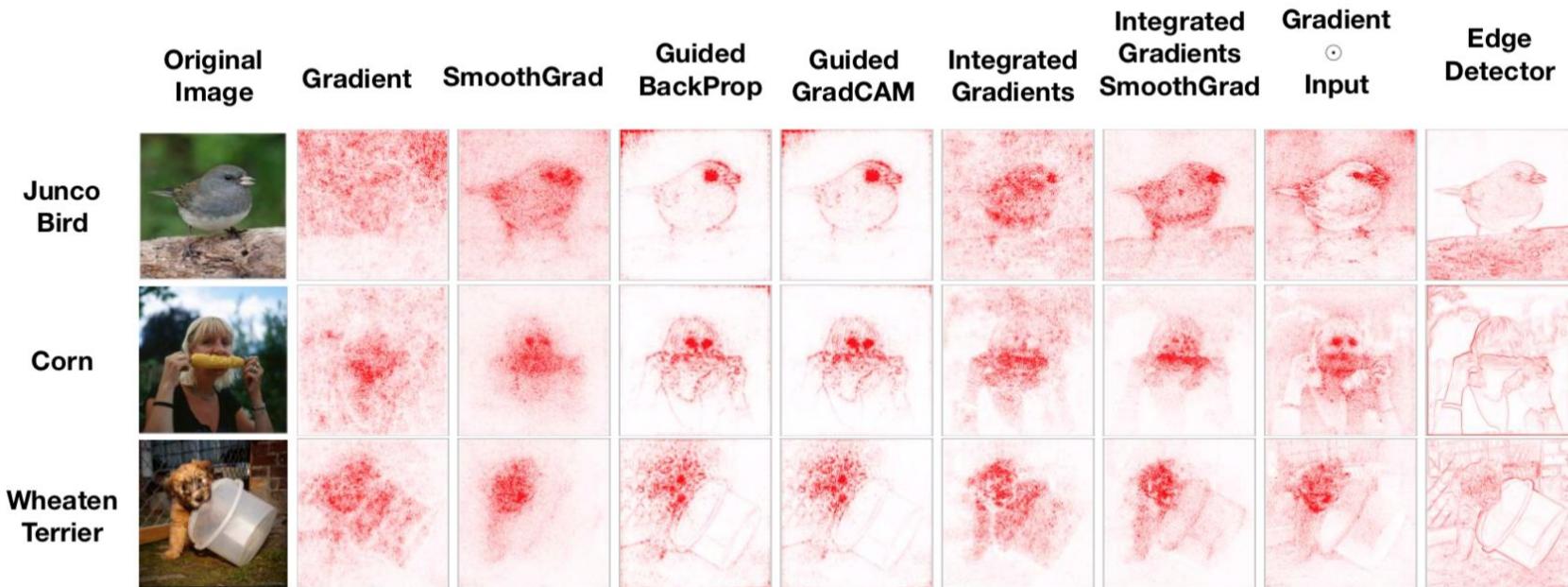


GAN DISSECTION: VISUALIZING AND UNDERSTANDING GENERATIVE ADVERSARIAL NETWORKS

Романов Игнат

Философия



Чуть конкретнее

- Хотим понять что выучивает GAN какие-то пиксельные паттерны или мы можем найти какие-то внутренние переменные, отвечающие за объекты различимые человеком?
- Что порождает артефакты?

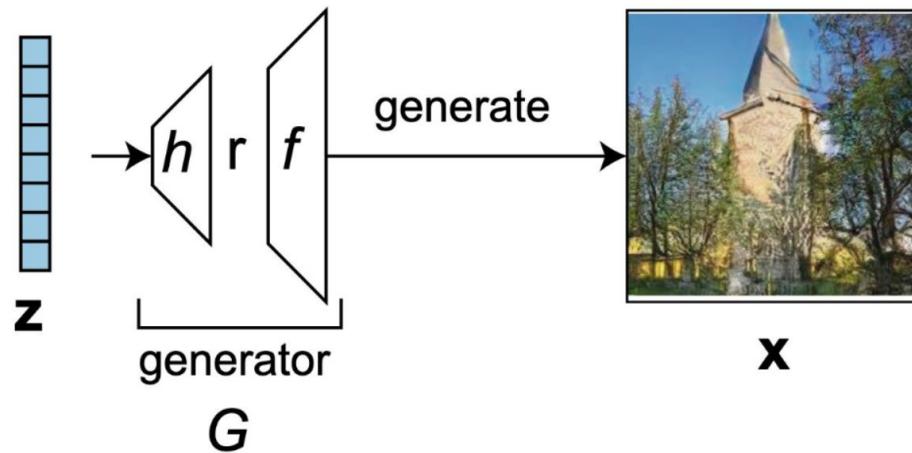


(f) Bedroom images with artifacts



(g) Ablating “artifact” units improves results

Определения



$$r = h(z) \text{ and } x = f(r) = f(h(z)) = G(z)$$

Определения

\mathbb{U} - множество всех юнитов (каналов)

\mathbb{P} - множество всех пикселей

U - подмножество юнитов

P - подмножество пикселей

C - множество всех концептов (классов объектов)

c - конкретный класс

План

1) Dissection

Хотим понять представления каких классов вообще присутствуют в гане и найти юниты хорошо скоррелированные с данным классом

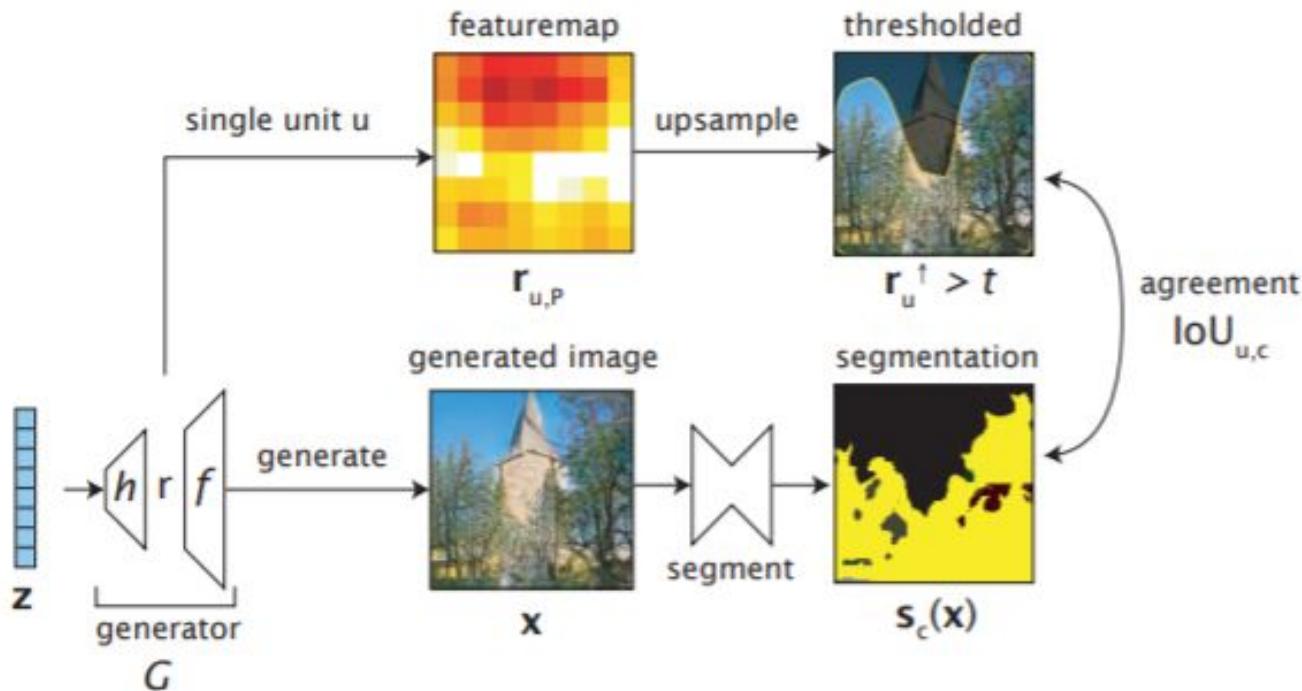
2) Intervention

Выделяем конкретные части изображения влияющие на генерацию объектов

Dissection



Dissection



Back to ПСМО

1) Взаимная информация

$$I(X;Y) = D_{\text{KL}}(P_{(X,Y)} \| P_X \otimes P_Y)$$

2) Совместная энтропия

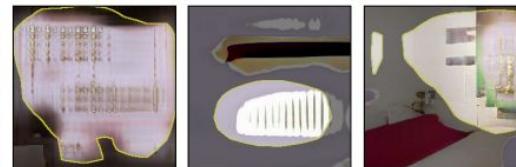
$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(x, y) \log_2 [P(x, y)]$$

Intersection over Union

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

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

Dissection



(e) Identify GAN units that cause artifacts



(f) Bedroom images with artifacts



(g) Ablating “artifact” units improves results

Intervention

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 :

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

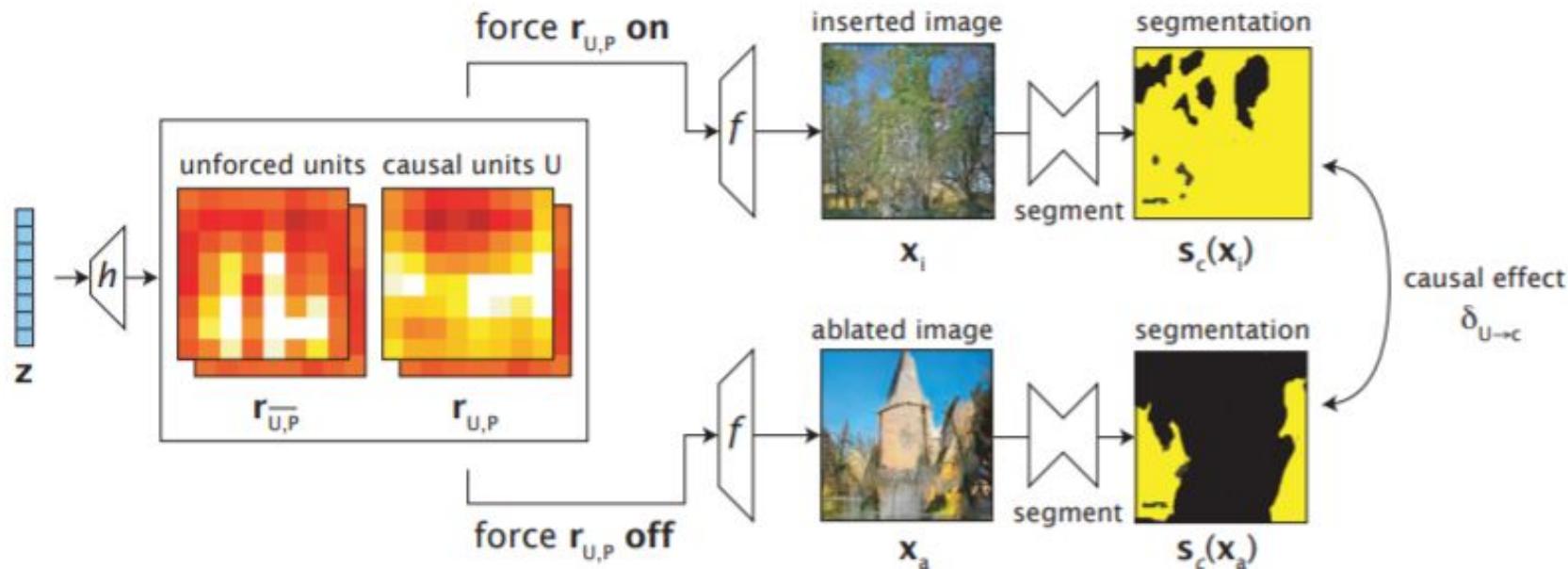
Image with U inserted at pixels P :

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

Average Causal Effect

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

Intervention



Intervention

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

Image with partial ablation at pixels P :	$\mathbf{x}'_a = f((\mathbf{1} - \boldsymbol{\alpha}) \odot \mathbf{r}_{\mathbb{U},P}, \mathbf{r}_{\mathbb{U},\bar{P}})$
Image with partial insertion at pixels P :	$\mathbf{x}'_i = f(\boldsymbol{\alpha} \odot \mathbf{k} + (\mathbf{1} - \boldsymbol{\alpha}) \odot \mathbf{r}_{\mathbb{U},P}, \mathbf{r}_{\mathbb{U},\bar{P}})$
Objective :	$\delta_{\boldsymbol{\alpha} \rightarrow c} = \mathbb{E}_{\mathbf{z},P} [\mathbf{s}_c(\mathbf{x}'_i)] - \mathbb{E}_{\mathbf{z},P} [\mathbf{s}_c(\mathbf{x}'_a)],$

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

То что осталось за кадром

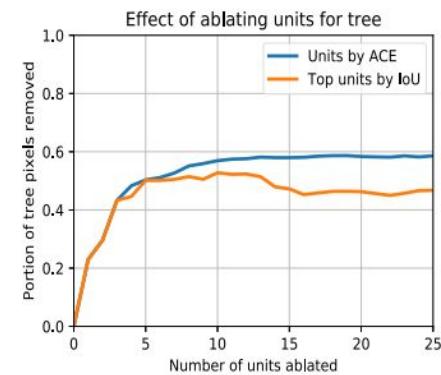
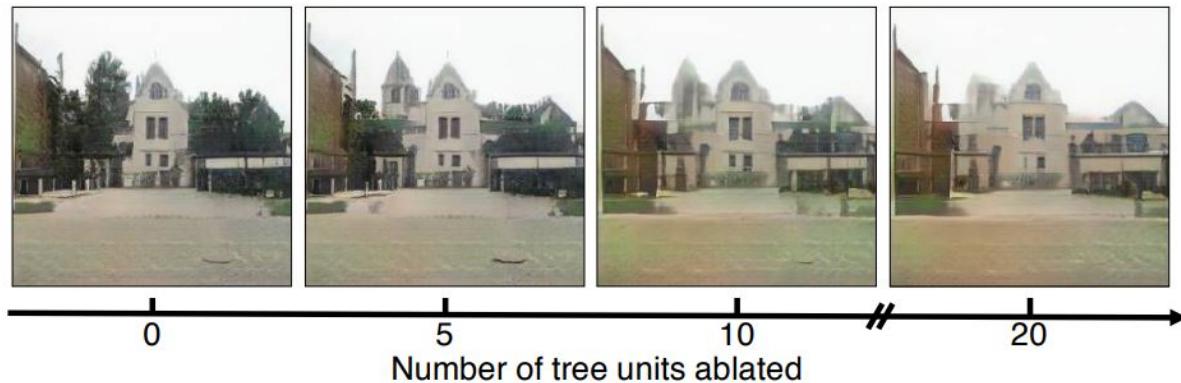
- 1) Инициализация альфа

$$\alpha_u = \frac{\text{IoU}_{u,c}}{\max_v \text{IoU}_{v,c}}$$

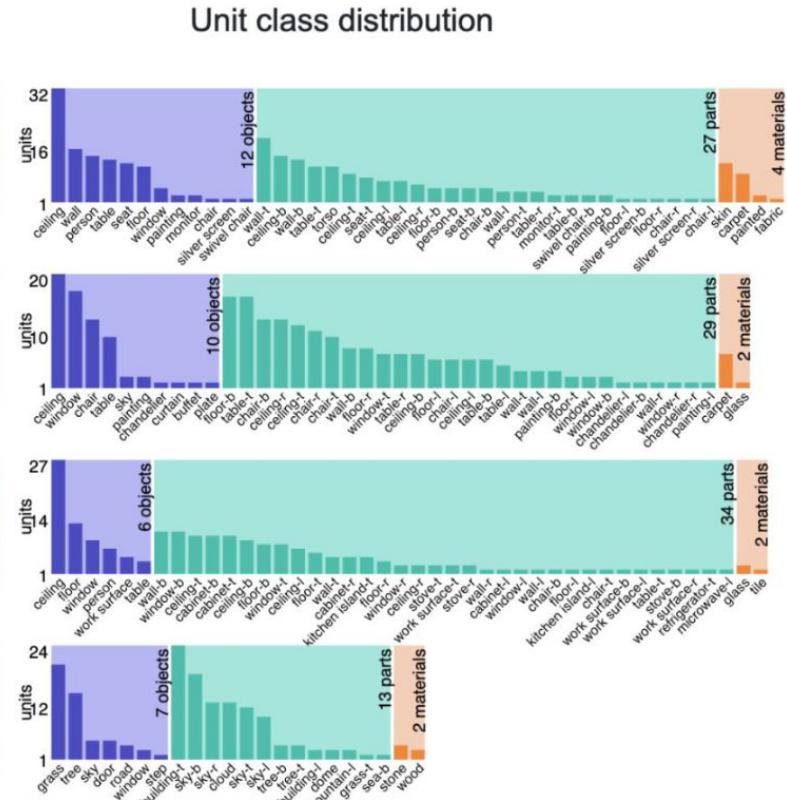
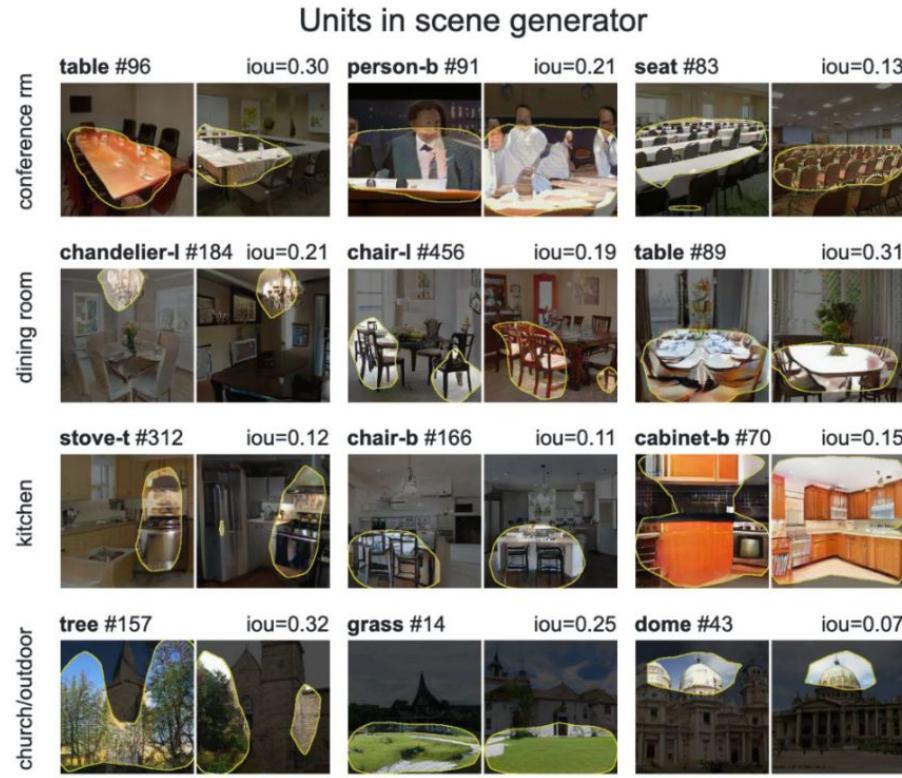
- 2) И все-таки Р

Будем сэмплить из масок найденных на первом шаге и сегментаций оригинальных картинок

Intervention



Интерпретируемые юниты для разных данных



Юниты для разных слоев

layer1
512 units total
0 object units
2 part units
0 material units



layer4
512 units total
86 object units
149 part units
10 material units



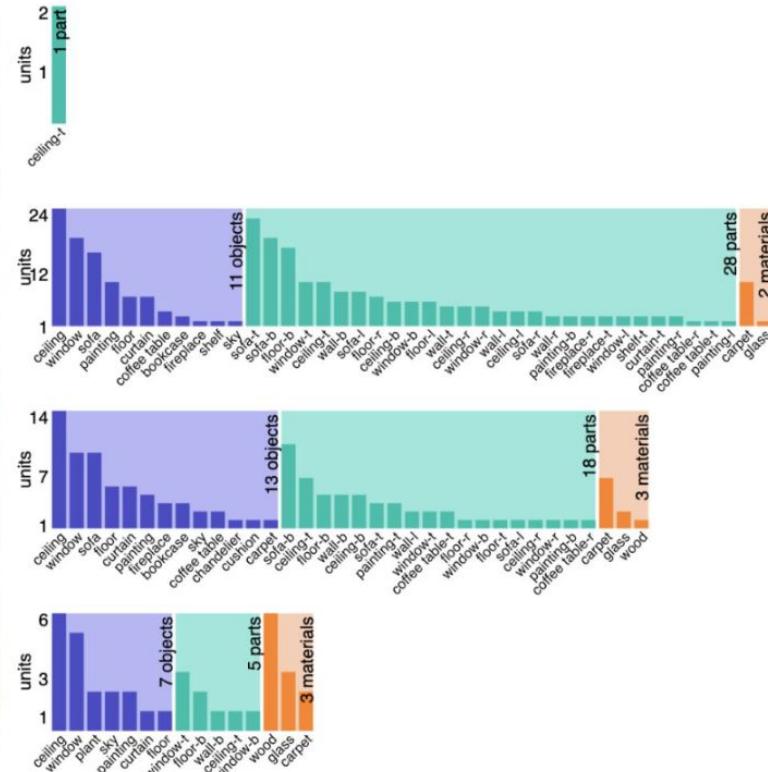
layer7
256 units total
59 object units
48 part units
9 material units



layer10
128 units total
19 object units
8 part units
11 material units



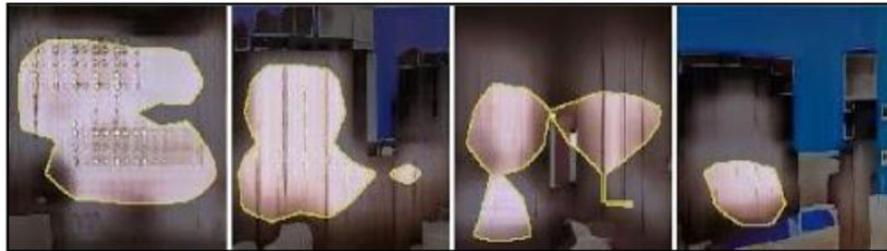
Unit class distribution



Юниты для различных вариантов обучения

interpretable units	SWD	Best "bed" unit	Best "window" unit	Unit class distribution
base prog GAN 512 units total 74 object units 84 part units 9 material units	167 units 7.60	bed layer4 #253 iou=0.18 	window layer4 #142 iou=0.19 	
+batch stddev 512 units total 55 object units 128 part units 6 material units	189 units 6.48	bed layer4 #88 iou=0.11 	window layer4 #422 iou=0.25 	
+pixelwise norm 512 units total 82 object units 128 part units 16 material units	226 units 4.01	bed layer4 #129 iou=0.29 	window layer4 #494 iou=0.26 	

Удаление артефактов



Example artifact-causing
units



Bedroom images with artifacts



Ablating "artifact" units improves
results

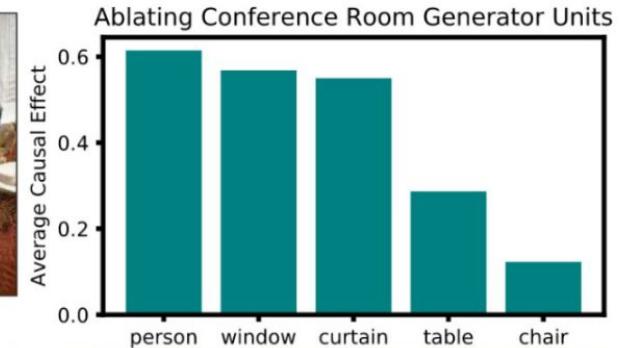
Удаление объектов



ablate person units



ablate curtain units



ablate window units

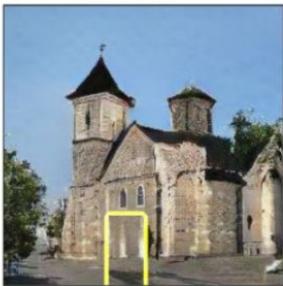


ablate table units



ablate chair units

Вставка объектов



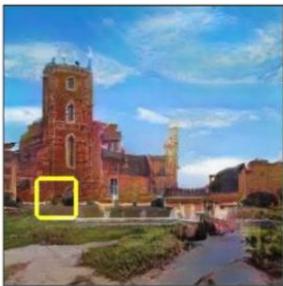
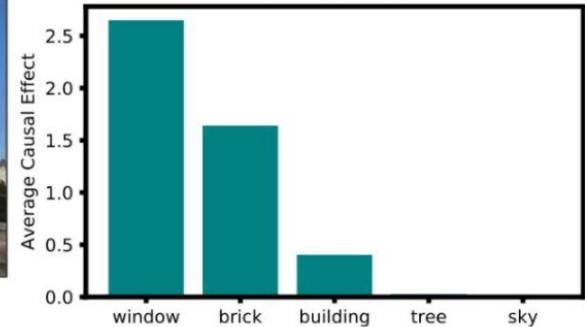
(a)



(b)



Where Can a Door Go?



(c)



(d)



(e)



Итоги кратко

- GAN -ы содержат подмножества нейронов, отвечающие за генерацию понятных человеку классов
- GAN -ы могут учить композицию
- Артефакты также могут генерироваться выделенным подмножеством нейронов