

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)
Dept. of Computer Science and Engineering,
Indian Institute of Technology, Jodhpur, India

Conference paper at ICLR 2019

THE PROBLEM

- Ever Imagined How GAN can give so realistic FAKE Images?.....
- To render a beautiful scene, What does a GAN need to know?
- And SOMETIMES....What causes the mistakes?

THE GOAL

GOAL: To analyze how objects such as trees **are encoded** by the **internal representations** of a GAN generator

$G: z \rightarrow x.$

Restaurant



Living room



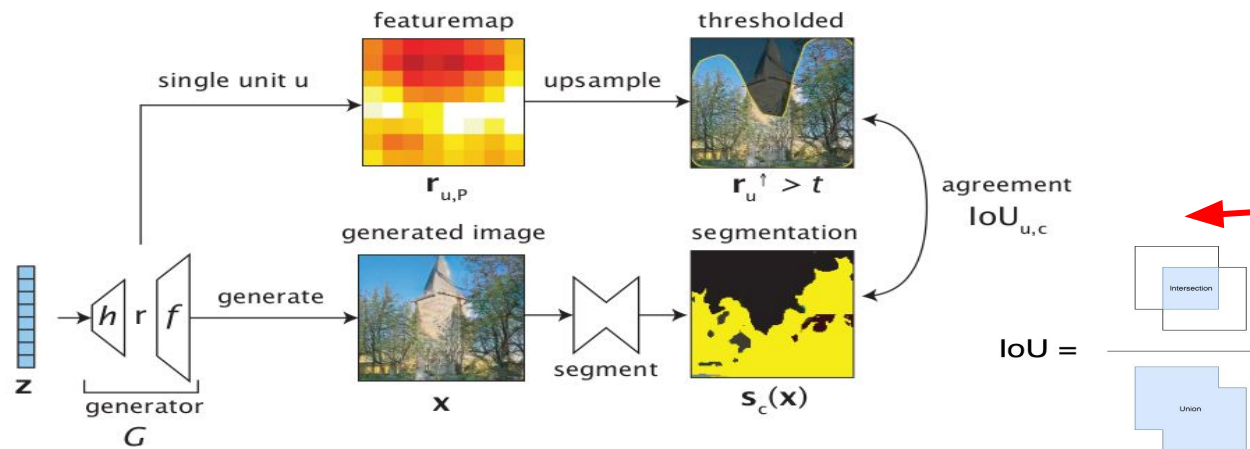
Church



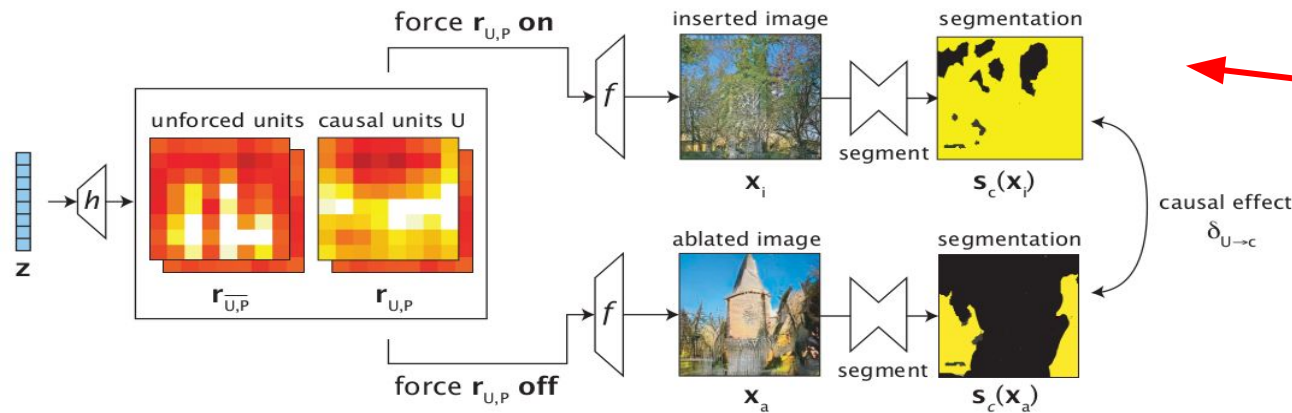
256x256 images synthesized by a Progressive GAN [Karras, et al 2017]

Bedroom





How
units correlate to an
object class?



Measuring the
relationship between
representation units and
trees in the output using
INTERVENTION

RESULTS

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



Figure 1 displays four bar charts showing the unit class distribution for different models. The y-axis represents the number of units (log scale). The x-axis lists various classes. The charts are color-coded: purple for 'objects', teal for 'parts', and orange for 'materials'.

- Top Chart (Model 1):** Shows a distribution with 1 part (ceiling-t) and 11 objects. The 'objects' class is the most frequent, with approximately 24 units.
- Second Chart (Model 2):** Shows a distribution with 13 objects, 18 parts, and 3 materials. The 'objects' class is the most frequent, with approximately 14 units.
- Third Chart (Model 3):** Shows a distribution with 7 objects, 5 parts, and 3 materials. The 'objects' class is the most frequent, with approximately 6 units.
- Bottom Chart (Model 4):** Shows a distribution with 1 part (ceiling-t) and 11 objects. The 'objects' class is the most frequent, with approximately 24 units.

