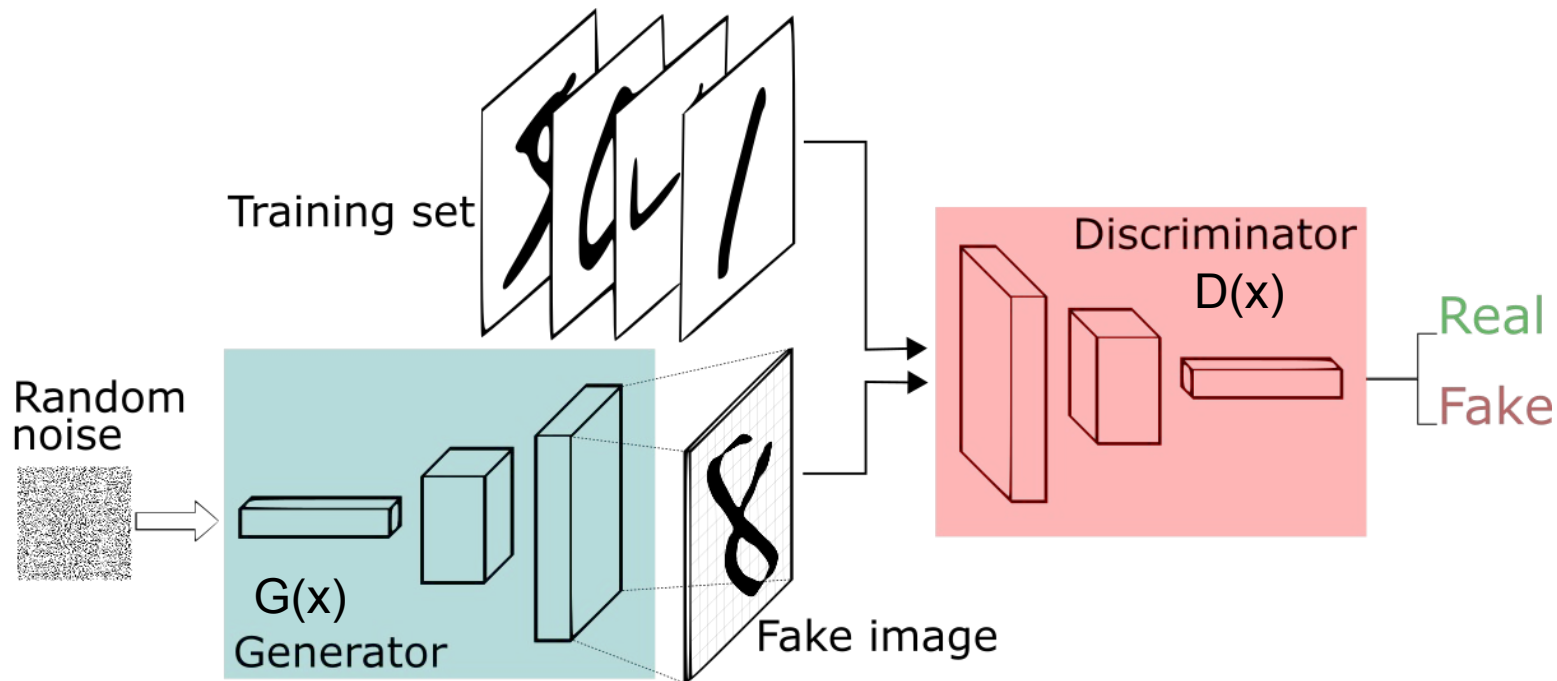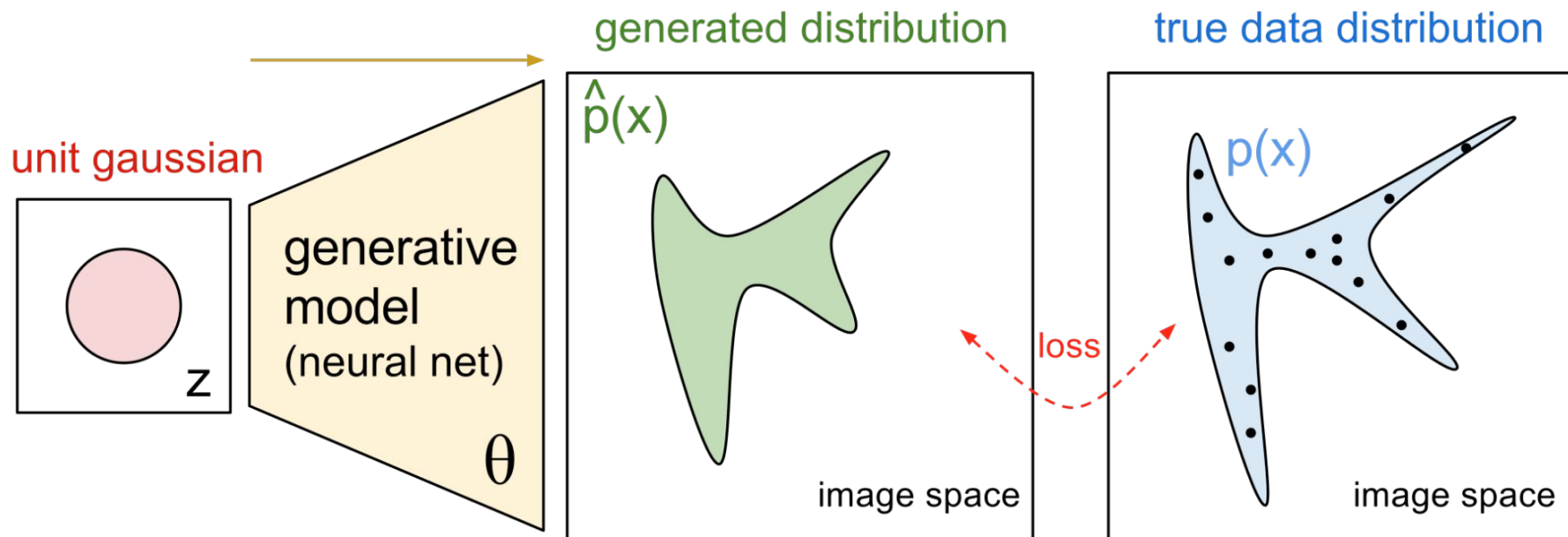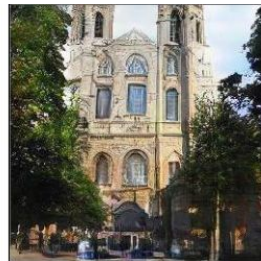# GAN Dissection

Polina Guseva, BAMI181
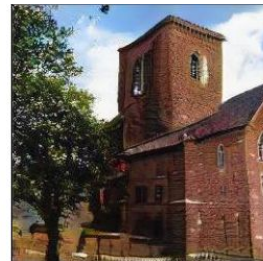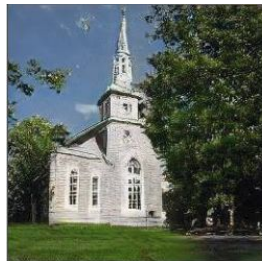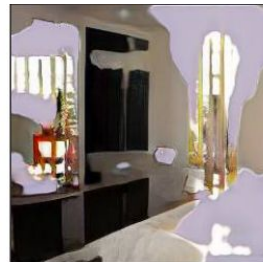
# GANs. Recap

# GANs. Recap

# Motivation

- What knowledge does GAN need to learn?

- What causes the mistakes?

- **Does GAN contain internal variables that correspond to the objects that humans perceive?**

# Goal

**Explain how an image can be generated by a network**

# Definitions



$\mathbf{z} \in \mathbb{R}^{|z|}$      Latent vector

$\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$      Generated image

$G : \mathbf{z} \to \mathbf{x}$      Generator

$\mathbf{r} = h(\mathbf{z})$      Representation

$$\mathbf{x} = f(\mathbf{r}) = G(\mathbf{z})$$

# Definitions

$U$     A set of units (channels)

$P$     A set of pixels in featuremap
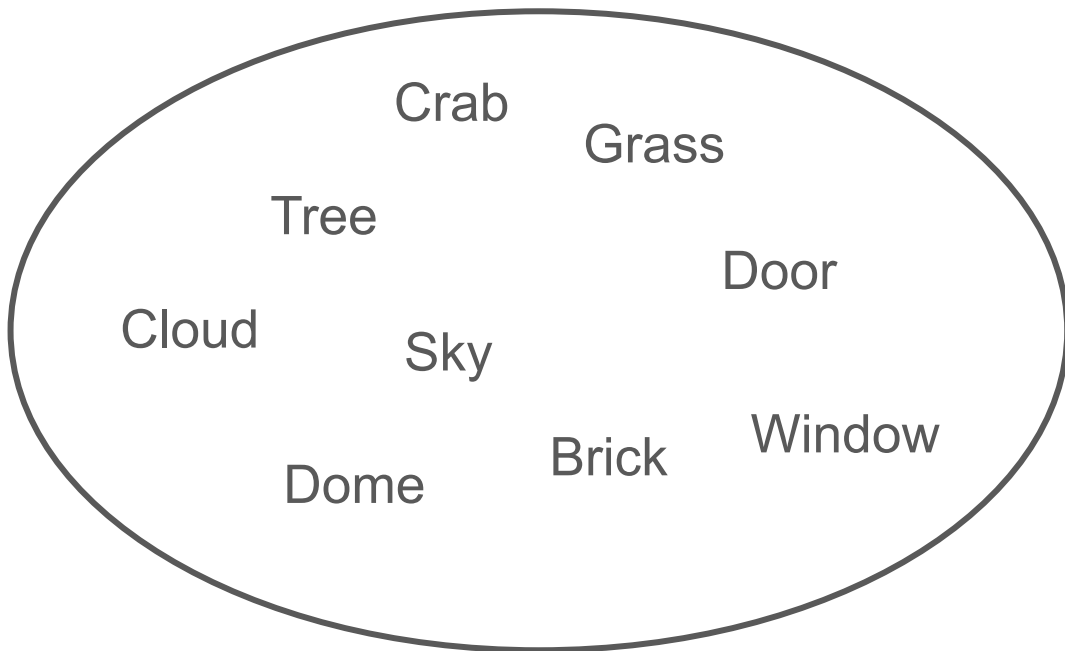
$\mathbb{U}$     All units

$\mathbb{P}$     All pixels



Convnet Filter

One Feature Map

All Feature Maps

# Definitions

$\mathcal{C}$      Universe of concepts

$c \in \mathcal{C}$      Concept

# Task

Factor representation **r** at location P into components

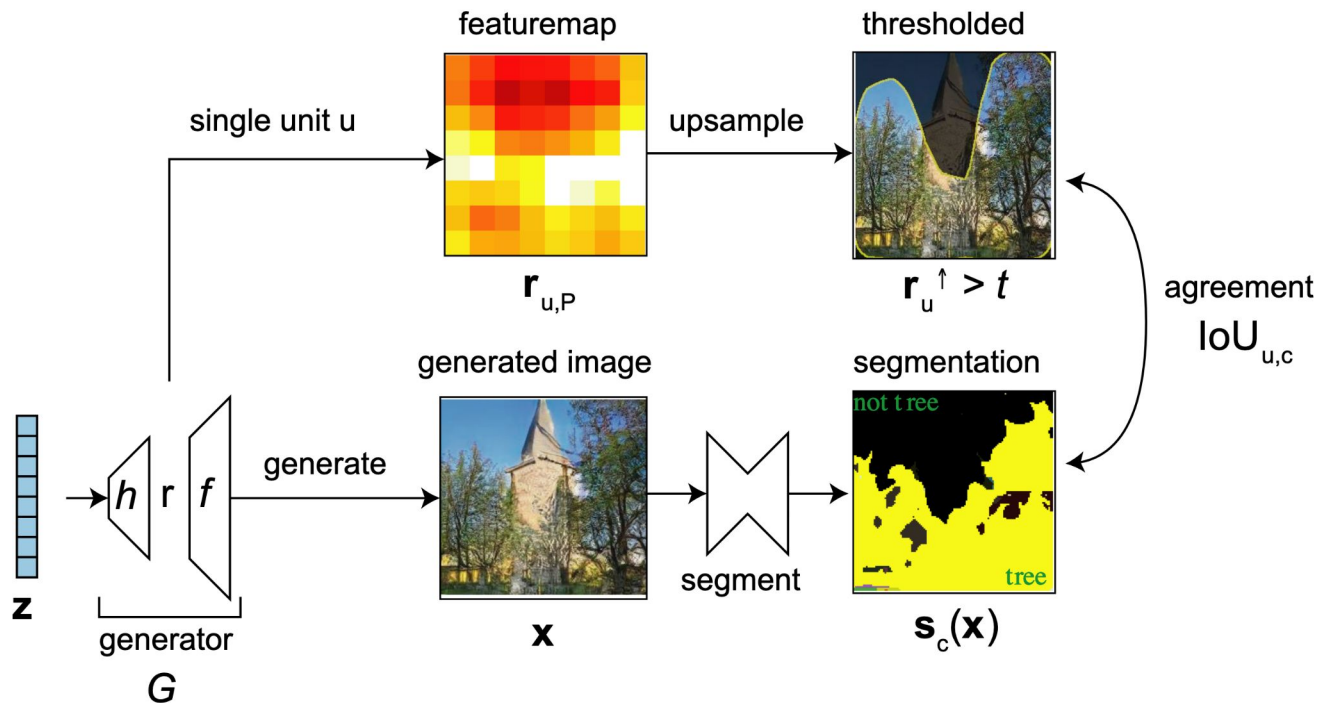$$\mathbf{r}_{\mathbb{U},P} = \left( \mathbf{r}_{U,P}, \mathbf{r}_{\overline{U},P} \right)$$

such that generation of object *c* is dependent on the units in first components and is insensitive to units in second component.
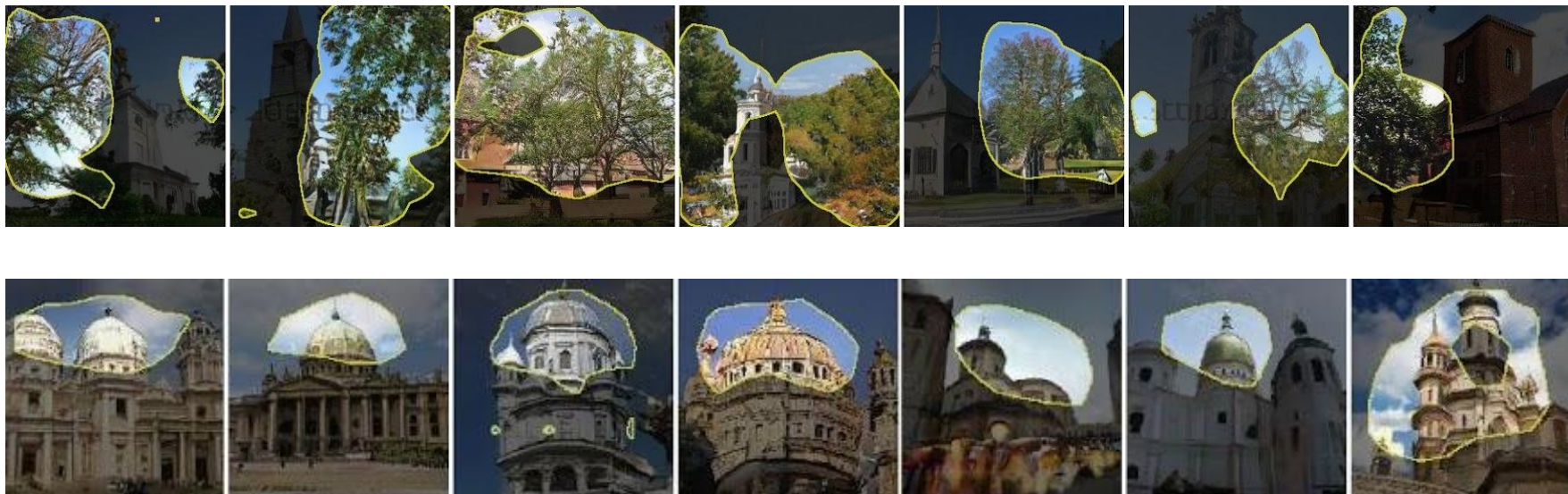
# Proposed method

Two phases:

- Dissection. Select classes with explicit representations.
- Intervention. Identify causal sets of units.

# Method. Dissection

# Method. Dissection

# Method. Dissection

$$\text{IoU}_{u,c} \equiv \frac{\mathbb{E}_{\mathbf{z}} \left| (\mathbf{r}_{u,\mathbb{P}}^{\uparrow} > t_{u,c}) \wedge \mathbf{s}_c(\mathbf{x}) \right|}{\mathbb{E}_{\mathbf{z}} \left| (\mathbf{r}_{u,\mathbb{P}}^{\uparrow} > t_{u,c}) \vee \mathbf{s}_c(\mathbf{x}) \right|}$$
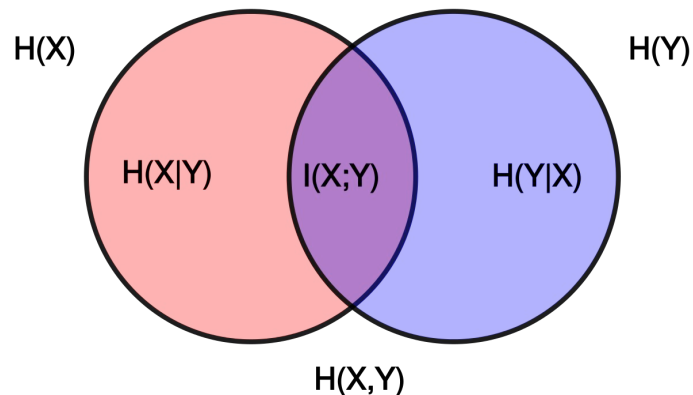
# Method. Dissection

Q. How to select threshold?

A. Maximize **information quality ratio**

$$t_{u,c} = \arg\max_t \frac{\mathbf{I}(\mathbf{r}_{u,\mathbb{P}}^{\uparrow} > t; \mathbf{s}_c(\mathbf{x}))}{\mathbf{H}(\mathbf{r}_{u,\mathbb{P}}^{\uparrow} > t, \mathbf{s}_c(\mathbf{x}))}$$

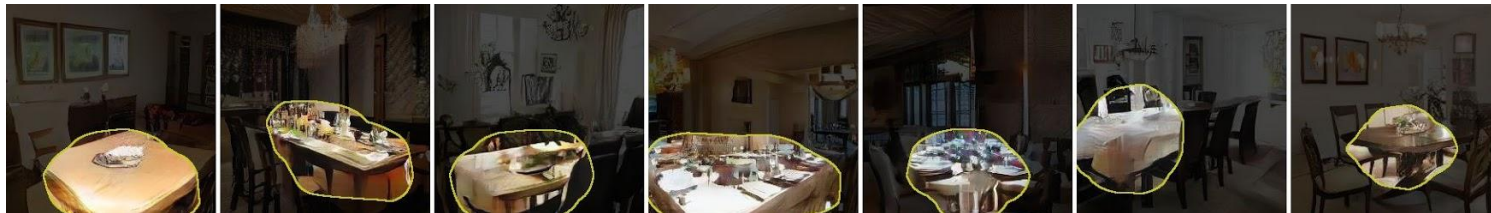$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$



H(X)  H(Y)

H(X|Y)  I(X;Y)  H(Y|X)
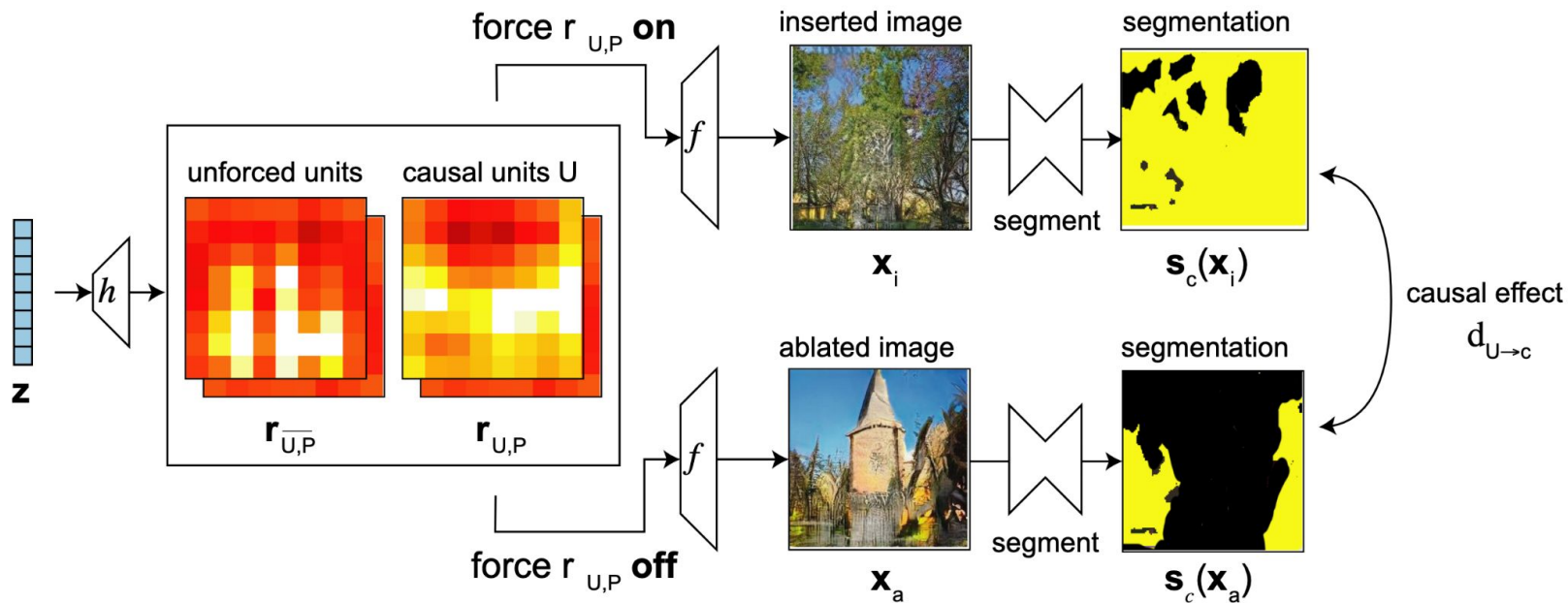
H(X,Y)

# Method. Dissection

Unit №65

IoU=0.34



Unit №37

IoU=0.29

# Method. Intervention

# Method. Intervention

Original image :
$$\mathbf{x} = G(\mathbf{z}) \equiv f(\mathbf{r}) \equiv f(\mathbf{r}_{\mathrm{U,P}}, \mathbf{r}_{\overline{\mathrm{U,P}}})$$

Image with U ablated at pixels P :
$$\mathbf{x}_a = f(\mathbf{0}, \mathbf{r}_{\overline{\mathrm{U,P}}})$$

Image with U inserted at pixels P :
$$\mathbf{x}_i = f(\mathbf{k}, \mathbf{r}_{\overline{\mathrm{U,P}}})$$

# Method. Intervention

An object is caused by $U$ if the object appears in $x_i$ and disappears from $x_a$.

The measure is average causal effect (ACE)

$$\delta_{\mathbf{U} \to c} \equiv \mathbb{E}_{\mathbf{z}, \mathbf{P}}\big[\mathbf{s}_c(\mathbf{x}_i)\big] - \mathbb{E}_{\mathbf{z}, \mathbf{P}}\big[\mathbf{s}_c(\mathbf{x}_a)\big]$$

# Method. Intervention

$r$ contains $d$ units

$\alpha \in [0,1]^d$ — continious intervention

$\alpha_u$ — degree of intervention for $u$

$$\mathbf{x}_a = f(\mathbf{0}, \mathbf{r}_{\overline{U,P}})$$

$$\mathbf{x}_i = f(\mathbf{k}, \mathbf{r}_{\overline{U,P}})$$

$$\delta_{U \to c} \equiv \mathbb{E}_{\mathbf{z},P}[\mathbf{s}_c(\mathbf{x}_i)] - \mathbb{E}_{\mathbf{z},P}[\mathbf{s}_c(\mathbf{x}_a)]$$
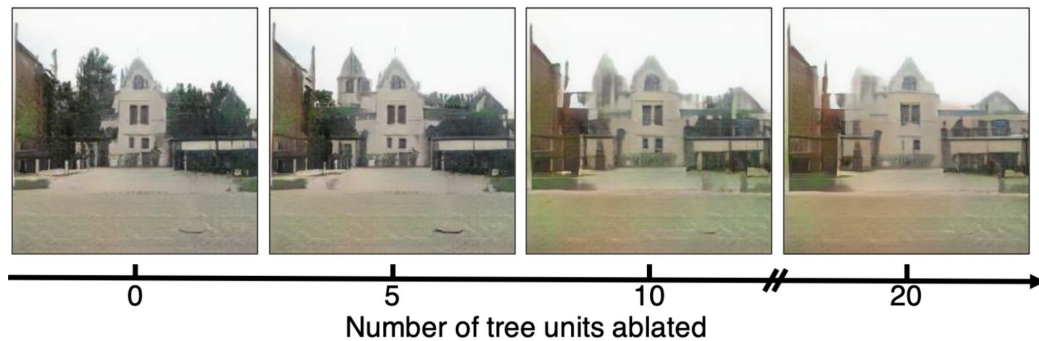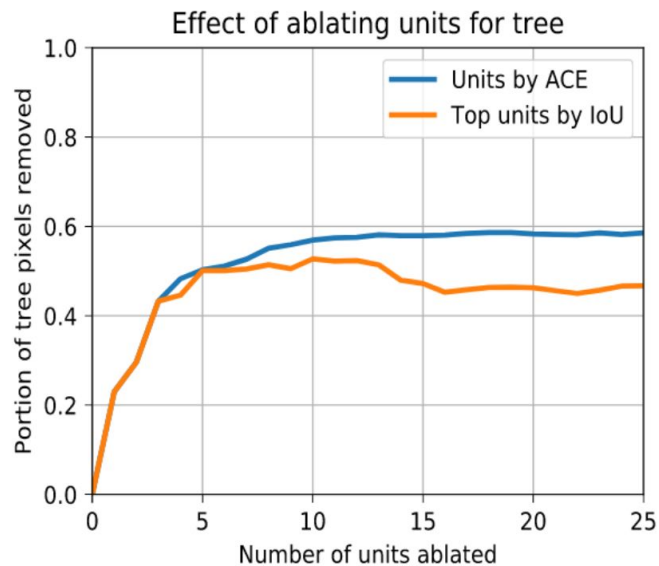
$\longrightarrow$

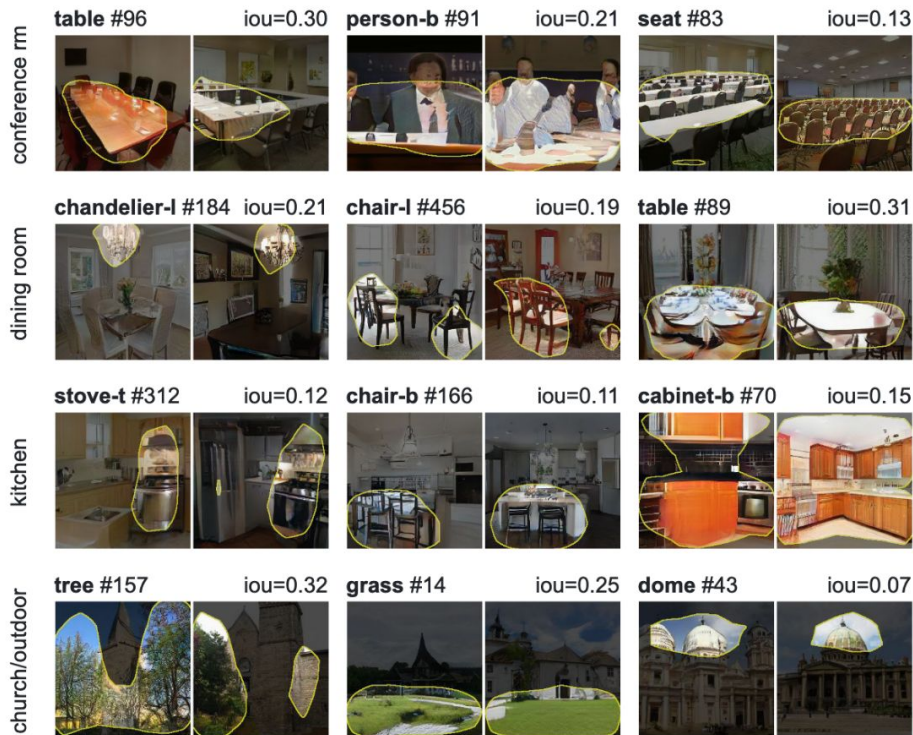$$\mathbf{x}'_a = f((\mathbf{1} - \boldsymbol{\alpha}) \odot \mathbf{r}_{U,P},\ \mathbf{r}_{U,\bar{P}})$$

$$\mathbf{x}'_i = f(\boldsymbol{\alpha} \odot \mathbf{k} + (\mathbf{1} - \boldsymbol{\alpha}) \odot \mathbf{r}_{U,P},\ \mathbf{r}_{U,\bar{P}})$$

$$\delta_{\boldsymbol{\alpha} \to c} = \mathbb{E}_{\mathbf{z},P}\left[\mathbf{s}_c(\mathbf{x}'_i)\right] - \mathbb{E}_{\mathbf{z},P}\left[\mathbf{s}_c(\mathbf{x}'_a)\right],$$

$$\boldsymbol{\alpha}^* = \arg\min_{\boldsymbol{\alpha}}(-\delta_{\boldsymbol{\alpha} \to c} + \lambda||\boldsymbol{\alpha}||_2)$$

# Method. Intervention



Effect of ablating units for tree

Portion of tree pixels removed vs. Number of units ablated

Units by ACE

Top units by IoU
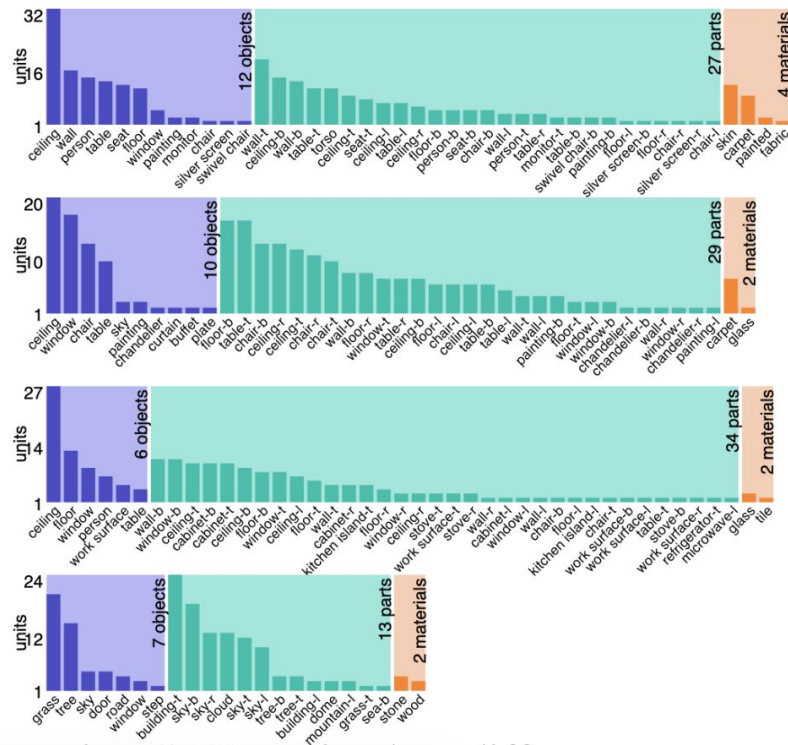
Number of tree units ablated

# Results. Interpretable units for different scene types



Units in scene generator

Unit class distribution
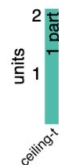
# Results. Units for different network layers

Units in layer

Unit class distribution

**layer1**
512 units total

0 object units
2 part units
0 material units

**ceiling-t** layer1 #457    iou=0.10

**ceiling-t** layer1 #194    iou=0.07

**layer4**
512 units total

86 object units
149 part units
10 material units

**sofa** layer4 #37    iou=0.28

**fireplace** layer4 #23    iou=0.15

**layer7**
256 units total

59 object units
48 part units
9 material units

**painting** layer7 #15    iou=0.23

**coffee table-t** #247    iou=0.07

**layer10**
128 units total

19 object units
8 part units
11 material units

**carpet** layer10 #53    iou=0.14

**glass** layer10 #126    iou=0.21

# Results. Units for various networks

# Results. Debugging GANs
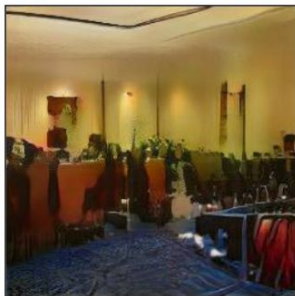


Example artifact-causing units

Bedroom images with artifacts

Ablating "artifact" units improves results

# Results. Erasing objects
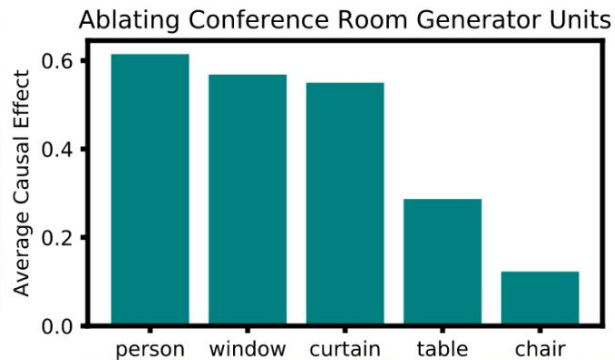


ablate person units

ablate curtain units

Ablating Conference Room Generator Units

ablate window units

ablate table units

ablate chair units

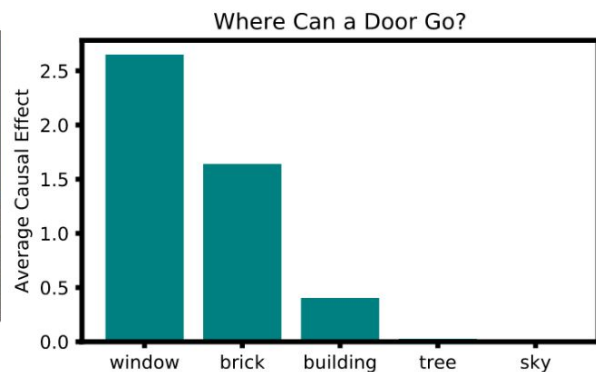# Results. Inserting objects



(a)

(b)

Where Can a Door Go?

(c)

(d)

(e)

# Summary

- GANs have **sets of neurons** that **explicitly control object generation**
- Suggested that GANs learned some **aspects of composition**
- Some **artifacts** may be triggered by **specific sets of neurons** — easy fix

# References

1. David Bau et al. GAN Dissection: Visualizing and Understanding Generative Adversarial Networks. 2019
2. David Bau et al. Understanding the role of individual units in a deep neural network. 2020
3. OpenAI Blog. Generative models. 2016
4. David Bau, Bolei Zhou et al. Network dissection: Quantifying interpretability of deep visual representations. 2017