

GAN DISSECTION: VISUALIZING AND UNDERSTANDING GENERATIVE ADVERSARIAL NETWORKS

David Bau, Jun-Yan Zhu, Joshua B. Tenenbaum,
William T. Freeman, Antonio Torralba
[davidbau@csail.mit.edu , junyanz@csail.mit.edu ,
jbt@csail.mit.edu, billf@csail.mit.edu,
torralba@csail.mit.edu]

Massachusetts Institute of Technology, Boston

Hendrik Strobelt
hendrik.strobelt@ibm.com

IBM Research, Cambridge MA

Bolei Zhou
bzhou@ie.cuhk.edu.hk

The Chinese University of Hong Kong

Presentation by: Tejas Gaikwad
MT19AI021
Indian Institute of Technology, Jodhpur, India

Conference paper at ICLR 2019

EVER IMAGINED HOW GAN CAN GIVE SO REALISTIC FAKE IMAGES?.....

Restaurant



Living room



Church

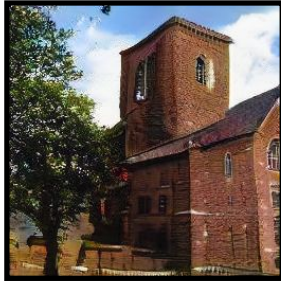


256x256 images synthesized by a Progressive GAN [Karras, et al 2017]

Credits : [arXiv:1811.10597](https://arxiv.org/abs/1811.10597)

TO RENDER A BEAUTIFUL SCENE, WHAT DOES A GAN NEED TO KNOW?

Church



AND SOMETIMES.... WHAT CAUSES THE MISTAKES?

Bedroom



BUT BEFORE THOSE QUESTIONS...

- What are Generative Adversarial Networks (GAN's) ?
- What makes them so “interesting” ?

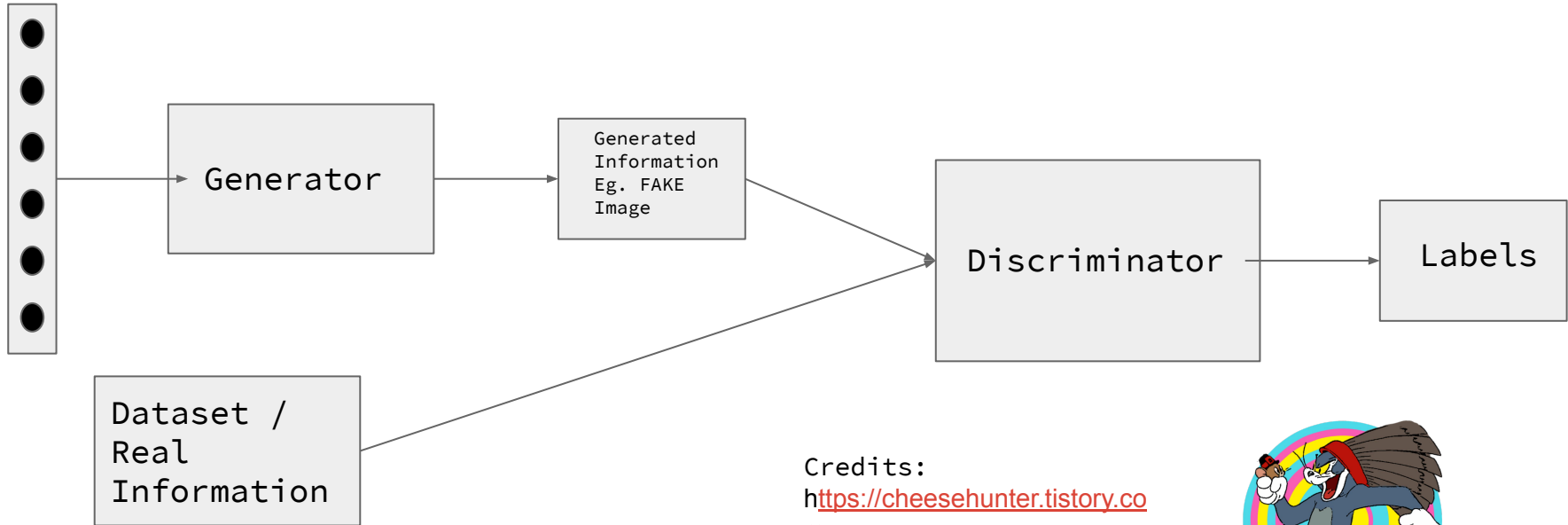


GAN'S

- These are the networks that belong to the **Generative Model**, introduced by **Ian Goodfellow in the year 2014**.
- GAN's are used to create new images which are not in the dataset, but looks natural.
- A GAME based approach is used to train the model.
- It consists of 2 main blocks named as a **Discriminator** and a **Generator**. Discriminator simply try to discriminate between generated images and images from the dataset. Whereas, generator generates the FAKE image (generated using random noise) and try to match the natural images in the dataset.

BLOCK DIAGRAM FOR GAN

Random Noise Vector



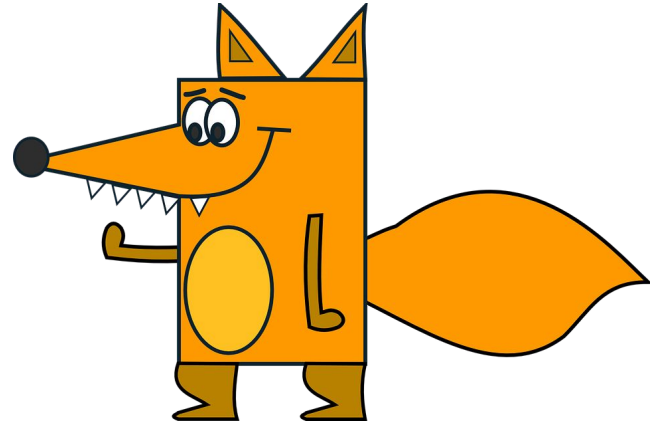
Credits:
<https://cheesehunter.tistory.com>



LITERATURE SURVEY

- **GAN(Generative Adversarial Networks) (Goodfellow et al., 2014):**
Introduced GAN's
- **Visualizing deep neural networks:**
Visualizations for CNNs (Zeiler & Fergus, 2014) and RNNs (Karpathy et al., 2016; Strobel et al., 2018), by locating and reconstructing salient image features (Simonyan et al., 2014; Mahendran & Vedaldi, 2015) or by mining patches that maximize hidden layers' activations (Zeiler & Fergus, 2014), or we can synthesize input images to invert a feature layer (Dosovitskiy & Brox, 2016)
- **Explaining the decisions of deep neural networks:**
Explains individual network decisions using informative heatmaps (Zhou et al., 2018b; 2016; Selvaraju et al., 2017) or modified backpropagation (Simonyan et al., 2014; Bach et al., 2015; Sundararajan et al., 2017)

THE PROPOSED MODEL



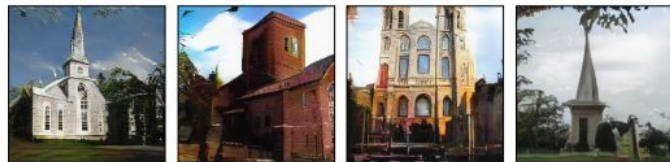
Goal is to analyze how objects such as trees are encoded by the internal representations of a GAN generator $G: z \rightarrow x$.



(a) Generate images of churches



(b) Identify GAN units that match trees



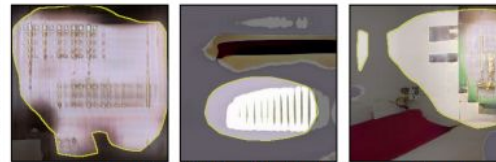
(c) Ablating units removes trees



(d) Activating units adds trees



(e) Identify GAN units that cause artifacts

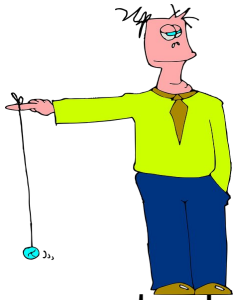


(f) Bedroom images with artifacts



(g) Ablating "artifact" units improves results

STEPS...



They've presented an analytic framework to visualize and understand GANs at the unit-, object-, and scene-level

- First step is identify a group of interpretable units that are related to Semantic Classes.
- These units' featuremaps closely match the semantic segmentation of a particular object class (e.g., trees)
- Then, intervene in units in the network to cause a type of object to disappear or appear
- And Finally, study the contextual relationships by observing where we can insert the object

NOW L'IL MATHEMATICAL EXPLANATION

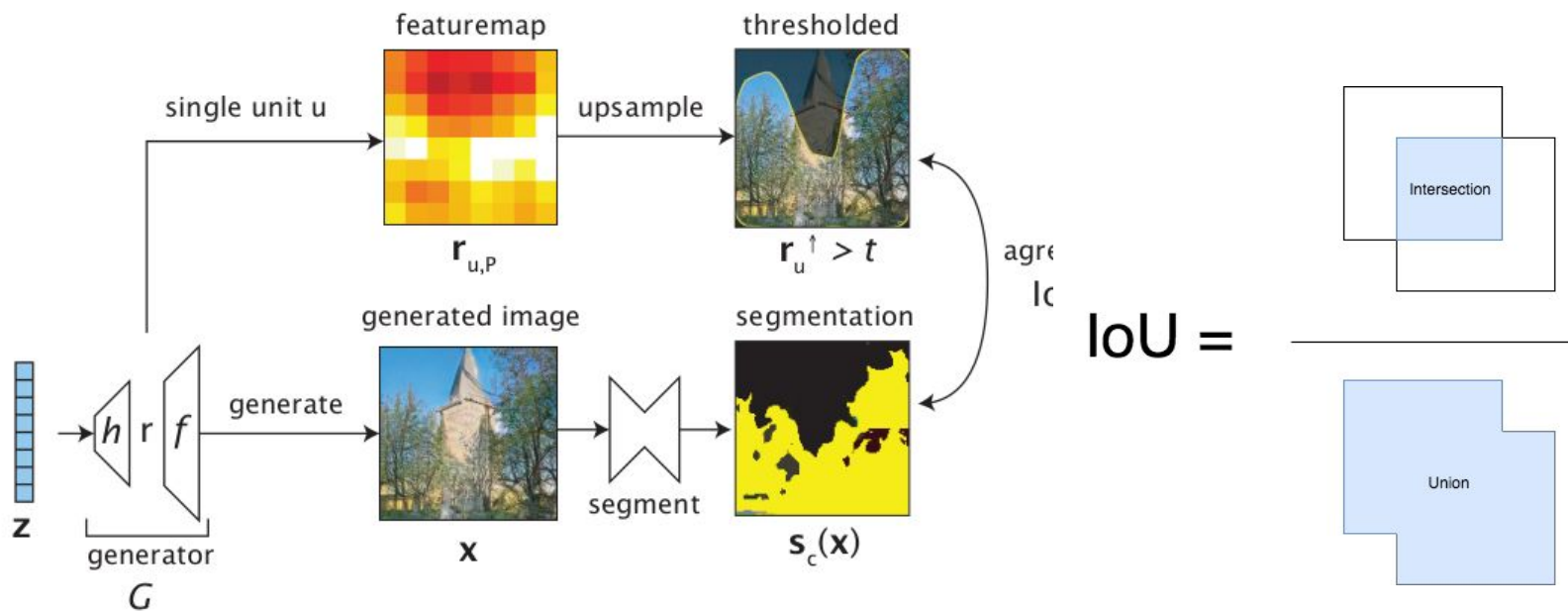
- They have **analysed the internal GAN representations** by decomposing the featuremap r at a layer into positions $P \subset \mathbf{P}$ and unit channels $u \in \mathbf{U}$
- To identify a unit u with semantic behavior, they have upsampled and thresholded the unit, and measured how well it matches an object class c in the image x as identified by a supervised semantic segmentation network $S_c(x)$
- Each unit is a little segmentation solution. A standard way to measure segmentation accuracy is IoU(Intersection over Union), which is given as

$$\text{IoU}_{u,c} \equiv \frac{\mathbb{E}_{\mathbf{z}} \left| (\mathbf{r}_{u,\mathbf{P}}^{\uparrow} > t_{u,c}) \wedge s_c(\mathbf{x}) \right|}{\mathbb{E}_{\mathbf{z}} \left| (\mathbf{r}_{u,\mathbf{P}}^{\uparrow} > t_{u,c}) \vee s_c(\mathbf{x}) \right|} \quad \text{Where the threshold } t \text{ is given as follows}$$

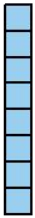
$$t_{u,c} = \arg \max_t \frac{I(\mathbf{r}_{u,\mathbf{P}}^{\uparrow} > t; s_c(\mathbf{x}))}{H(\mathbf{r}_{u,\mathbf{P}}^{\uparrow} > t, s_c(\mathbf{x}))}$$

The threshold $t_{u,c}$ is chosen to maximize the information quality ratio, that is, the **portion of the joint entropy H which is mutual information I**

HOW UNITS CORRELATE TO AN OBJECT CLASS?



WHICH UNITS CORRELATE TO AN OBJECT CLASS?



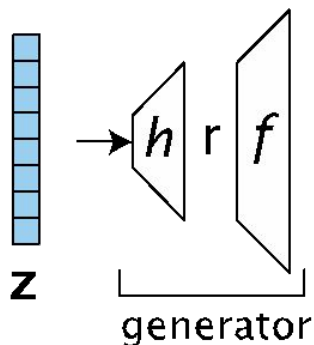
Z

WHICH UNITS CORRELATE TO AN OBJECT CLASS?

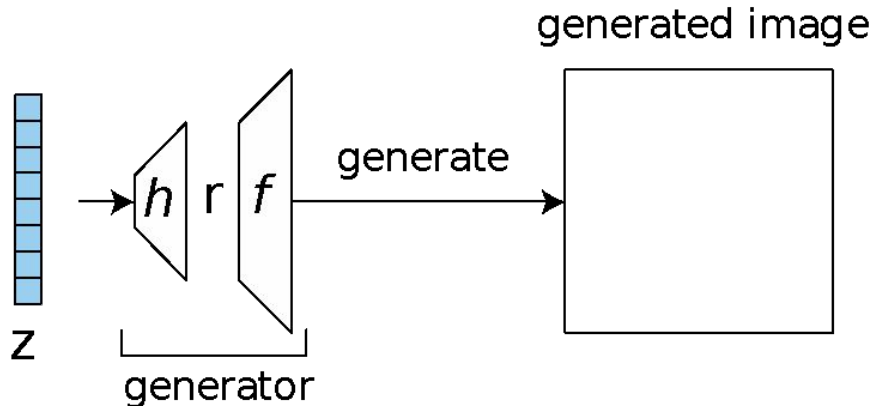
r : the current layer

h : the first half

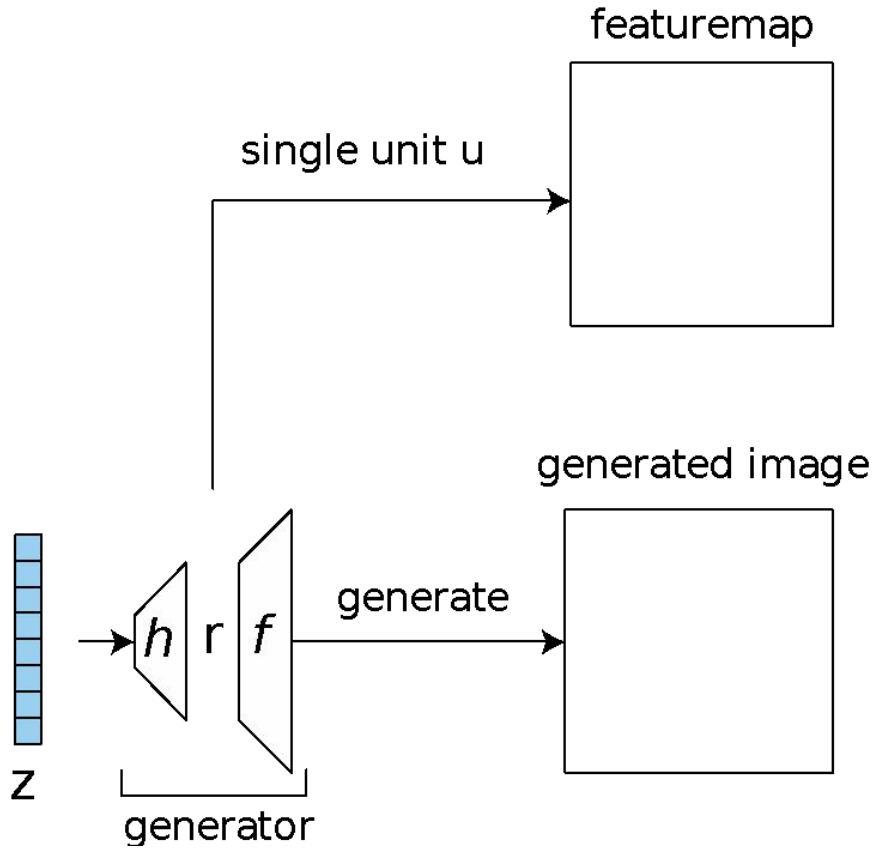
f : the second half



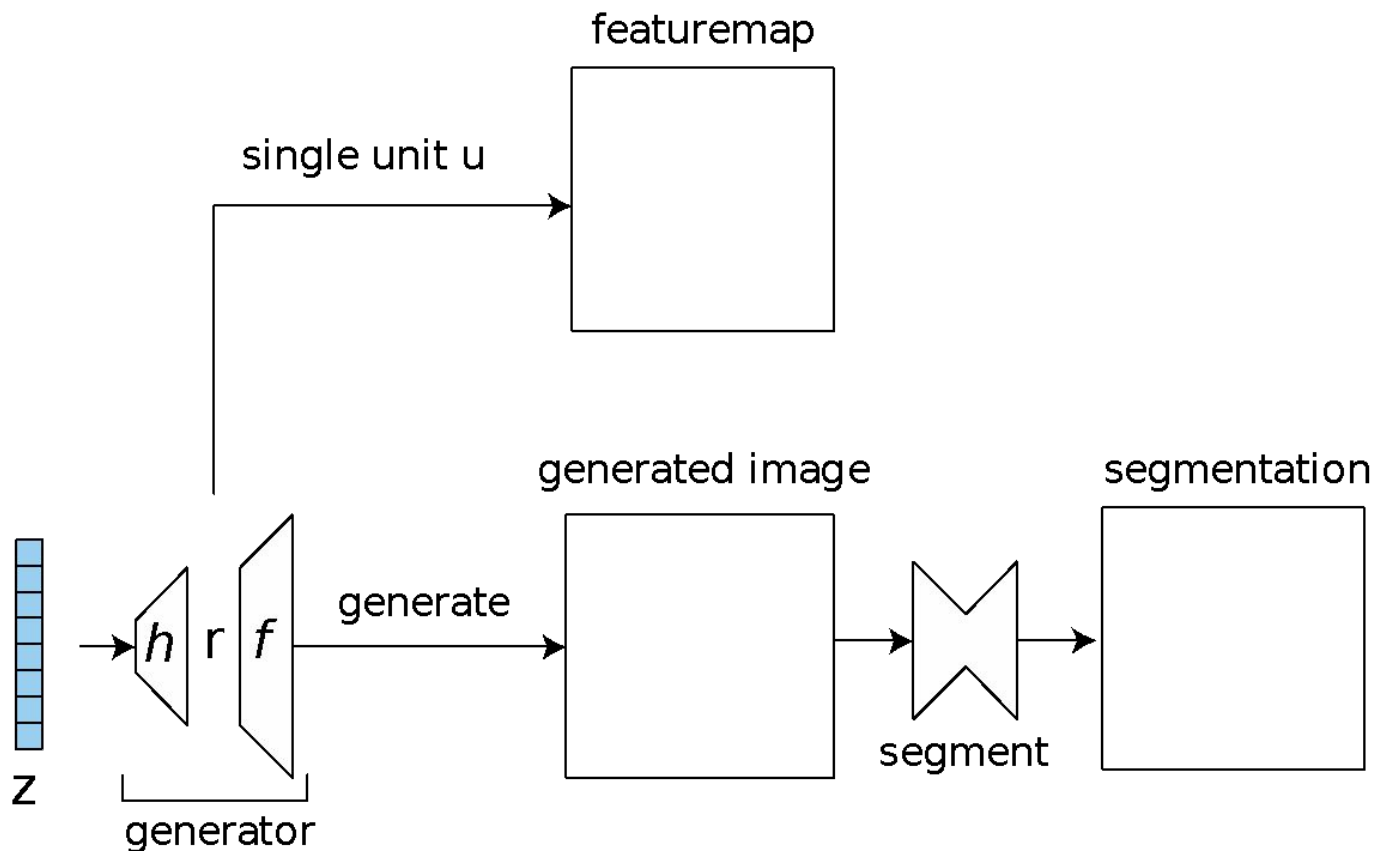
WHICH UNITS CORRELATE TO AN OBJECT CLASS?



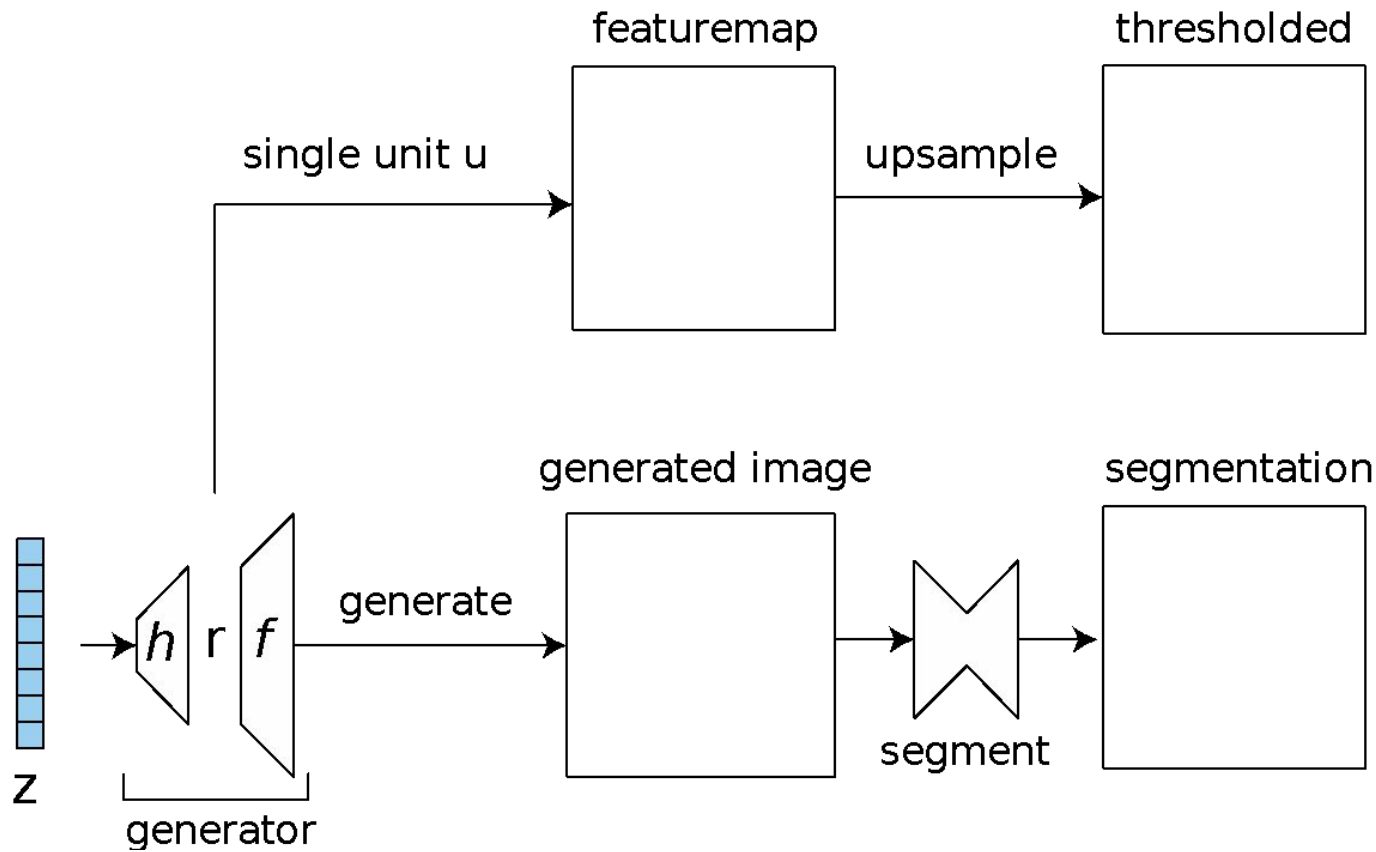
WHICH UNITS CORRELATE TO AN OBJECT CLASS?



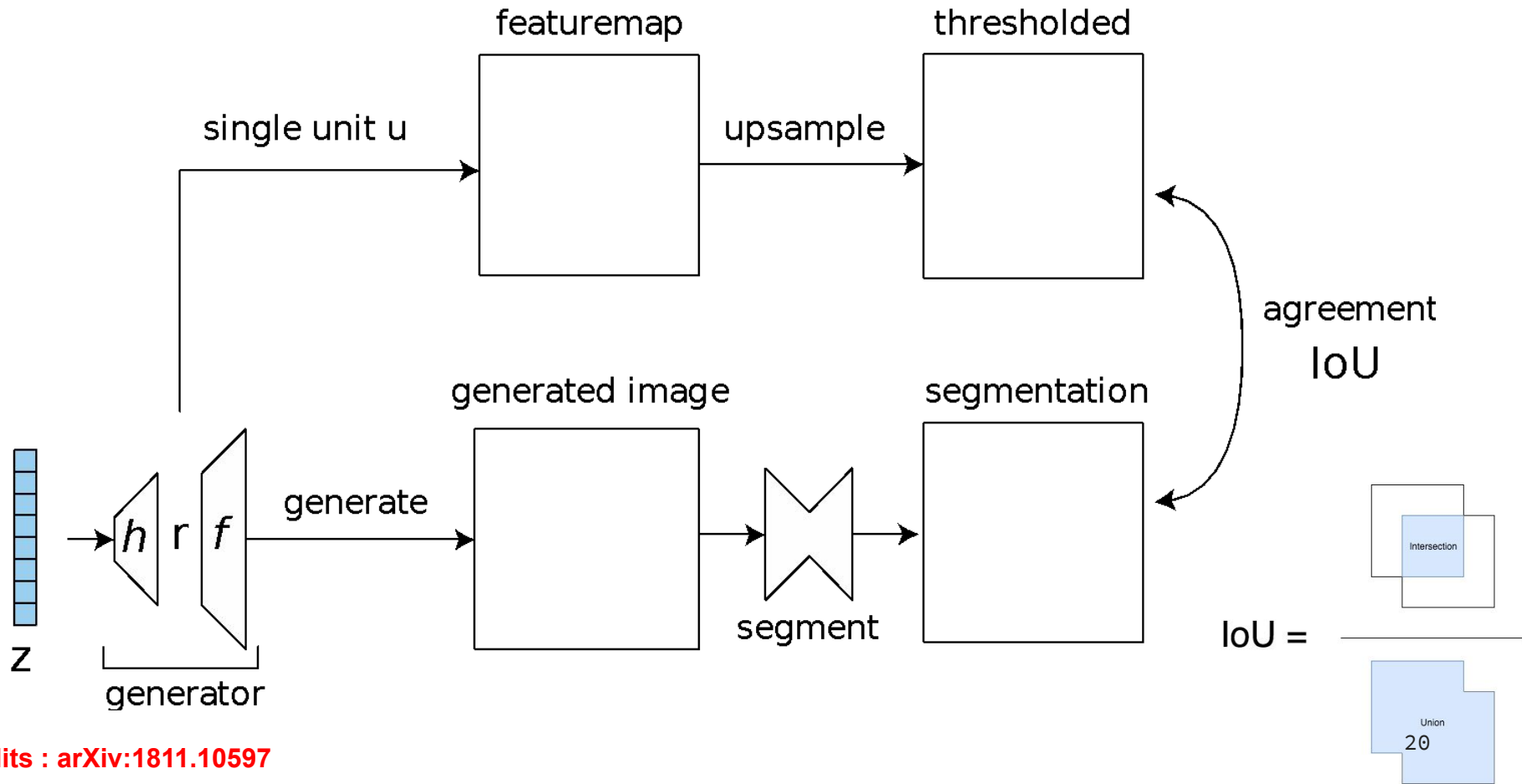
WHICH UNITS CORRELATE TO AN OBJECT CLASS?



WHICH UNITS CORRELATE TO AN OBJECT CLASS?



WHICH UNITS CORRELATE TO AN OBJECT CLASS?



WHICH UNITS CORRELATE TO AN OBJECT CLASS?

Church samples



Unit:
Tree



Unit:
Dome

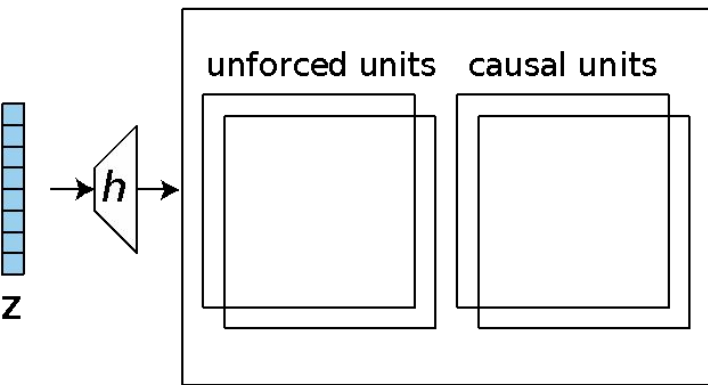


WHICH UNITS CAUSE AN OBJECT CLASS?

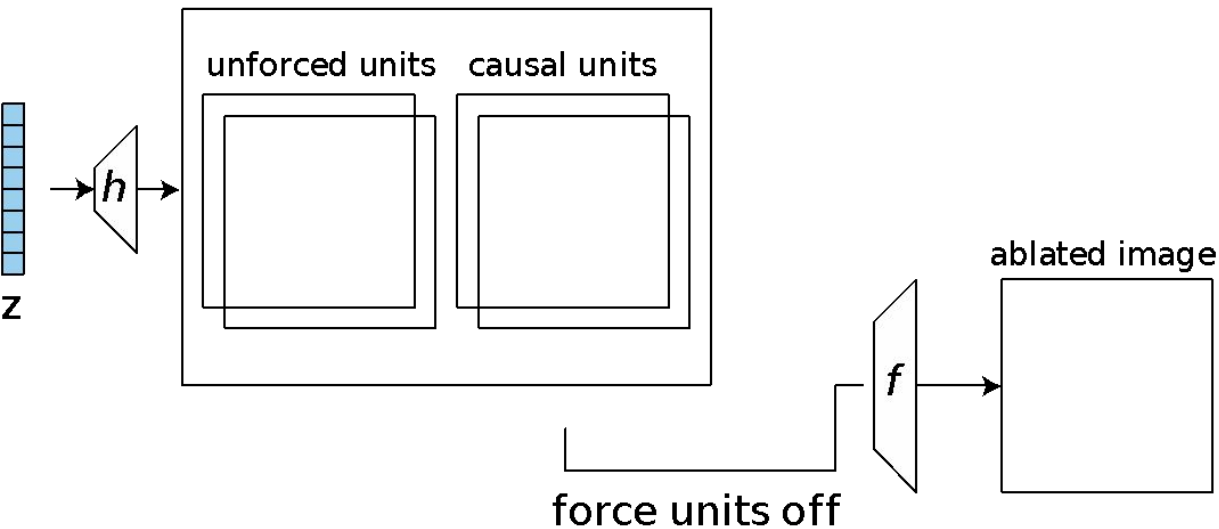


Z

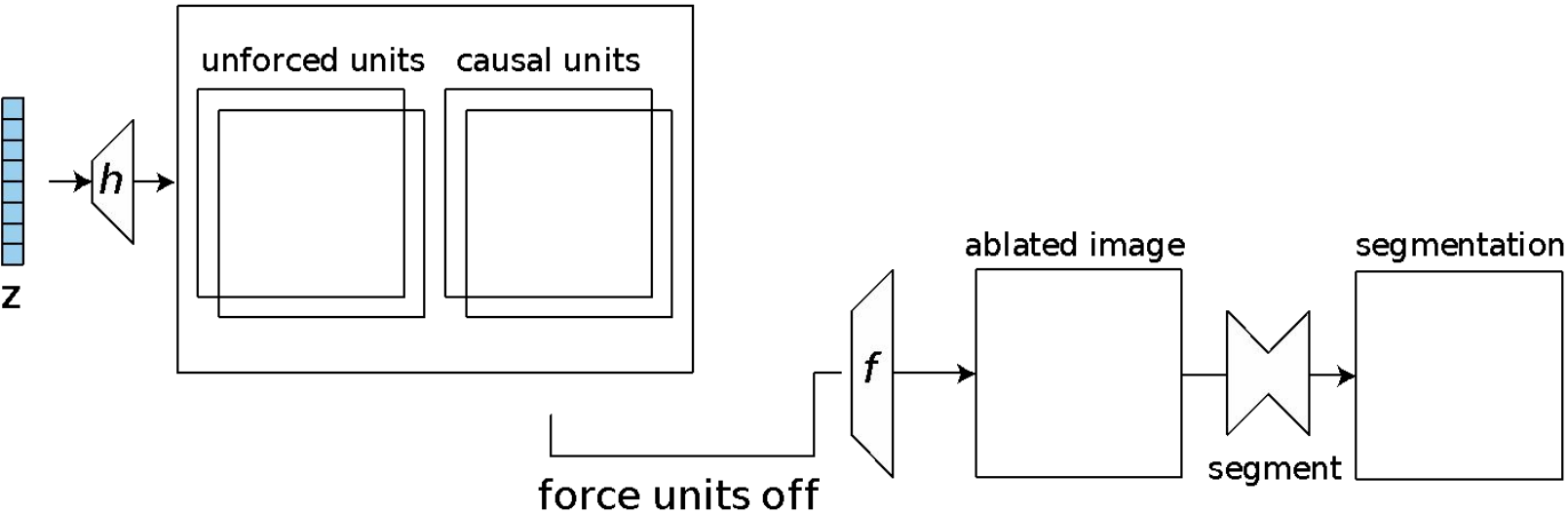
WHICH UNITS CAUSE AN OBJECT CLASS?



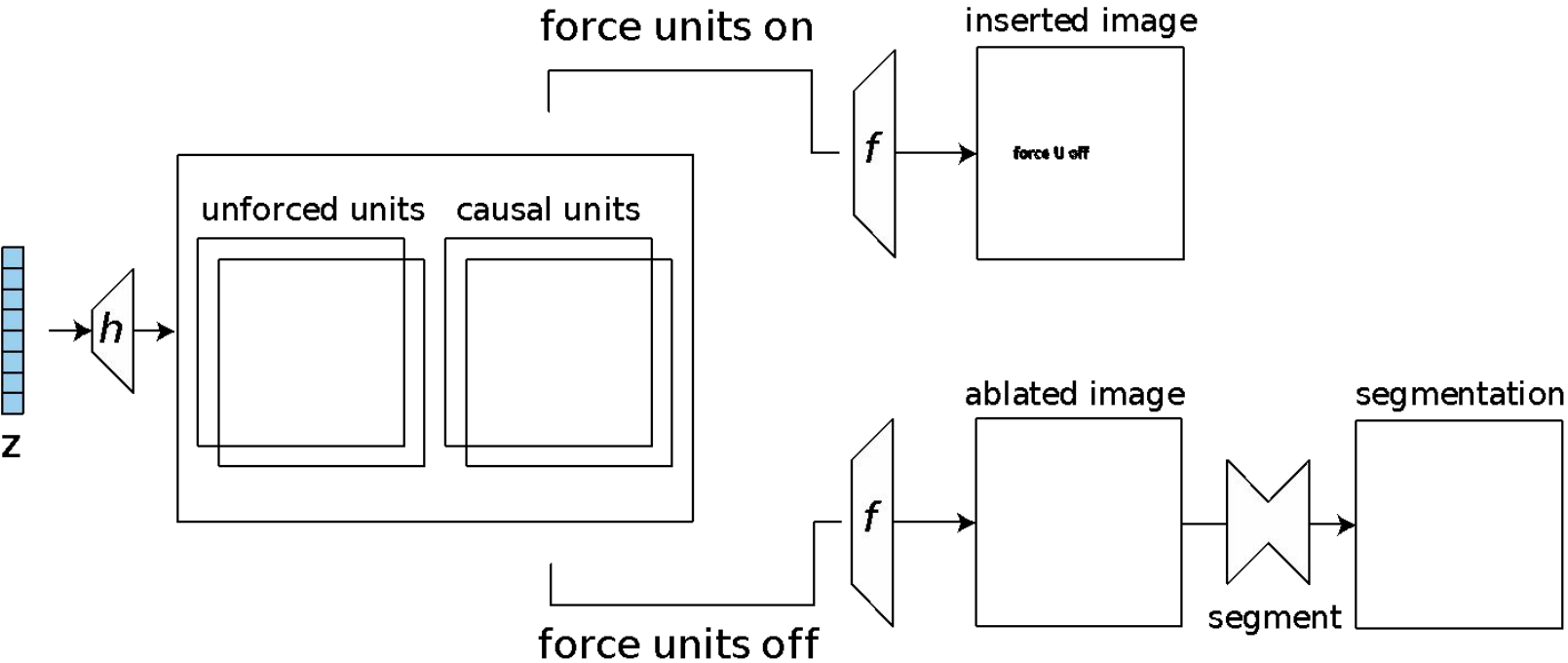
WHICH UNITS CAUSE AN OBJECT CLASS?



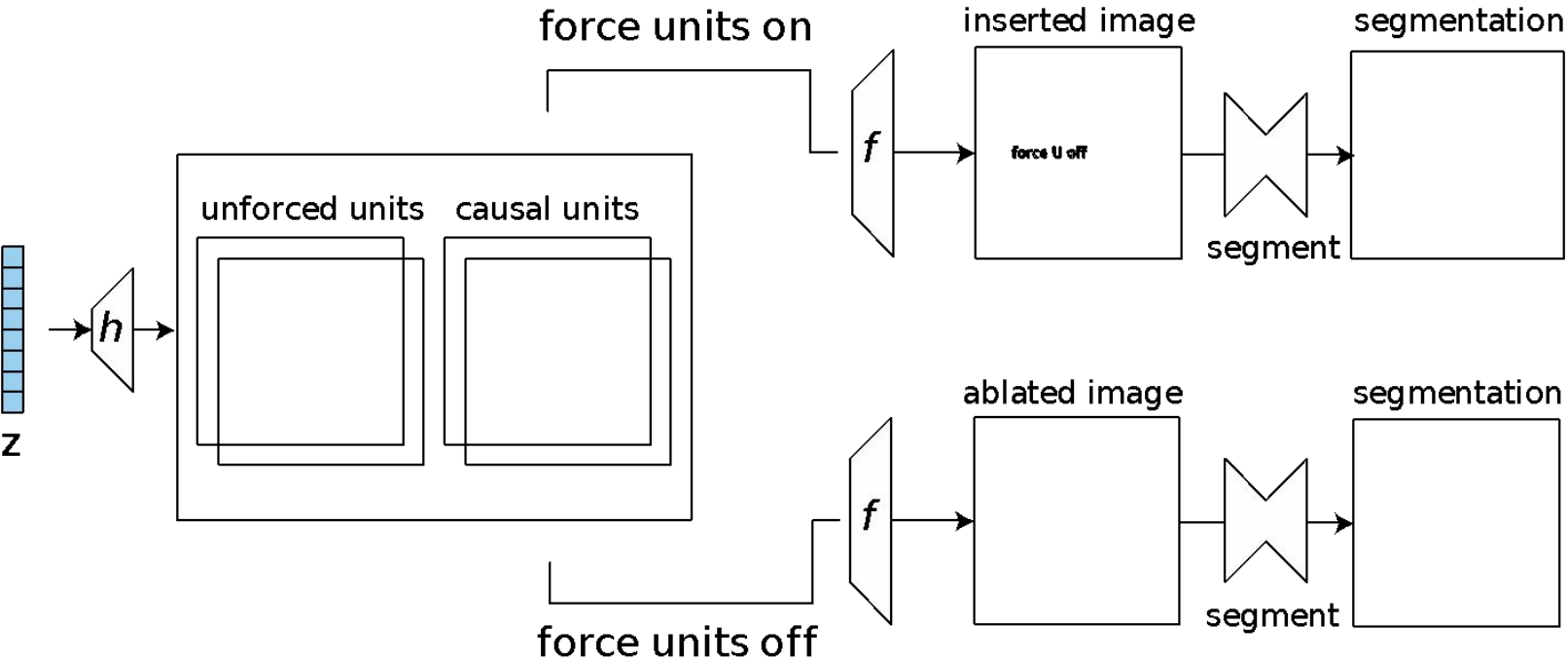
WHICH UNITS CAUSE AN OBJECT CLASS?



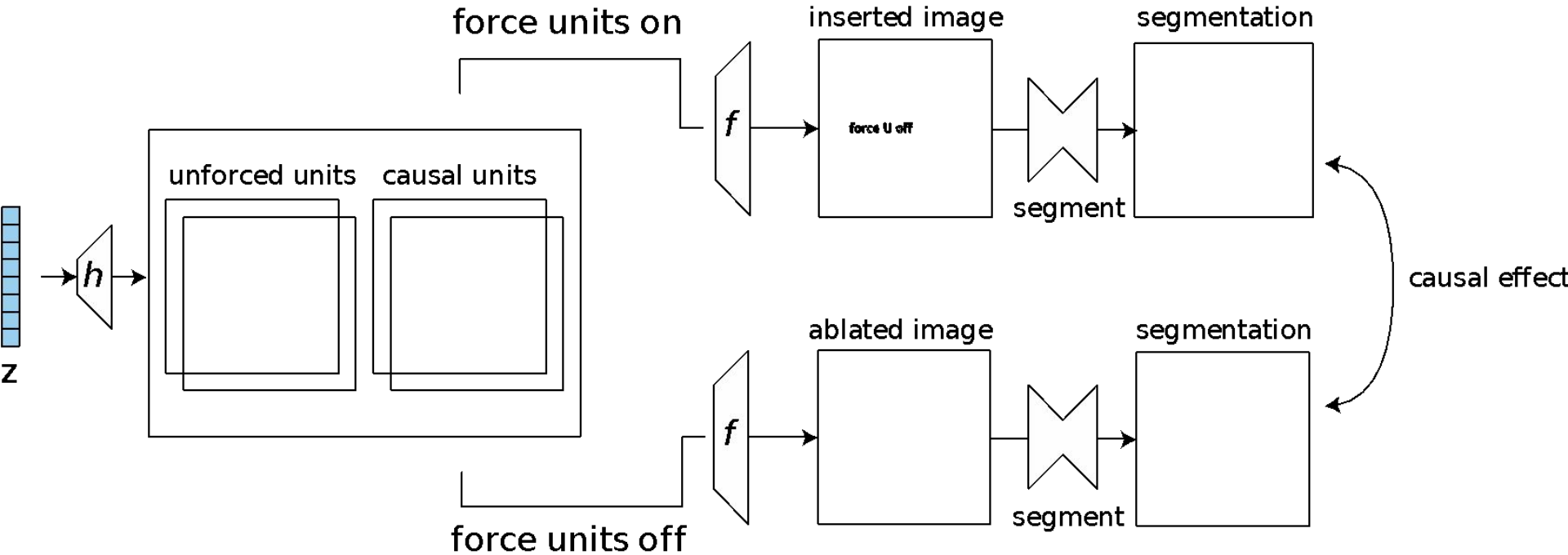
WHICH UNITS CAUSE AN OBJECT CLASS?



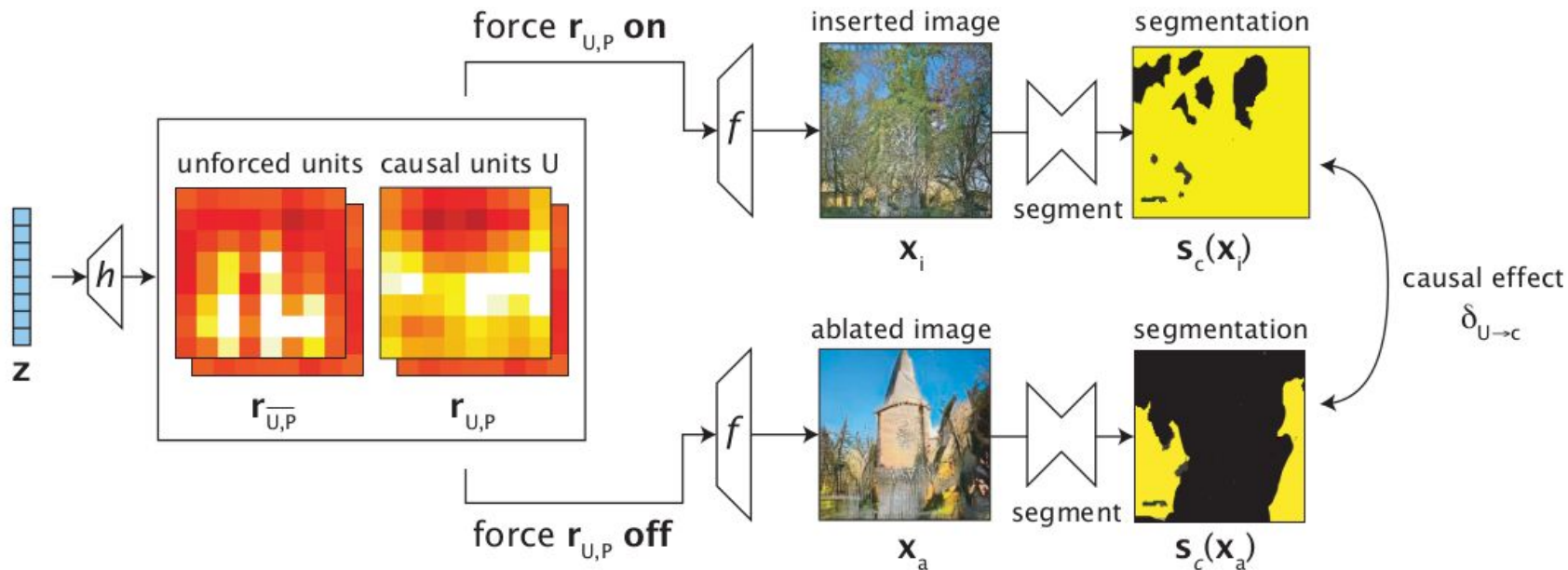
WHICH UNITS CAUSE AN OBJECT CLASS?



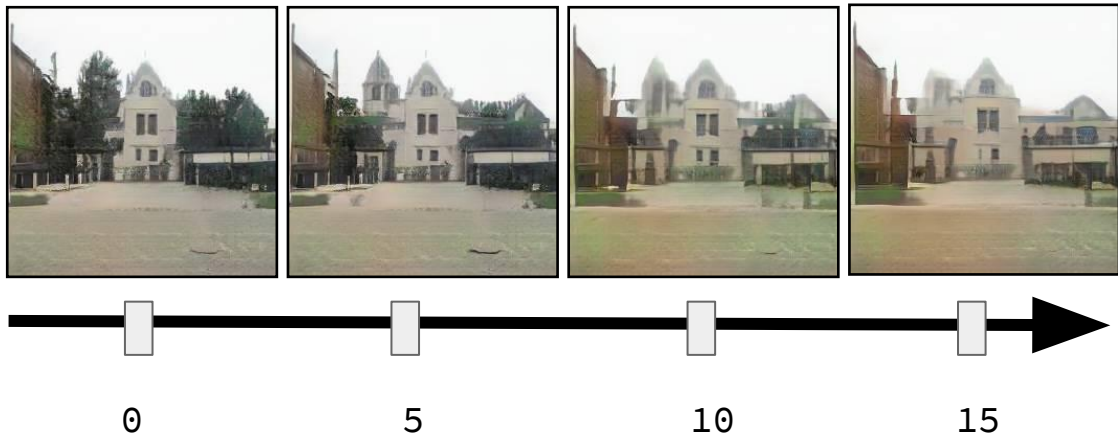
WHICH UNITS CAUSE AN OBJECT CLASS?



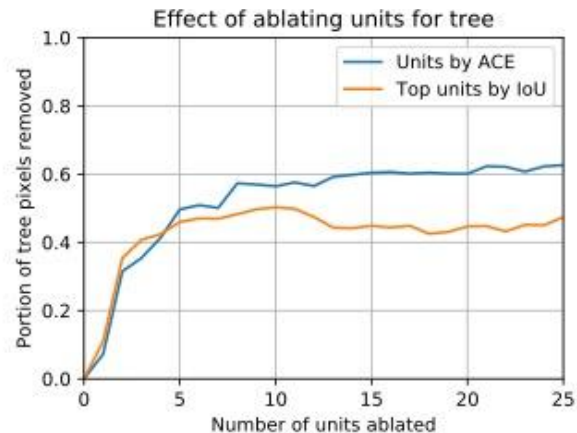
MEASURING THE RELATIONSHIP BETWEEN REPRESENTATION UNITS AND TREES IN THE OUTPUT USING INTERVENTION



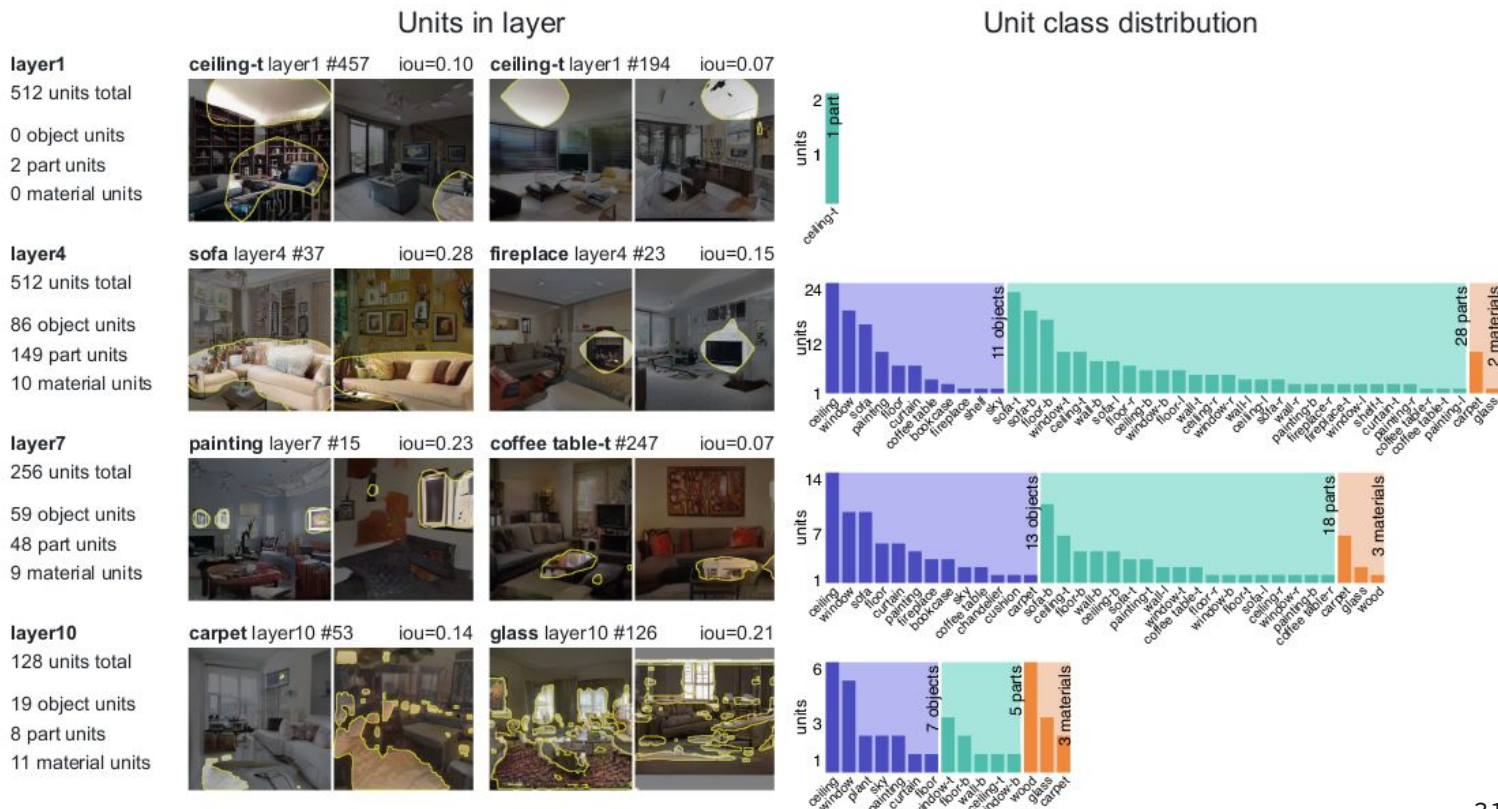
REMOVING OR ADDING UNITS

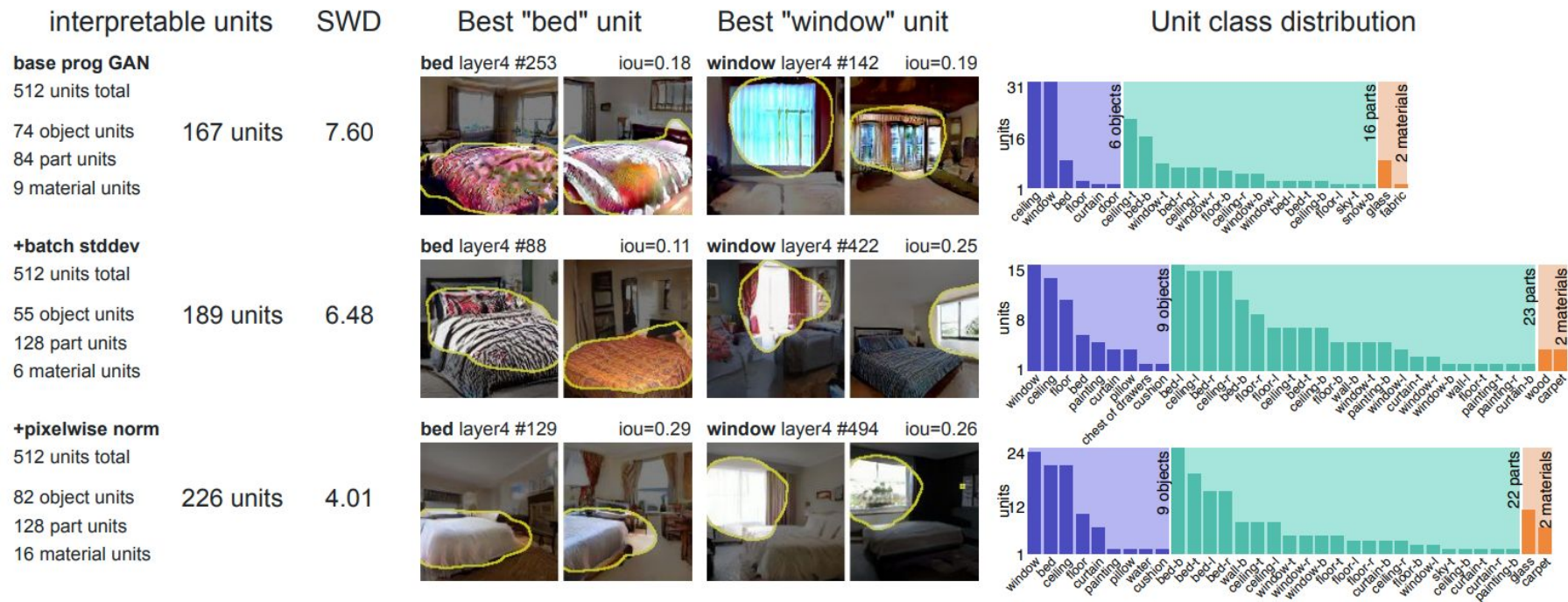


Number of tree units ablated



GAN DISSECTION: COMPARING DATASETS , DATASET USED: LSUN

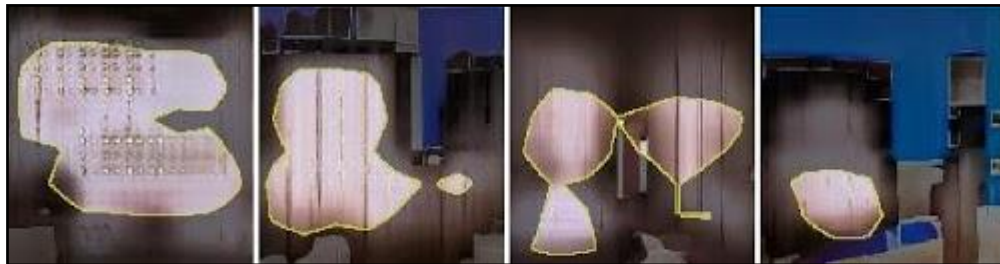




Comparing layer4 representations learned by different training variations. Sliced Wasserstein Distance (SWD) is a GAN quality metric suggested by Karras et al. (2018): lower SWD indicates more realistic image statistics

DEBUGGING AND IMPROVING GANS

Unit #63



Bedroom images with artifacts

Unit #231



Example artifact-causing
units



Ablating “artifact” units improves
results

OBJECT-SCENE RELATIONSHIP



ablate person



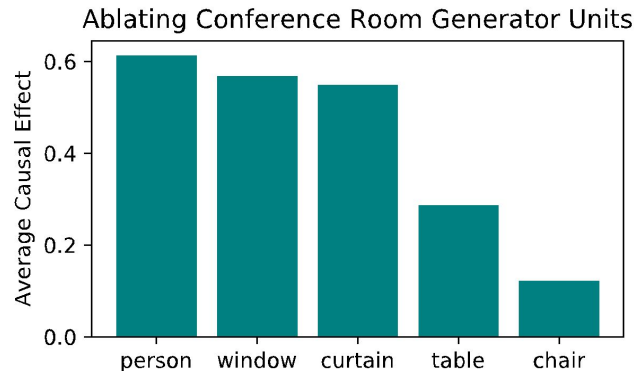
ablate curtain



ablate window

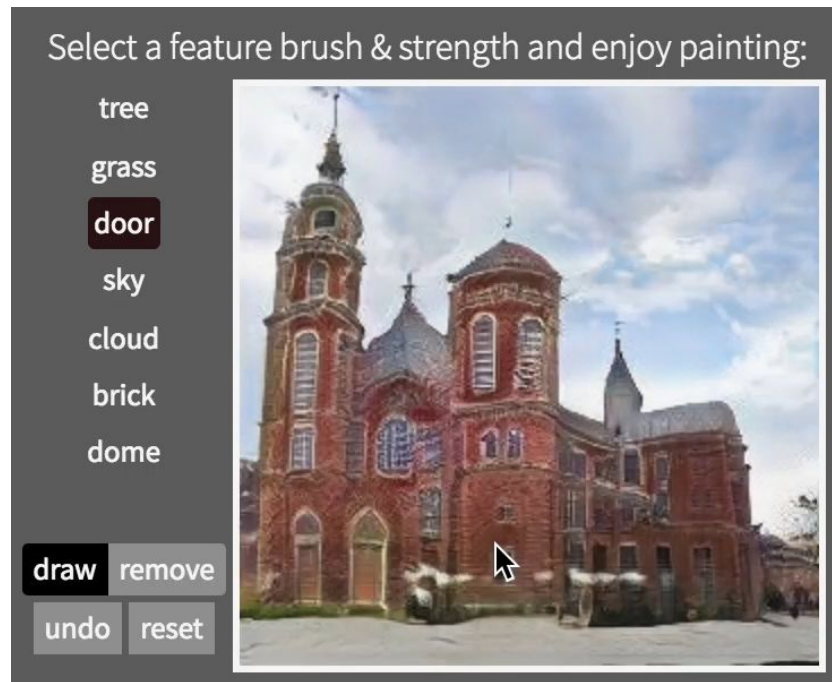


ablate table
units

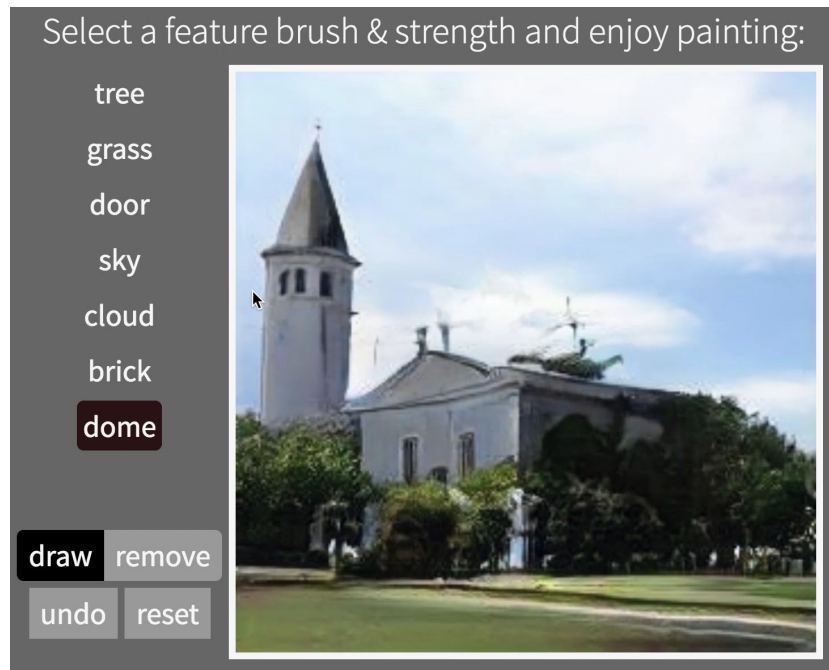


ablate chair
units

PAINT WITH GANS



PAINT WITH GANS




PAINT WITH GANS

Select a feature brush & strength and enjoy painting:

- tree
- grass
- door**
- sky
- cloud
- brick
- dome

draw remove

undo reset



THANK YOU!



Credits:

https://wallpapersprinted.com/wallpaper/2/minions_14.html