

Synthesis: Image Generation and Hole Filling

Danna Gurari

University of Colorado Boulder
Fall 2021



Review

- Last lecture topic:
 - Synthesis: style transfer
- Assignments
 - Final project outline due tonight
 - Final project presentation due in three weeks
 - Peer evaluation due in three weeks
 - Final project report due in four weeks
- Questions?

Today's Topics

- Problem
- Applications
- Image generation methods
- Hole filling methods
- Evaluation approaches

Today's Topics

- Problem
- Applications
- Image generation methods
- Hole filling methods
- Evaluation approaches

Synthesis (With and Without Context)

Generation & Alteration:



Hole Filling
(constrained generation):

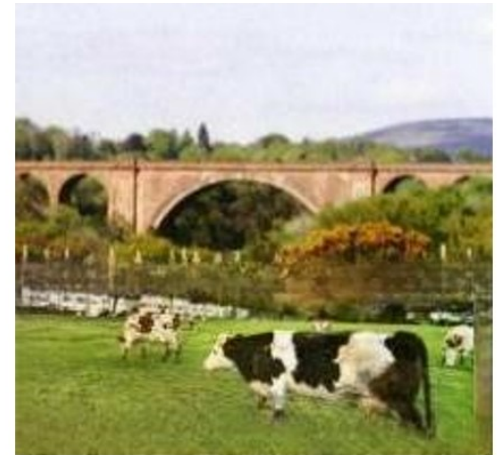


Image source: <https://medium.com/image-recreation-a-method-to-make-learning-of-gan/image-generation-text-to-image-d7c4210ecb90>

Today's Topics

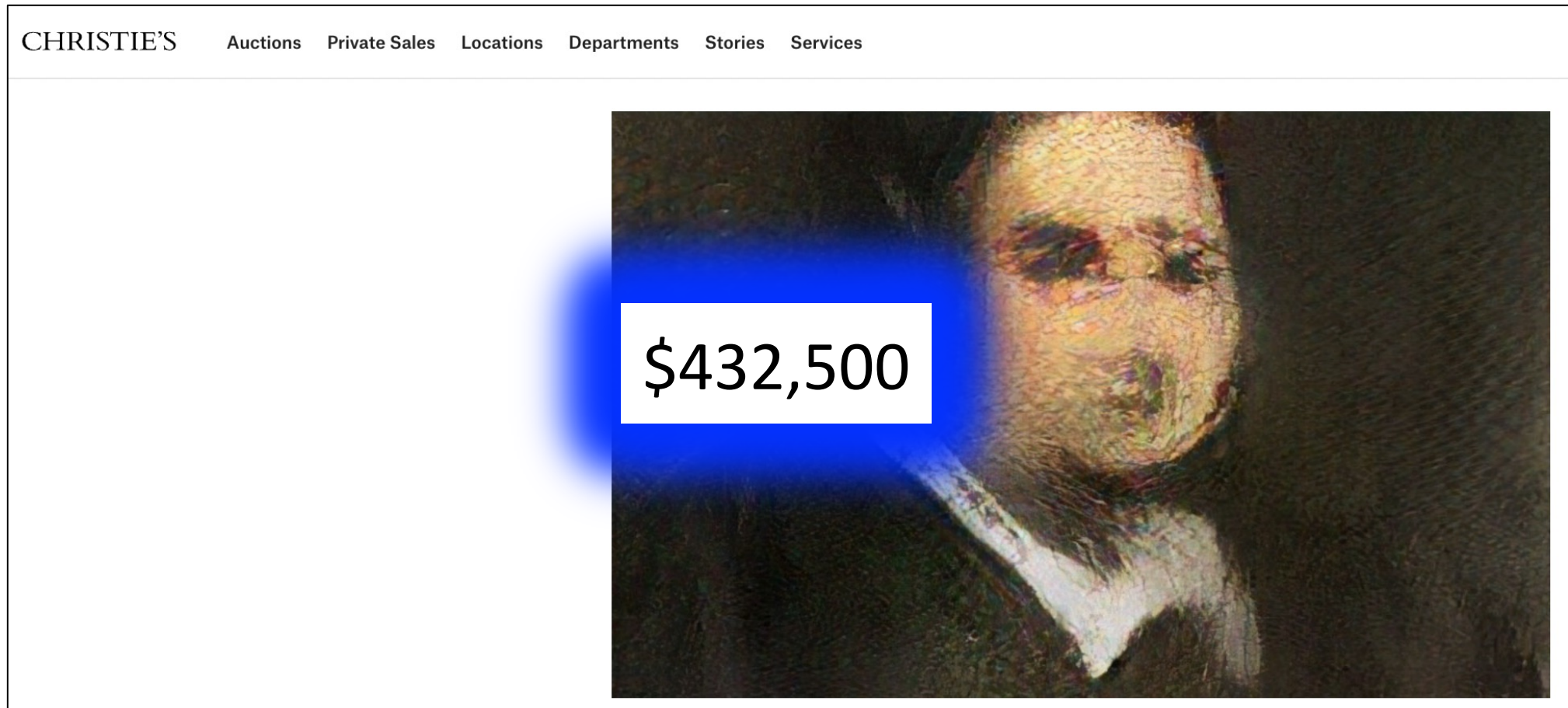
- Problem
- **Applications**
- Image generation methods
- Hole filling methods
- Evaluation approaches

Refinement (e.g., enhance, avoid payment, and/or rewrite history)

- Damaged regions from camera
- Blurred areas
- Watermarks
- Undesired content e.g., remove ex-partner
→



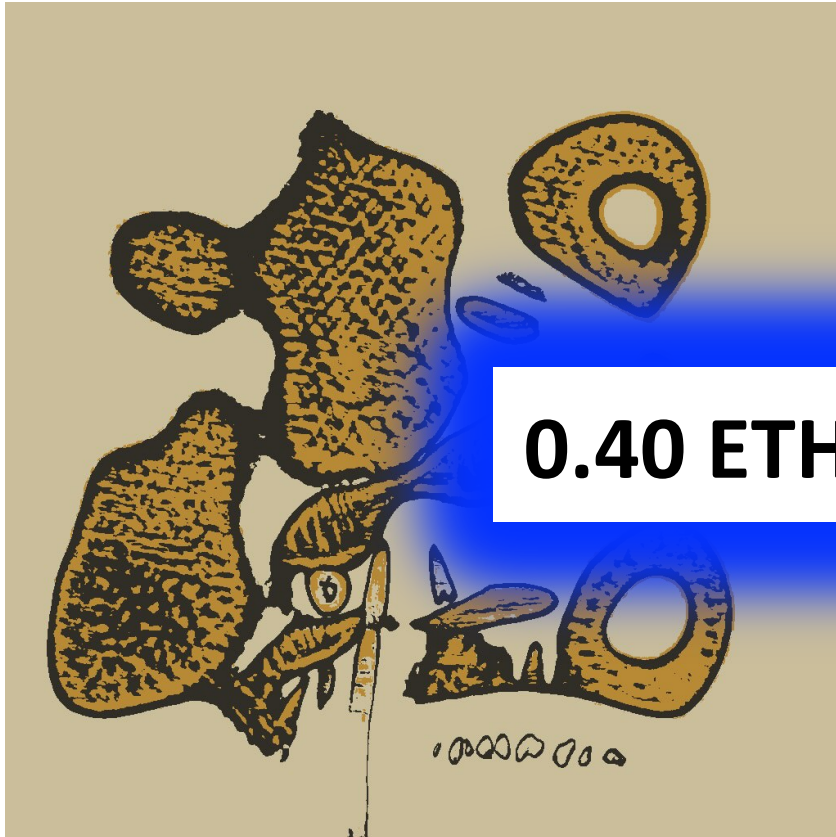
Commercial Art




How much do you think this sold for?

<https://www.christies.com/features/A-collaboration-between-two-artists-one-human-one-a-machine-9332-1.aspx>

Commercial Art



 @deviparikh

Attending to details


↓ Artwork information

0.40 ETH (\$1,700.89)

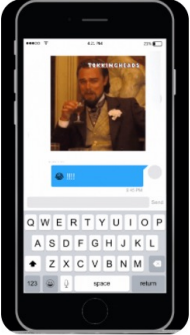

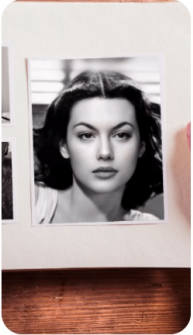
...an analog art, and part human digital
...this piece is generated using a neural generative model
(AI) trained on the artist's physical sketchbook. The
generation has been re-rendered in a new palette. The
physical sketchbook had sketches, doodles, paper cutting,
mandalas, zentangles -- often with a brief hand-written
description.

How much do you think this sold for?


Entertainment

 Company

Animate Old Photos Bring your Art to Life Prank your Friends



Puppet any Avatar from just an Image





Social Media

The image shows a screenshot of an Instagram profile for the user 'lilmiquela'. The profile is verified and has 1,146 posts, 3.1 million followers, and 1,911 accounts following. The bio identifies the user as 'Miquela', a 19-year-old robot living in LA, and includes the hashtag #BlackLivesMatter. A bio link is provided: smarturl.it/MiquelaTikTok?iqid=m.ig. Below the bio, there are seven story highlights with circular icons and labels: 'LIFE RN', 'EVOLVING...', 'FITS', 'MAY', 'APRIL', 'MARCH', and 'FEBRUAR...'. At the bottom, the navigation bar shows 'POSTS' selected, with 'REELS', 'VIDEOS', and 'TAGGED' also visible. Three post thumbnails are shown at the bottom: a group of five young people, a woman and a man smiling together, and a woman and a man sitting together.

Instagram

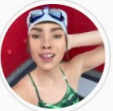


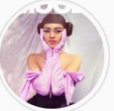
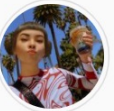

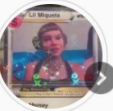
Search

Log In Sign Up




 **lilmiquela**  [Follow](#)

1,146 posts 3.1m followers 1,911 following

Miquela
#BlackLivesMatter
🤖 19-year-old Robot living in LA 💕
(still figuring the rest out, one post at a time)
Check out my new video 📺👇
smarturl.it/MiquelaTikTok?iqid=m.ig

 LIFE RN 🍷
 EVOLVING...
 FITS
 MAY 💜
 APRIL 🍷
 MARCH 🍷
 FEBRUAR...

POSTS REELS VIDEOS TAGGED

News



e.g., face re-enactment

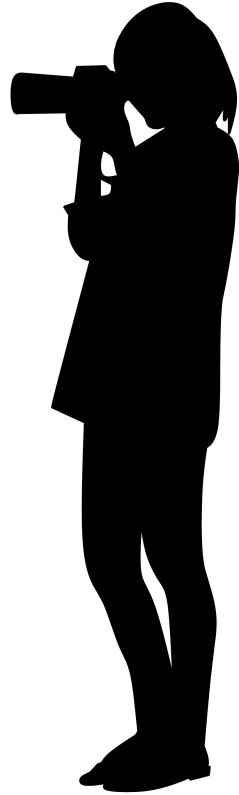
Demo: <https://www.youtube.com/watch?v=ttGUIwfTYvg>

Training Data Creation



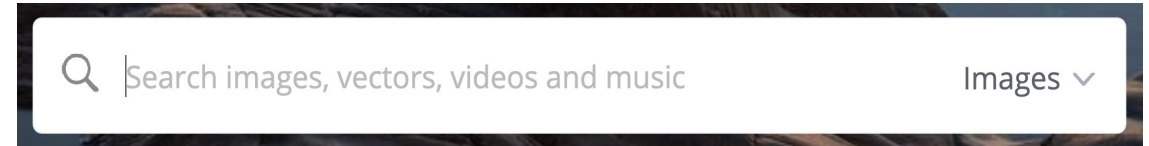
Improved Messaging via Visual Content

- Marketing
- Artwork
- Presentations
- Blogs
- Websites



Potential sources:

Photographer
(self or hired)



shutterstock



Pexels



Adobe Stock

BIGSTOCK™



dreamstime®

Google
Images



Unsplash
Photos for everyone

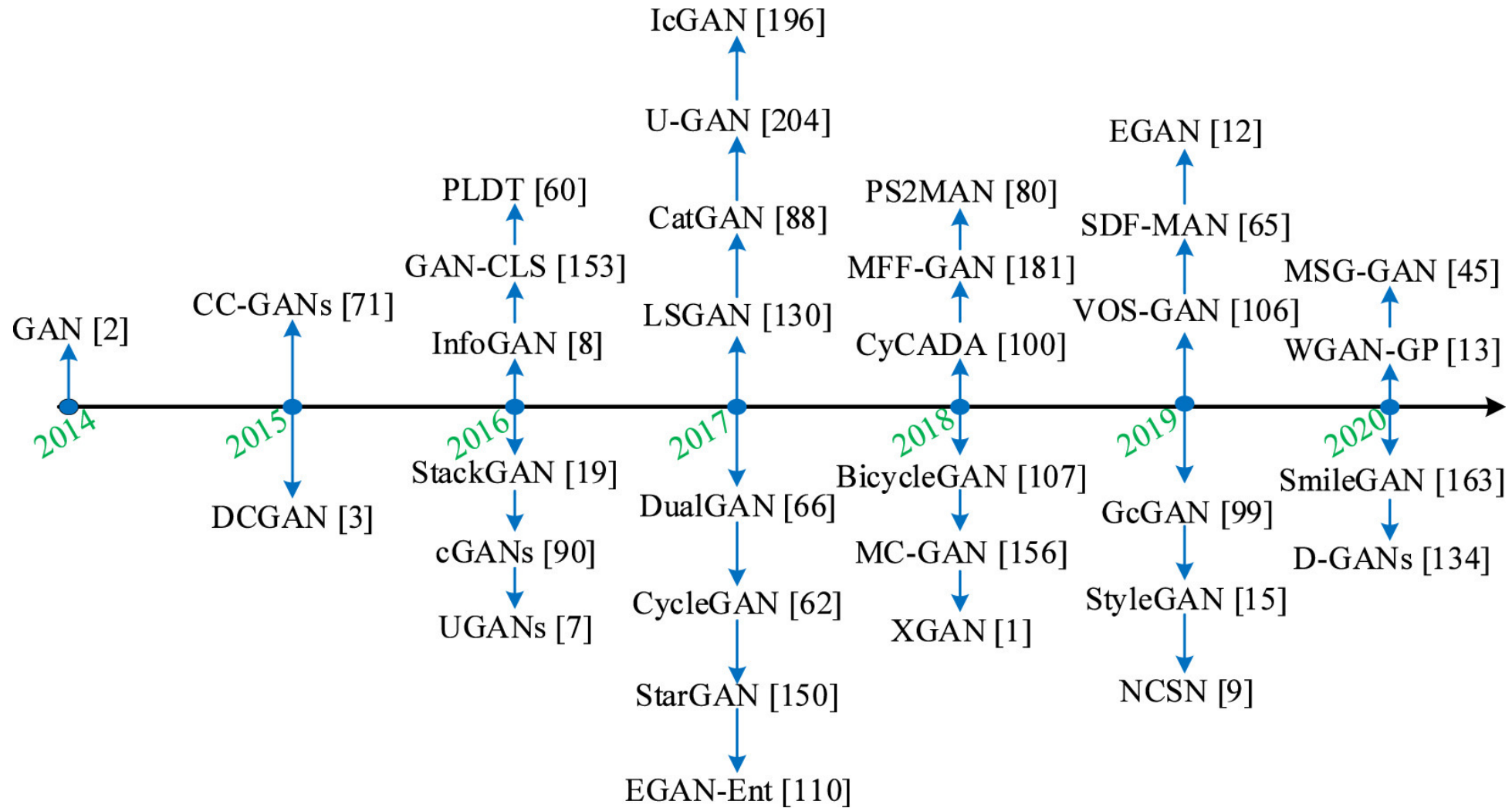
Stock photos

What are other possible applications of image synthesis, for good and harm?

Today's Topics

- Problem
- Applications
- **Image generation methods**
- Hole filling methods
- Evaluation approaches

Generative adversarial networks



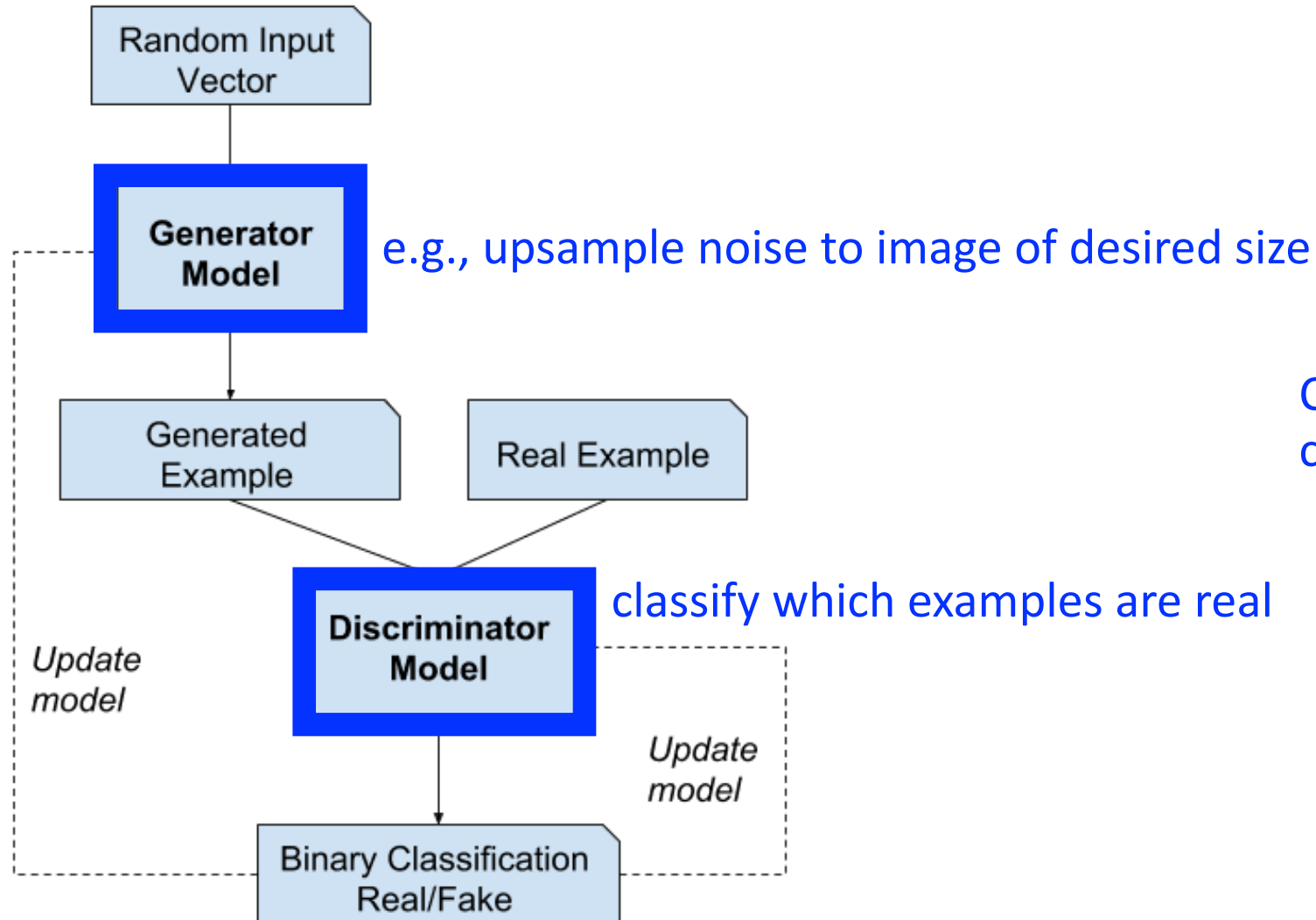
Generative adversarial networks

- Generative adversarial networks (GANs)
- Deep convolutional generative adversarial networks (DCGANs)
- GIRAFFE

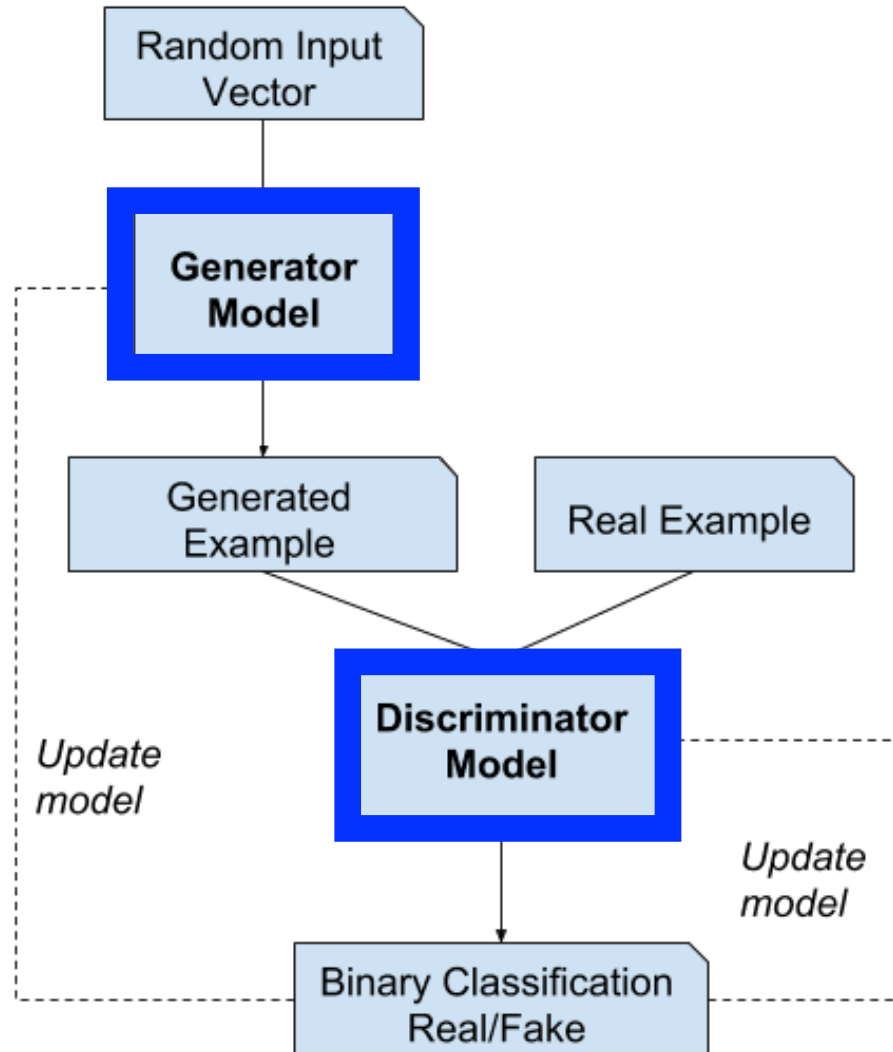
Generative adversarial networks

- Generative adversarial networks (GANs)
- Deep convolutional generative adversarial networks (DCGANs)
- GIRAFFE

GAN: Basic Architecture



GAN: Training



The two models are iteratively trained separately

- Train discriminator using fake and real images
- Train generator using just fake images and penalize it when the discriminator recognizes images are fake

GAN: Discriminator Loss Function

Discriminator tries to minimize classification error

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

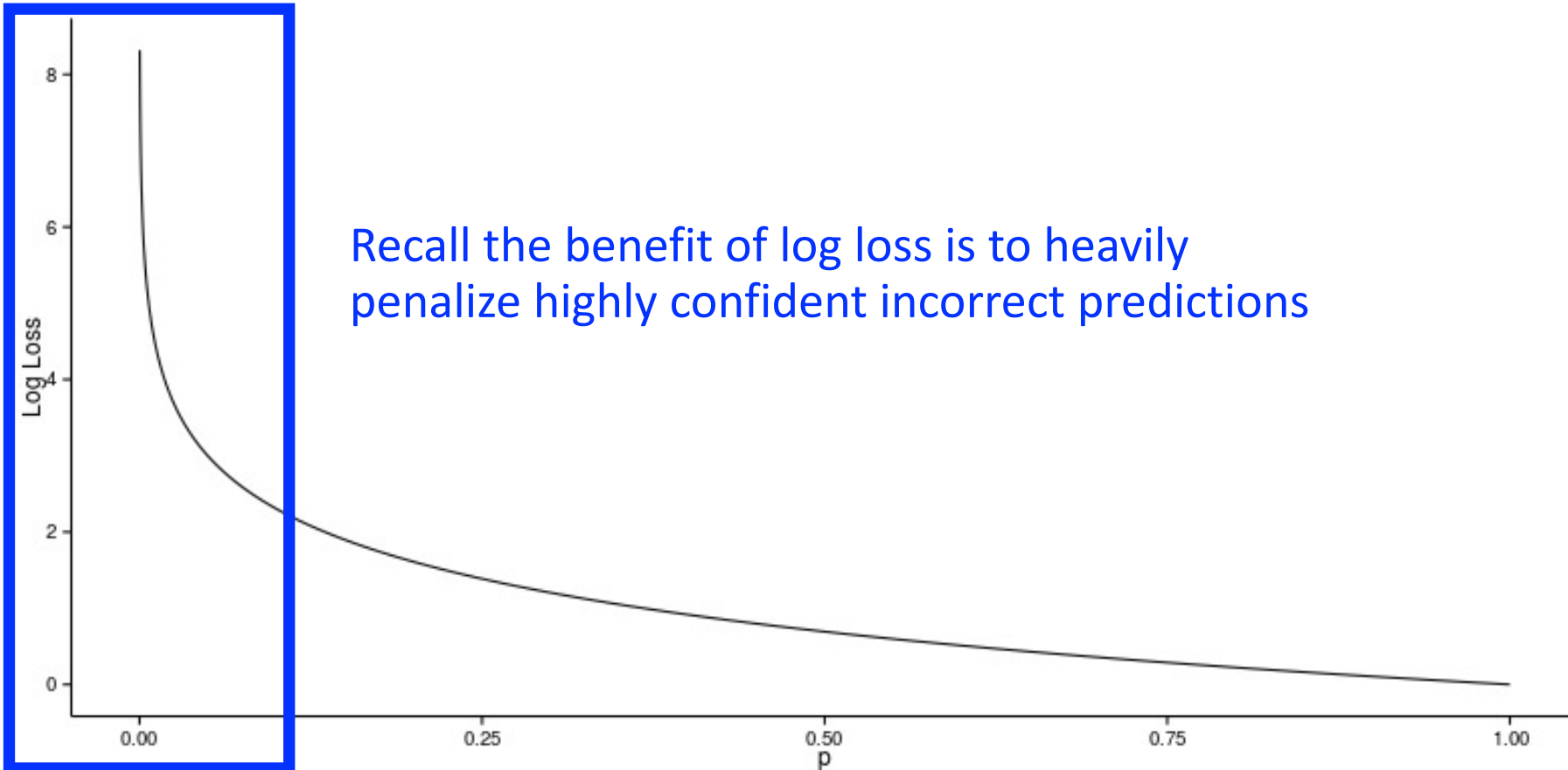
Discriminator wants a value of 1 for real images

Discriminator wants a value of 0 for fake images

Real image

Input noise

GAN: Discriminator Loss Function



GAN: Generator Loss Function

Generator tries to maximize classification error

$$J^{(G)} = -J^{(D)}$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Want the discriminator to mistakenly arrive at a value of 1 for fake images

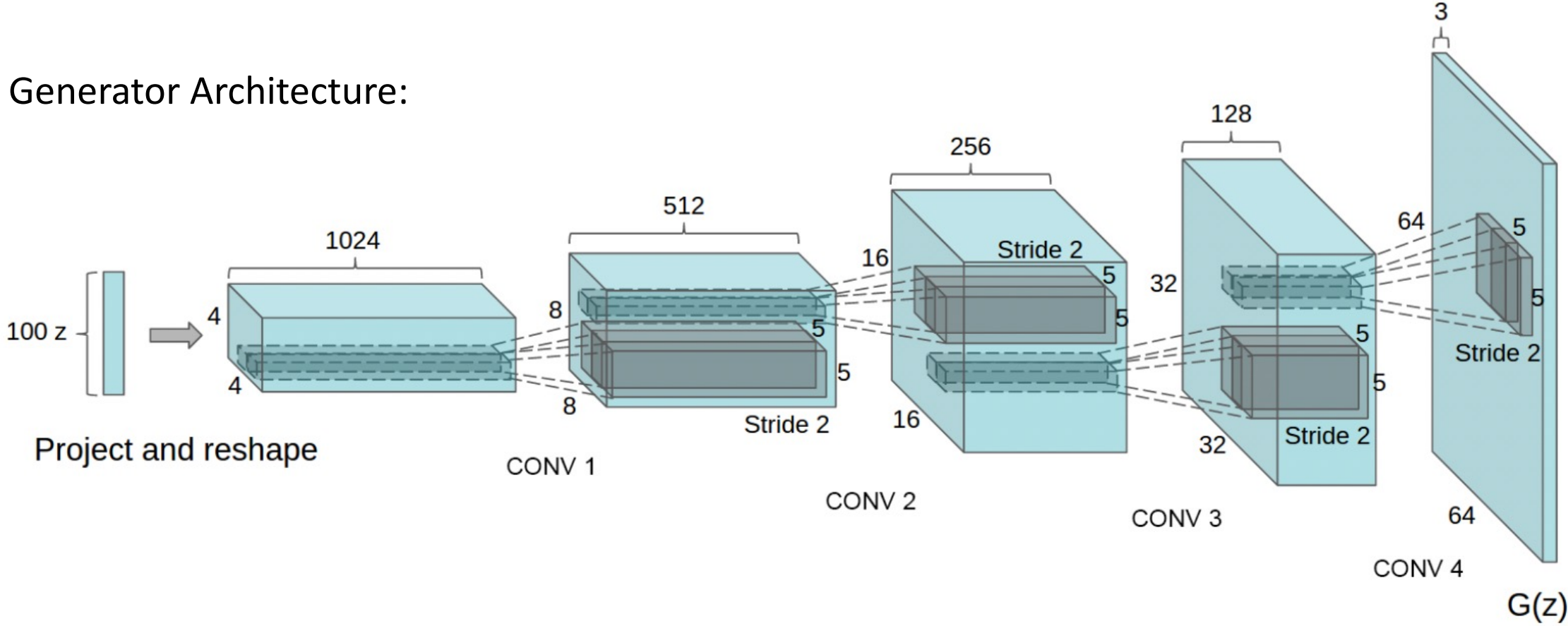
Input noise

Generative adversarial networks

- Generative adversarial networks (GANs)
- Deep convolutional generative adversarial networks (DCGANs)
- GIRAFFE

DGANs: GANs that Use Convolutional Layers

Generator Architecture:



What is the resolution of the image generated by this network?

DGANs: Qualitative Results



Bedrooms generated by observing over 3M bedroom images

DGANs: Qualitative Results



What objects does it learn to generate?

DGANs: Qualitative Results



What objects may it not have learned to generate?

DGANs: Qualitative Results



Faces generated by observing over 3M images of 10K people

DGANs: Qualitative Results



What does it generate poorly or not all?

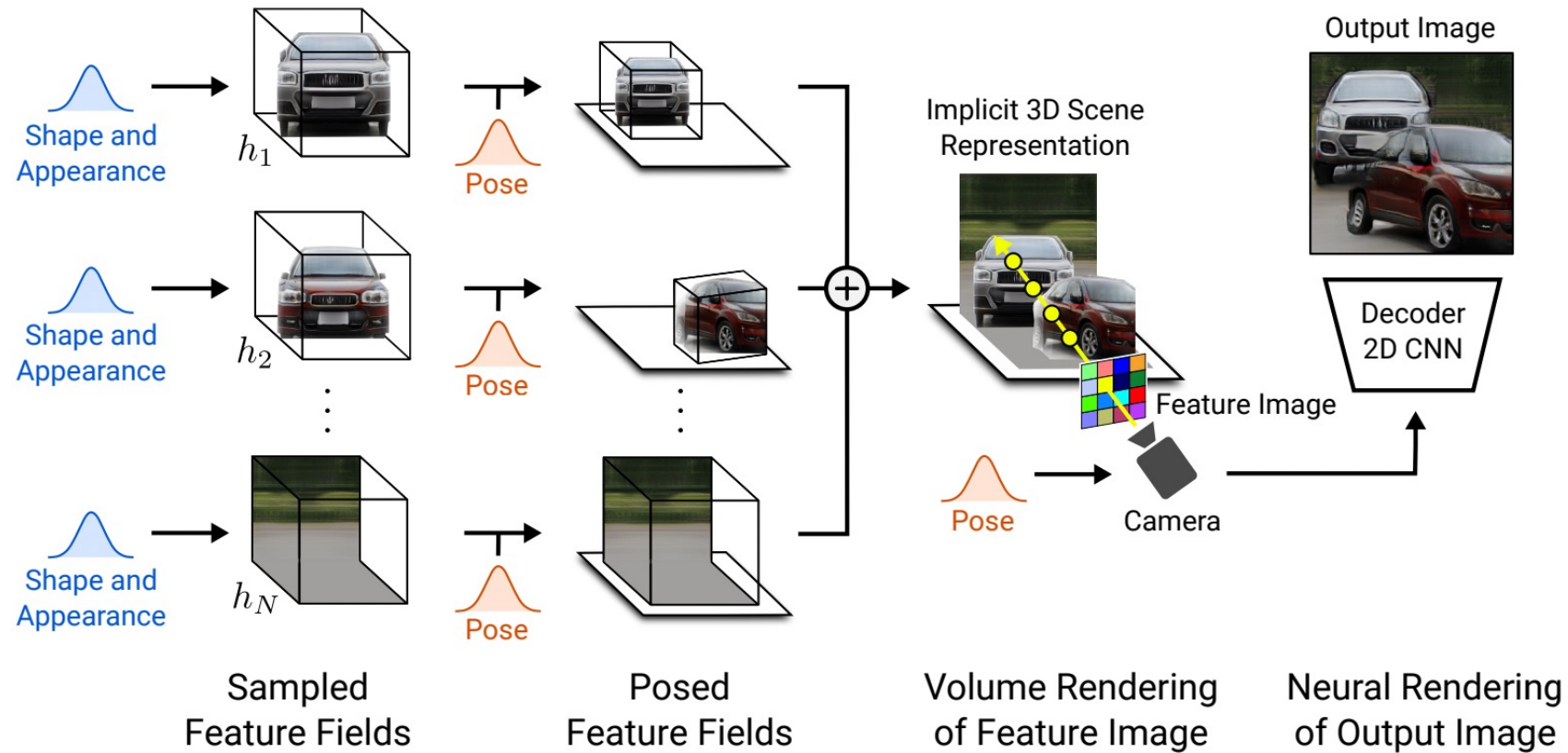
DGANs: Limitation

Cannot control what is generated

Generative adversarial networks

- Generative adversarial networks (GANs)
- Deep convolutional generative adversarial networks (DCGANs)
- **GIRAFFE**

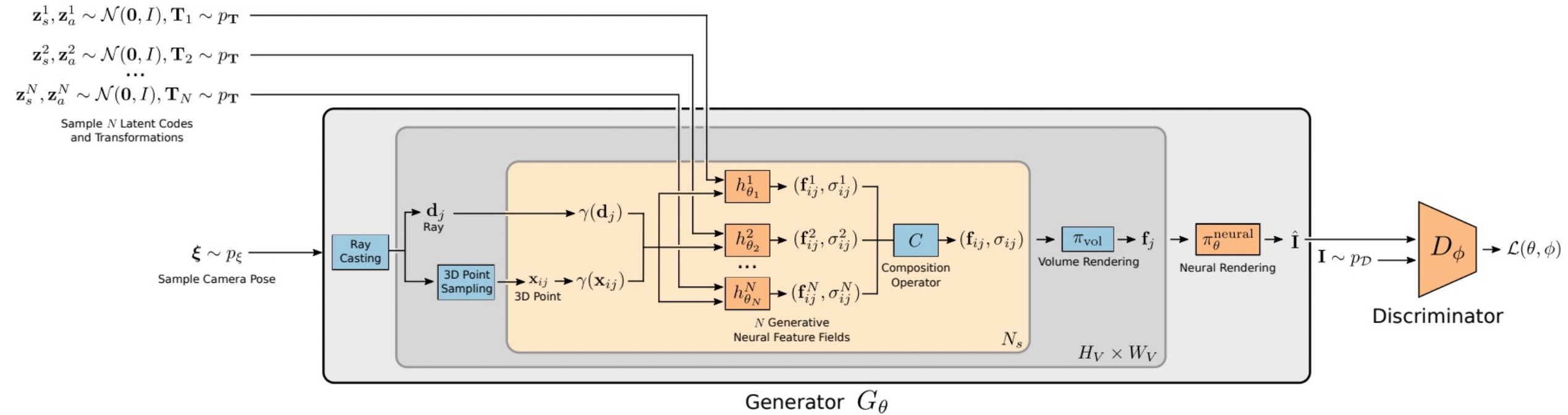
GIRAFFE: Idea



Key idea: control what is synthesized using a 3D scene representation in the generator

(Recognized with Best Paper Award)

GIRAFFE: Architecture

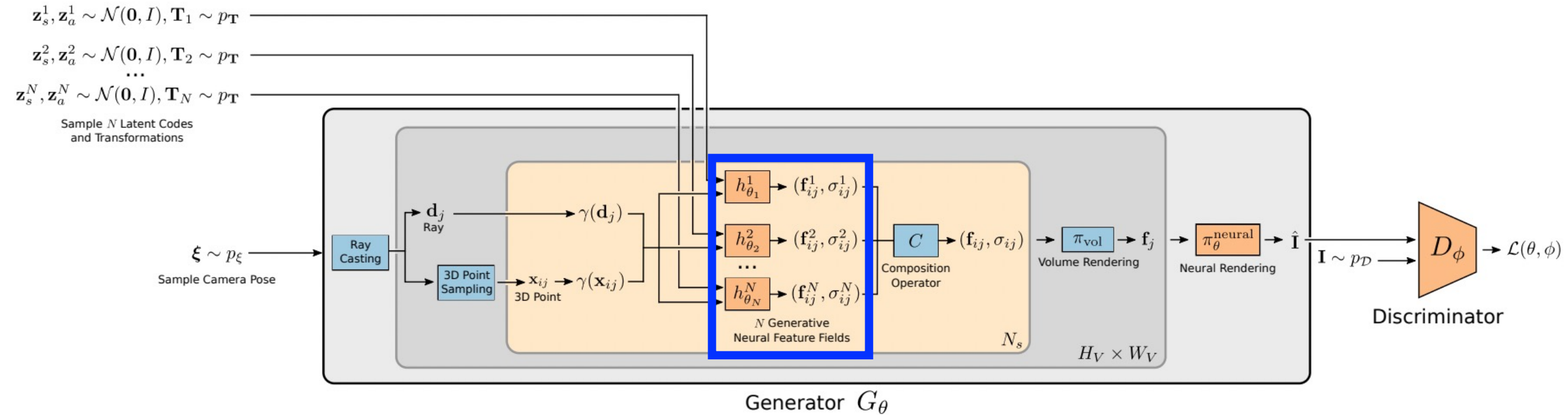


Key idea: control what is synthesized using a 3D scene representation in the generator

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

GIRAFFE: Architecture

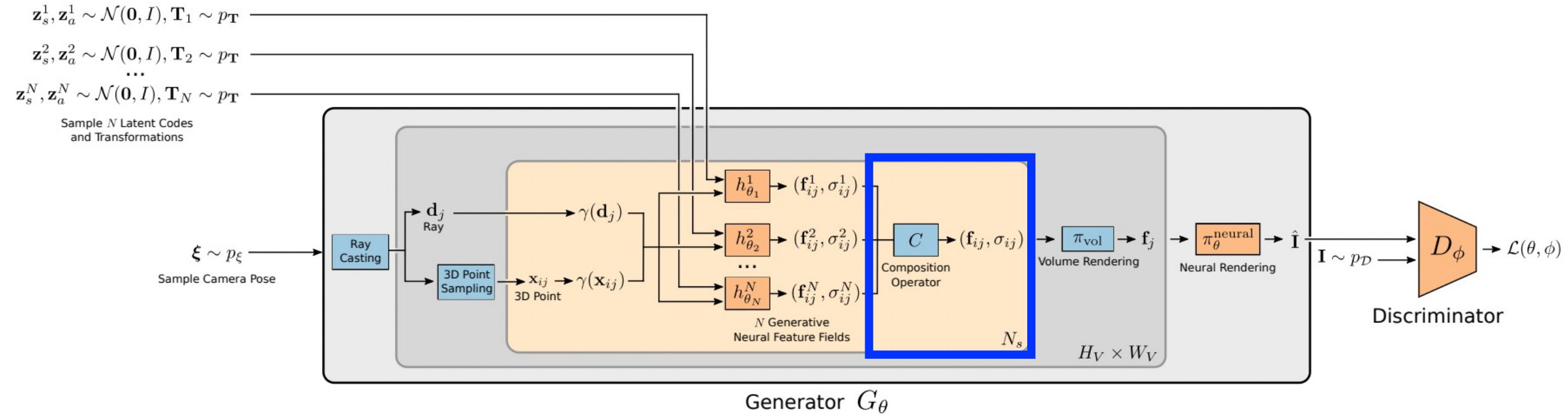


The $N-1$ objects and background to appear in the scene are not only represented separately but also separate from their shape and appearance

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

GIRAFFE: Architecture

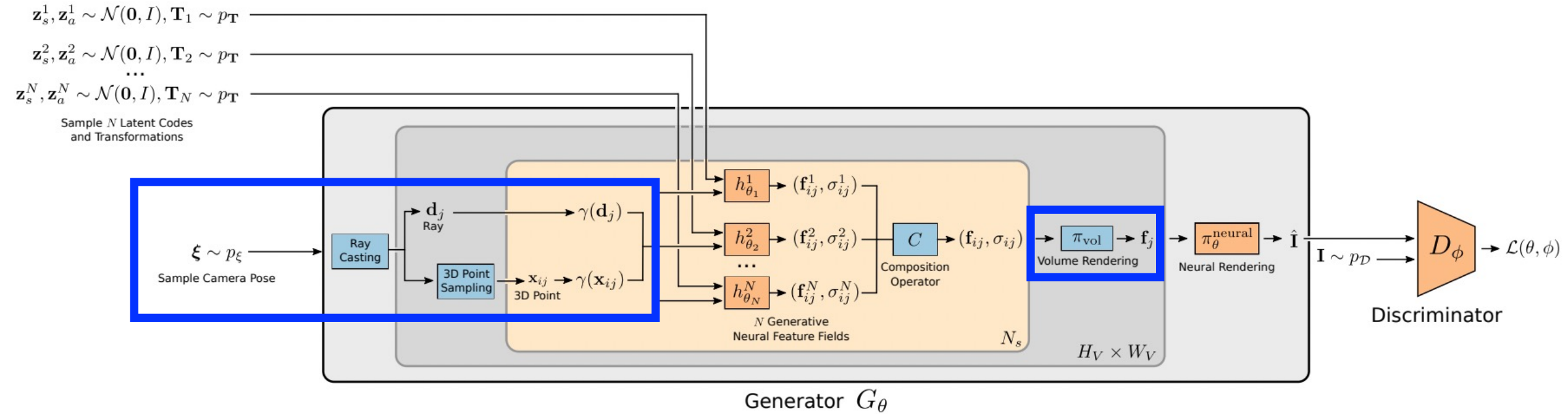


These N entities are then incorporated into a scene representation

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

GIRAFFE: Architecture

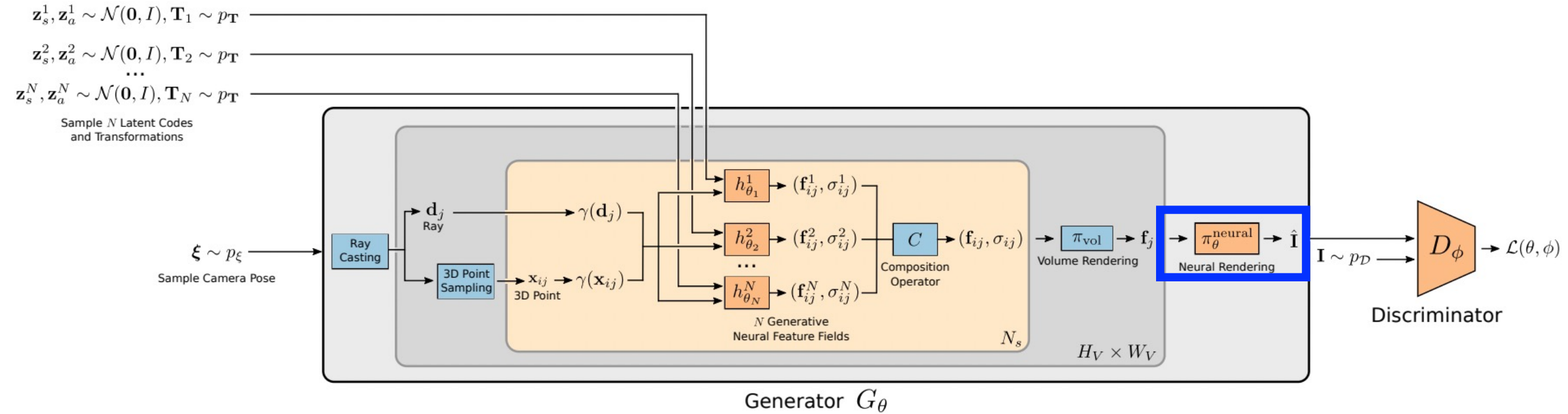


Knowledge about the camera pose is used to render a high dimensional feature vector for each pixel

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

GIRAFFE: Architecture



The final 2D image is then rendered

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

GIRAFFE: Qualitative Results



(a) Object Rotation

(b) Camera Elevation

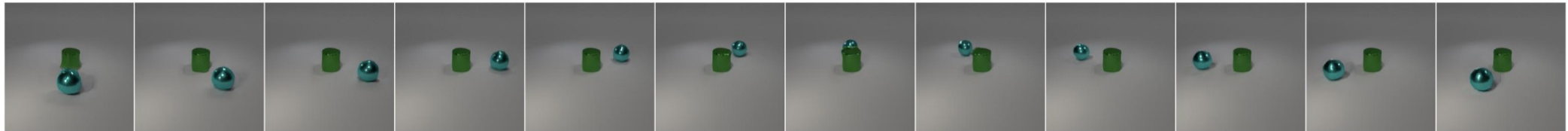


(c) Object Appearance



(d) Depth Translation

(e) Horizontal Translation



(f) Circular Translation of One Object Around Another Object

Can control synthesized results!

(Recognized with Best Paper Award)

GIRAFFE: Qualitative Results



Can control synthesized results!

(Recognized with Best Paper Award)

Niemeyer and Geiger. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR 2021.

Generative adversarial networks

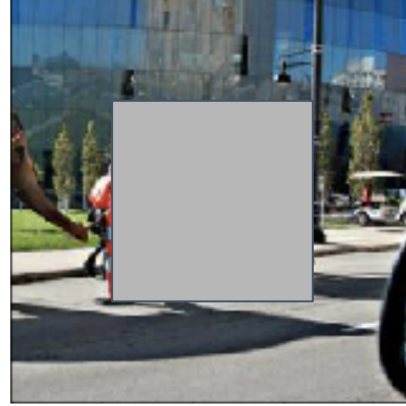
- Generative adversarial networks (GANs)
- Deep convolutional generative adversarial networks (DCGANs)
- GIRAFFE

Today's Topics

- Problem
- Applications
- Image generation methods
- **Hole filling methods**
- Evaluation approaches

Key Challenge

- What might fit into this hole?



- Many items may plausibly fit into the hole:



- Challenge: have up to 1 known ground truth region per hole

Methods

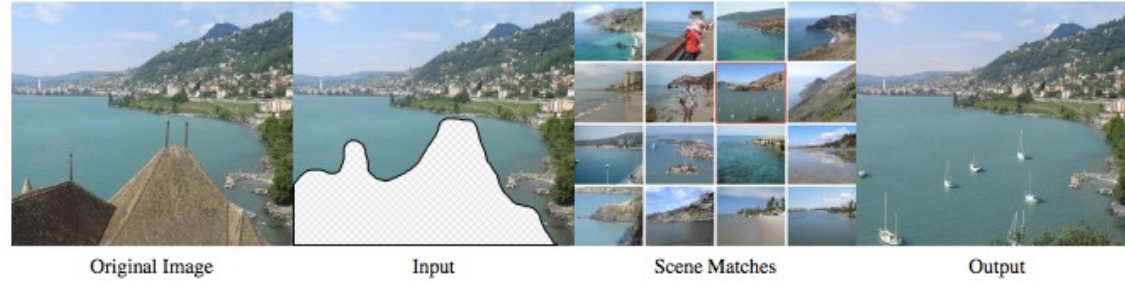
- Before deep learning era: cut-paste from nearest neighbors
- Context encoder
- Guided image inpainting

Methods

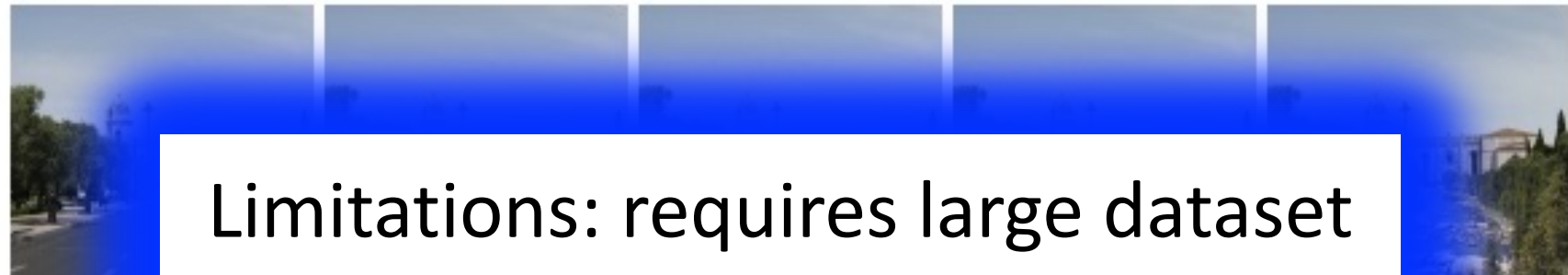
- Before deep learning era: cut-paste from nearest neighbors
- Context encoder
- Guided image inpainting

Cut-Paste from Nearest Neighbors

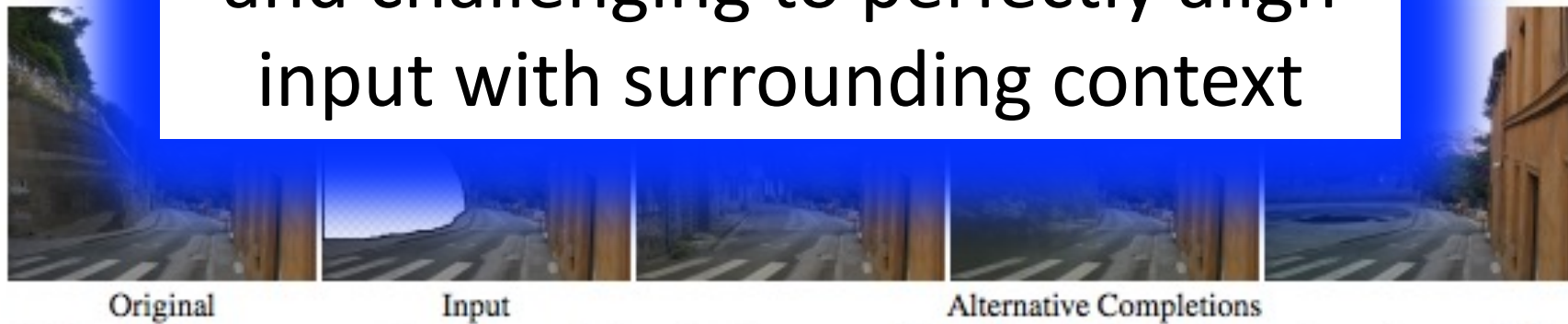
Idea:



Example:



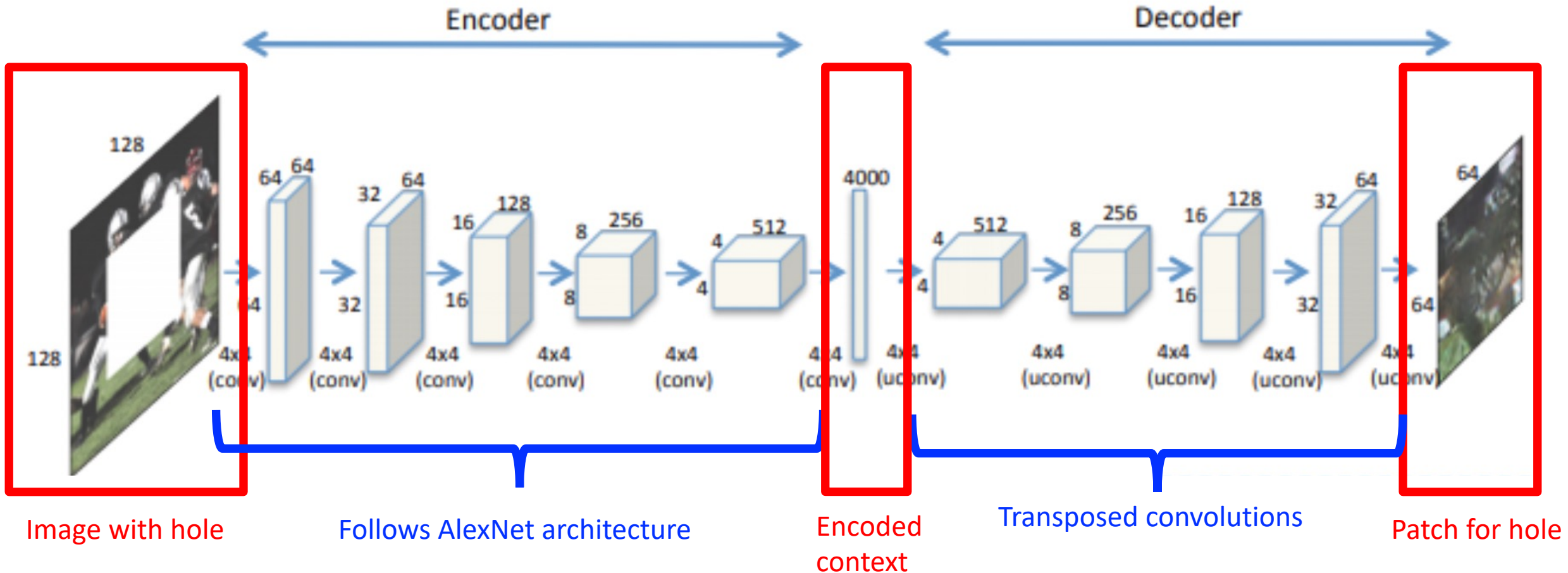
Example:



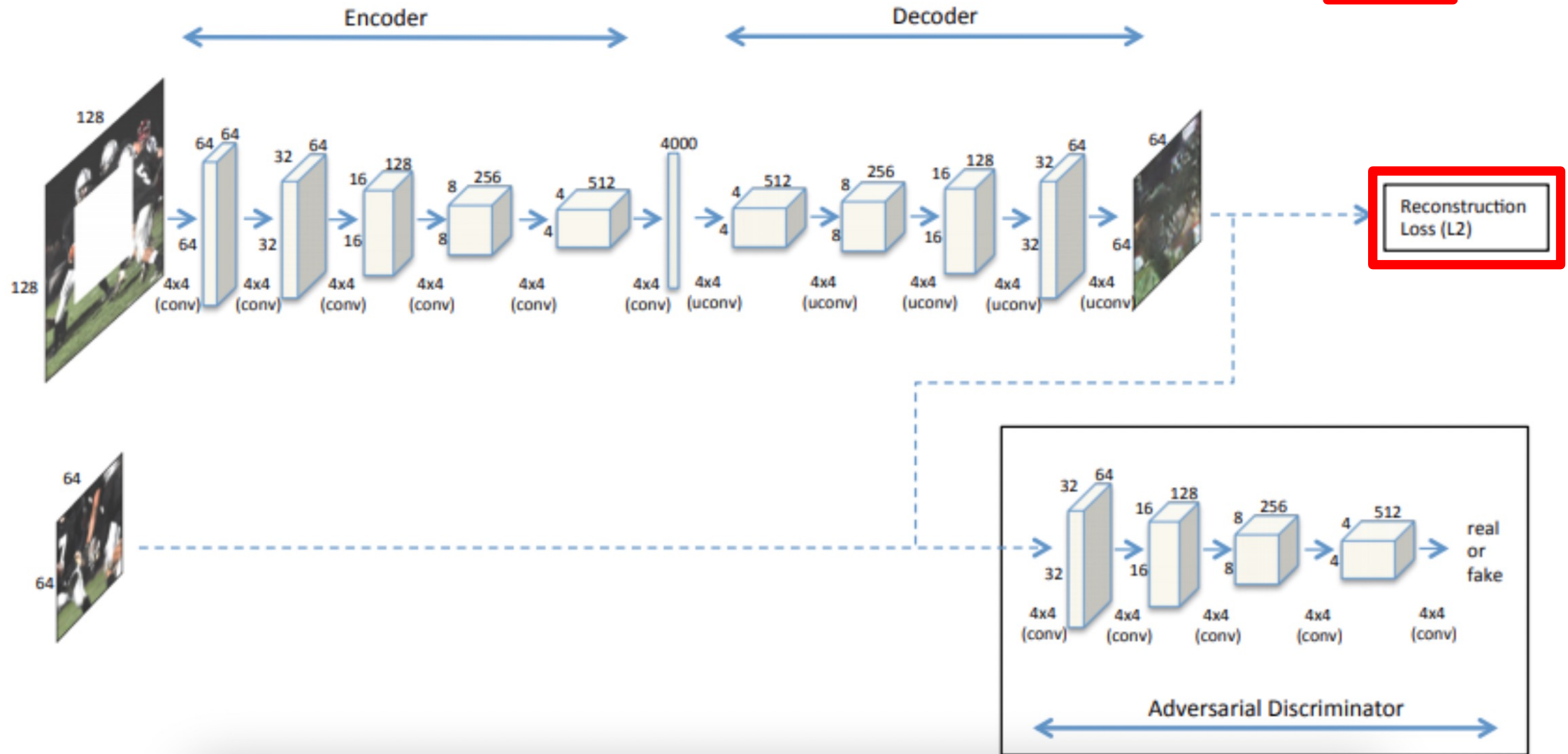
Methods

- Before deep learning era: cut-paste from nearest neighbors
- **Context encoder**
- Guided image inpainting

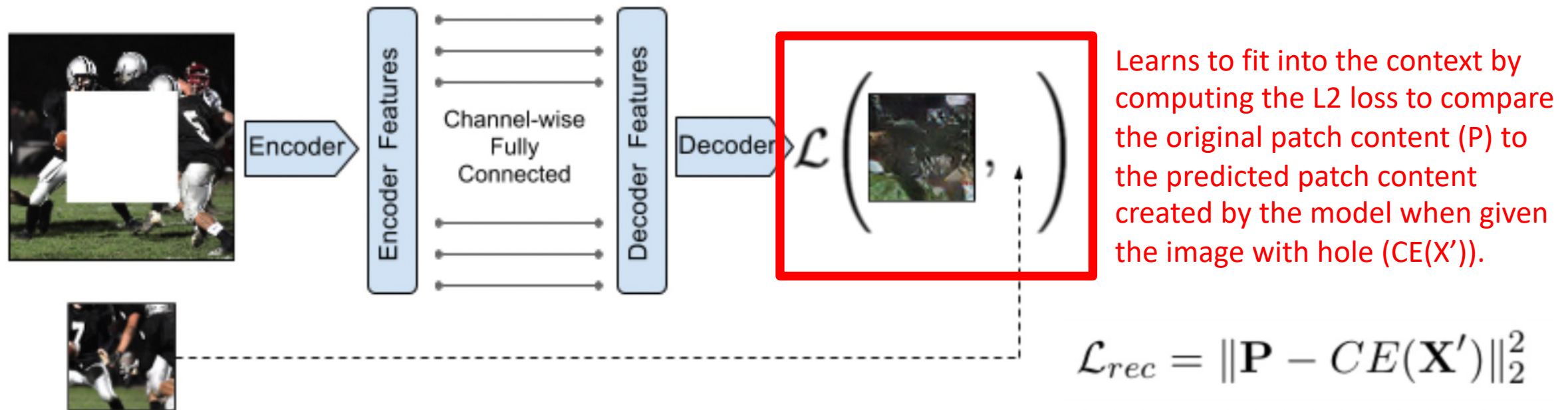
Architecture



Training: Loss Functions ($\mathcal{L} = \lambda_{adv}\mathcal{L}_{adv} + \lambda_{rec}\mathcal{L}_{rec}$)



Training: Reconstruction Loss (i.e., Self-Supervised Learning Approach)



Training: Reconstruction Loss (i.e., Self-Supervised Learning Approach)



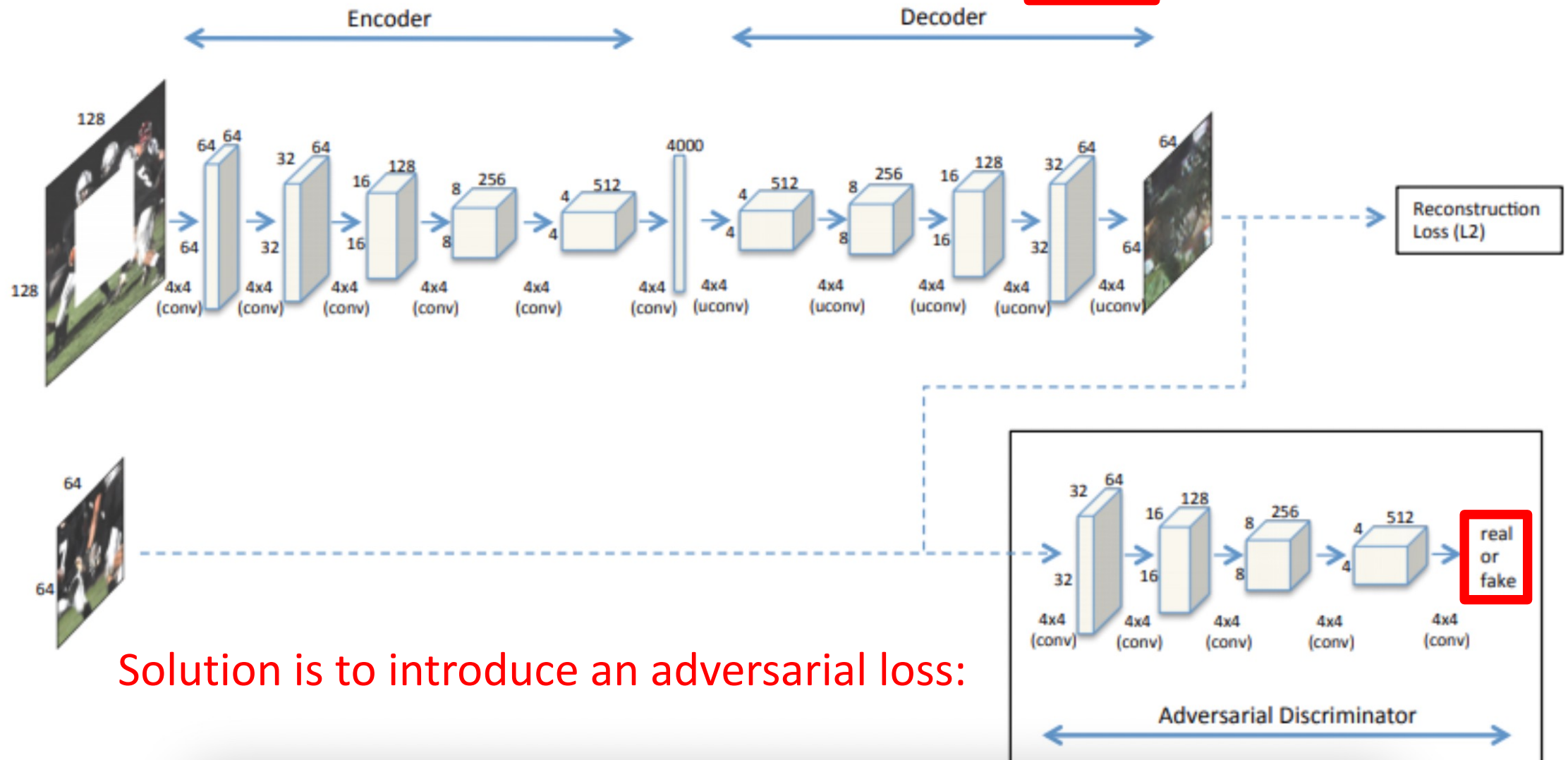
(a) Input context



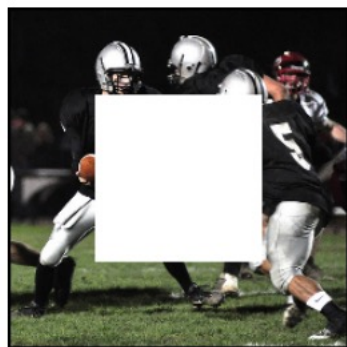
(c) Context Encoder
(L_2 loss)

Why might training with this loss function alone lead to blurry results?
- It averages the multiple plausible inpaintings for a hole

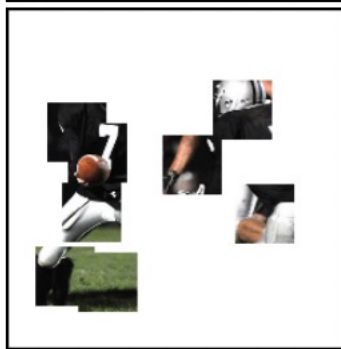
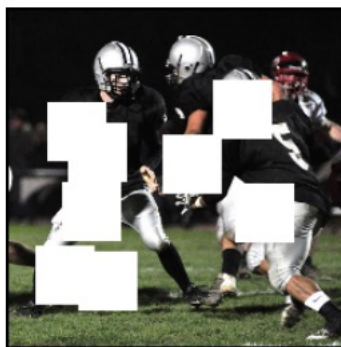
Training: Loss Functions ($\mathcal{L} = \lambda_{adv} \mathcal{L}_{adv} + \lambda_{rec} \mathcal{L}_{rec}$)



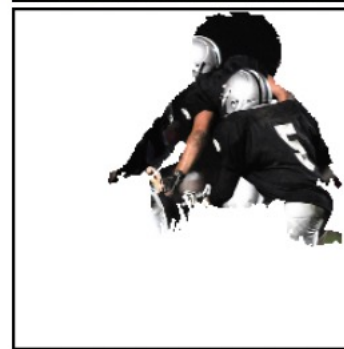
Training: Datasets



(a) Central region



(b) Random block



(c) Random region

Training completed on ImageNet (all 1.2M and a 100K subset) for three hole types

Results Demo

https://www.cs.cmu.edu/~dpathak/context_encoder/

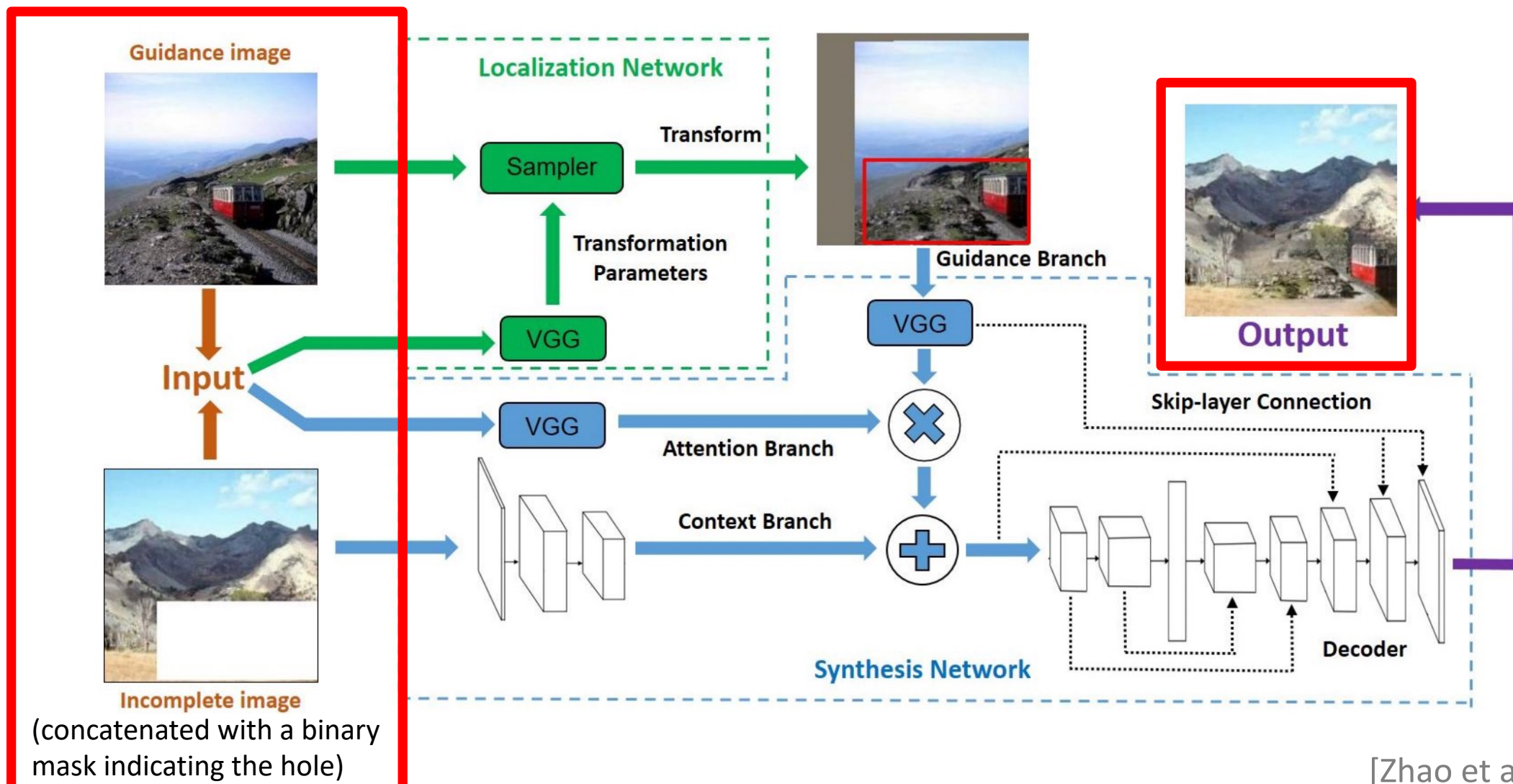
Key Limitation

Users cannot control what
content to insert in the hole

Methods

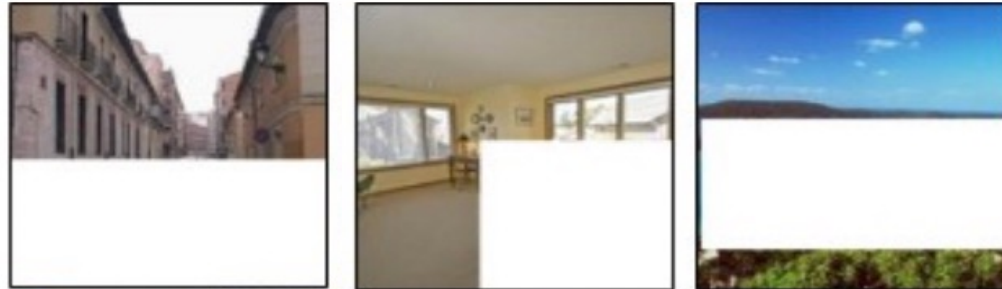
- Before deep learning era: cut-paste from nearest neighbors
- Context encoder
- Guided image inpainting

Architecture



Examples of Input and Output

Input Image:



Input user content:

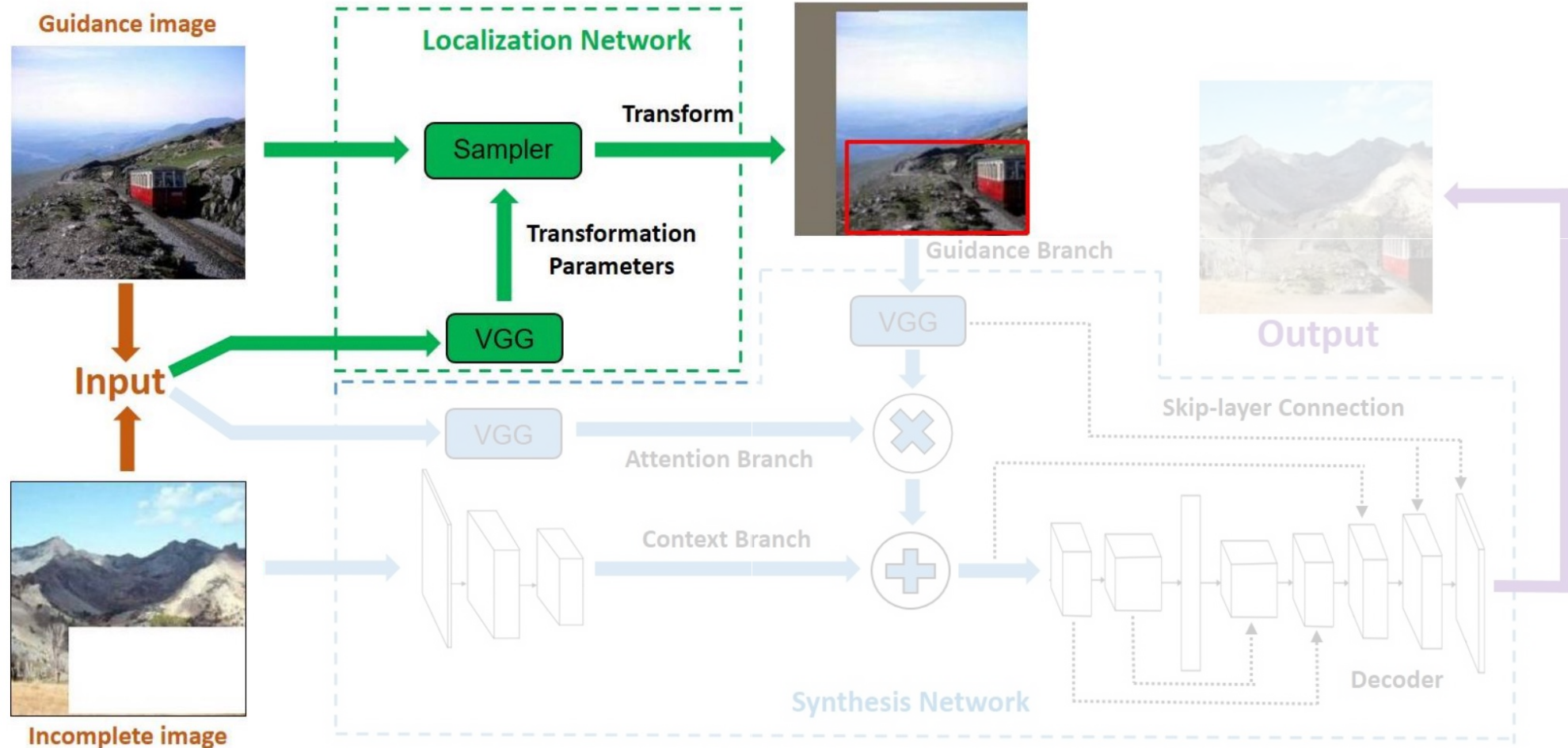


Result:



Architecture

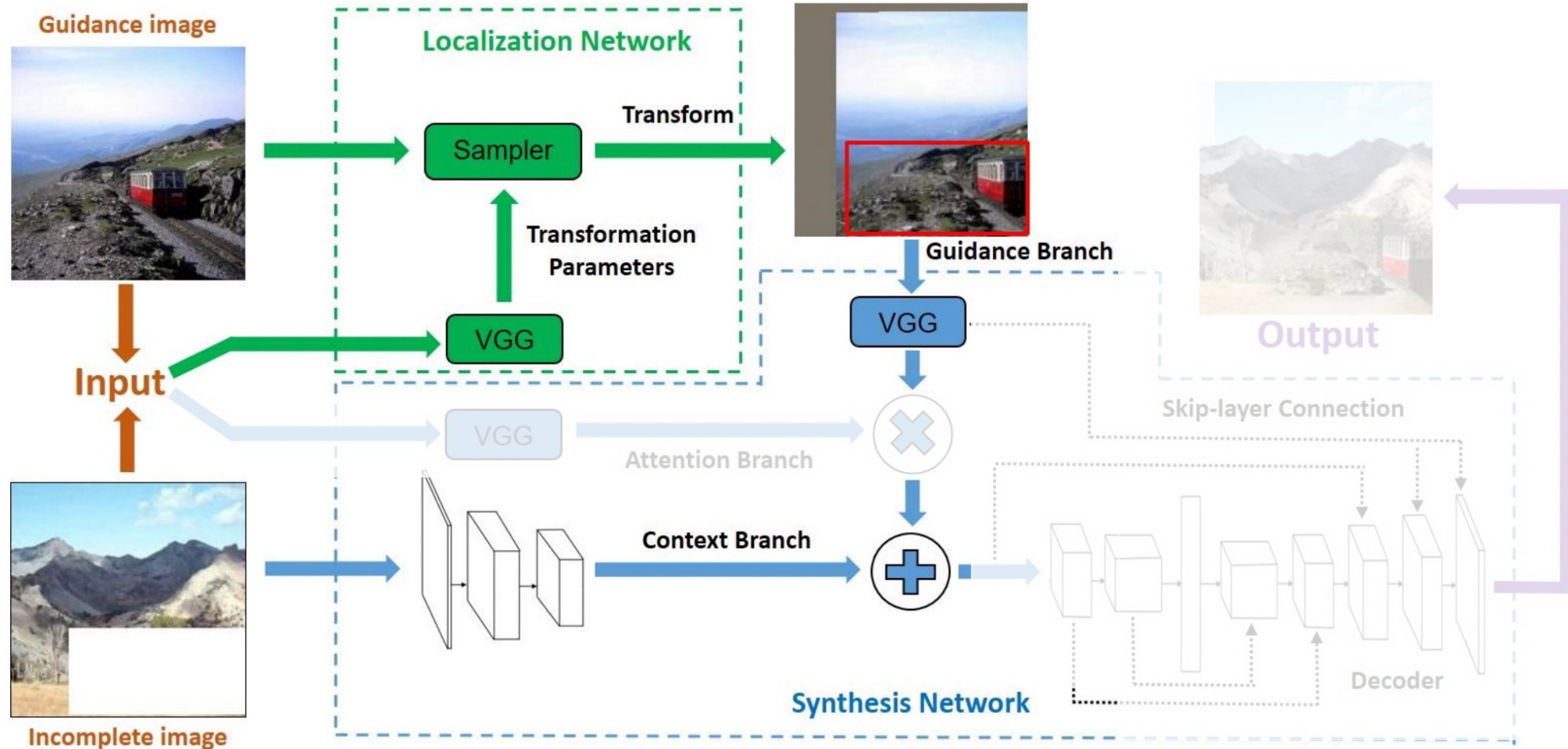
Transforms the guidance image with (predicted) transformation parameters to locate which patch to align with the hole of the incomplete image



(concatenated with a binary mask indicating the hole)

Architecture

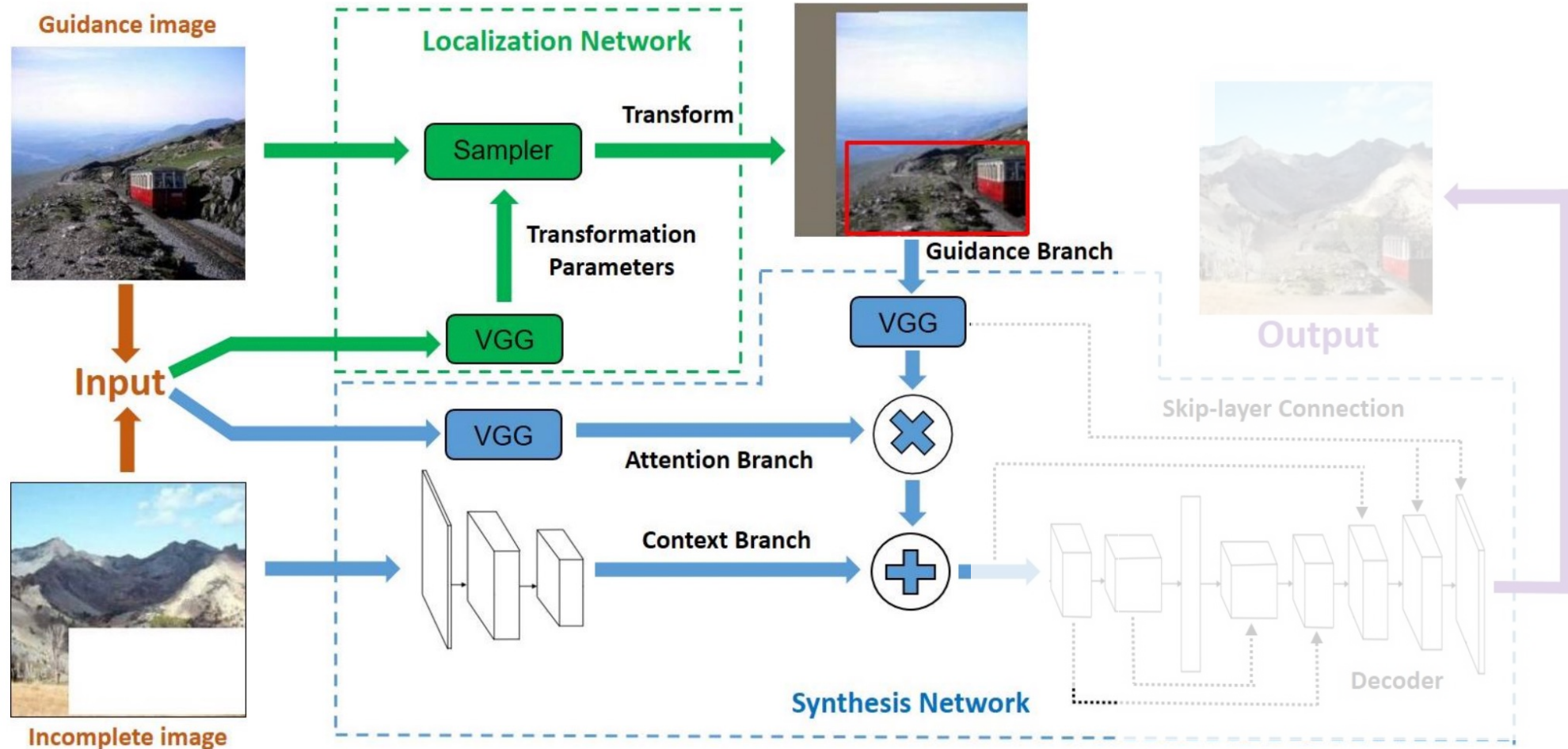
High-level features describing the aligned guidance image fused with features describing the incomplete image



(concatenated with a binary mask indicating the hole)

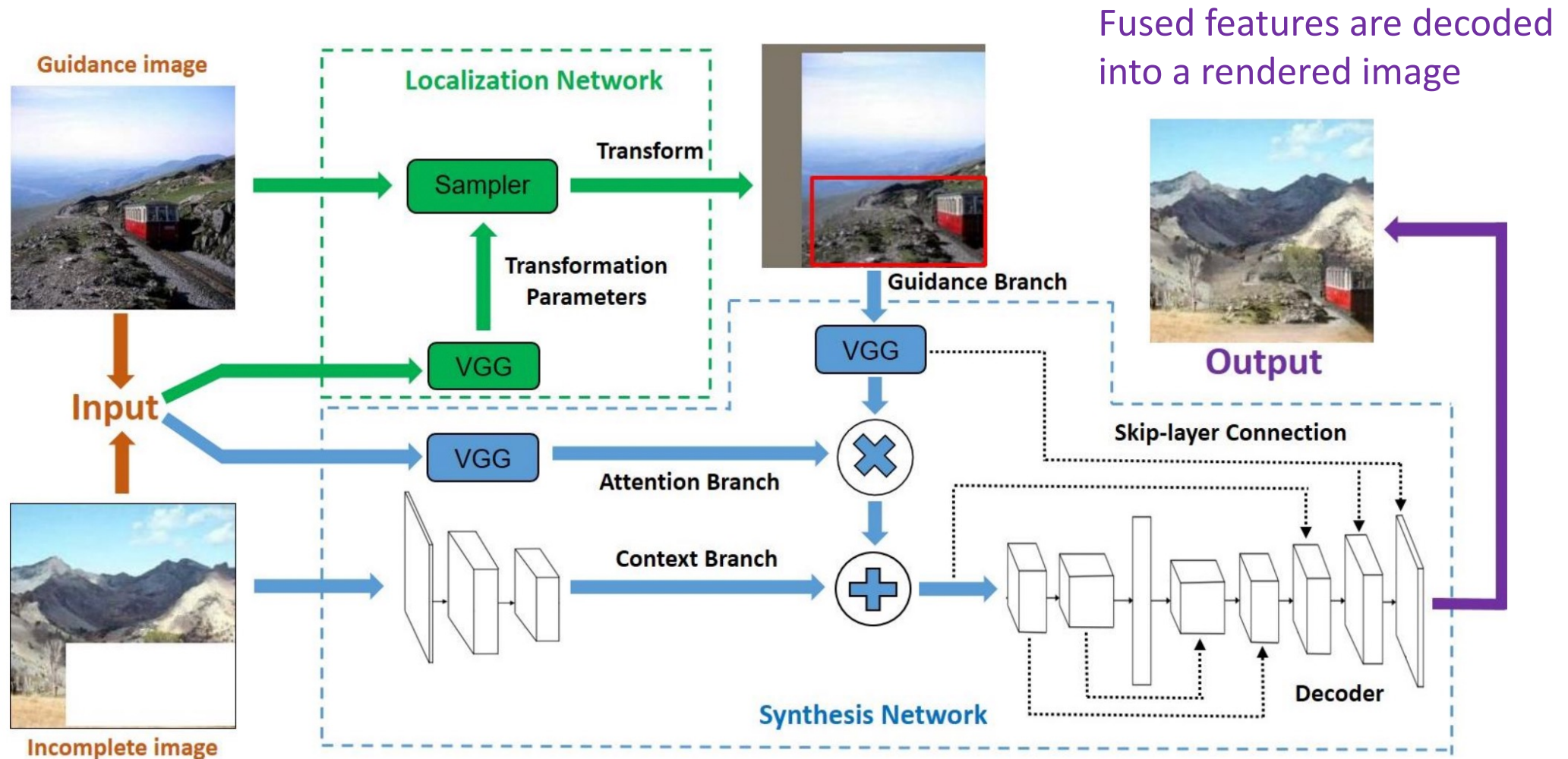
Architecture

Attention branch indicates to what extent regions in the guidance image patch are inconsistent with the context of the incomplete image to indicate whether the decoder should synthesize new content



(concatenated with a binary mask indicating the hole)

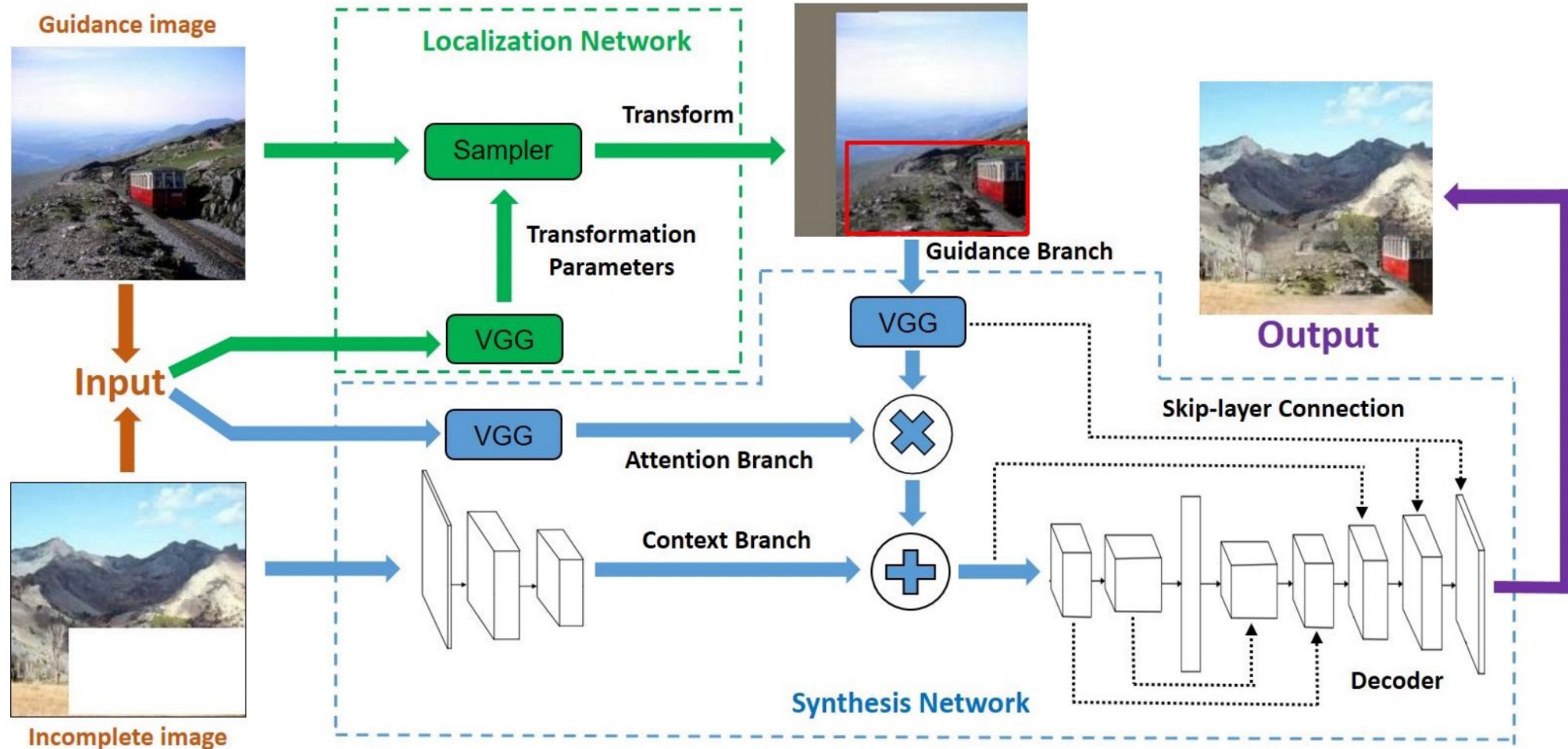
Architecture



(concatenated with a binary mask indicating the hole)

Fused features are decoded into a rendered image

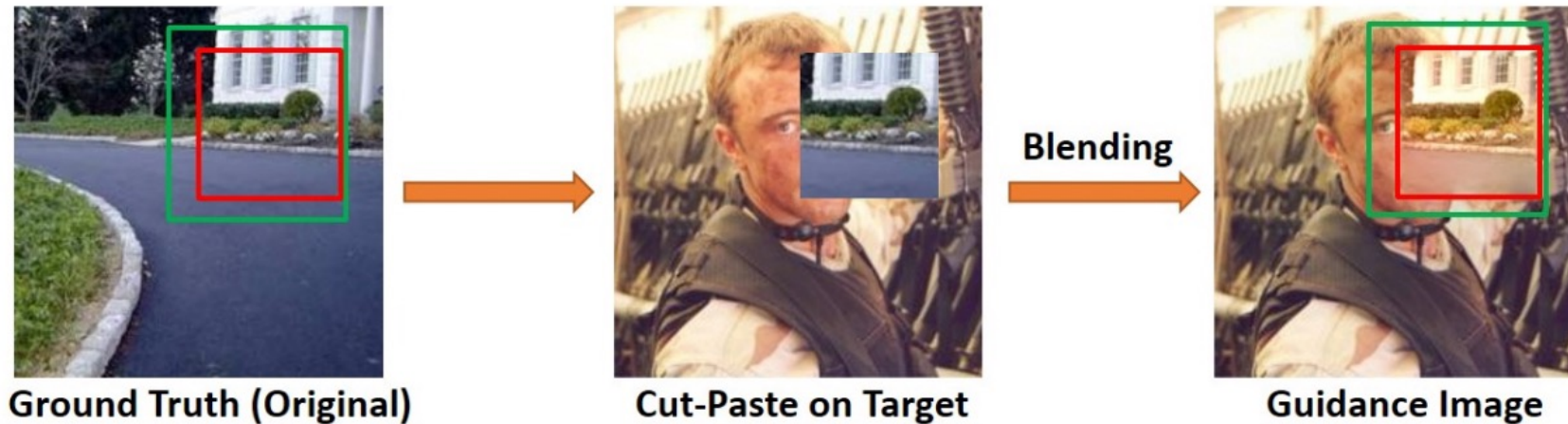
Key Challenge: How to Collect Training Data



(concatenated with a binary mask indicating the hole)

Key Challenge: How to Collect Training Data

Synthesize training data such that we know the true inpainting (i.e., self-supervised learning)

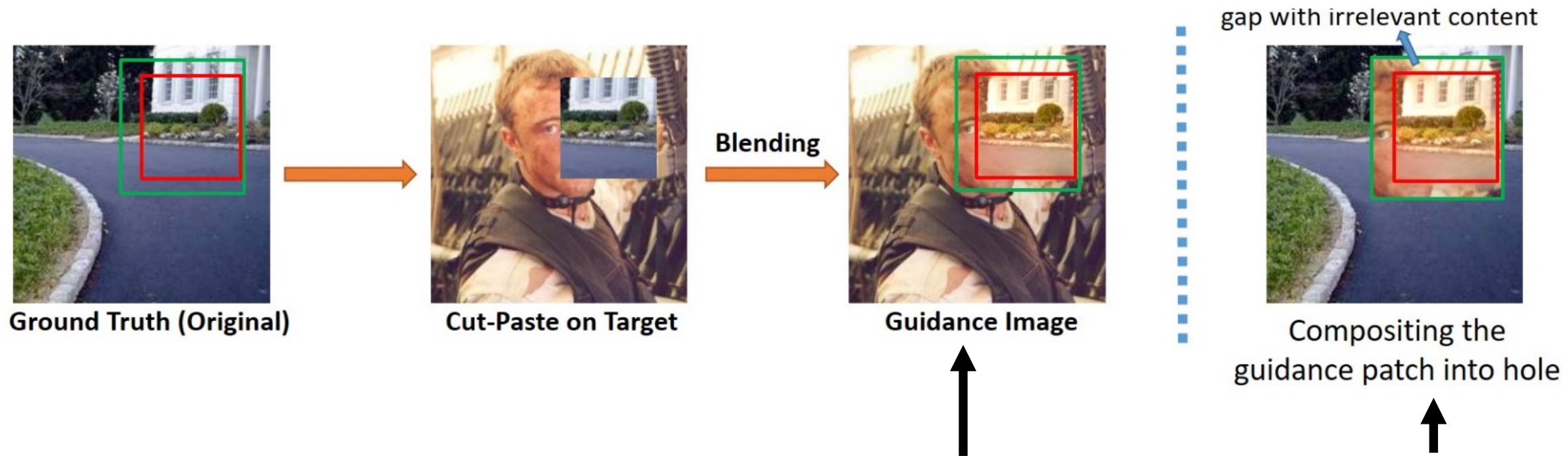


Hole to remove from the original image
Patch of the original image to corrupt

Guidance image contains **original content** that is corrupted and content irrelevant to the original image (gap between **red** and **green** regions)

Key Challenge: How to Collect Training Data

Synthesize training data such that we know the true inpainting (i.e., self-supervised learning)



The localization network must find the **patch** to use from the guidance image while the synthesis network must then recover the original content, including by synthesizing new content in the gap containing irrelevant content

What is a Key Limitation?



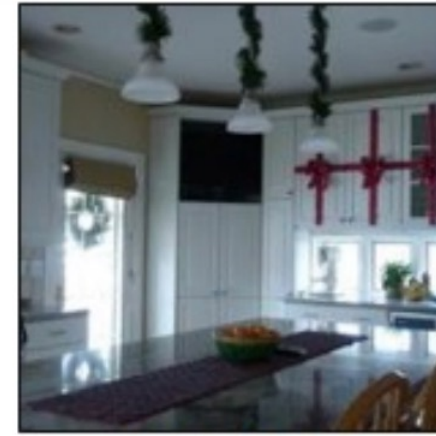
Original Image



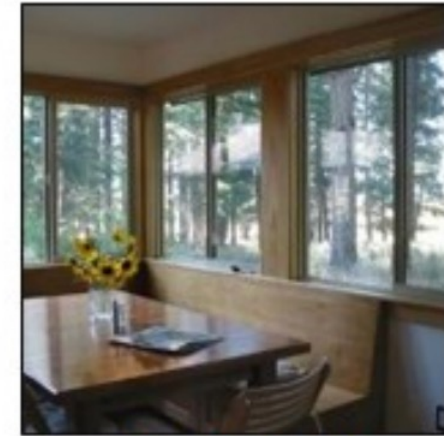
Guidance #1



Guidance #2



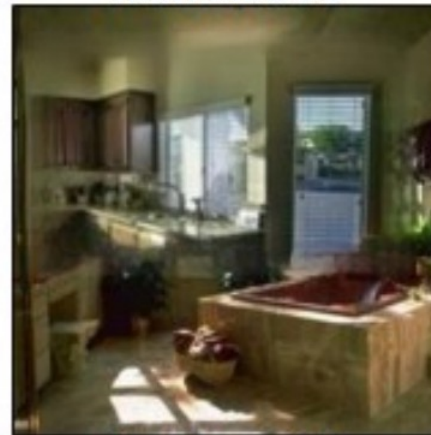
Guidance #3



Guidance #4



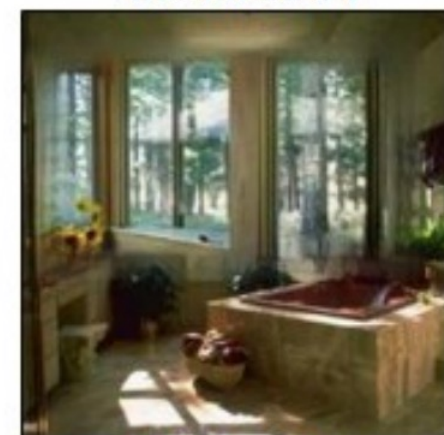
Synthesis #1



Synthesis #2



Synthesis #3



Synthesis #4

Does not account for lighting!

Methods

- Before deep learning era: cut-paste from nearest neighbors
- Context encoder
- Guided image inpainting

Today's Topics

- Problem
- Applications
- Image generation methods
- Hole filling methods
- Evaluation approaches

Experiments

Method	HR[10]	PB[23]	DH[24]	CAF[13]	CE[9]	IM[25]	GLCIC[15]
Retrieval (a)	76%	76%	71%	70%	70%	67%	66%
Retrieval (b)	71%	73%	72%	70%	73%	67%	70%

Crowd workers rate which generated images look more realistic between the method and baselines

(chance score is 50%)

Experiments

Method	NI	CE[9]	HR[10]	CAF[13]	PB[23]	DH[24]	GLCIC[15]	IM[25]	Ours
Retrieval (a)	97.7%	10.0%	14.0%	31.0%	18.0%	23.0%	14.0%	23.0%	33.0%
Retrieval (b)	97.7%	22.0%	15.0%	16.0%	20.0%	22.0%	12.0%	27.0%	36.0%

Crowd workers indicate if images look realistic independently for the method and baselines when also shown real images

Today's Topics

- Problem
- Applications
- Image generation methods
- Hole filling methods
- Evaluation approaches

The image features a dark gray background with a central, soft, circular white glow. This glow is framed by a white film strip border with rectangular sprocket holes. The text "The End" is centered within the glow in a white, elegant, cursive script font.

The End