

# Understanding Invariance via Feedforward Inversion of Discriminatively Trained Classifiers

Piotr Teterwak<sup>1</sup>, Chiyuan Zhang<sup>2</sup>, Dilip Krishnan<sup>2</sup>,  
Michael C. Mozer<sup>2</sup>



1 Boston University



2 Google Research

# What is preserved in the logit (pre-softmax) outputs?

- A trained neural network can filter all information from the input other than class label

# What is preserved in the logit (pre-softmax) outputs?

- A trained neural network can filter all information from the input other than class label
- Prior work has shown that *some* information can be preserved in the logits, the most successful of which is Dosovitskiy and Brox(2016, NeurIPS)

# What is preserved in the logit (pre-softmax) outputs?

- A trained neural network can filter all information from the input other than class label
- Prior work has shown that *some* information can be preserved in the logits, the most successful of which is Dosovitskiy and Brox(2016, NeurIPS)

Original



Reconstruction



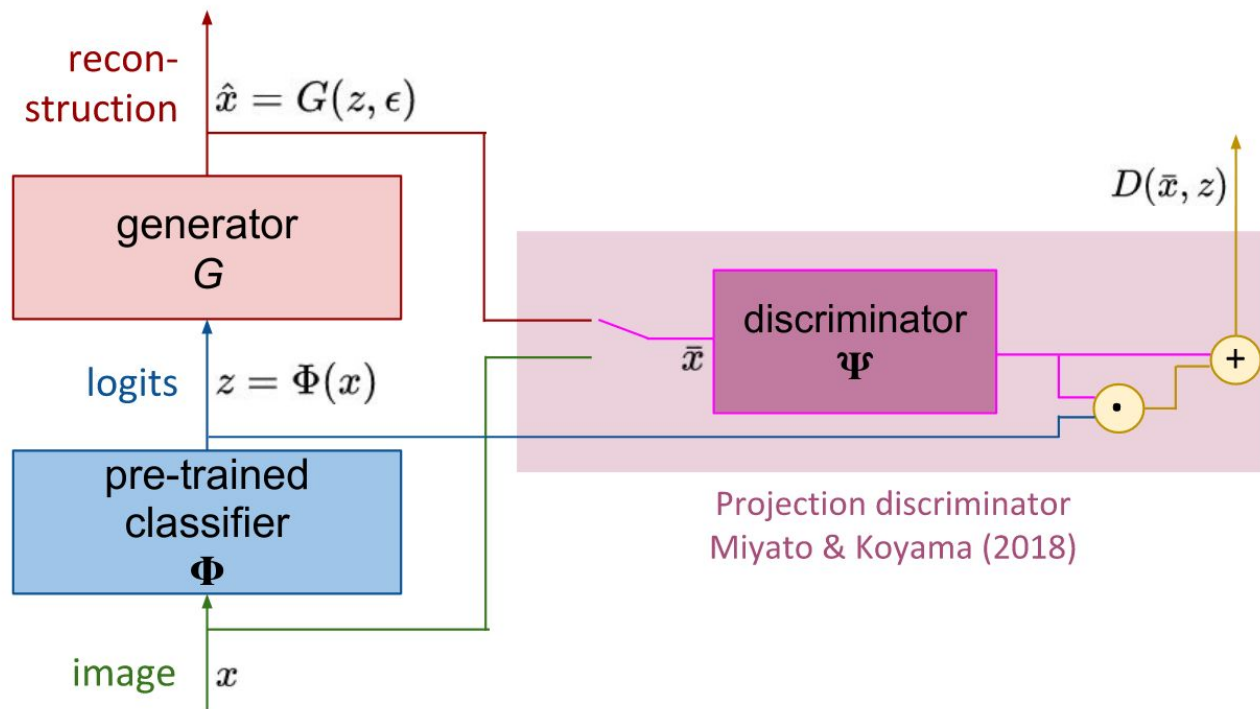
# What is preserved in (pre-softmax) outputs of neural nets?

- The field was overdue to revisit this question with:
  - Modern classifiers
  - Better methods for Image synthesis

# What is preserved in (pre-softmax) outputs of neural nets?

- The field was overdue to revisit this question with:
  - Modern classifiers
  - Better methods for Image synthesis
- We show:
  - Surprising fidelity in reconstruction
  - Information preserved in logits depends on optimization method and architecture
  - Robust models contain information in the logits which supports better reconstructions of both shapes and textures relative to standard models

# Method



# Results

original



original



original

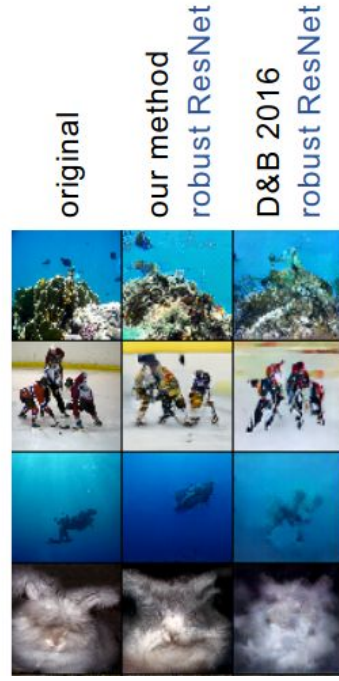




# Results



# Results

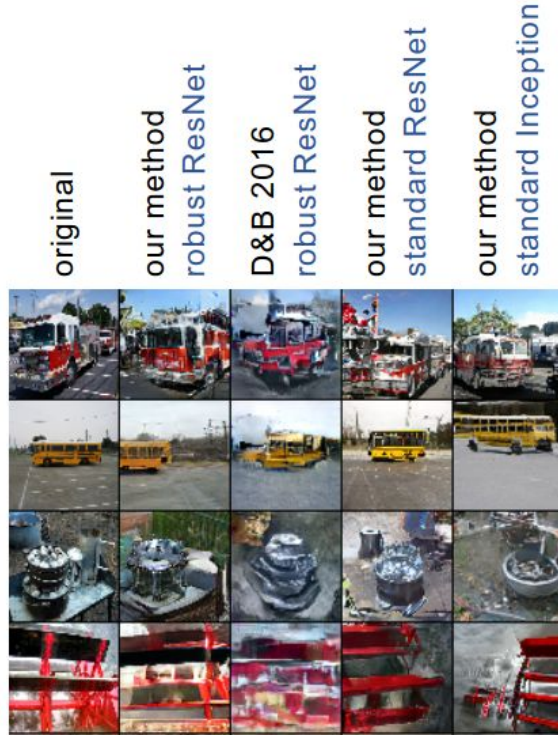


# Results





# Results



# Invariances: Constant Logit, Varying GAN Noise Input



## And More!

- Manipulating logits with scales/shifts/perturbations
- Logit interpolations between images
- Reconstructions of adversarially attacked images
- Reconstructions of incorrectly classified images
- Reconstructions of OOD data
- Find the paper here: <https://arxiv.org/abs/2103.07470>

Thank you!