

GAN DISSECTION

Аъзам Бехруз

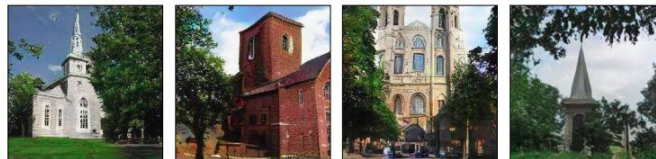
22.02.2022 - палиндром

Интересные вопросы

- Существуют ли во внутренностях GAN понятия об объектах (деревья, двери и т.д.)?
- Есть ли внутри переменные отвечающие за деревья и двери?
- Если да, то каким образом представлено их взаимодействие?

Возможные применения

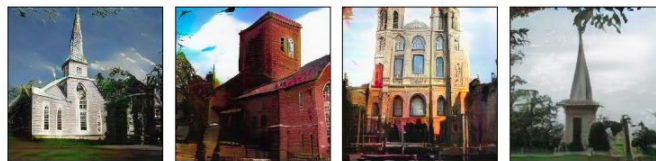
- Манипуляции с объектами
- Удаление артефактов



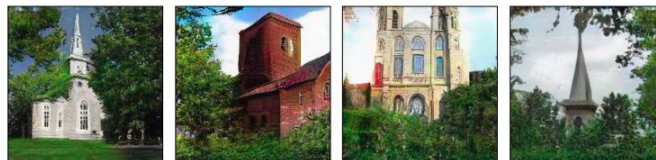
(a) Generate images of churches



(b) Identify GAN units that match trees



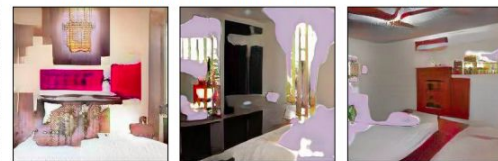
(c) Ablating units removes trees



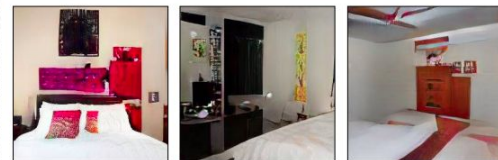
(d) Activating units adds trees



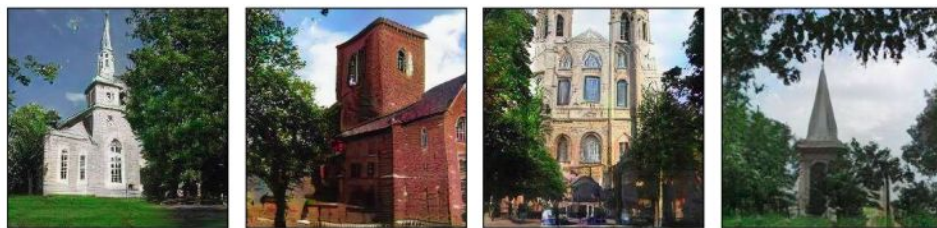
(e) Identify GAN units that cause artifacts



(f) Bedroom images with artifacts



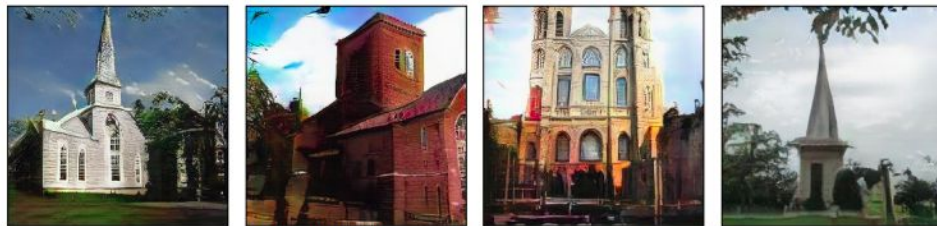
(g) Ablating "artifact" units improves results



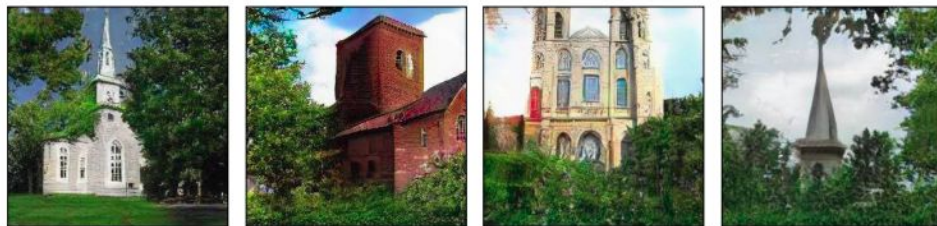
(a) Generate images of churches



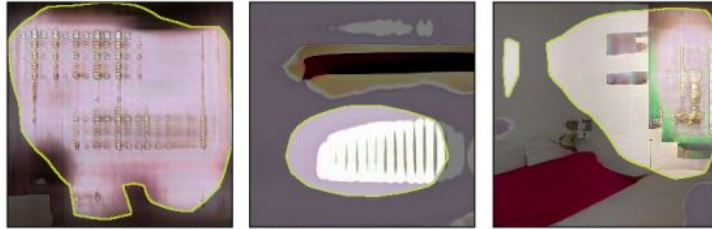
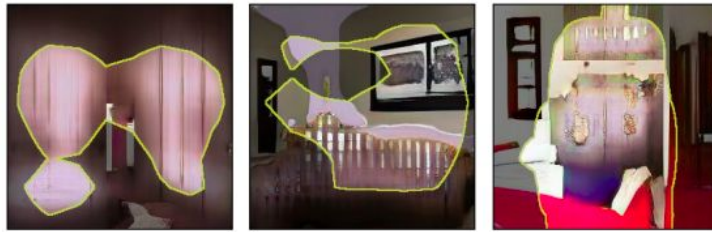
(b) Identify GAN units that match trees



(c) Ablating units removes trees



(d) Activating units adds trees



(e) Identify GAN units that cause artifacts

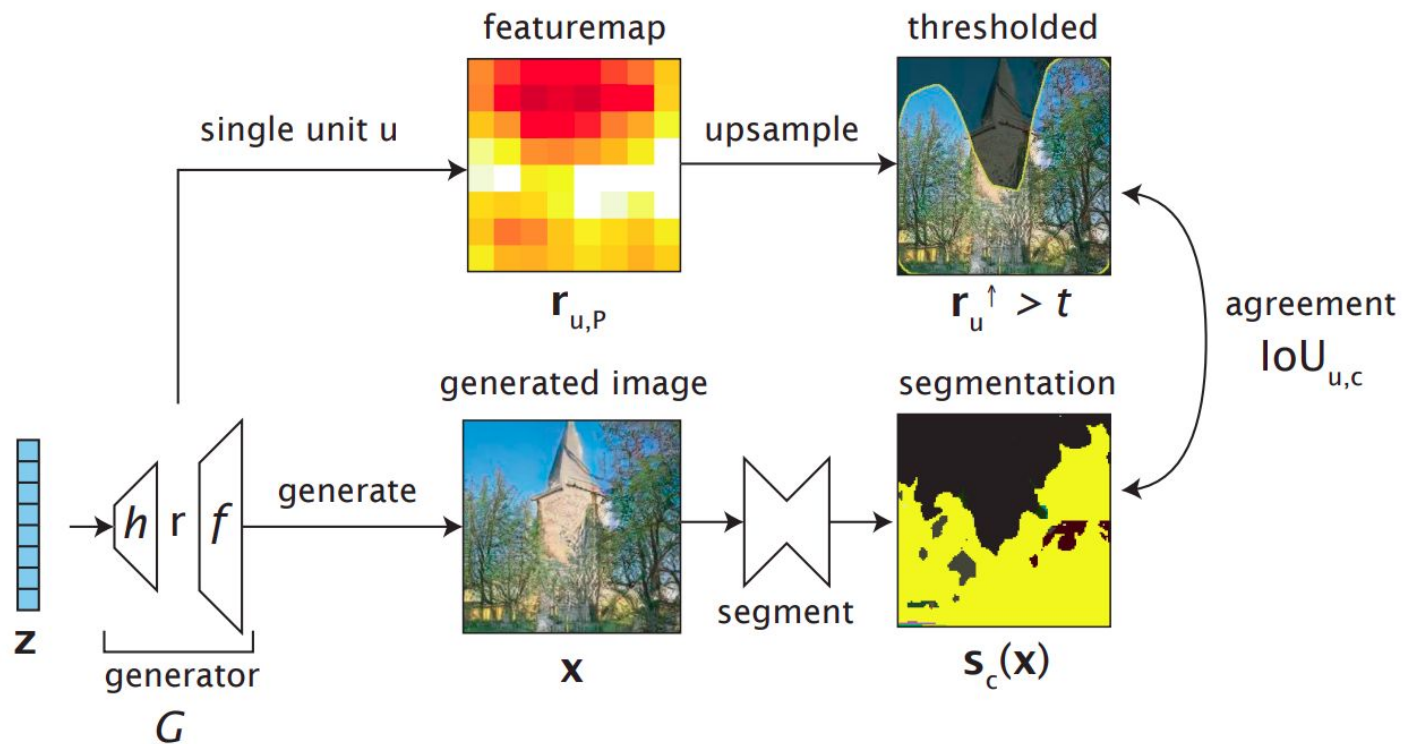


(f) Bedroom images with artifacts



Выявление и вмешательство

Выявление



Пространственная согласованность

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

Мы присваиваем каждому слою класс с которым у него наибольшая IoU

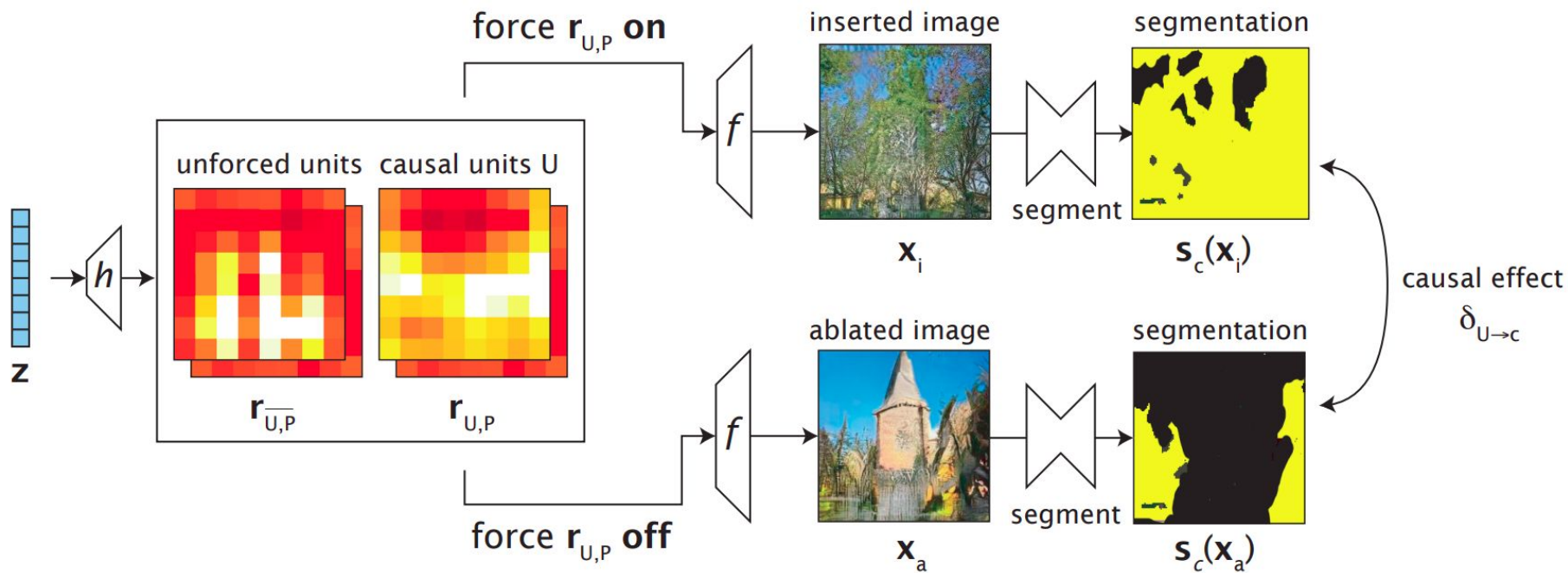


Thresholding unit #65 layer 3 of a dining room generator matches 'table' segmentations with IoU=0.34.



Thresholding unit #37 layer 4 of a living room generator matches 'sofa' segmentations with IoU=0.29.

Вмешательство



Подсчет каузального эффекта

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

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 :

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

Каузальный эффект:

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

**Комбинаторный
взрыв!**



Вмешательство

$$\boldsymbol{\alpha} \in [0, 1]^d$$

Image with partial ablation at pixels \mathbf{P} :

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

Image with partial insertion at pixels \mathbf{P} :

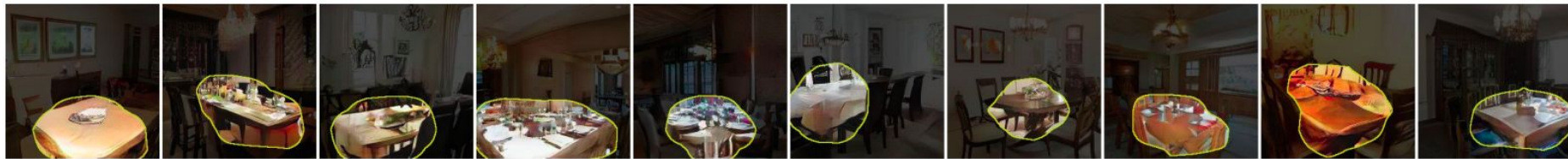
$$\mathbf{x}'_i = f(\boldsymbol{\alpha} \odot \mathbf{k} + (\mathbf{1} - \boldsymbol{\alpha}) \odot \mathbf{r}_{\mathbb{U}, \mathbf{P}}, \mathbf{r}_{\mathbb{U}, \bar{\mathbf{P}}})$$

Objective :

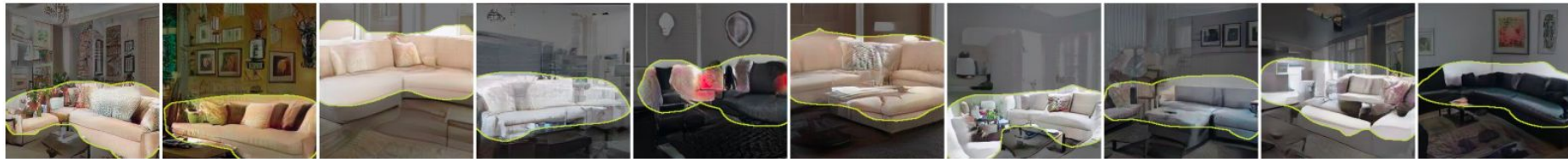
$$\delta_{\boldsymbol{\alpha} \rightarrow c} = \mathbb{E}_{\mathbf{z}, \mathbf{P}} [\mathbf{s}_c(\mathbf{x}'_i)] - \mathbb{E}_{\mathbf{z}, \mathbf{P}} [\mathbf{s}_c(\mathbf{x}'_a)],$$

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

Выводы



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 $\text{IoU}=0.29$.

GAN имеет такое же понятие о столах и диванах что и мы

- Слои с 4-7 коррелируют с семантическими объектами
- 10 слой и далее отвечают за группы пикселей (углы, материал)

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

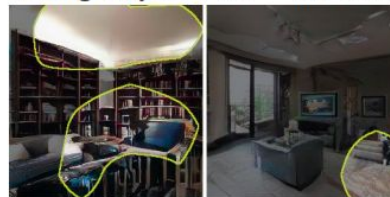
layer4
512 units total
86 object units
149 part units
10 material units

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

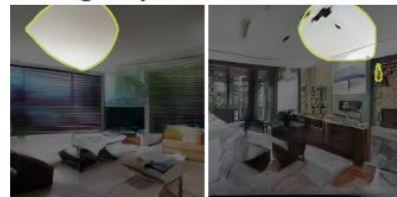
layer10
128 units total
19 object units
8 part units
11 material units

Units in layer

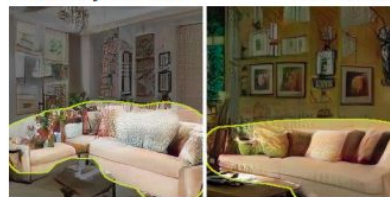
ceiling-t layer1 #457 iou=0.10



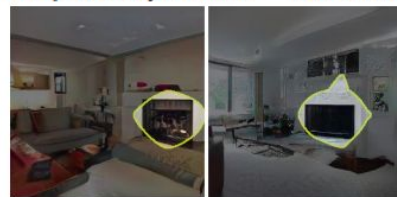
ceiling-t layer1 #194 iou=0.07



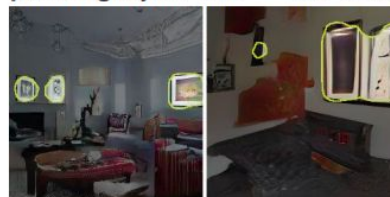
sofa layer4 #37 iou=0.28



fireplace layer4 #23 iou=0.15



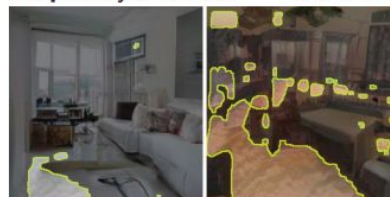
painting layer7 #15 iou=0.23



coffee table-t #247 iou=0.07



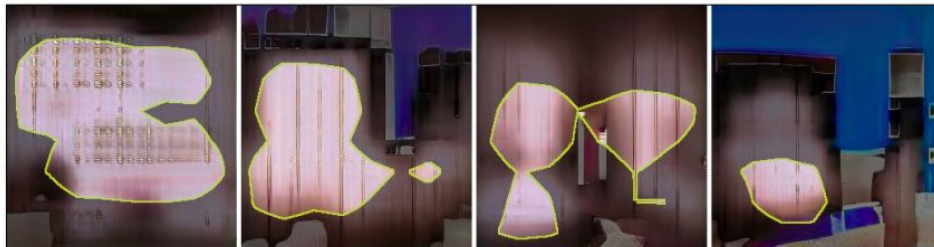
carpet layer10 #53 iou=0.14



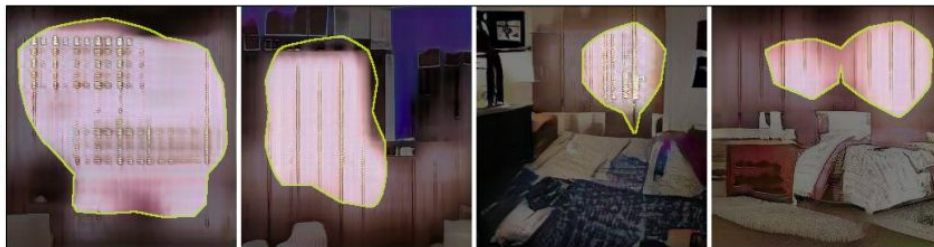
glass layer10 #126 iou=0.21



Unit #63



Unit #231



(a) Example artifact-causing units



(b) Bedroom images with artifacts



(c) Ablating “artifact” units improves results

Fréchet Inception Distance (FID)

original images	43.16
“artifacts” units ablated (ours)	27.14
random units ablated	43.17

Human preference score

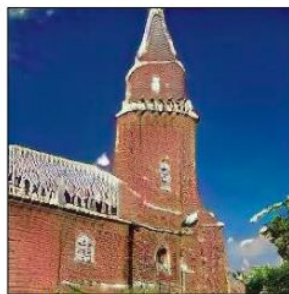
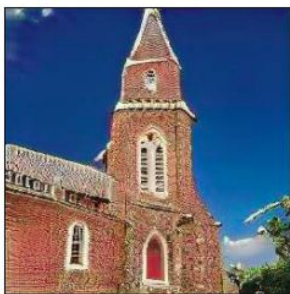
original images

“artifacts” units ablated (ours)	72.4%
random units ablated	49.9%

Можем удалять артефакты

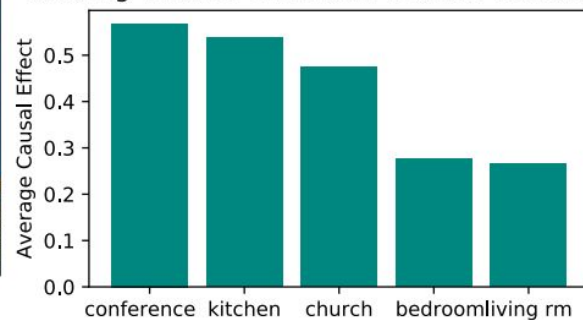


conference room



church

Ablating Window Units from Several Generators



kitchen

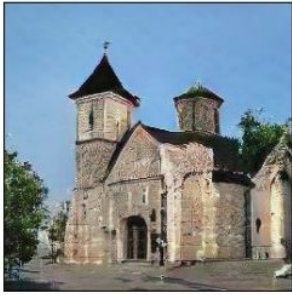
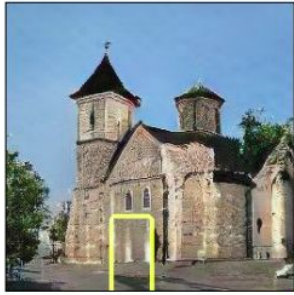


living room



bedroom

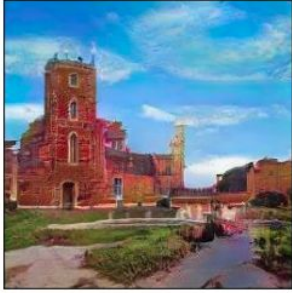
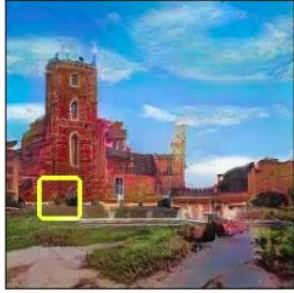
Не всегда можно просто взять и удалить окно.
Все зависит от контекста.



(a)



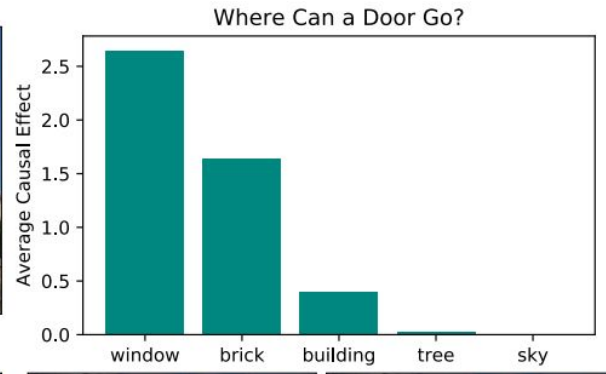
(b)



(c)



(d)



(e)

Можем вставлять двери куда хотим!

Understanding GANs



GANDissection

Interacting with GANs



GANPaint



Спасибо за внимание!