

# GAN DISSECTION: VISUALIZING AND UNDERSTANDING GENERATIVE ADVERSARIAL NETWORKS

**David Bau<sup>1,2</sup>, Jun-Yan Zhu<sup>1</sup>, Hendrik Strobelt<sup>2,5</sup>, Bolei Zhou<sup>3</sup>,  
Joshua B. Tenenbaum<sup>1</sup>, William T. Freeman<sup>1,4</sup>, Antonio Torralba<sup>1,2</sup>**

<sup>1</sup>Massachusetts Institute of Technology, <sup>2</sup>MIT-IBM Watson AI Lab,

<sup>3</sup>The Chinese University of Hong Kong, <sup>4</sup>Google Research, <sup>5</sup>IBM Research

ICLR 2019

Presenter: Jack Lanchantin

# Introduction

- $G: z \rightarrow x$ , where  $z \in \mathbb{R}^{|z|}$  and  $x \in \mathbb{R}^{H \times W \times 3}$
- Tensor  $r$  output from a particular layer of  $G$ :  $r = h(z)$  and  $x = f(r) = f(h(z)) = G(z)$
- $r$  certainly contains the information to deduce the presence of any visible class  $c$  in the image
- Question is *how* the information about  $c$  is encoded in  $r$

# Introduction

- In particular, we seek to understand whether  $r$  explicitly represents the concept  $c$  in some way where it is possible to factor  $r$  at locations  $P$  into components

$$\mathbf{r}_{\mathbb{U},P} = (\mathbf{r}_{U,P}, \mathbf{r}_{\bar{U},P})$$

where the generation of the object  $c$  at locations  $P$  depends mainly on the units  $r_{U,P}$  and is insensitive to the other units  $r_{\bar{U},P}$

- Refer to each channel of the featuremap as a unit;  $U$  denotes the set of unit indices of interest and denotes its complement
- we will write  $\mathbb{U}$  and  $\mathbb{P}$  to refer to the entire set of units and feature map pixels in  $r$

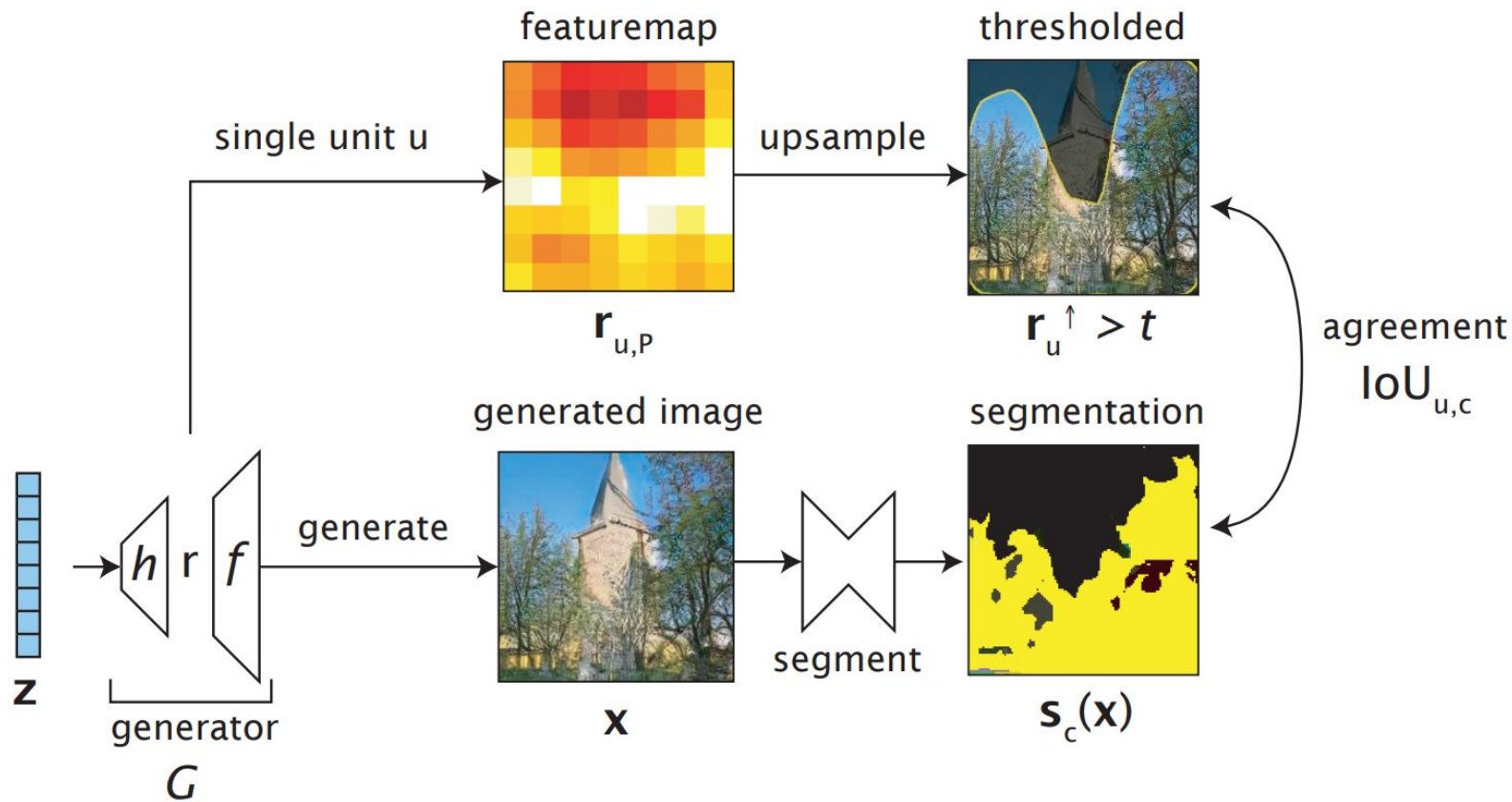
# Characterizing Units by Dissection

- Quantify the spatial agreement between the unit  $U$ 's thresholded featuremap and a concept  $c$ ' segmentation with the following intersection-over-union (IoU) measure:

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

where  $\wedge$  and  $\vee$  denote intersection and union,  $t_{u,c}$  is a fixed threshold, and  $\mathbf{s}_c(\mathbf{x})$  is a binary segmentation mask for concept  $c$

# Characterizing Units by Dissection



# Characterizing Units by Dissection



Thresholding unit #65 layer 3 of a dining room generator matches ‘table’ segmentations with  $\text{IoU}=0.34$ .



Thresholding unit #37 layer 4 of a living room generator matches ‘sofa’ segmentations with  $\text{IoU}=0.29$ .

Figure 3: Visualizing the activations of individual units in two GANs. 10 top activating images are shown, and  $\text{IoU}$  is measured over a sample of 1000 images. In each image, the unit feature is upsampled and thresholded as described in Eqn. 2.

# Measuring Causal Relationships Using Intervention

- Which of those units are actually responsible for triggering the rendering of that object?
  - Correlation != causation
  - Furthermore, any output will jointly depend on several parts of the representation (need to identify combinations)

# Measuring Causal Relationships Using Intervention

- Recall that  $\mathbf{r}_{U,P}$  denotes the feature map  $\mathbf{r}$  at unit  $U$  and location  $P$
- **Ablate** such unit by forcing  $\mathbf{r}_{U,P} = \mathbf{0}$ .
- **Insert** such unit by forcing  $\mathbf{r}_{U,P} = \mathbf{c}$ , where  $\mathbf{c}$  is a big constant.
- Decompose  $\mathbf{r}$  into two parts  $(\mathbf{r}_{U,P}, \mathbf{r}_{\overline{U,P}})$ , where  $\mathbf{r}_{\overline{U,P}}$  are unforced components of  $\mathbf{r}$

Original image :

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

Image with  $U$  ablated at pixels  $P$  :

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

Image with  $U$  inserted at pixels  $P$  :

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



# Measuring Causal Relationships Using Intervention

- Recall that  $\mathbf{r}_{U,P}$  denotes the feature map  $r$  at unit  $U$  and location  $P$
- **Ablate** such unit by forcing  $\mathbf{r}_{U,P} = 0$ .
- **Insert** such unit by forcing  $\mathbf{r}_{U,P} = \mathbf{c}$ , where  $\mathbf{c}$  is a big constant.
- Decompose  $r$  into two parts  $(\mathbf{r}_{U,P}, \mathbf{r}_{\overline{U,P}})$ , where  $\mathbf{r}_{\overline{U,P}}$  are unforced components of  $\mathbf{r}$

Original image :

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

Image with  $U$  ablated at pixels  $P$  :

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

Image with  $U$  inserted at pixels  $P$  :

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

➤ An object is caused by  $U$  if the object appears in  $\mathbf{x}_i$  and disappears from  $\mathbf{x}_a$

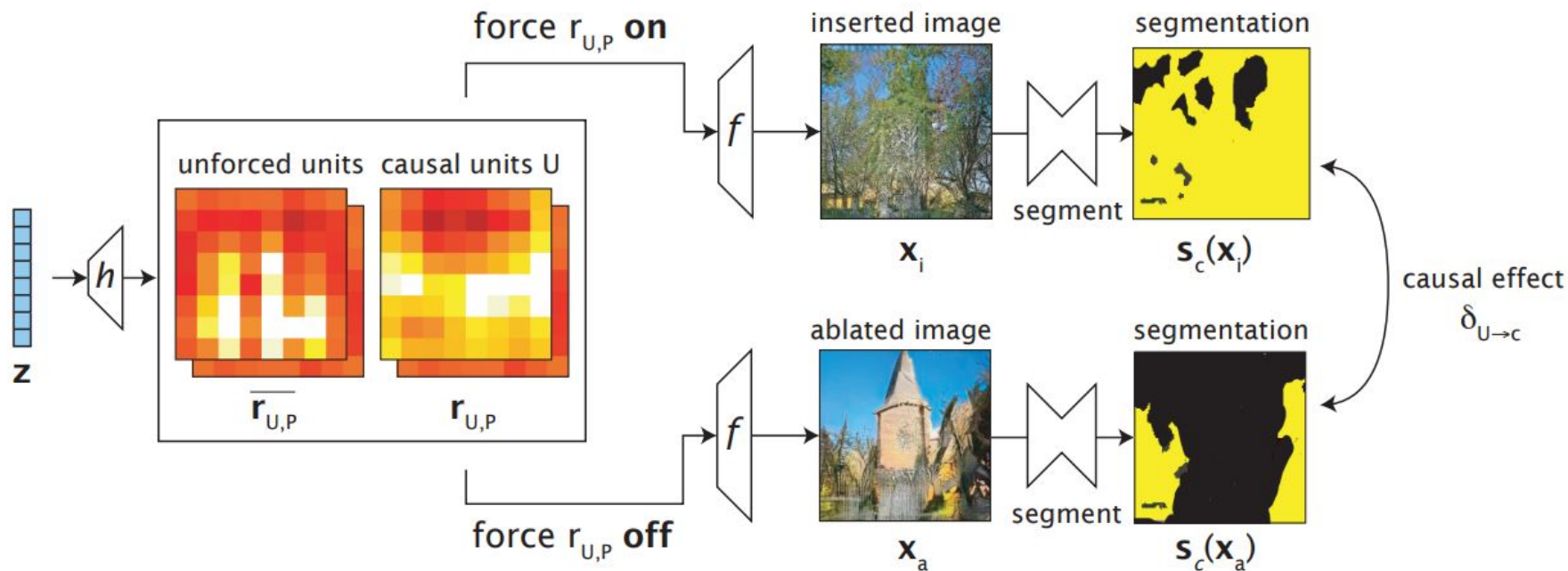
# Measuring Causal Relationships Using Intervention

- This causality can be quantified by comparing the presence of an object in  $\mathbf{x}_i$  and  $\mathbf{x}_a$  and averaging effects over all locations and images
- Define the average causal effect (ACE) of unit  $U$  on the generation of on class  $c$  as:

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

where  $\mathbf{s}_c(\mathbf{x})$  denotes a segmentation indicating the presence of class  $c$  in image  $\mathbf{x}$  at  $P$

# Measuring Causal Relationships Using Intervention



# Sets of Units with High Causal Effect

- Objects tend to depend on more than one unit.
- Thus we need to identify a *set of units*  $U$  that maximize the average causal effect  $\bar{\delta}_{U \rightarrow c}$  for a class  $c$

# Finding sets of units with high average causal effect

- Given a representation  $\mathbf{r}$  with  $d$  units, searching for a fixed-size set  $U$  with high  $\delta_{U \rightarrow c}$  requires  $\binom{d}{|U|}$  operations
- Instead, we optimize a continuous intervention  $\alpha \in [0, 1]^d$ , where each dimension  $\alpha_u$  indicates the degree of intervention for unit  $u$ .

# Finding sets of units with high average causal effect

- We maximize the following average causal effect formulation  $\delta_{\alpha \rightarrow c}$ :

Image with partial ablation at pixels  $P$  :  $\mathbf{x}'_a = f((\mathbf{1} - \alpha) \odot \mathbf{r}_{U,P}, \mathbf{r}_{U,\bar{P}})$

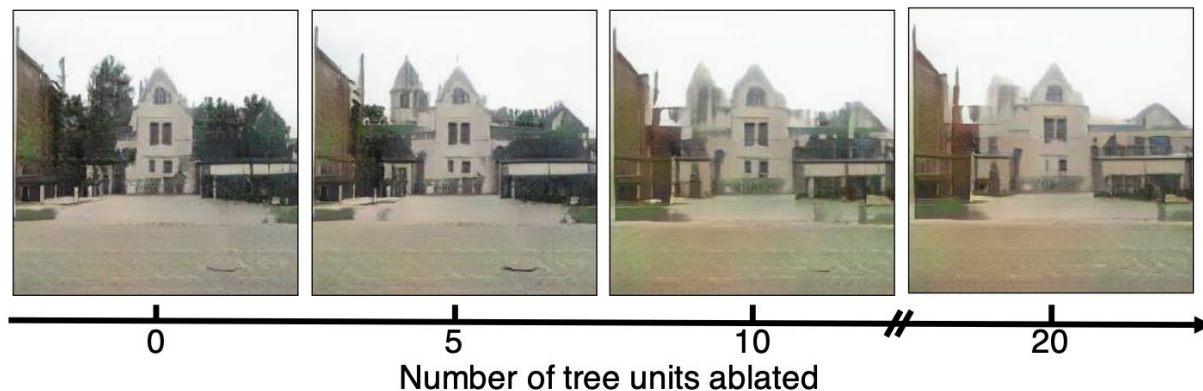
Image with partial insertion at pixels  $P$  :  $\mathbf{x}'_i = f(\alpha \odot \mathbf{c} + (\mathbf{1} - \alpha) \odot \mathbf{r}_{U,P}, \mathbf{r}_{U,\bar{P}})$

Objective :  $\delta_{\alpha \rightarrow c} = \mathbb{E}_{\mathbf{z},P} [\mathbf{s}_c(\mathbf{x}'_i)] - \mathbb{E}_{\mathbf{z},P} [\mathbf{s}_c(\mathbf{x}'_a)]$ ,

where  $\mathbf{r}_{U,P}$  denotes the all-channel featuremap at locations  $P$ ,  $\mathbf{r}_{U,\bar{P}}$  denotes the all-channel featuremap at other locations  $\bar{P}$ , and  $\odot$  applies a per-channel scaling vector  $\alpha$  to the featuremap  $\mathbf{r}_{U,P}$

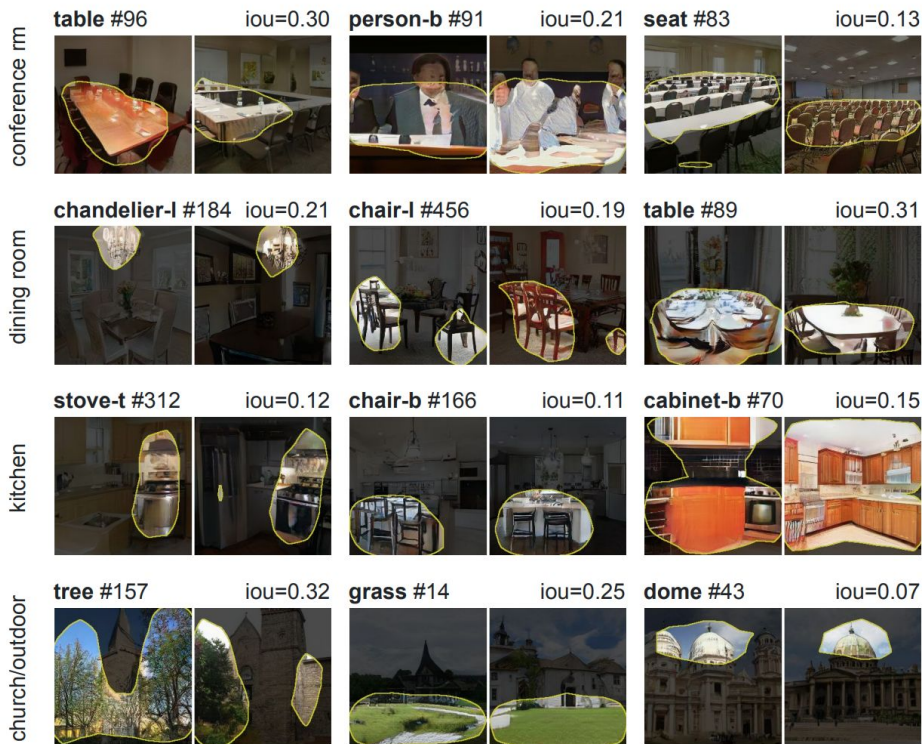
- $$\alpha^* = \arg \min_{\alpha} (-\delta_{\alpha \rightarrow c} + \lambda \|\alpha\|_2)$$

# Finding sets of units with high average causal effect

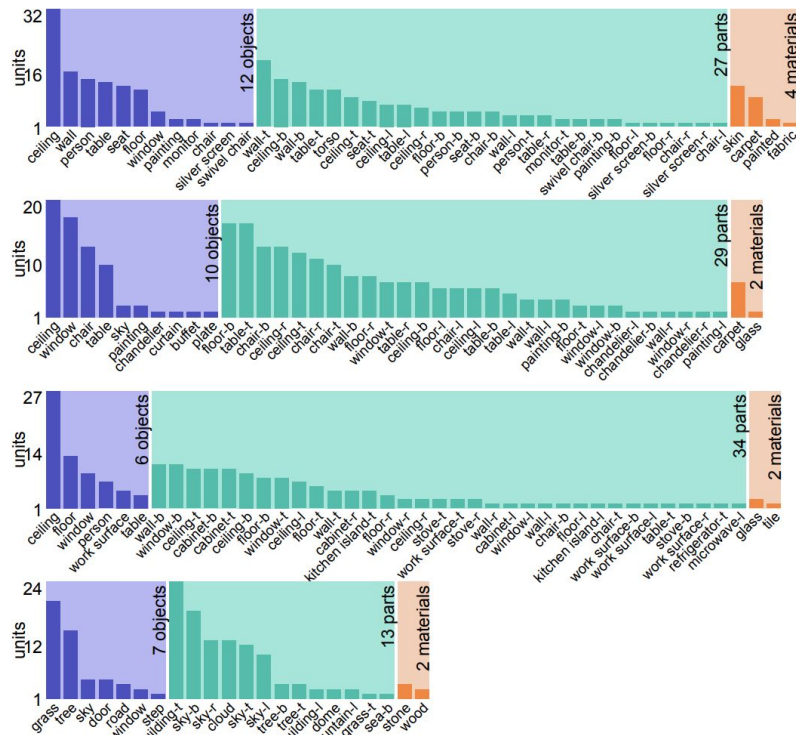


# Units that match objects (from layer 4 of trained CNN)

Units in scene generator



Unit class distribution



A unit is counted as a class predictor if it matches a supervised segmentation class with pixel accuracy > 0.75 and IoU > 0.05 when upsampled and thresholded.

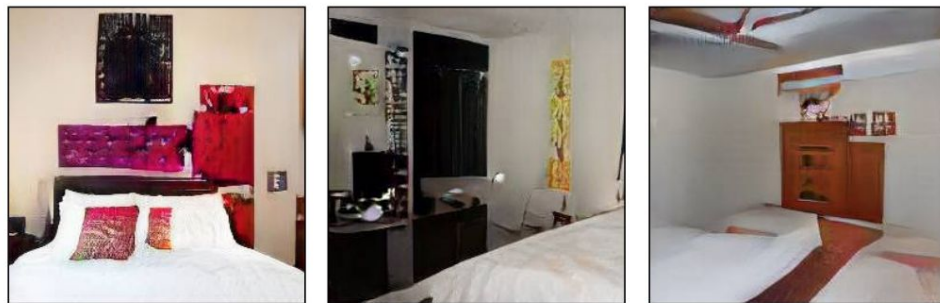




# Ablating Artifacts



(f) Bedroom images with artifacts



(g) Ablating “artifacts” units improves results

Ablate the 20 artifact-causing units out of 512 units in layer4.

# Ablating Artifacts

---

## Fréchet Inception Distance (FID)

---

original images	52.87
“artifacts” units ablated (ours)	<b>32.11</b>
random units ablated	52.27

---

---

## Human preference score

---

## original images

---

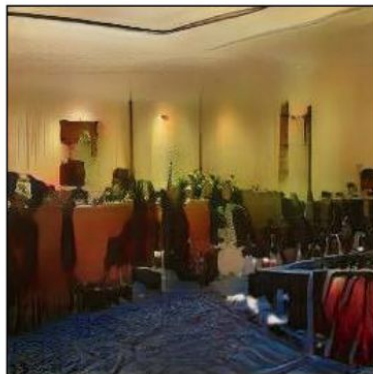
“artifacts” units ablated (ours)	<b>79.0%</b>
random units ablated	50.8%

---

# Ablating Objects



ablate person units



ablate curtain units



ablate window units

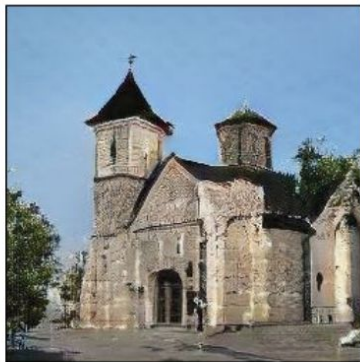
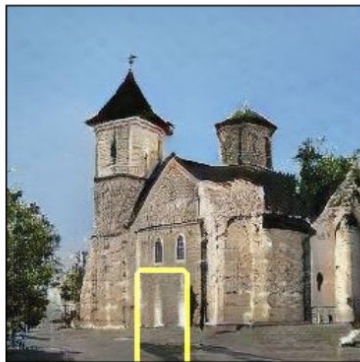


ablate table units





# Inserting Objects



(a)



(b)



(c)



(d)