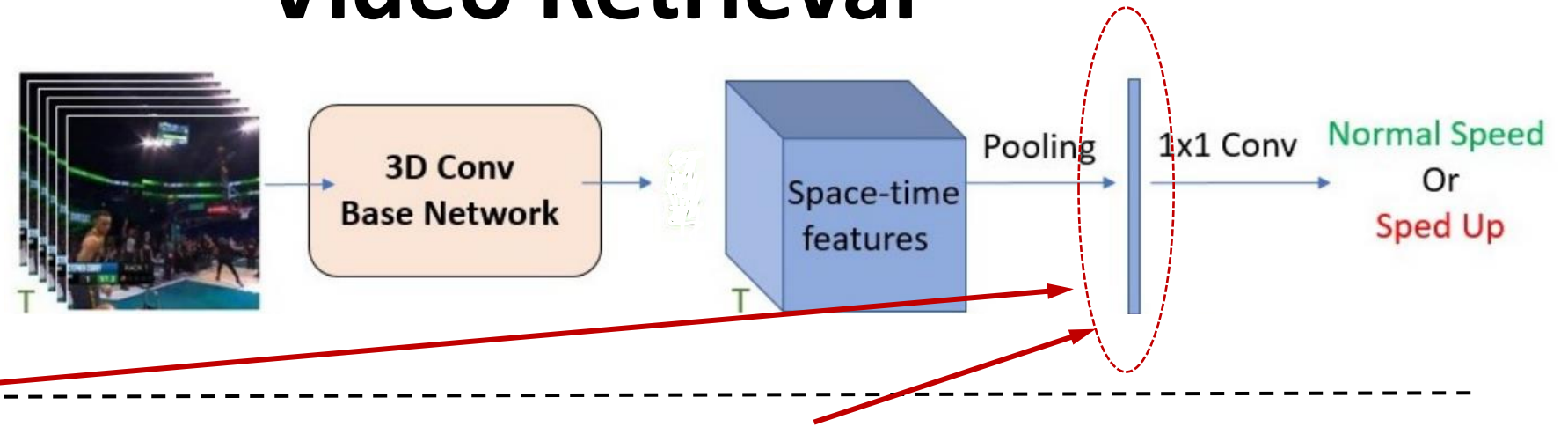
## Self Supervised Action Recognition

Method	Initialization Architecture	Supervised accuracy	
		UCF101	HMDB51
Random init	S3D-G	73.8	46.4
ImageNet inflated	S3D-G	86.6	57.7
Kinetics supervised	S3D-G	96.8	74.5
CubicPuzzle [19]	3D-ResNet18	65.8	33.7
Order [40]	R(2+1)D	72.4	30.9
DPC [13]	3D-ResNet34	75.7	35.7
AoT [38]	T-CAM	79.4	-
SpeedNet (Ours)	S3D-G	81.1	48.8
Random init	I3D	47.9	29.6
SpeedNet (Ours)	I3D	66.7	43.7



# Other self supervised tasks: Video Retrieval

Train SpeedNet



Query



Retrieved top-3 results (Within)



Query



Retrieved top-3 results (Across)





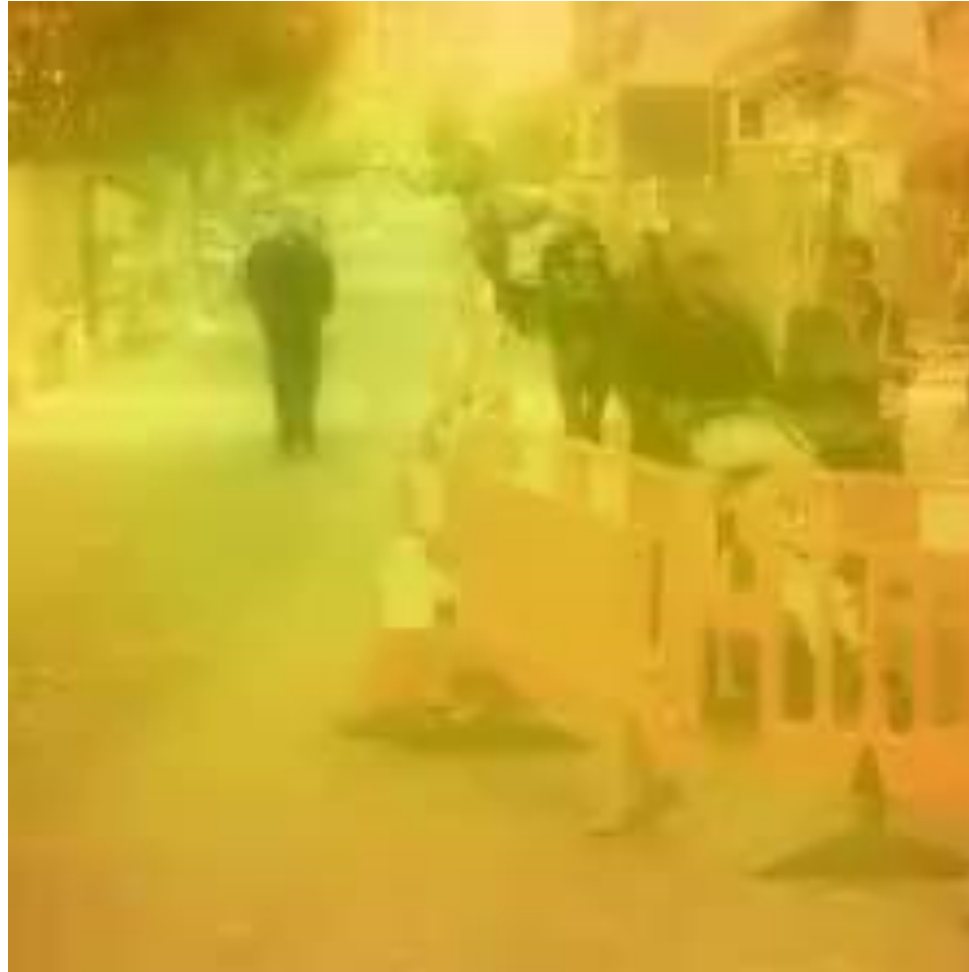
“Memory Eleven”: An artistic video by Bill Newsinger:  
[https://www.youtube.com/watch?v=djylS0Wi\\_lo](https://www.youtube.com/watch?v=djylS0Wi_lo)



# Spatio-Temporal Visualizations

blue/green =  
normal speed

yellow/orange =  
slowed down



## Manipulating Structure

- Multi-sample approaches
- Structural analogies
- Novel videos of similar structure
- Few shot anomaly detection

## Manipulating by Understanding Structure

- Speed up videos “gracefully” using “speed” as supervision
- Image classification and domain adaptation by reducing bias towards global statistics (CVPR 2021)

**Structure** is Key to **Image Understanding**

**Demonstrate** using **Structure Aware Manipulation**

### Next?

- 3D-aware structure manipulation
- Manipulating multiple objects from multiple scenes
- Functional relationships: A person riding a bike vs a person beside a bike

Thank You! Questions?