

Analyzing Artifacts in Discriminative and Generative Models

Richard Zhang (章睿嘉)

Research Scientist, Adobe SF
GAMES Webinar

Aug 2020



Example classifications



86.7



P(correct class)



69.2

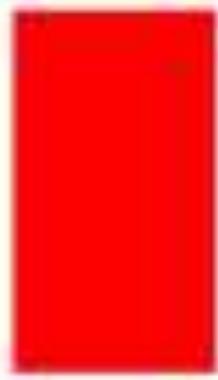


P(correct class)

Deep Networks are not Shift-Invariant



46.3



P(correct class)



18.0

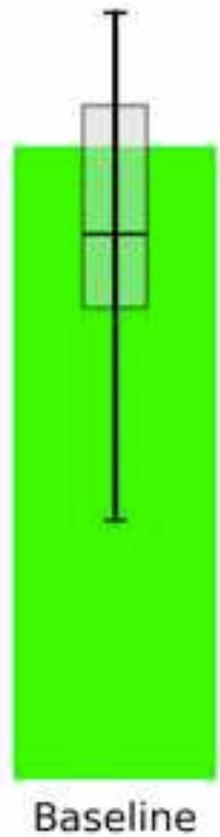


P(correct class)

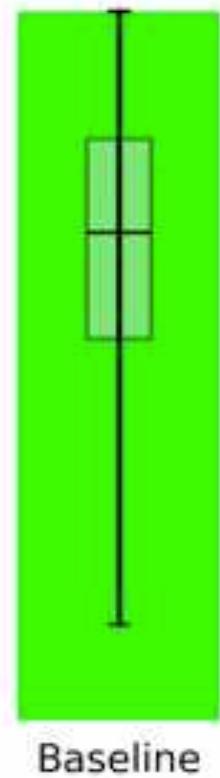
Deep Networks are not Shift-Invariant



75.5



84.5



c.f. Azulay and Weiss. Why do deep convolutional networks generalize so poorly to small image transformations? Arxiv, 2018. JMLR, 2019.
Engstrom, Tsipras, Schmidt, Madry. Exploring the Landscape of Spatial Robustness. Arxiv, 2017. ICML, 2019.

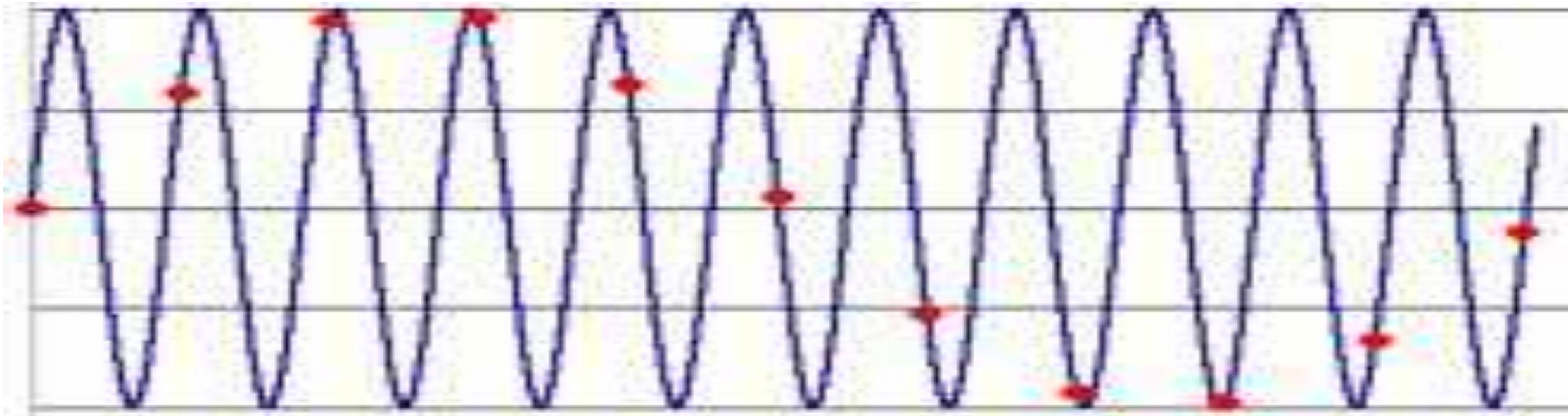
Why is shift-invariance lost?

“Convolutions are **shift-equivariant**”

“Pooling builds up shift-invariance”

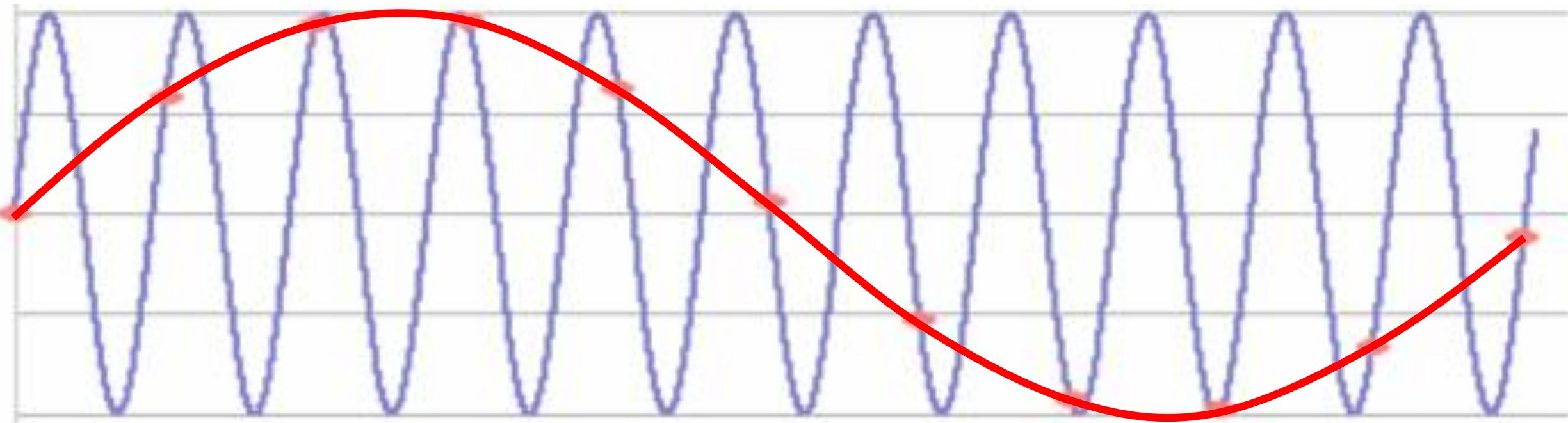
...but striding ignores Nyquist sampling theorem
and **aliases**

Nyquist-Shannon theorem



If the **sampling frequency** is less than twice the **underlying signal frequency**, the samples cannot faithfully represent the underlying signal...

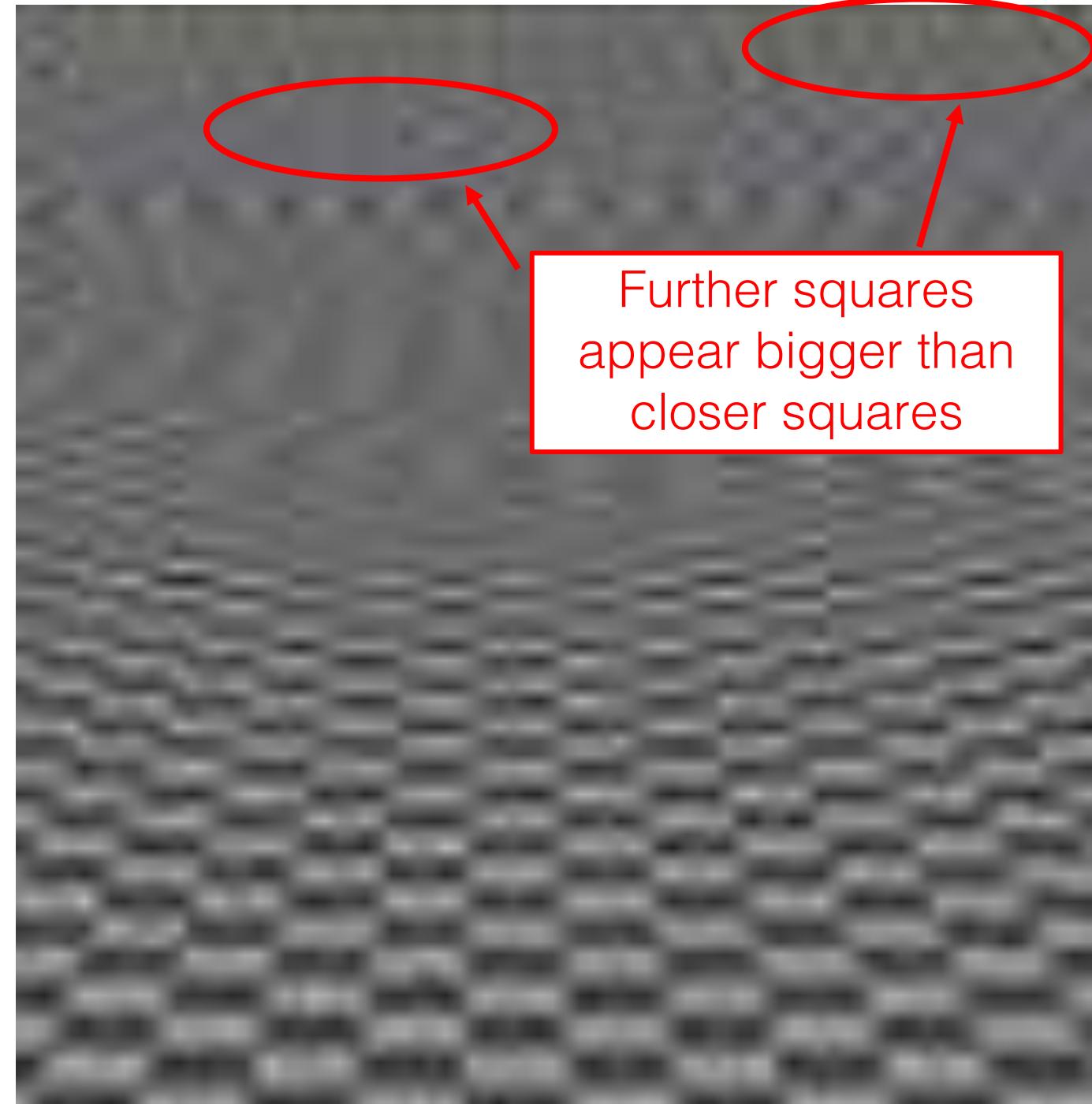
What is aliasing?



...worse, the samples will **misrepresent** the underlying signal

Rendering

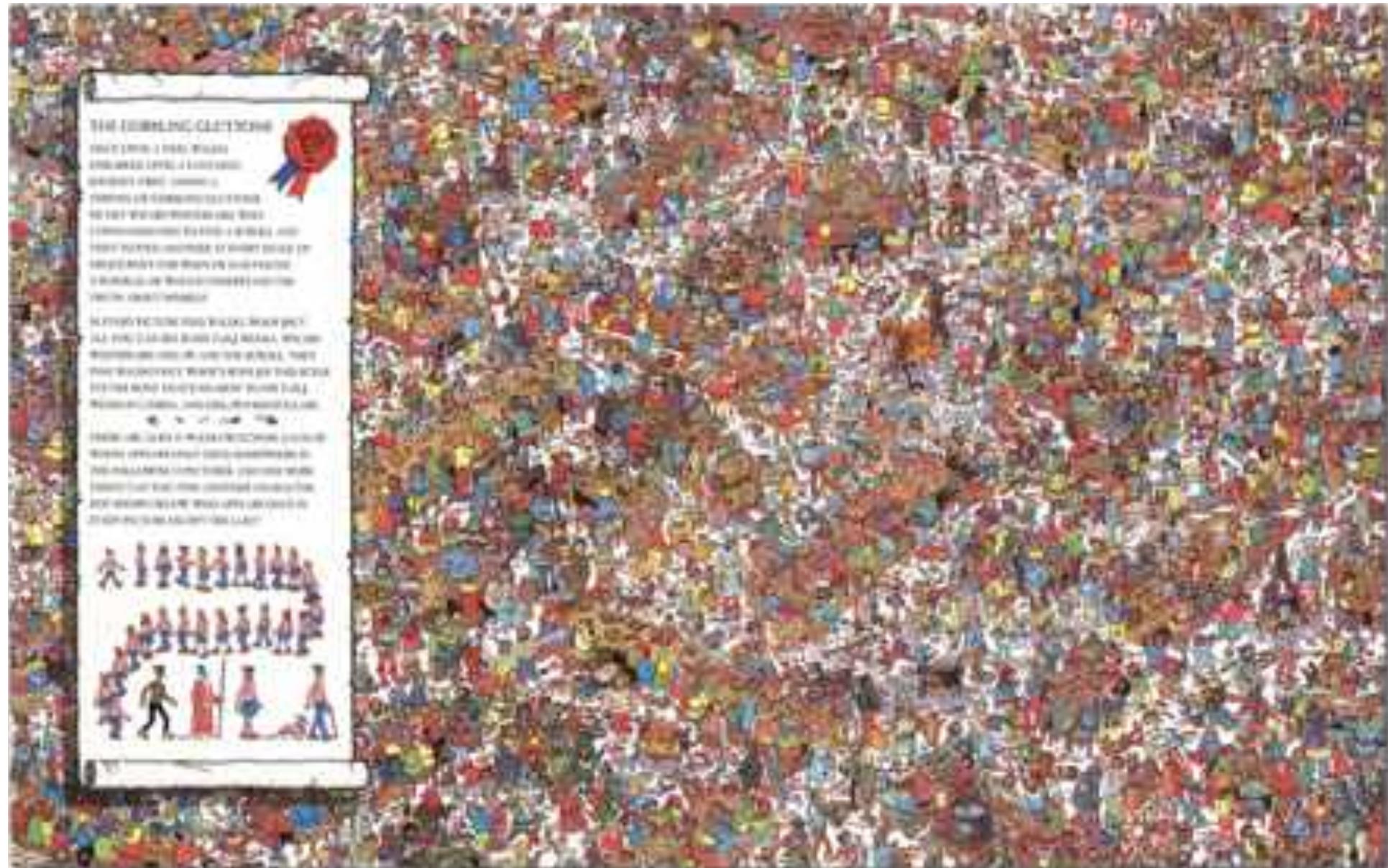
Further away,
decreased
sampling rate



Temporal aliasing

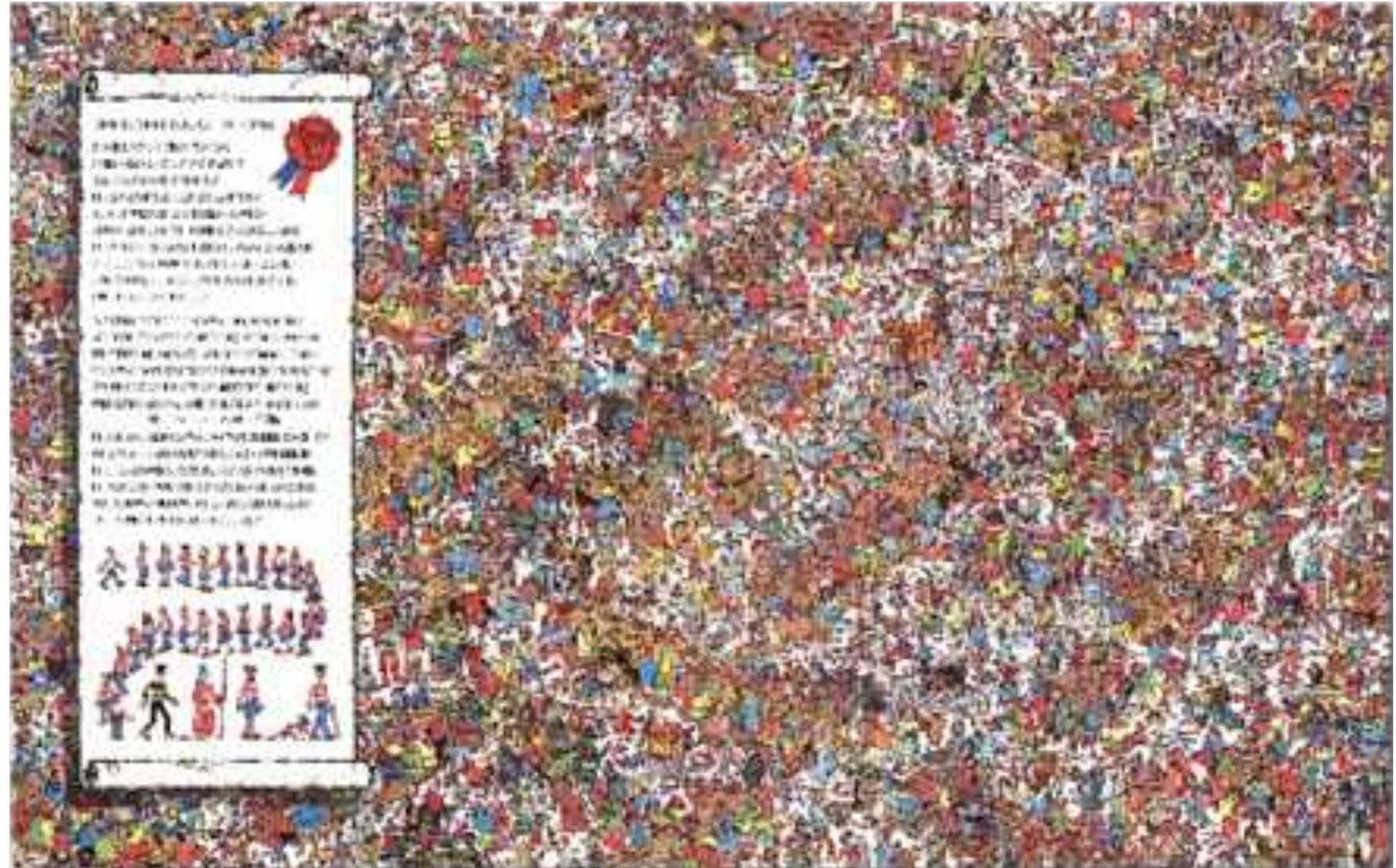


Image



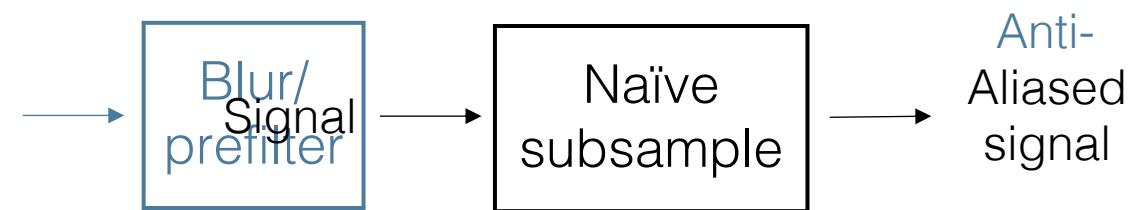
Adapted from Lectures on Digital Photography. Marc Levoy. <https://sites.google.com/site/marcllevylectures/schedule>

4x naive sub- sampling



Antialiasing

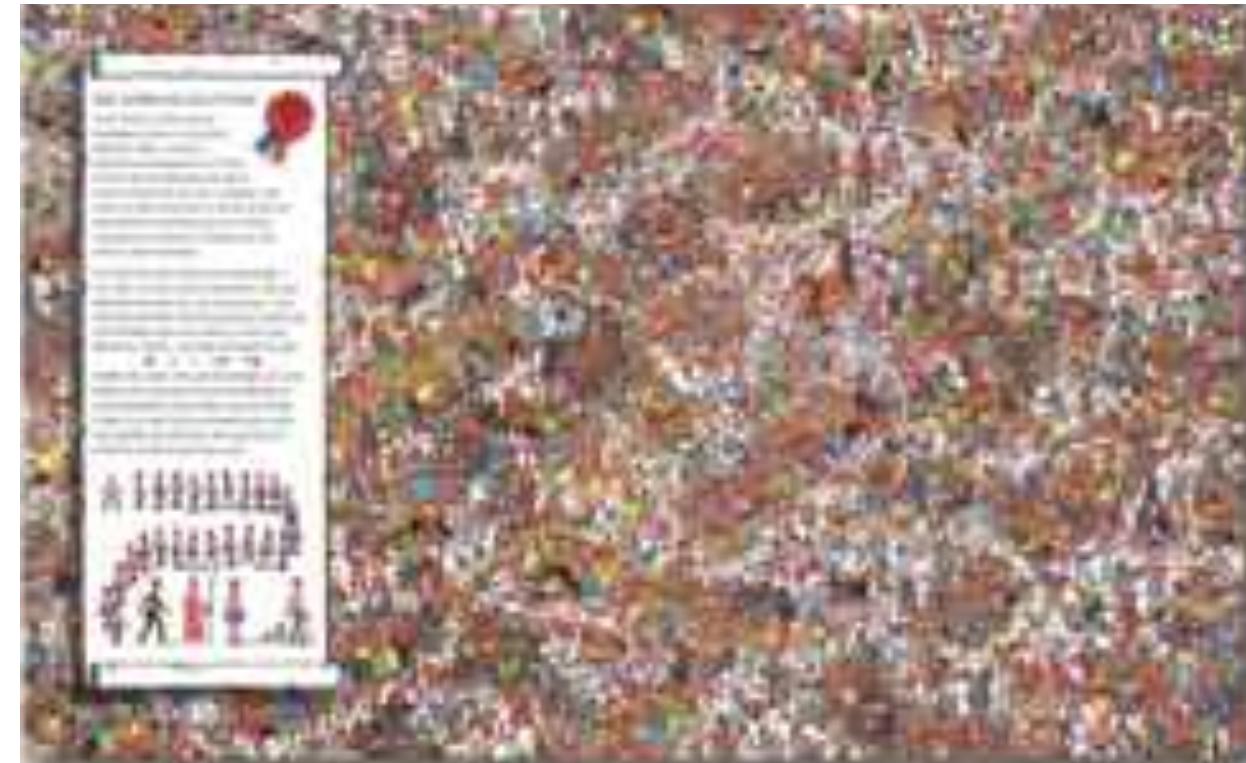
- Low sampling rate → signal cannot be represented
 - Just don't downsample? Expensive...
- Naïve subsampling → high freqs misrepresented
- Before subsampling, antialias by blurring → high freqs unrepresented
 - Better than aliasing!
 - Blur / prefilter / antialias / low-pass filter



Effect of prefiltering

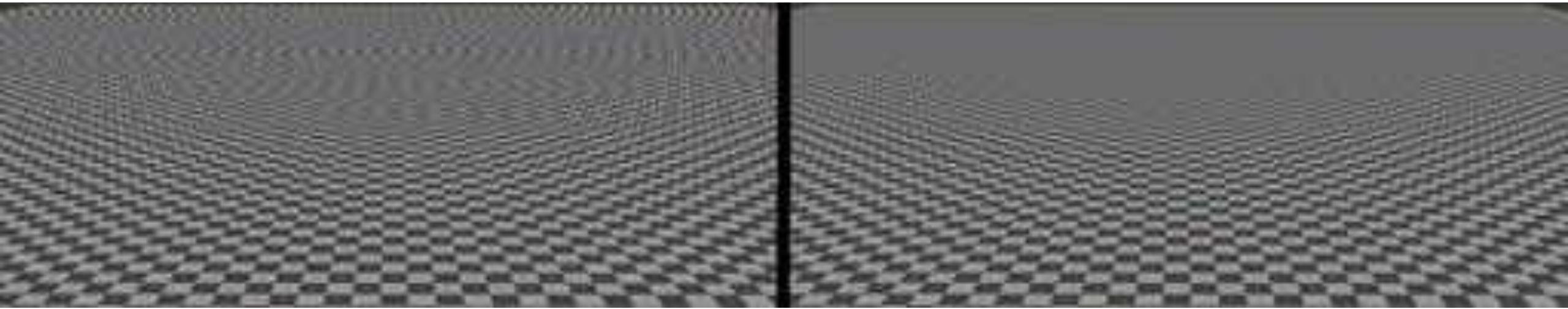


Naïvely subsampled



Prefiltered, then subsampled

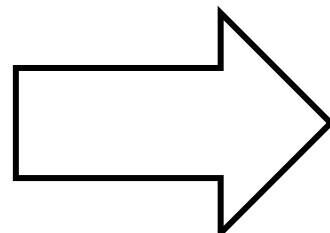
Effect of prefiltering



Naïvely subsampled

Prefiltered, then subsampled

Aliasing → Loss of shift-equivariance



Naive

Prefiltered

Shift to the right

Making Convolutional Networks Shift-Invariant Again

Richard Zhang

In ICML, 2019

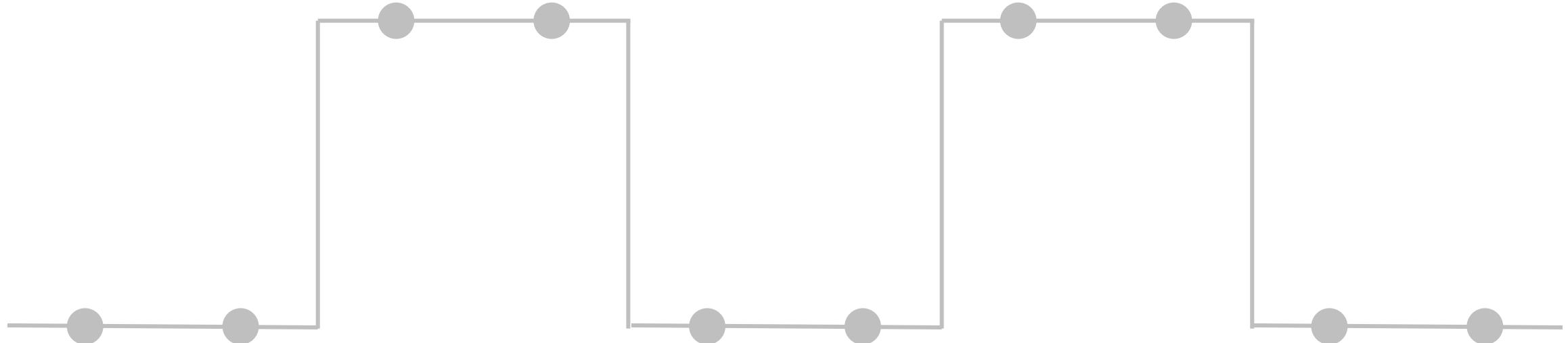


Alternative downsampling methods

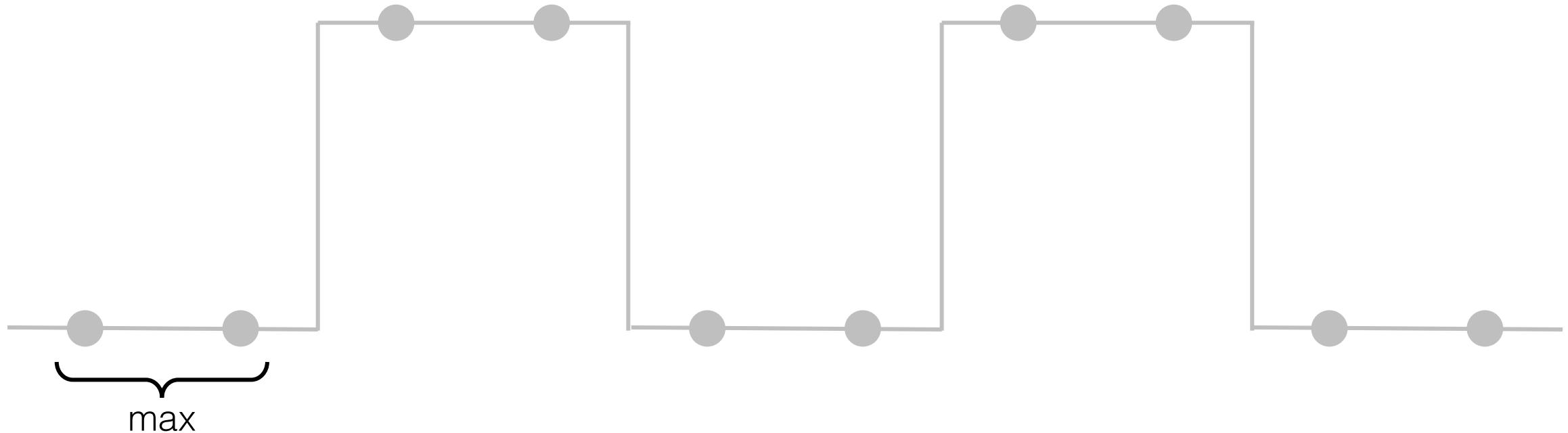
- Signal processing; image processing; graphics
 - Prefilter before subsampling
- Deep learning; computer vision
 - Szeliski vision book, Sec 2.3.1 “Sampling and aliasing”
 - Average pooling [LeNet 1989] does some antialiasing
 - Max-pooling gets better accuracy [Scherer 2010]

Antialiasing abandoned and forgotten

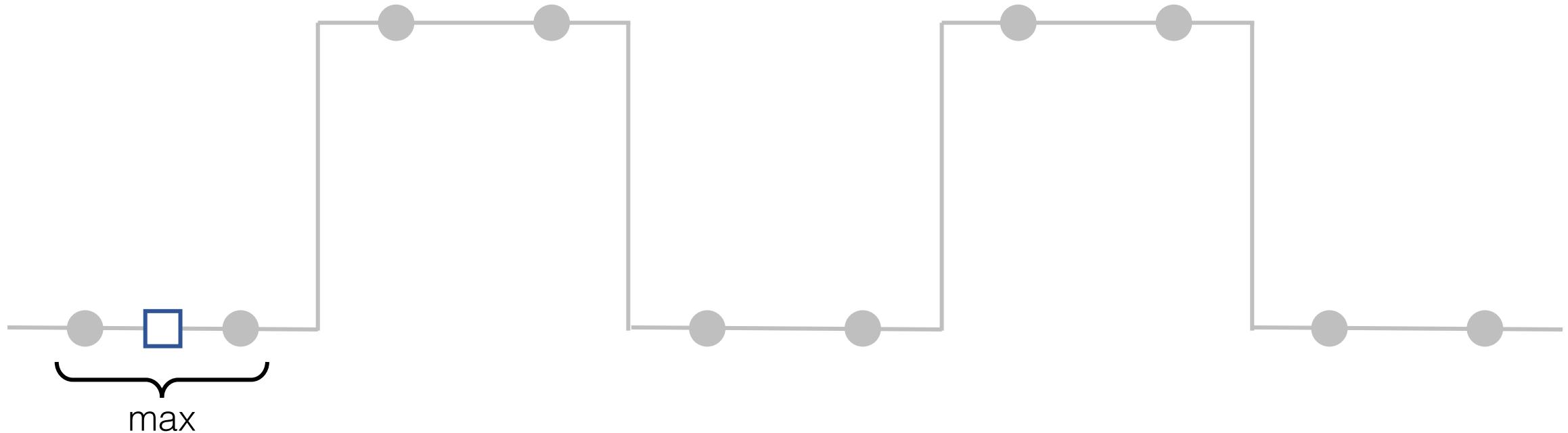
Re-examining Max-Pooling



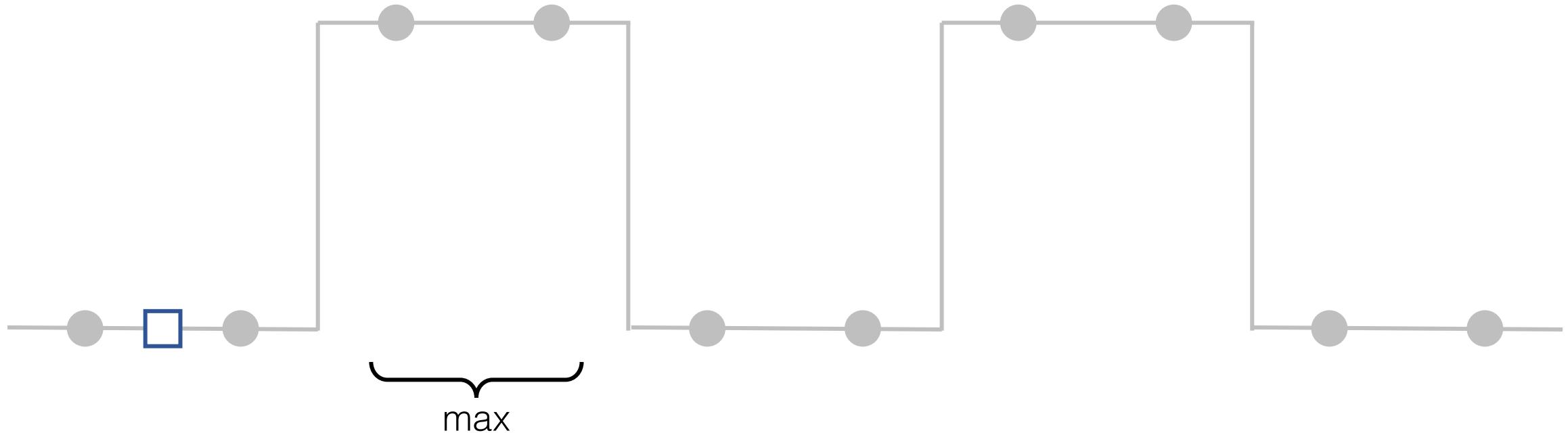
Re-examining Max-Pooling



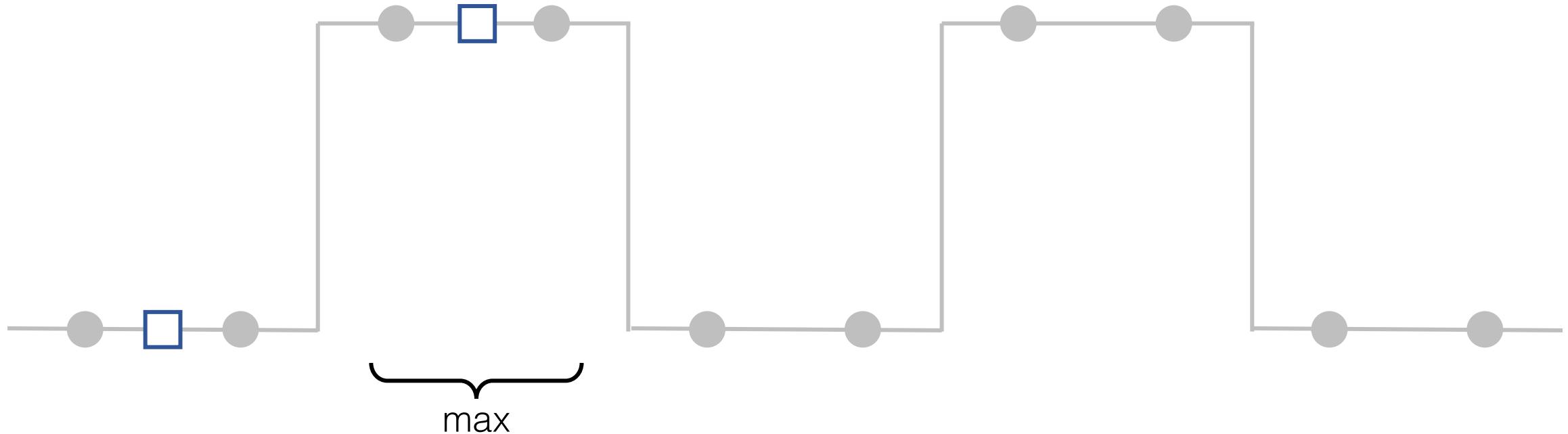
Re-examining Max-Pooling



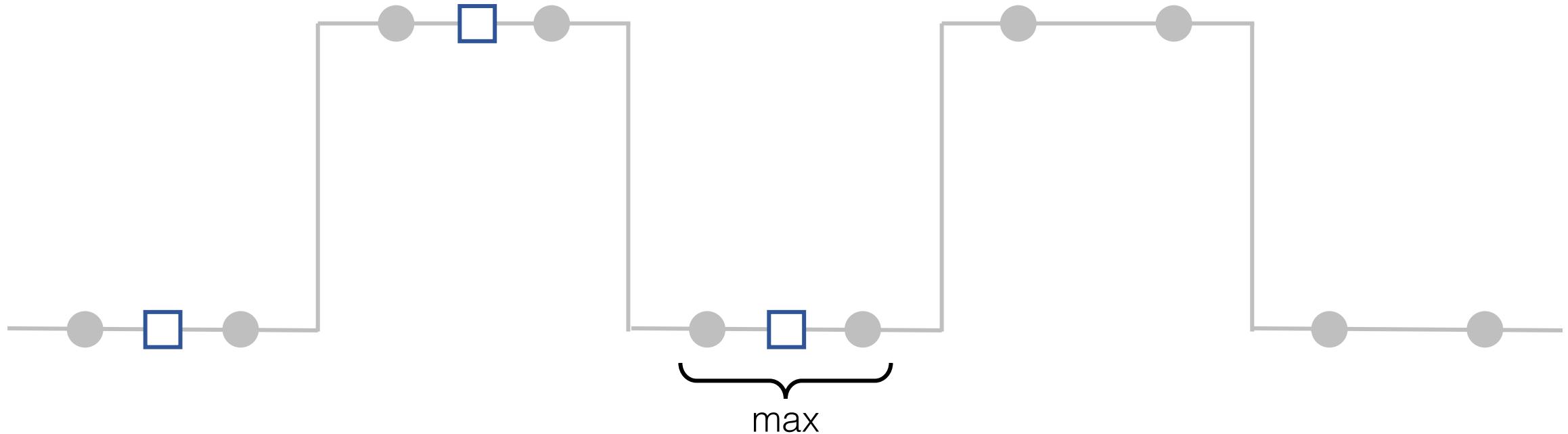
Re-examining Max-Pooling



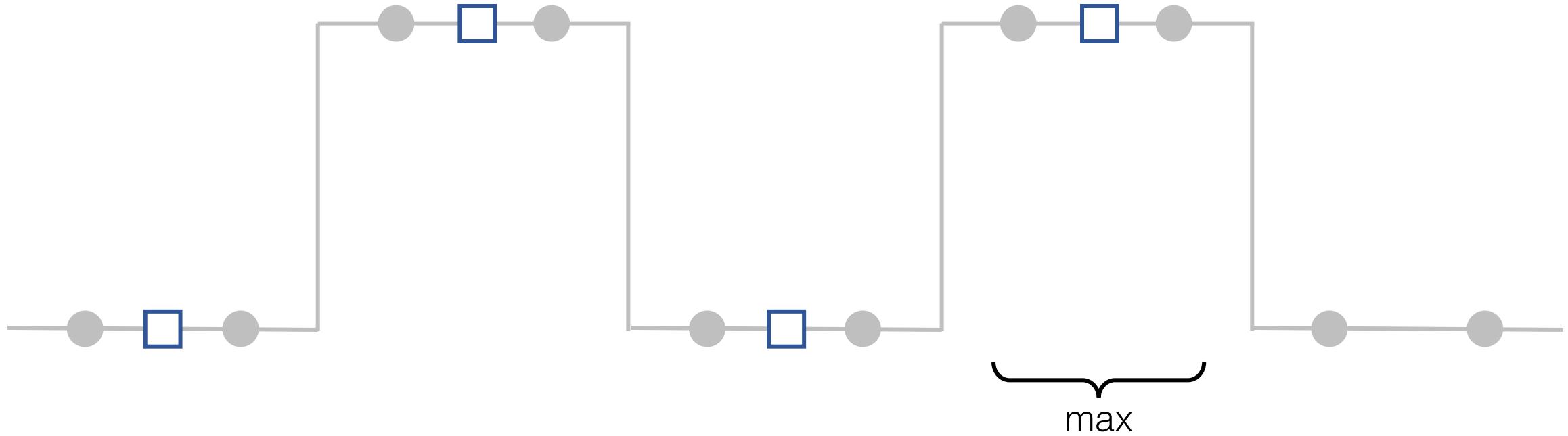
Re-examining Max-Pooling



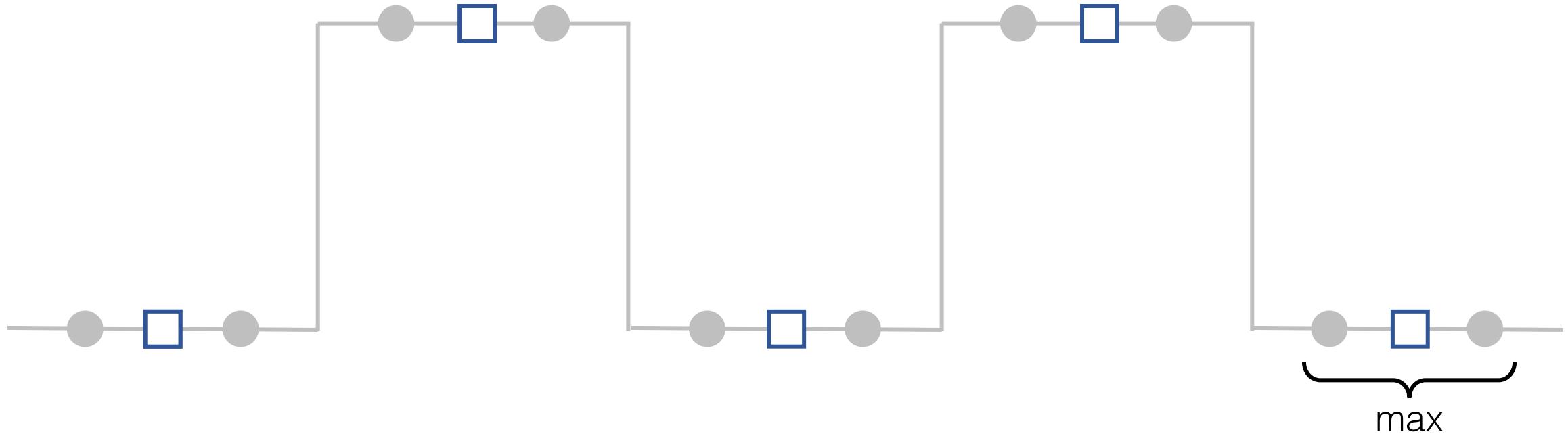
Re-examining Max-Pooling



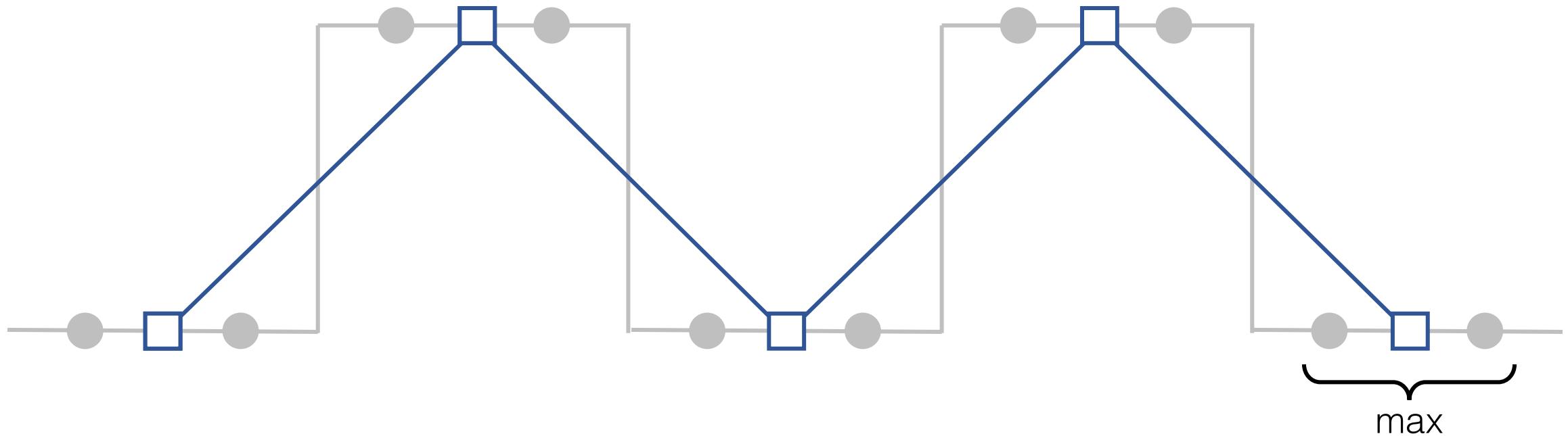
Re-examining Max-Pooling



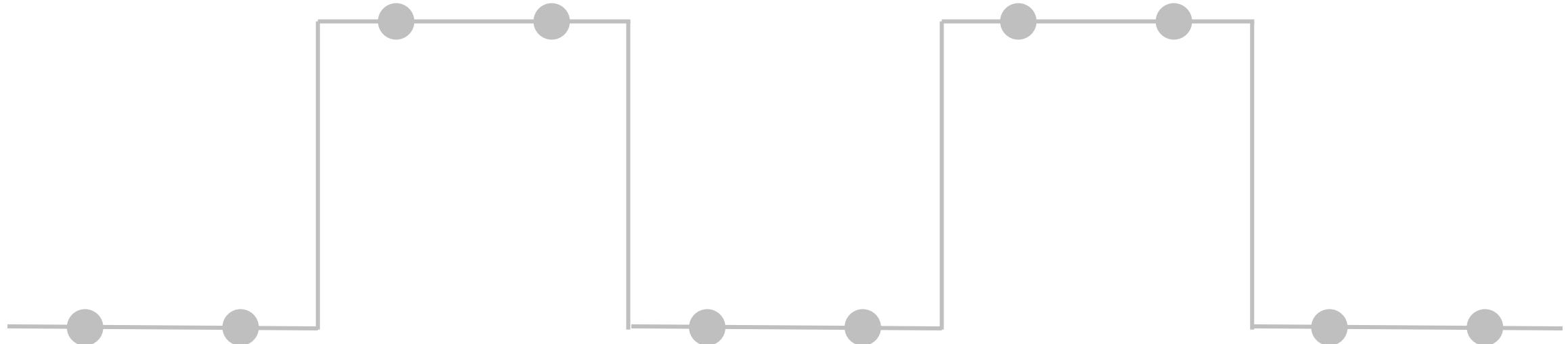
Re-examining Max-Pooling



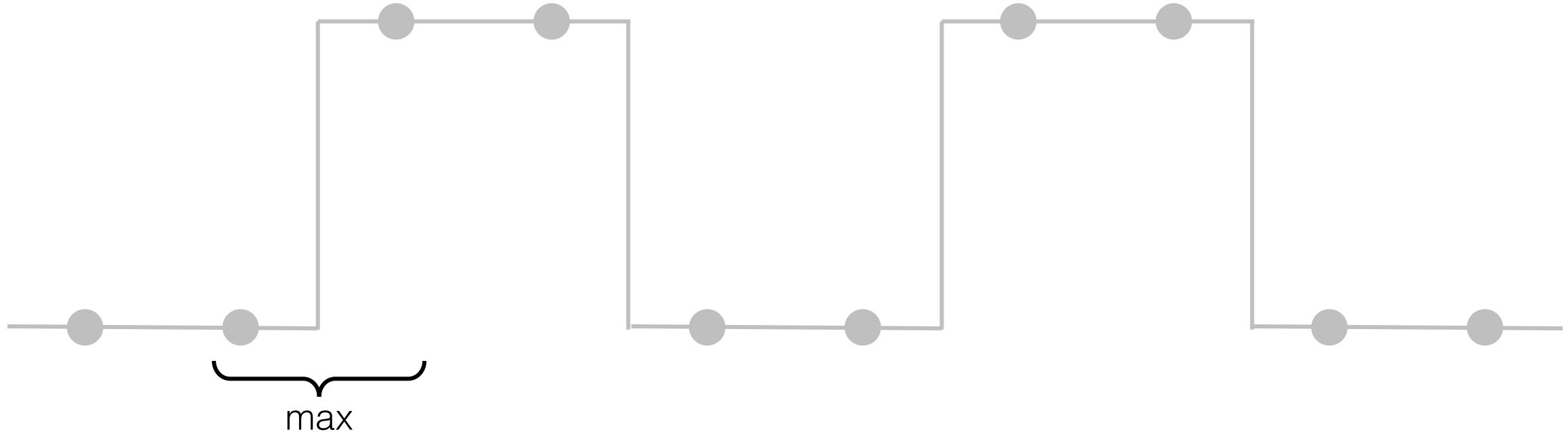
Re-examining Max-Pooling



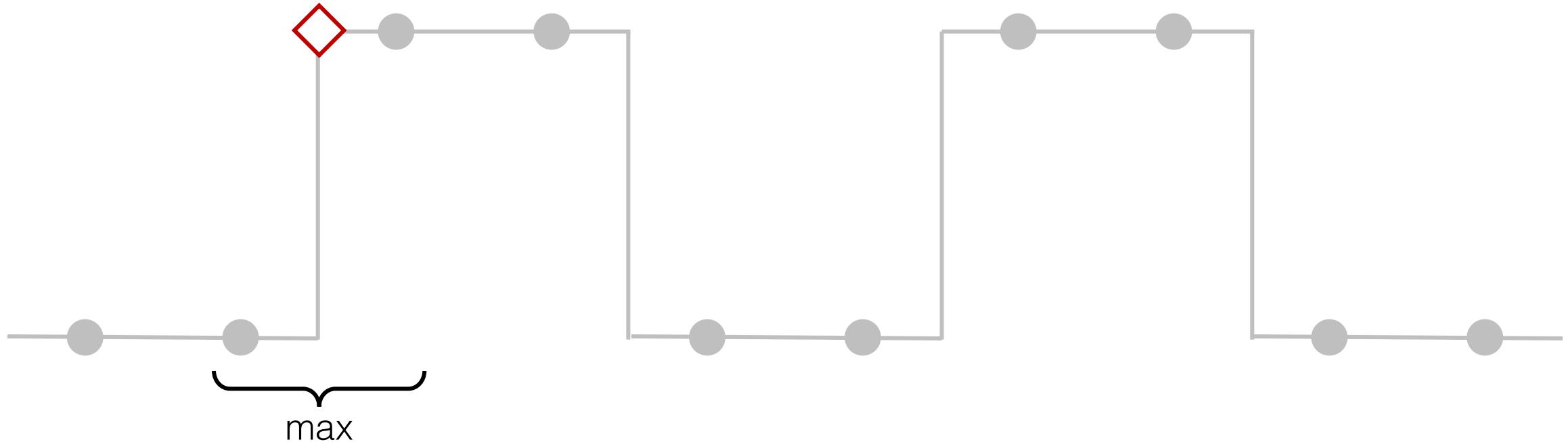
Re-examining Max-Pooling



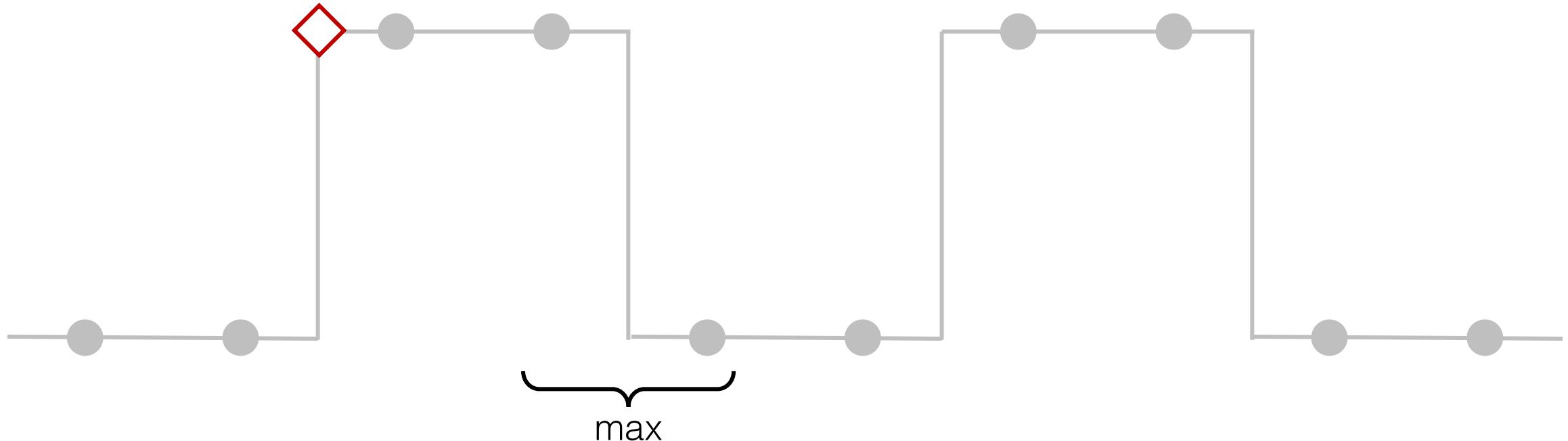
Re-examining Max-Pooling



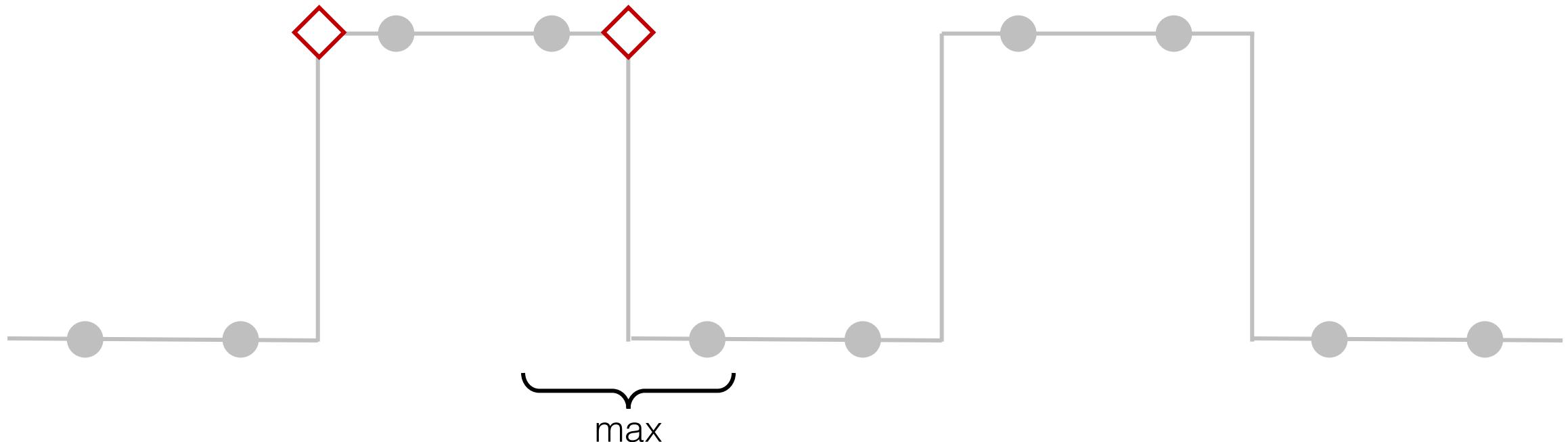
Re-examining Max-Pooling



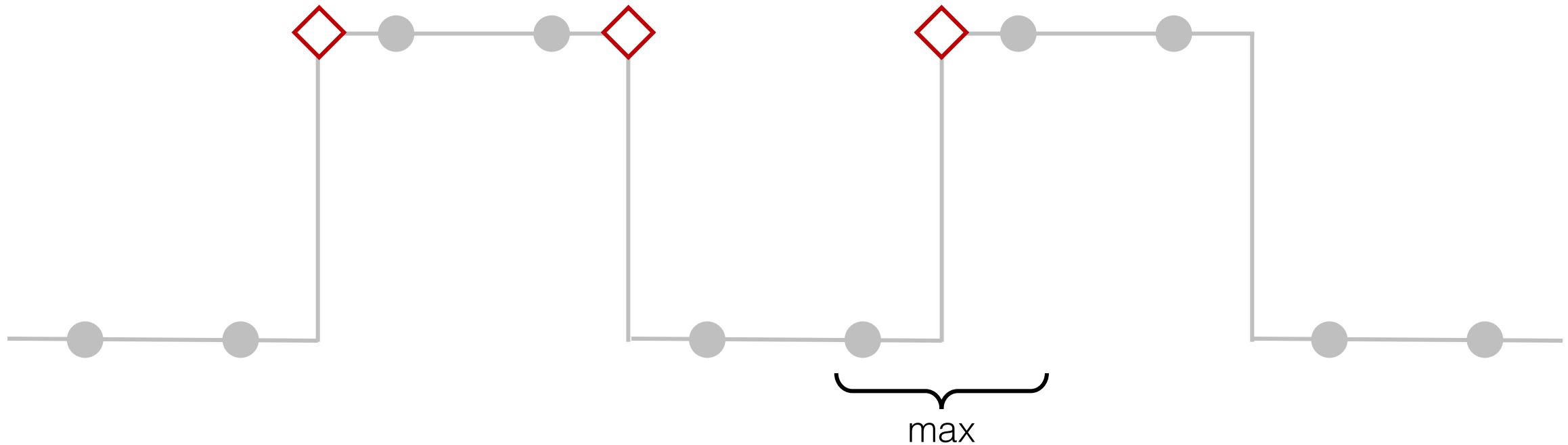
Re-examining Max-Pooling



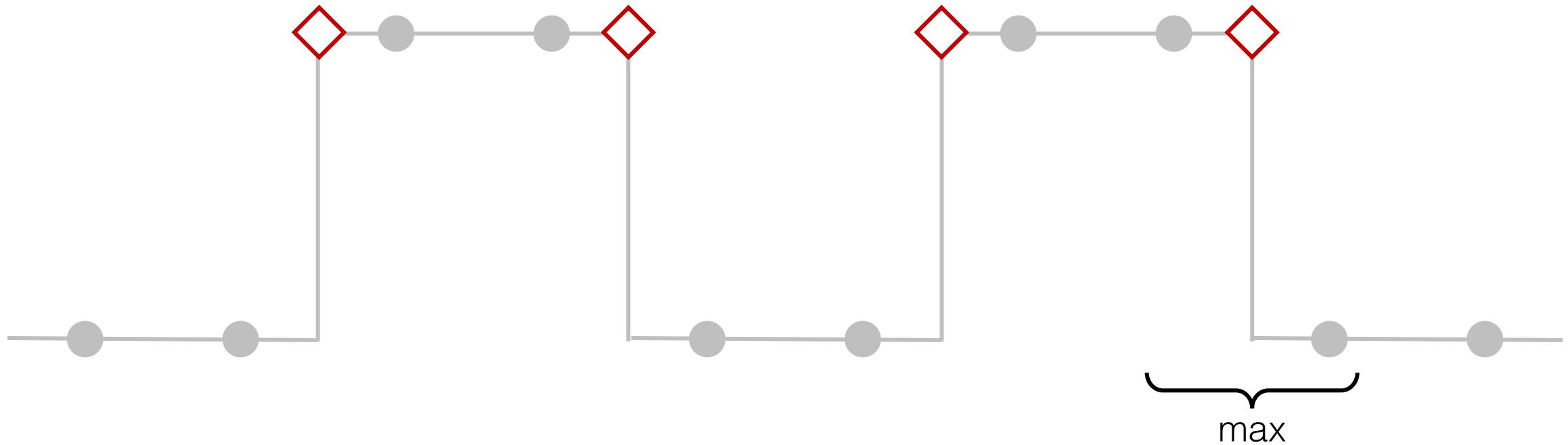
Re-examining Max-Pooling



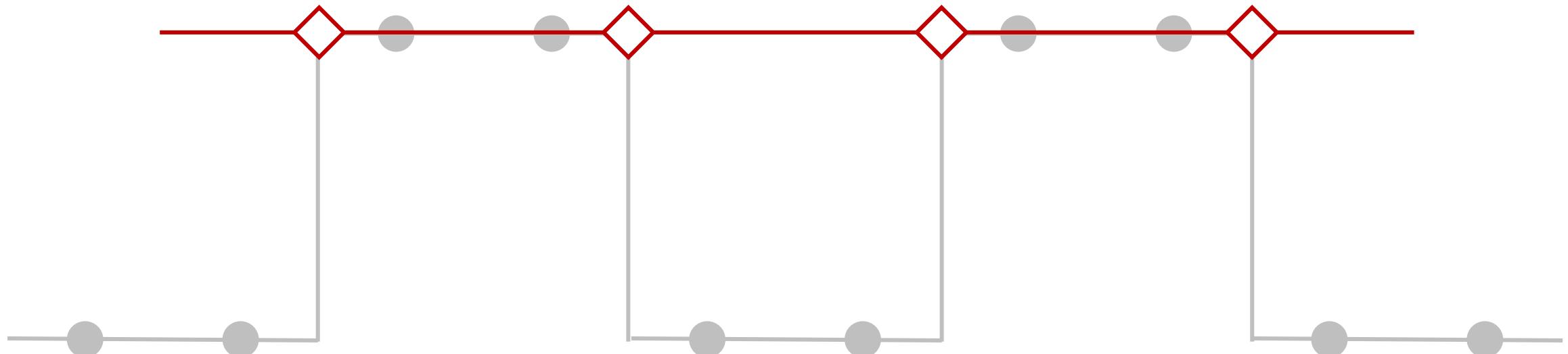
Re-examining Max-Pooling



Re-examining Max-Pooling



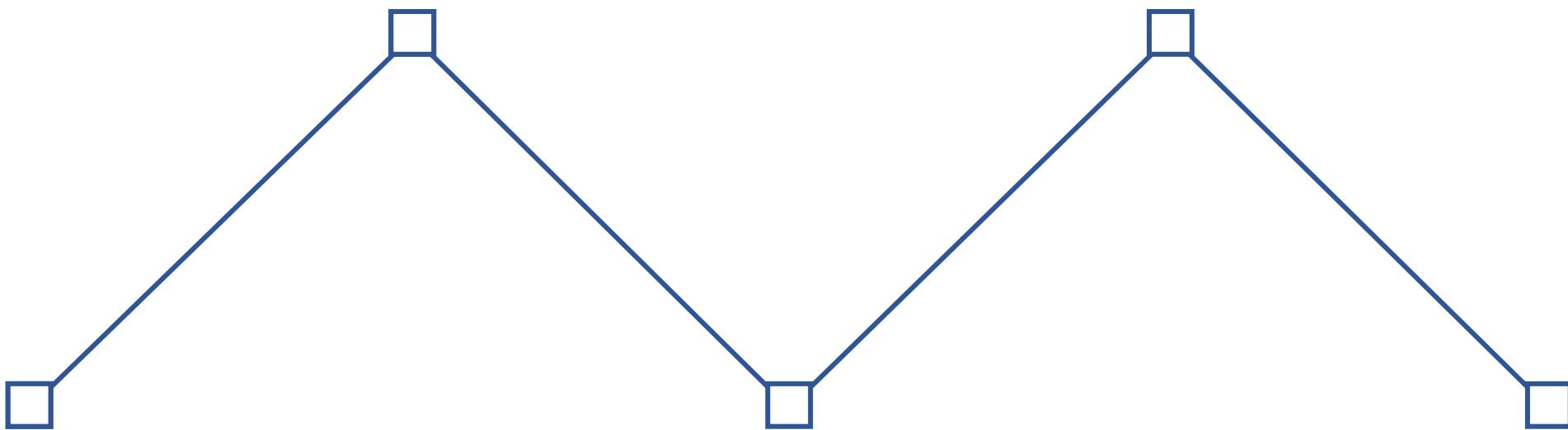
Re-examining Max-Pooling



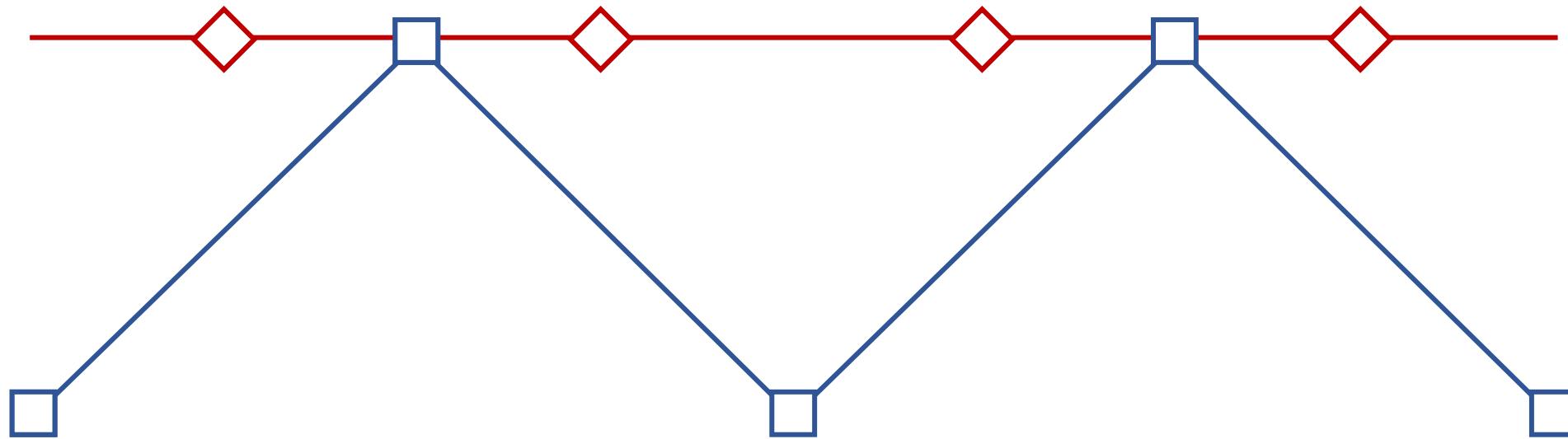
Re-examining Max-Pooling



Re-examining Max-Pooling



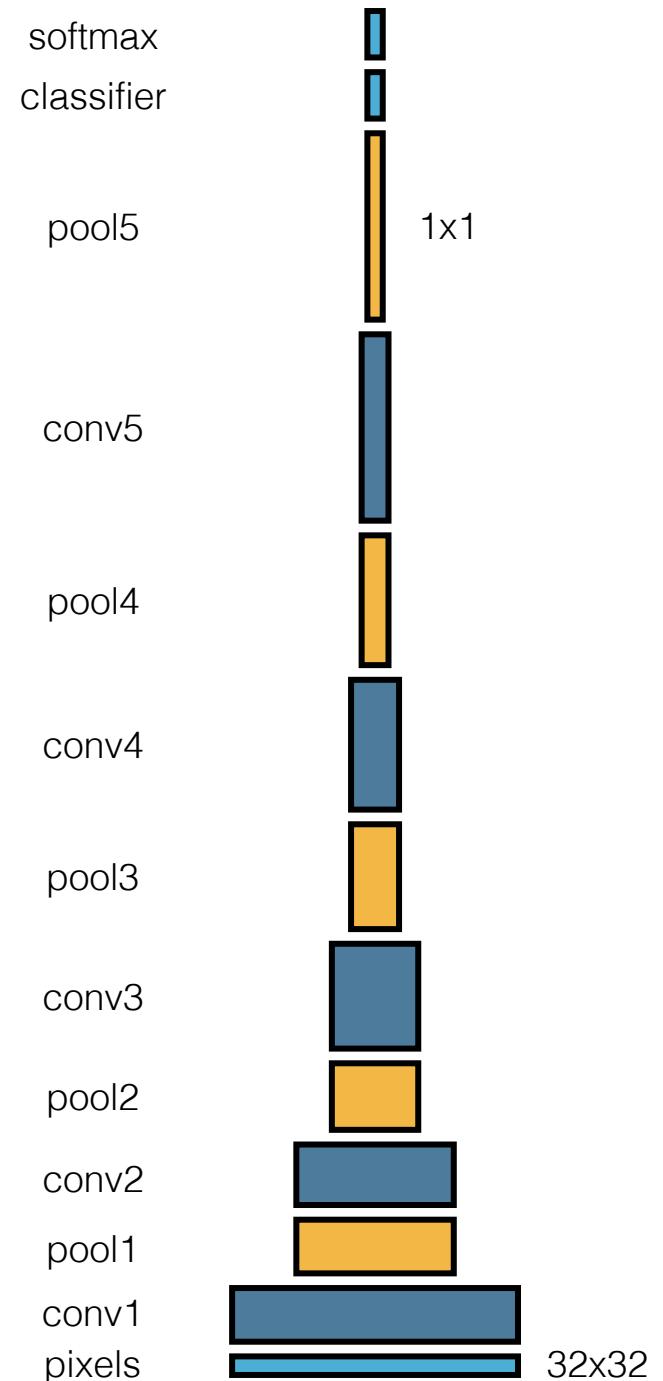
Re-examining Max-Pooling



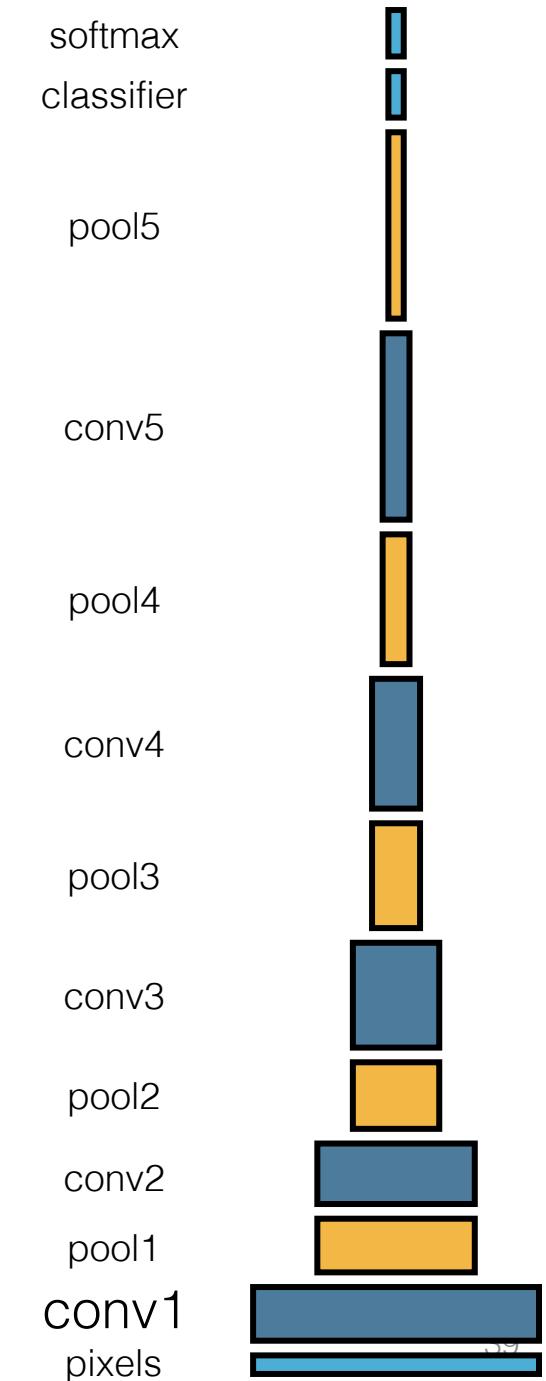
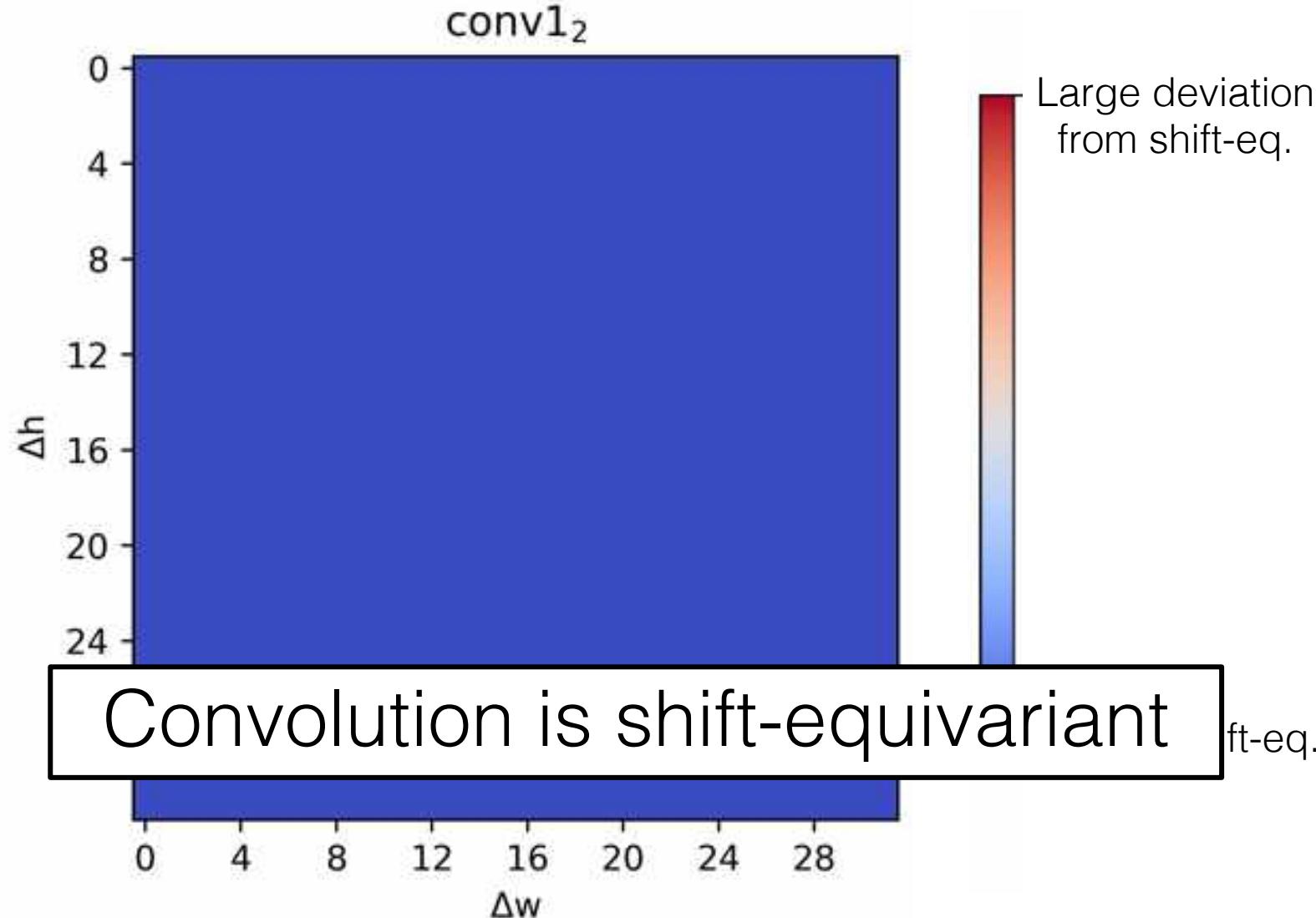
Max-pooling breaks shift-equivariance

Shift-equivariance Testbed

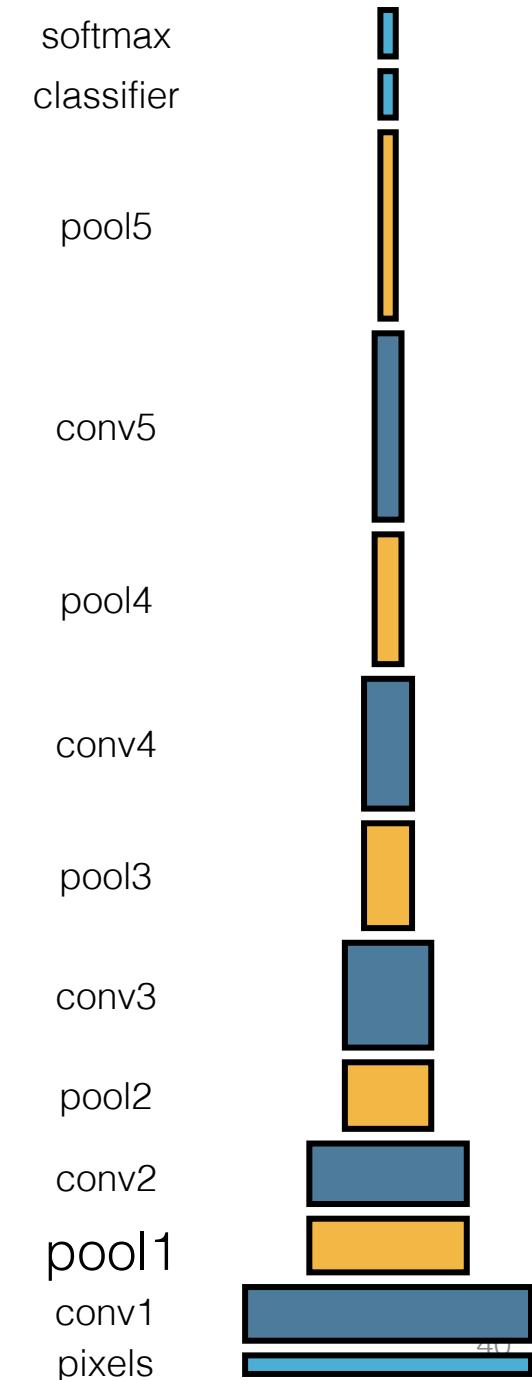
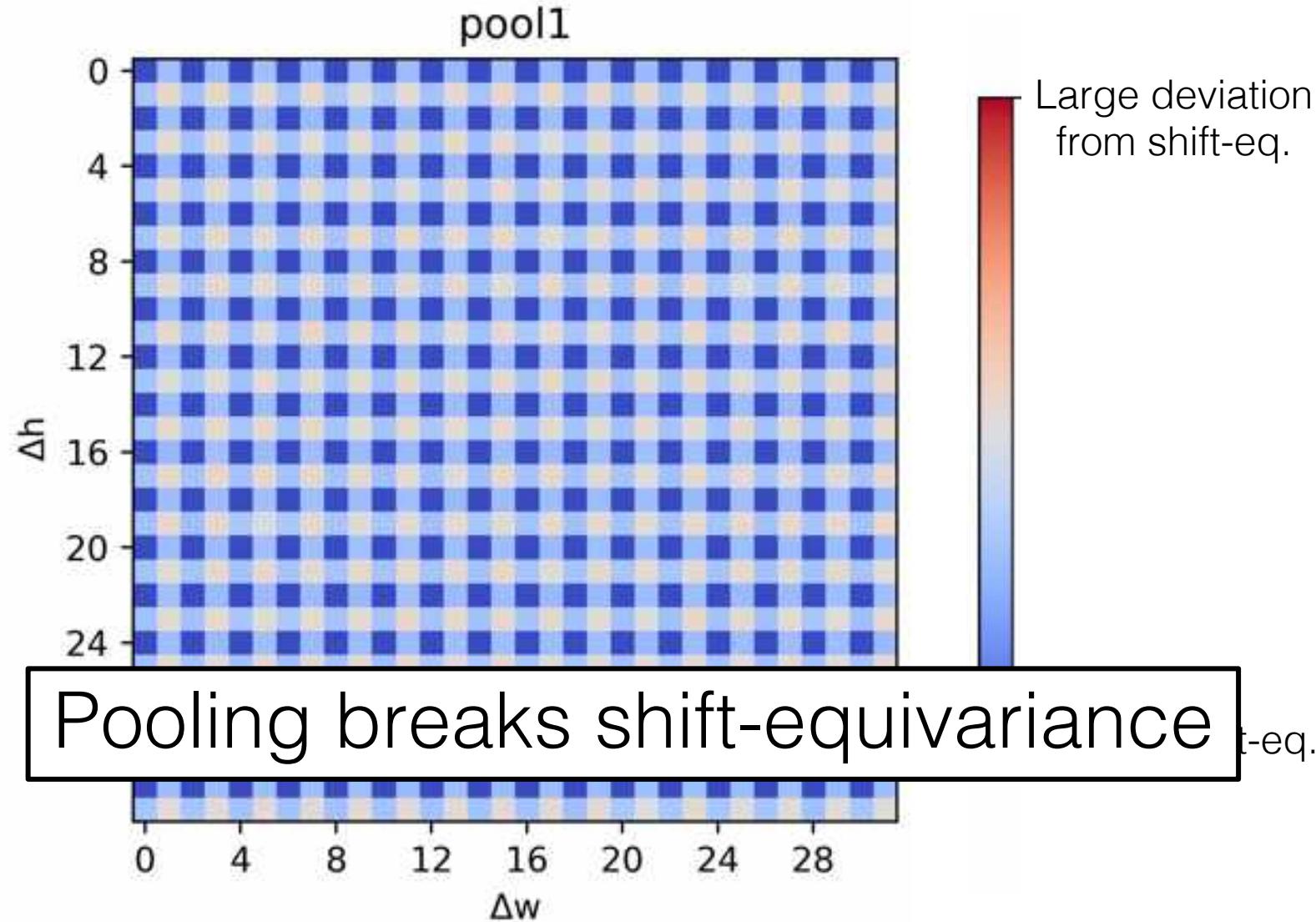
- CIFAR
- VGG network
 - 5 max-pools
- Test shift-equivariance condition
 - $\underline{F}(\underline{\text{Shift}_{\Delta h, \Delta w}}(X)) == \text{Shift}_{\Delta h, \Delta w}(F(X))$
- Circular convolution/shift



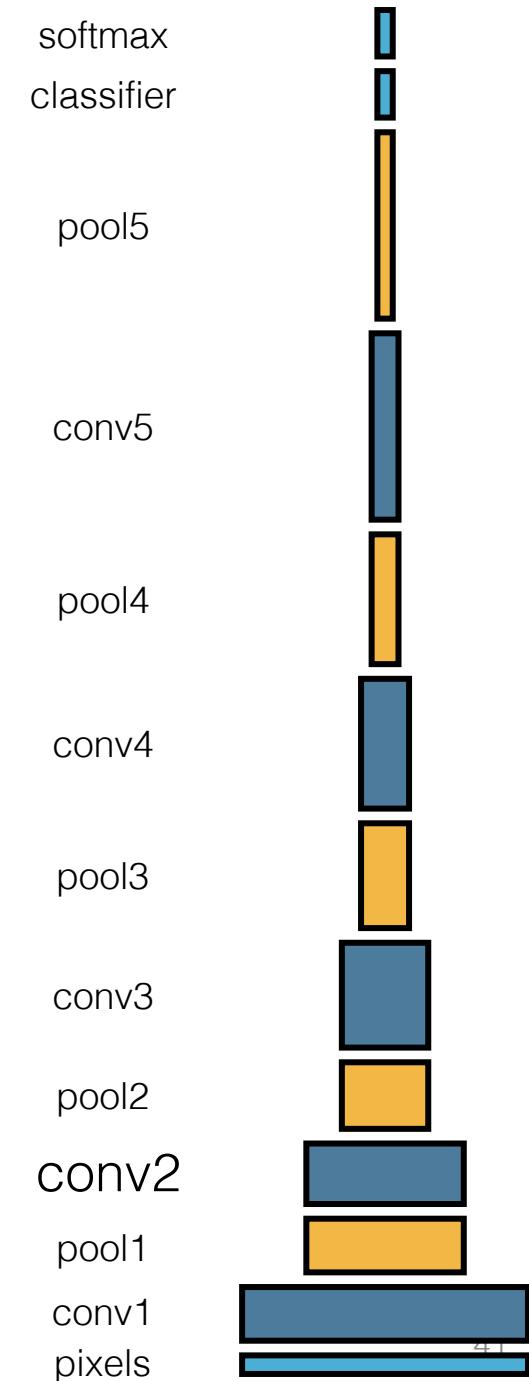
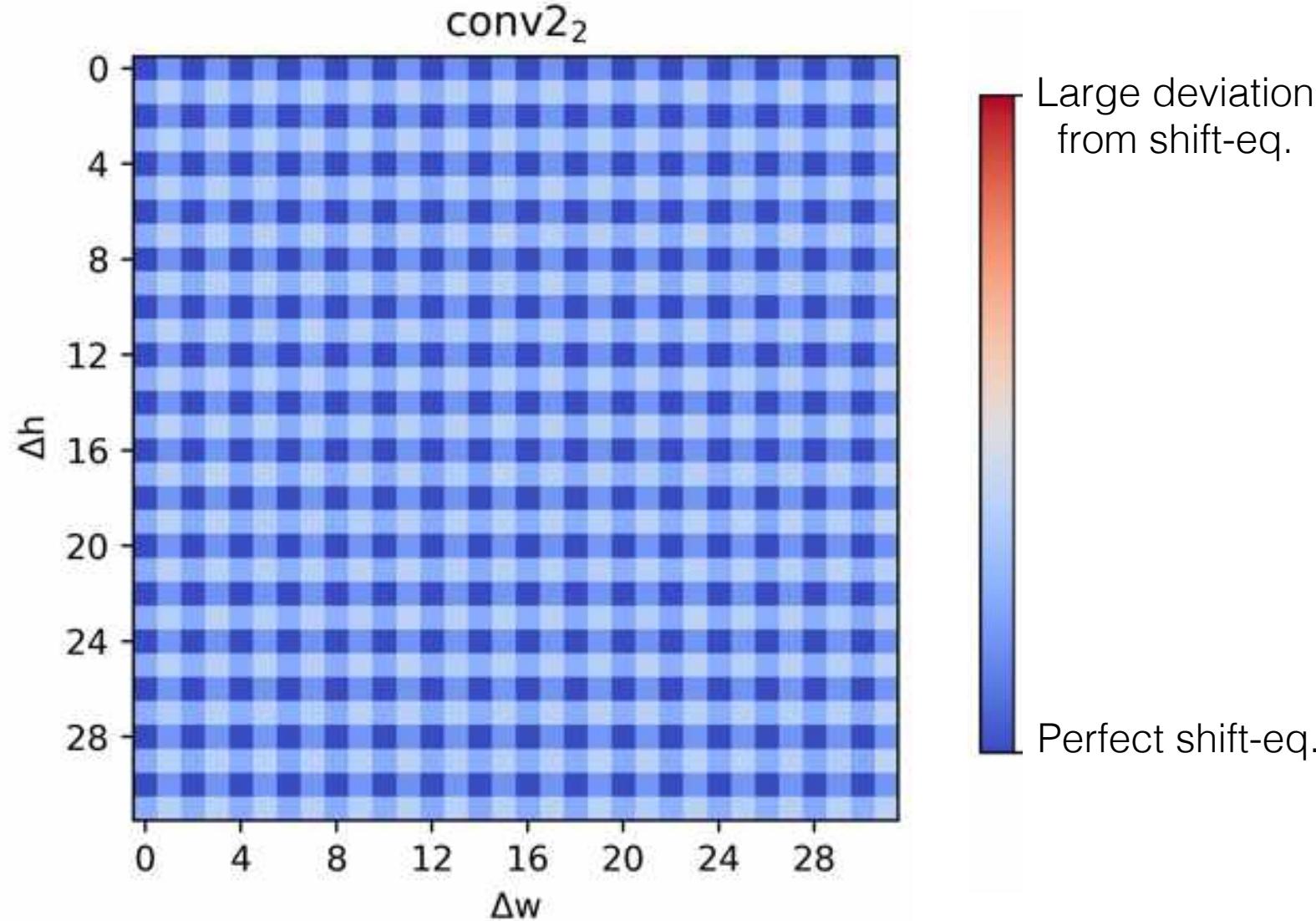
Shift-equivariance, per layer



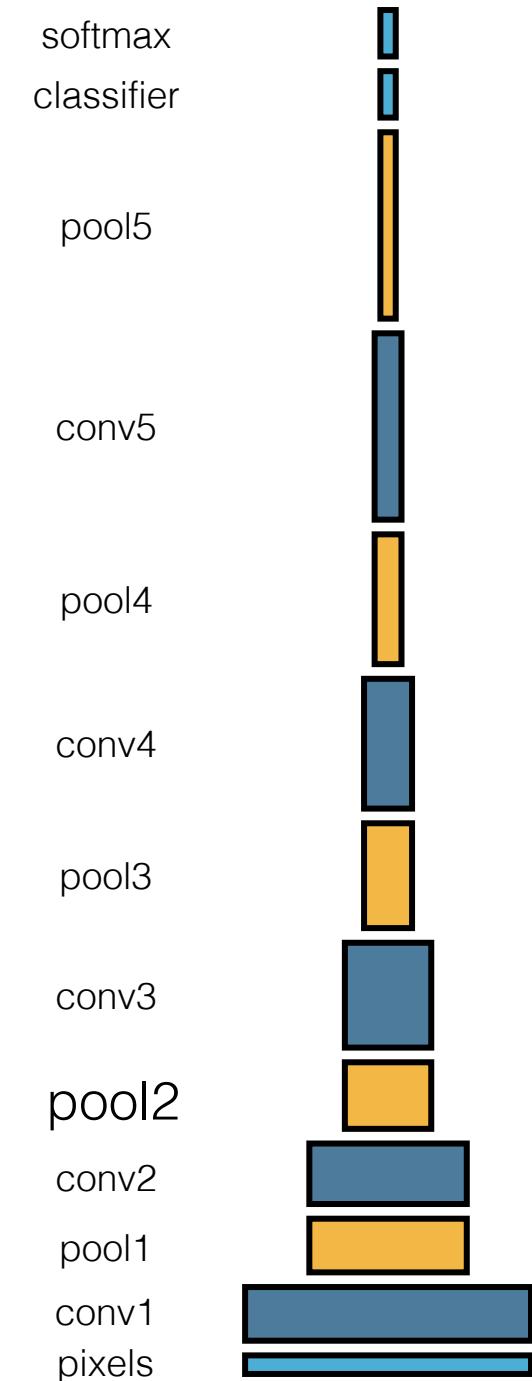
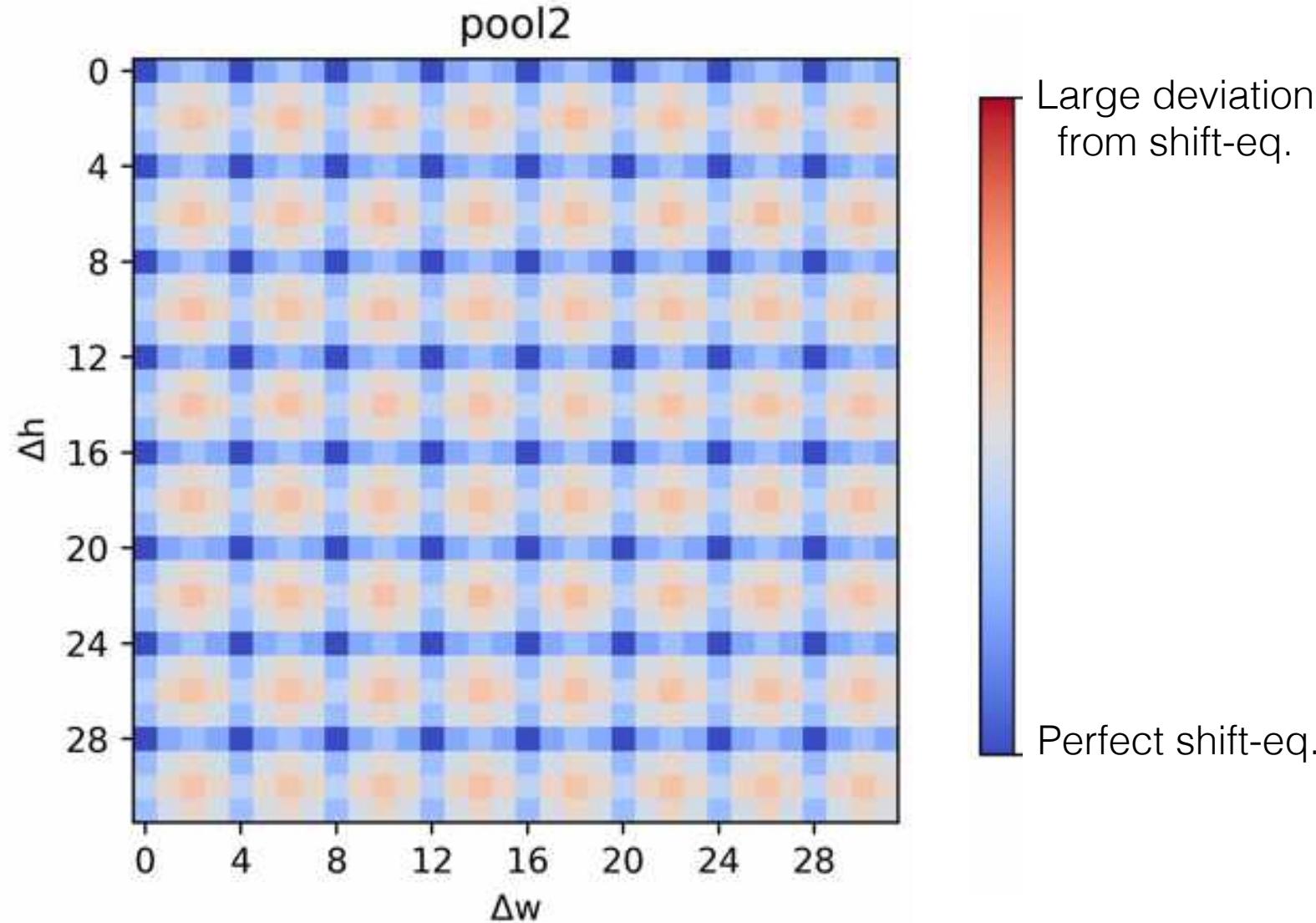
Shift-equivariance, per layer



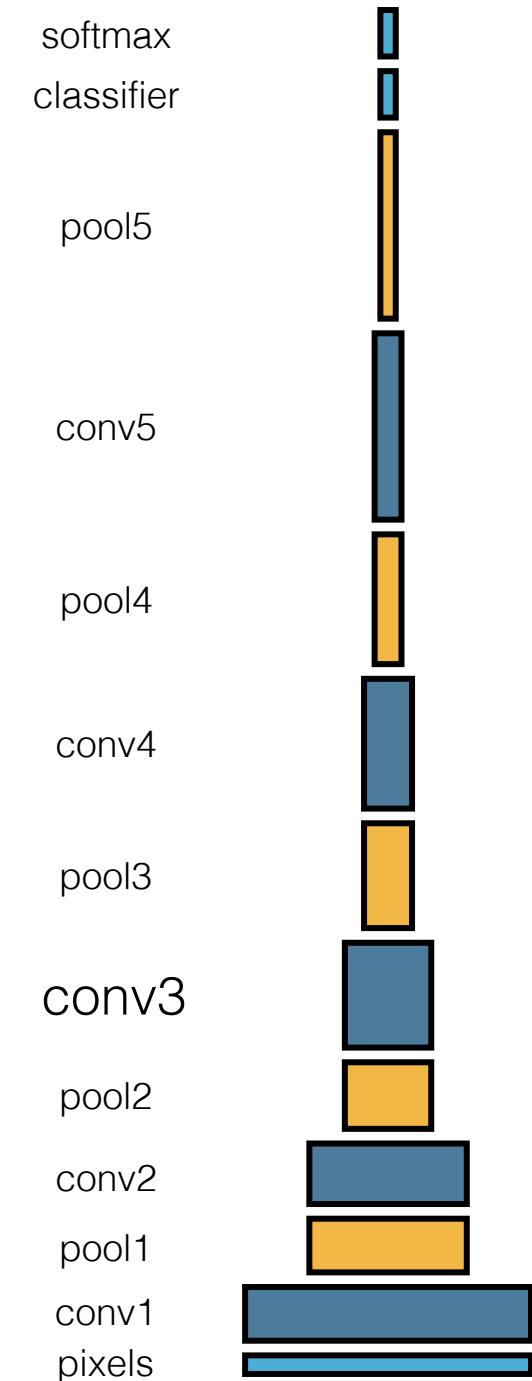
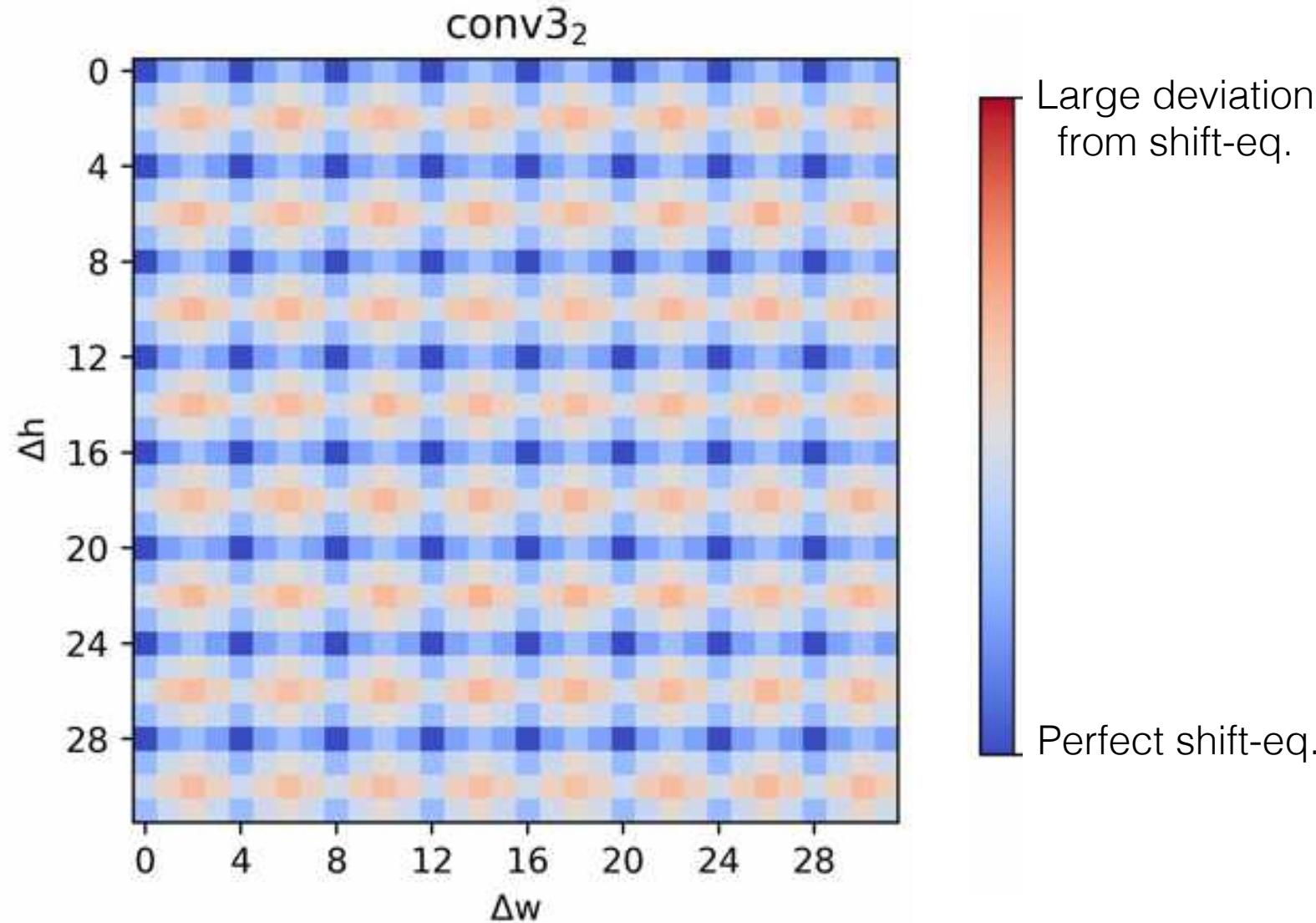
Shift-equivariance, per layer



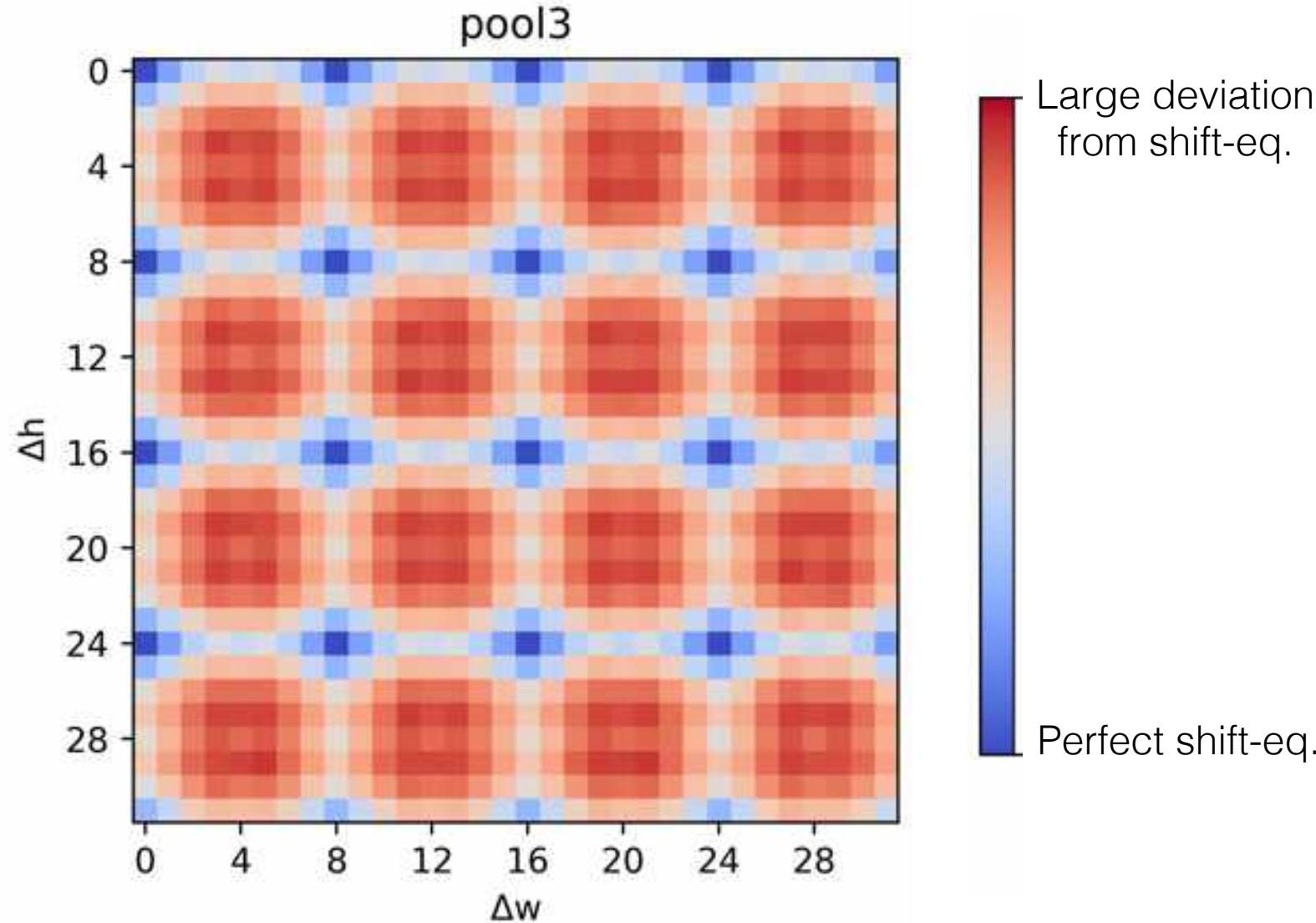
Shift-equivariance, per layer



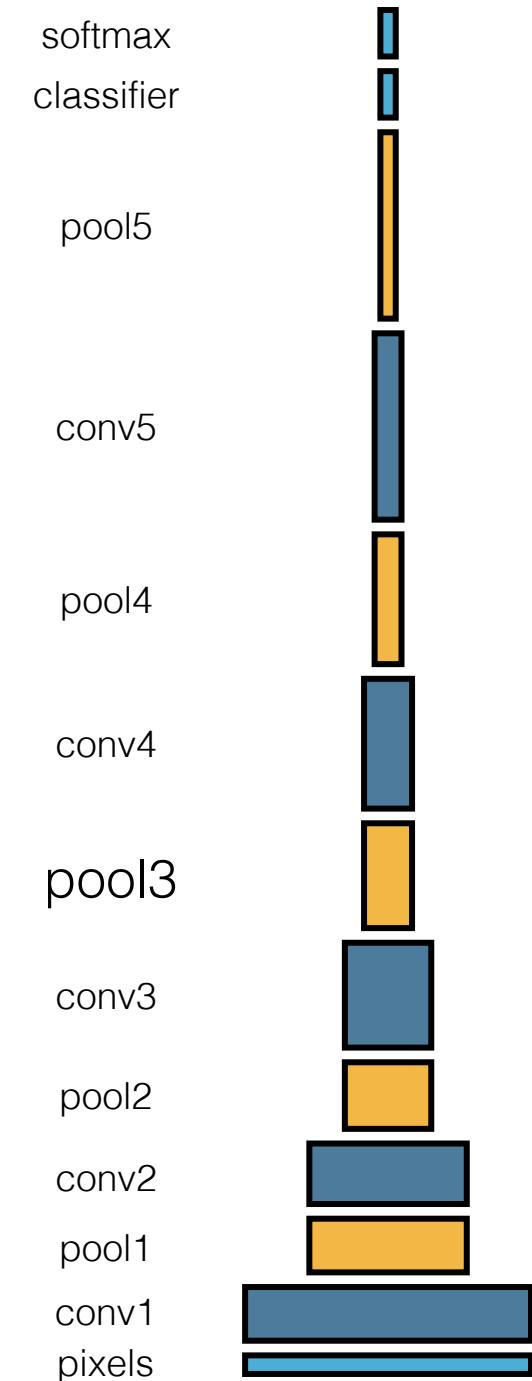
Shift-equivariance, per layer



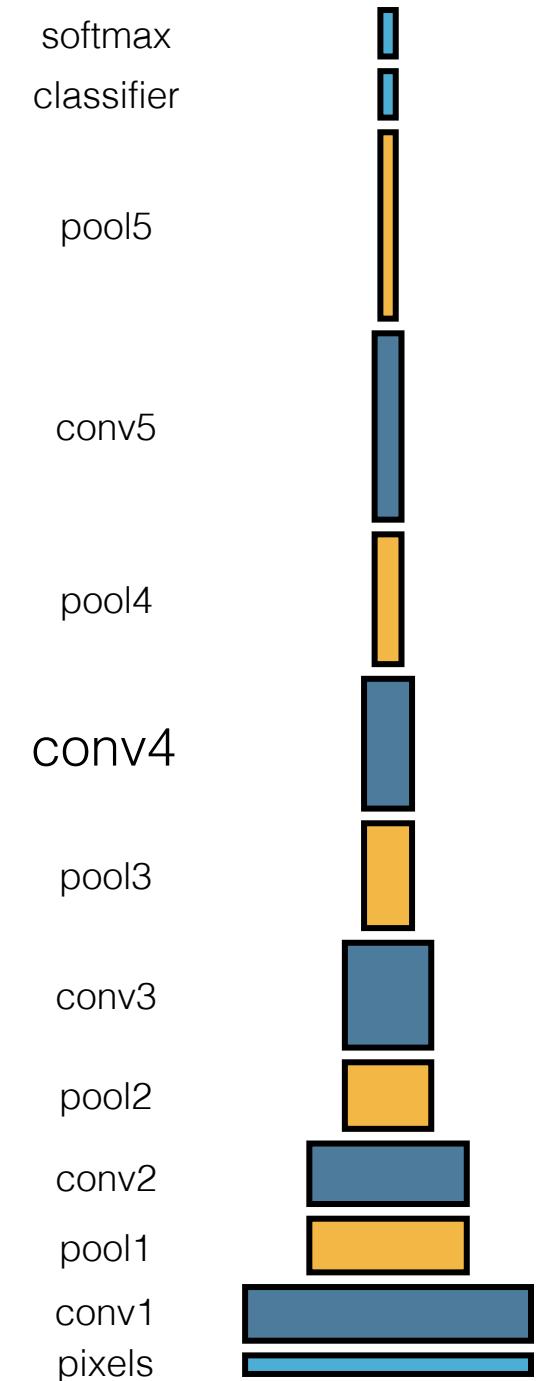
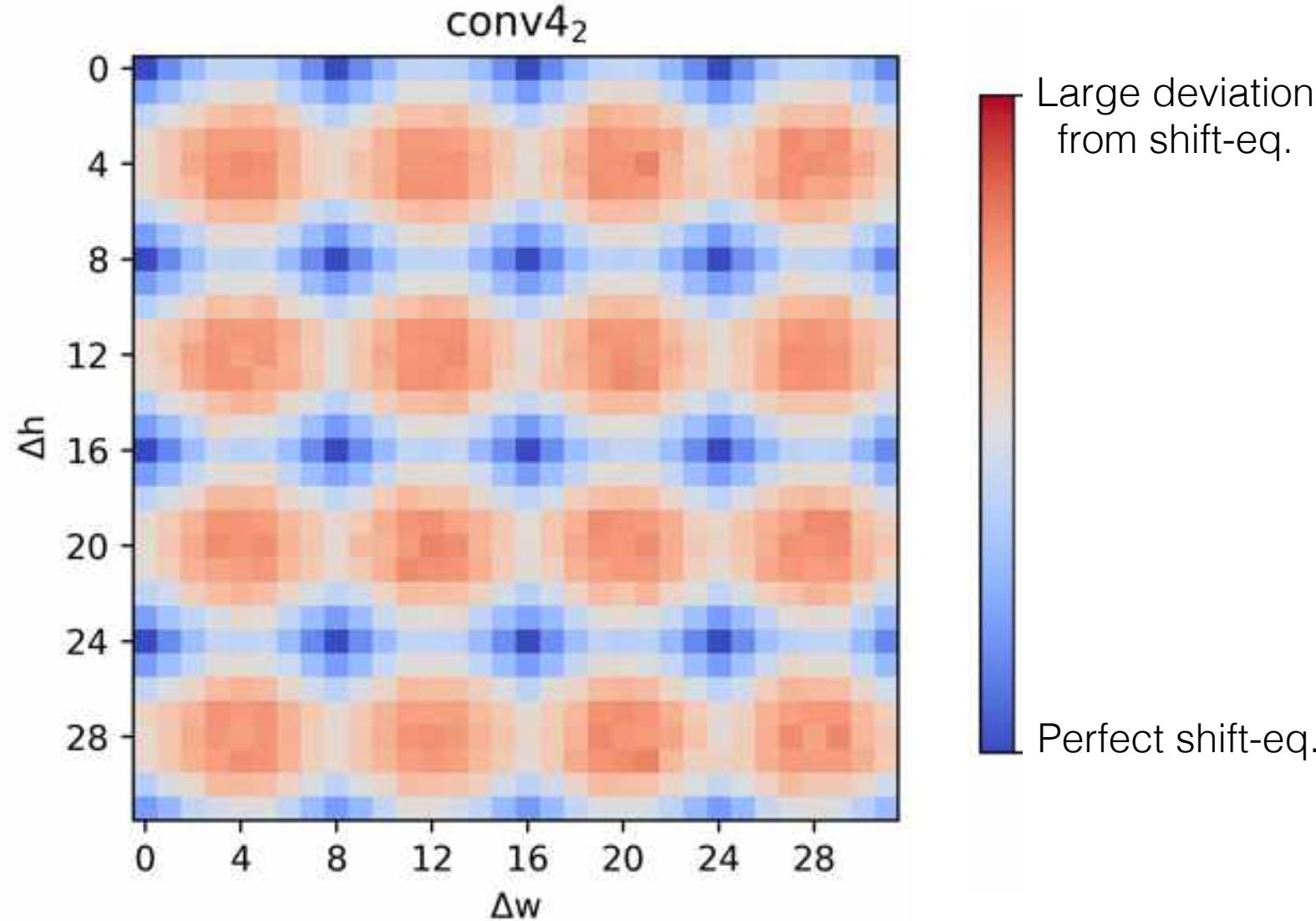
Shift-equivariance, per layer



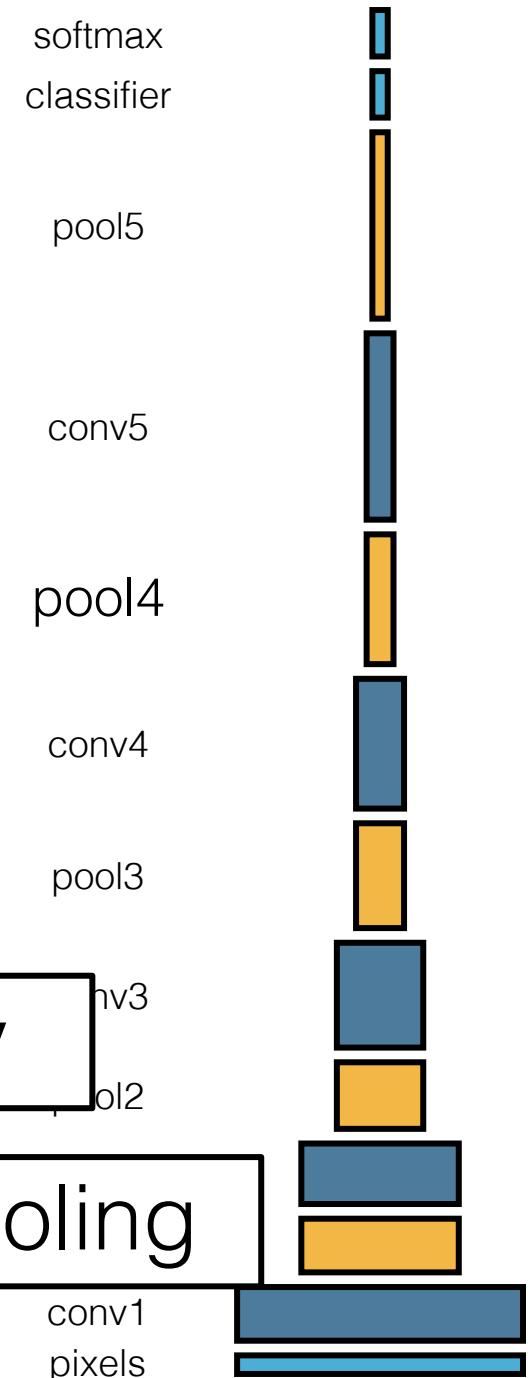
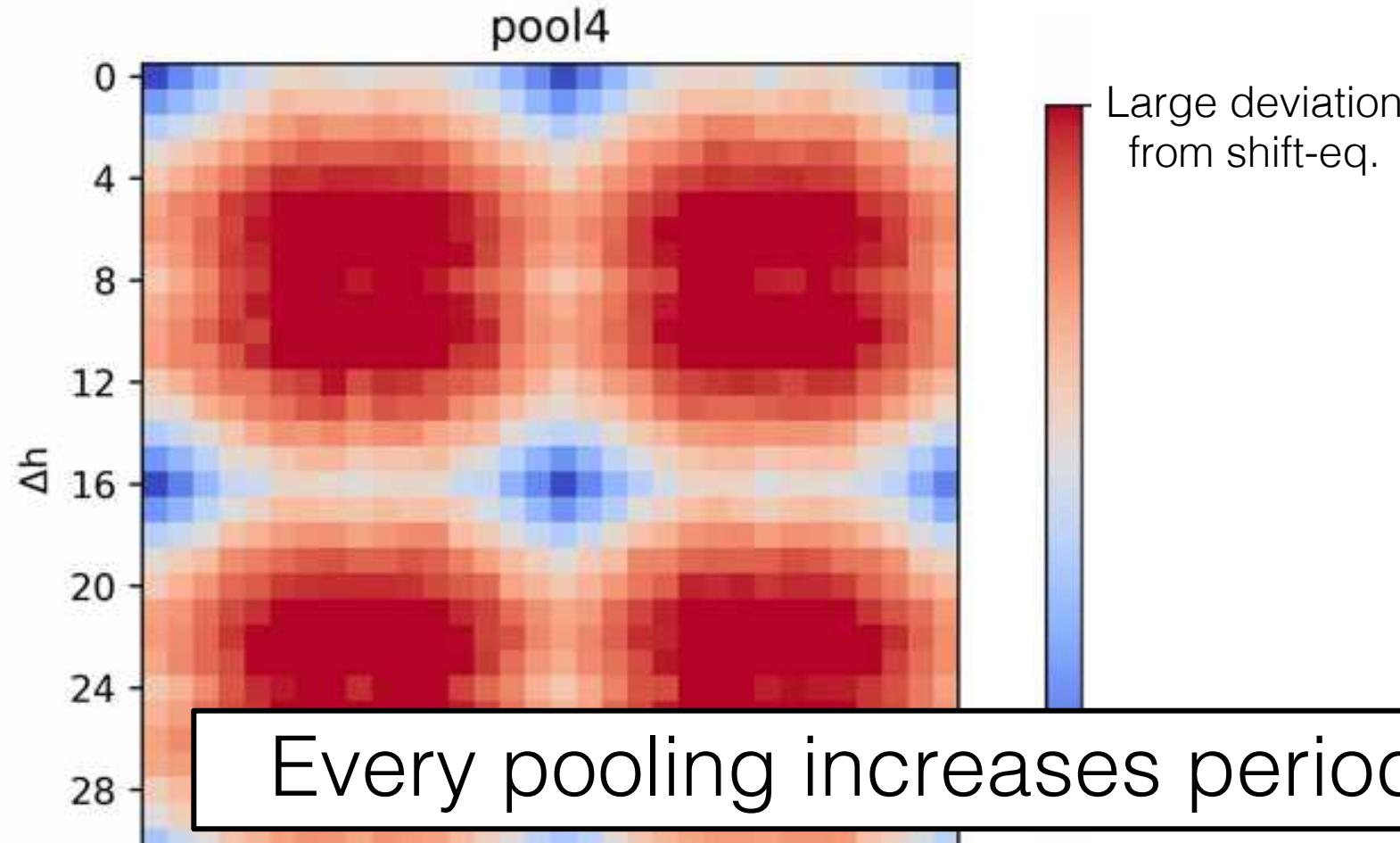
Large deviation
from shift-eq.
Perfect shift-eq.



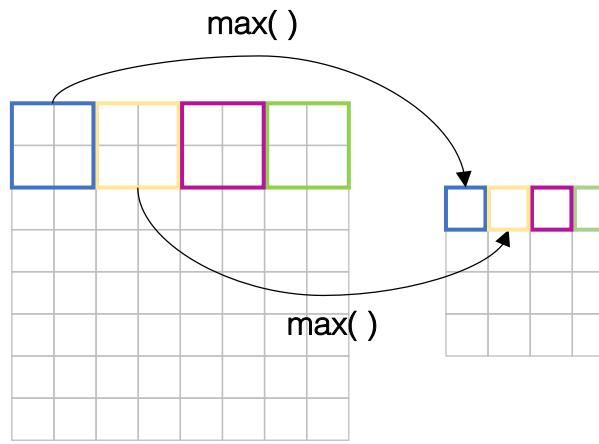
Shift-equivariance, per layer



Shift-equivariance, per layer

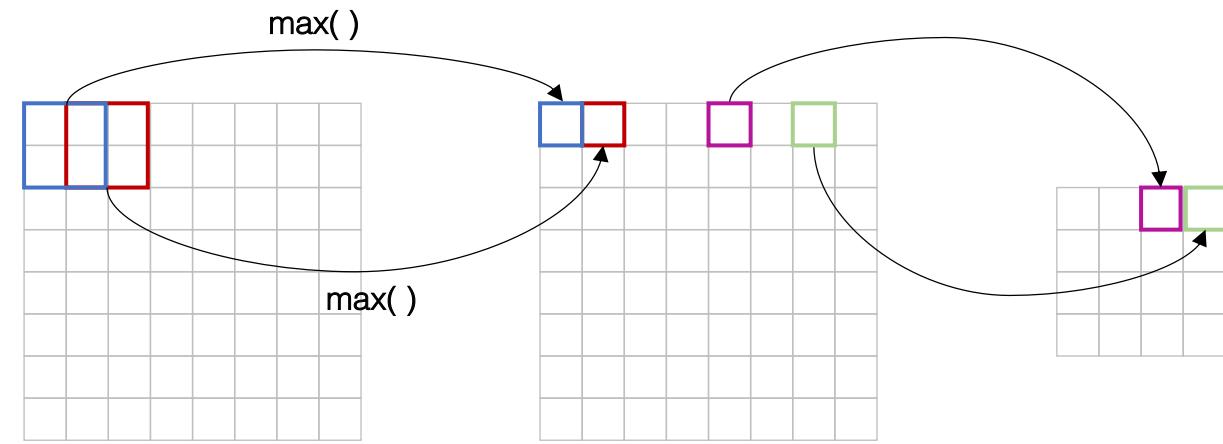


Goal: Reconcile antialiasing with max-pooling



Shift-equivariance lost; heavy aliasing

Strided-MaxPool

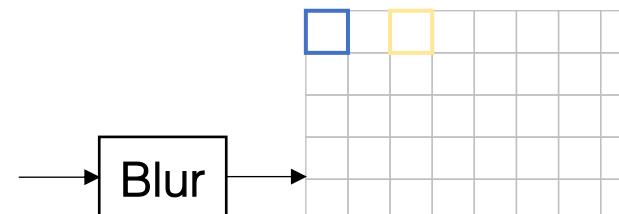


Max (densely)
Preserves shift-equivariance



Subsampling
Shift-eq. lost; heavy aliasing

Equivalent Interpretation



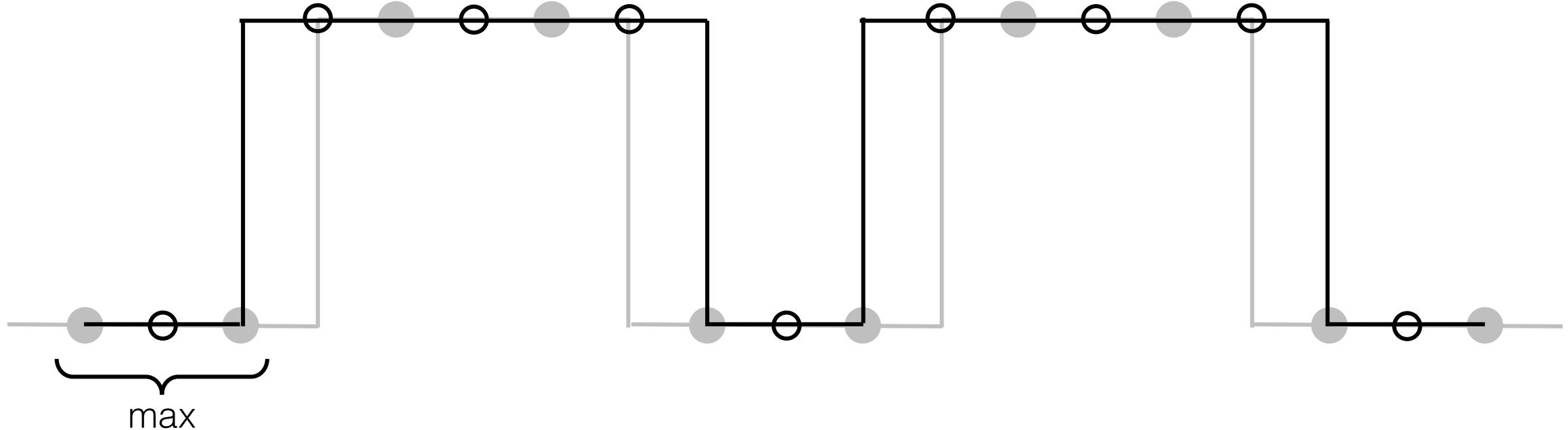
Evaluated together as “BlurPool”



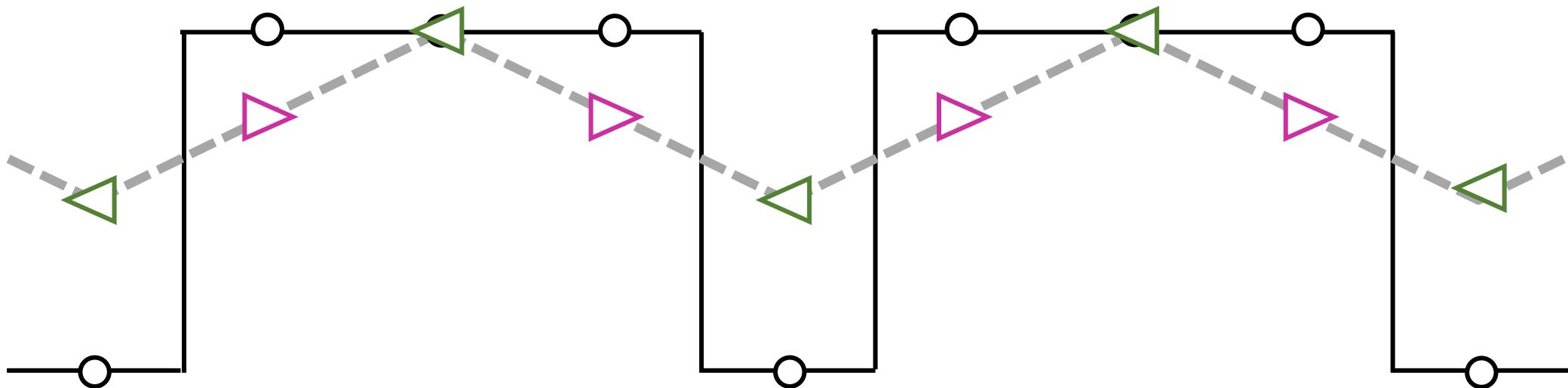
Blur
Preserves shift-eq.
Shift eq. lost, but reduced aliasing

Anti-aliased pooling (MaxBlurPool)

(1) Max (densely)

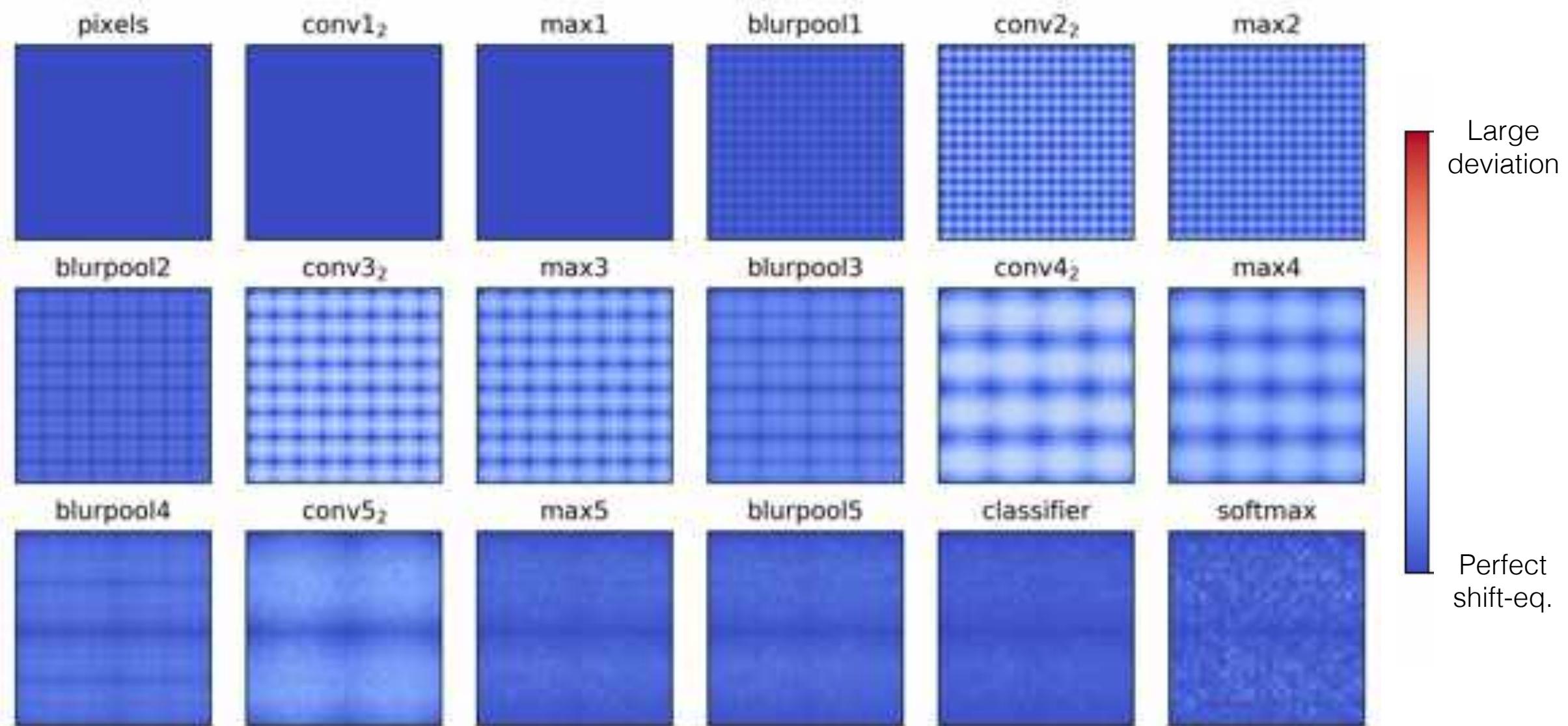


(2) Blur+Subsample

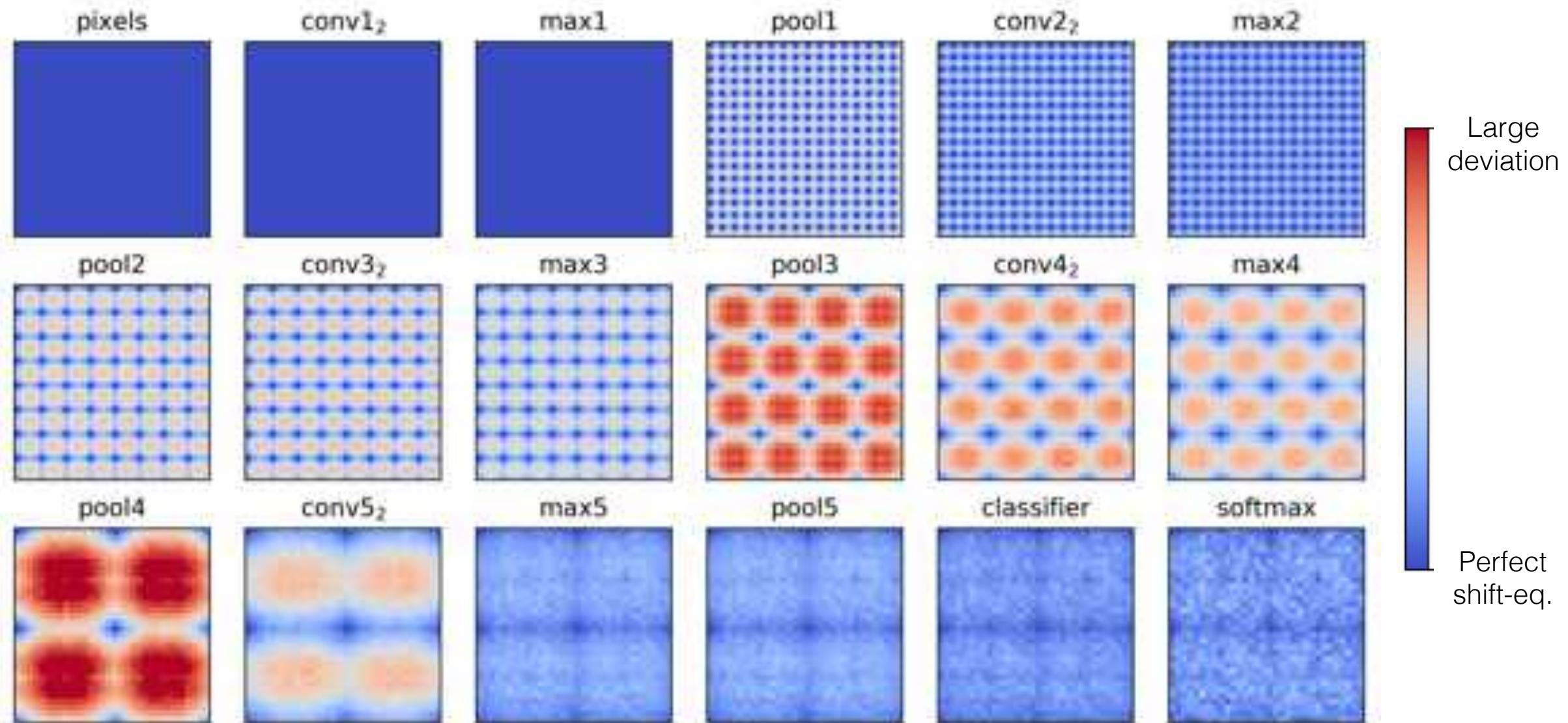


Shift-equivariance better preserved

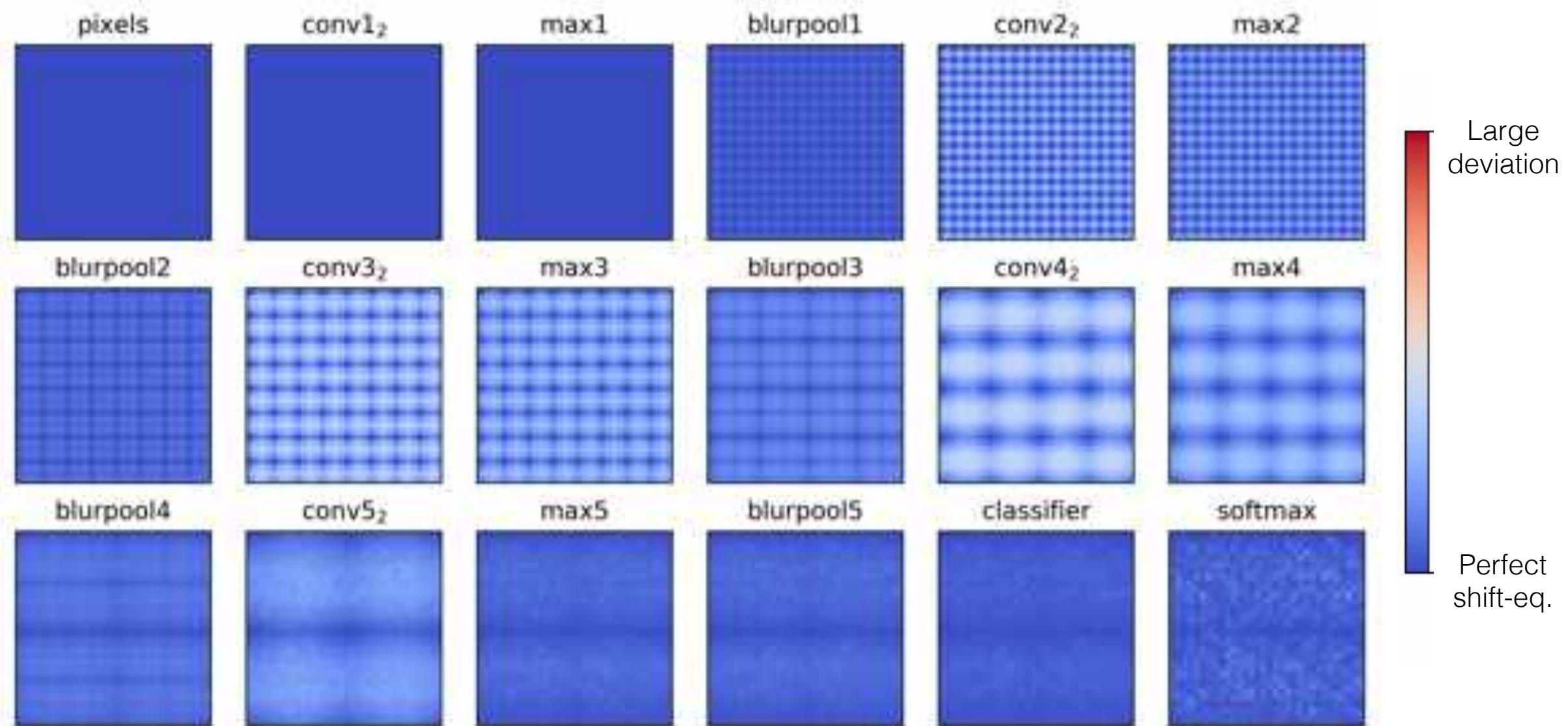
Anti-aliased



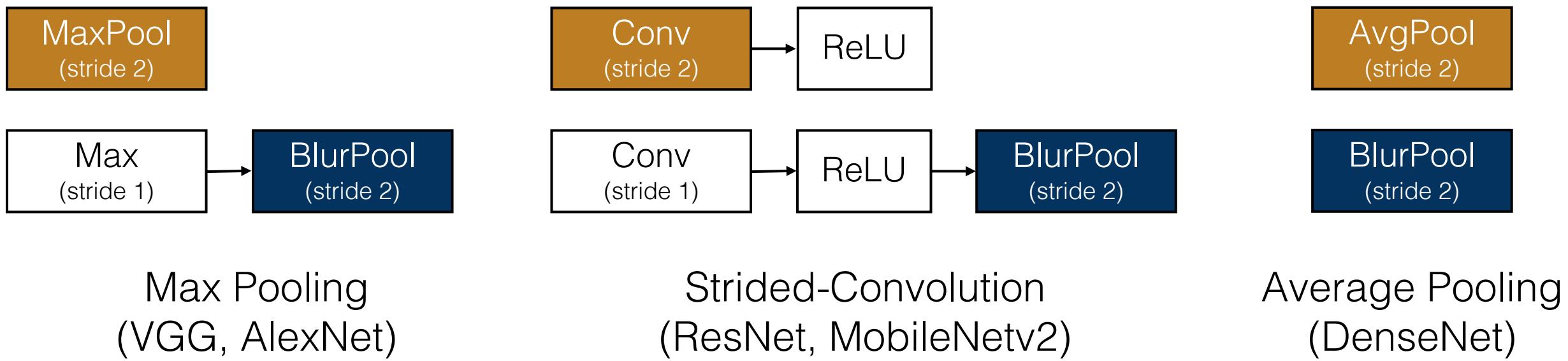
Baseline



Anti-aliased

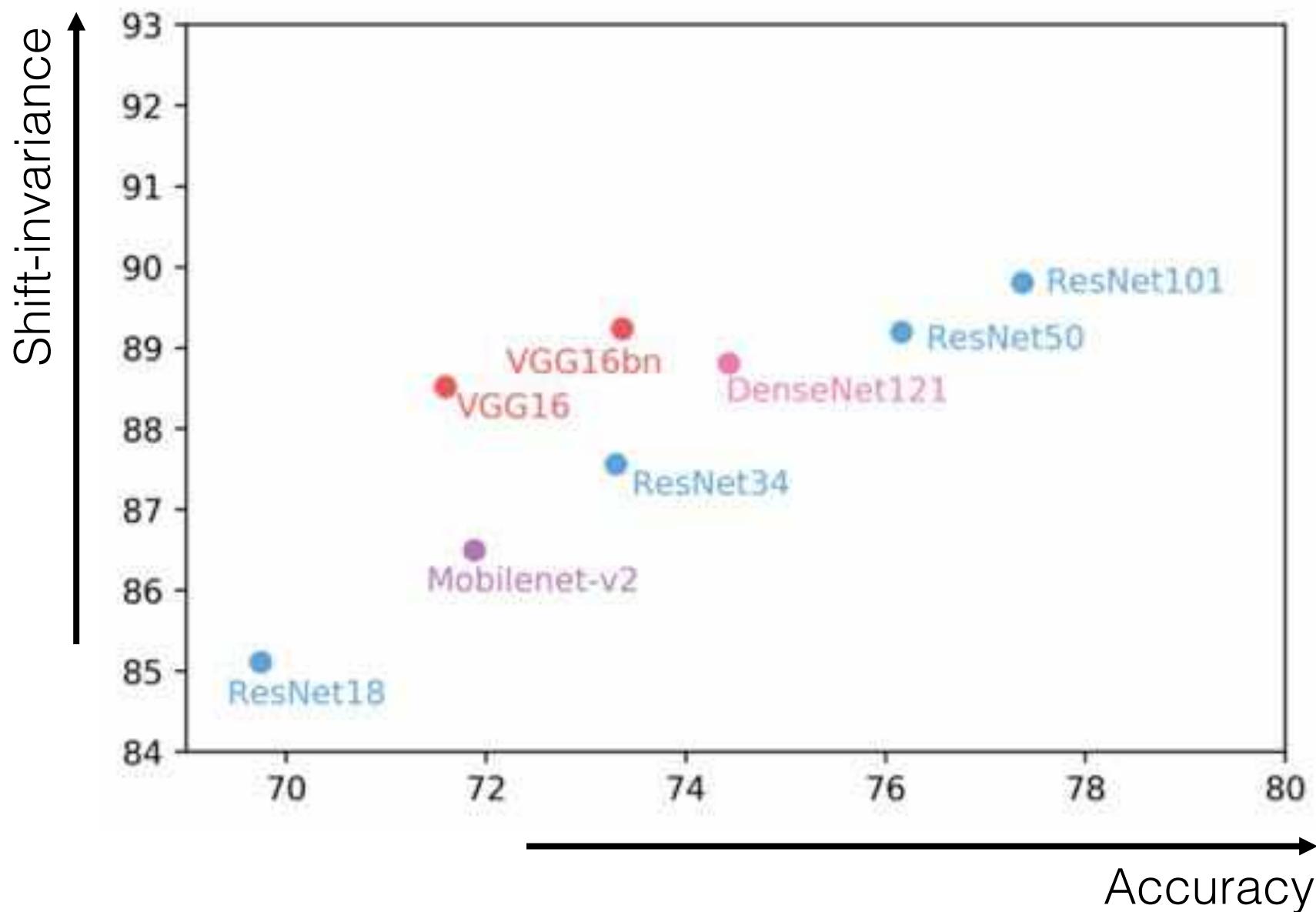


Other downsampling layers

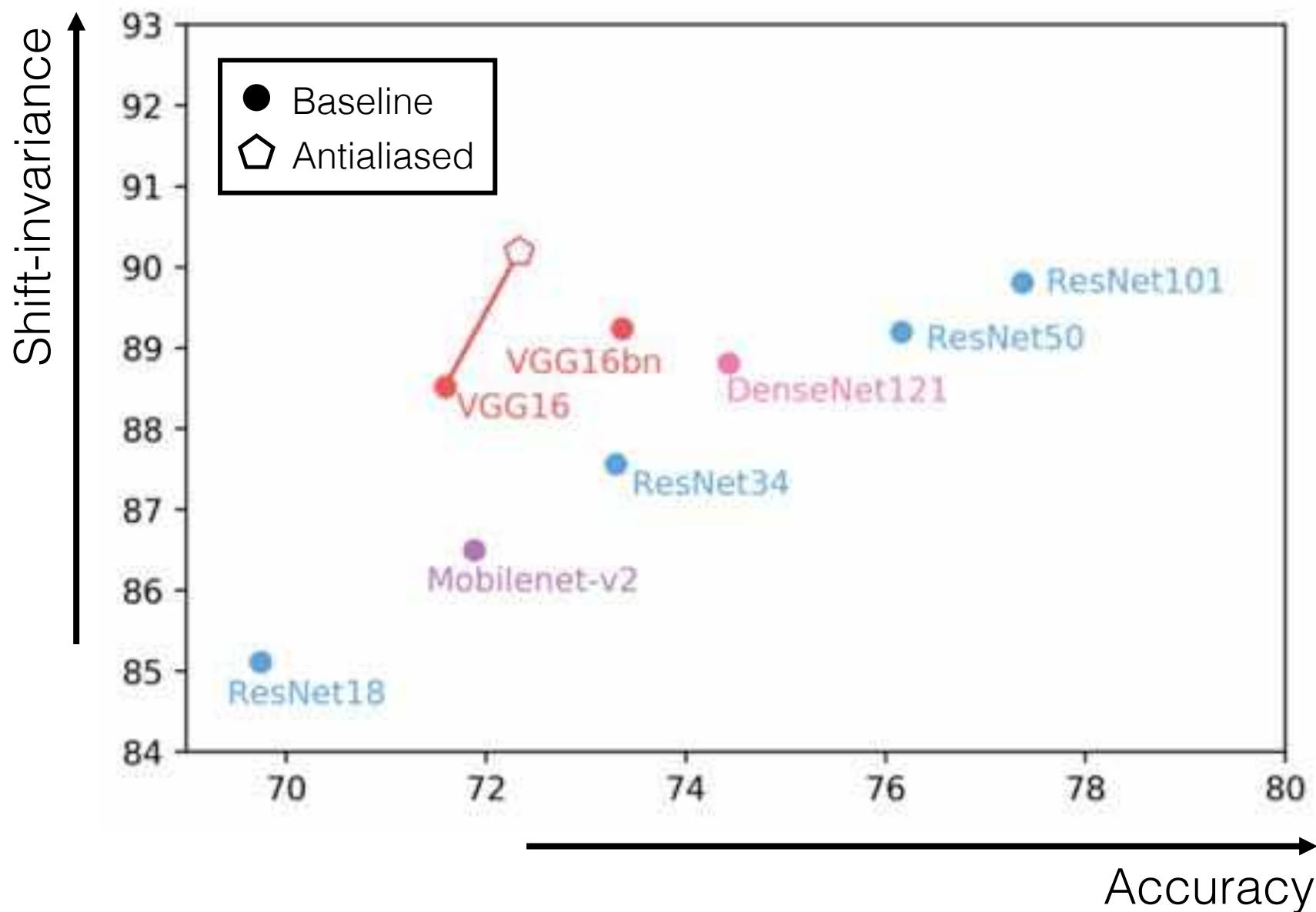


Antialias any off-the-shelf network

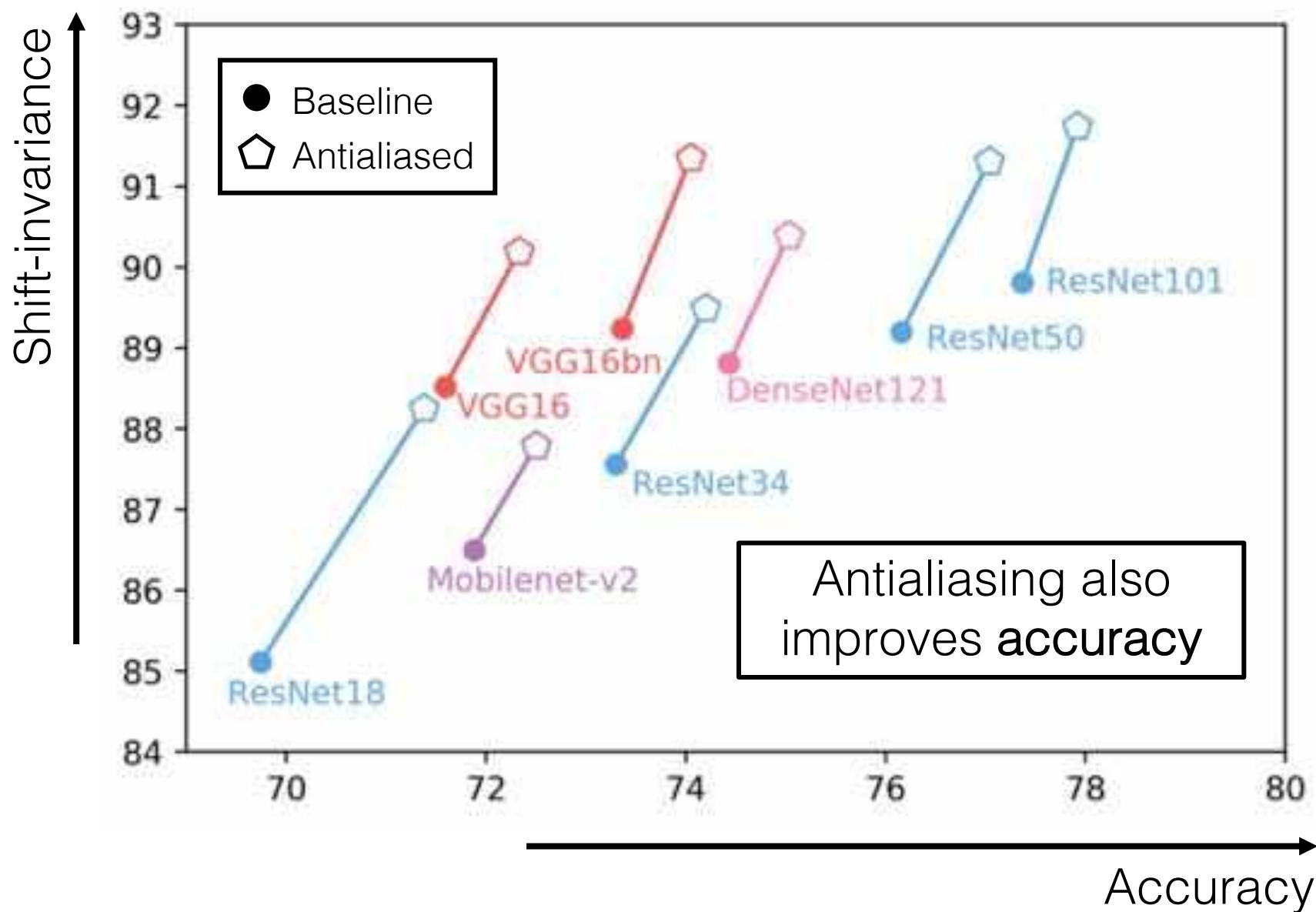
ImageNet



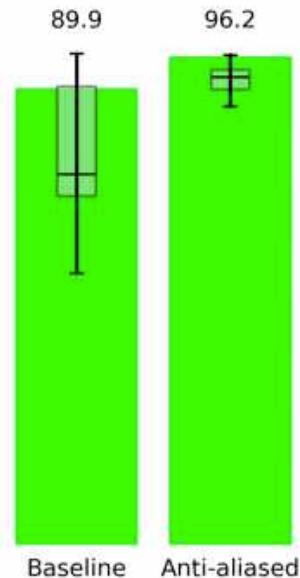
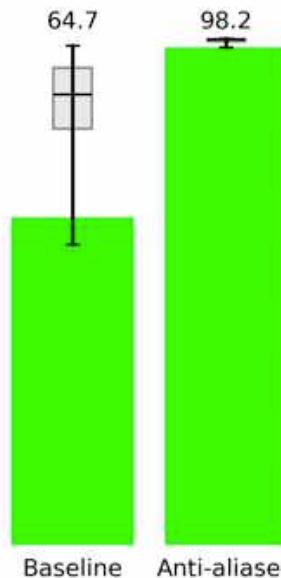
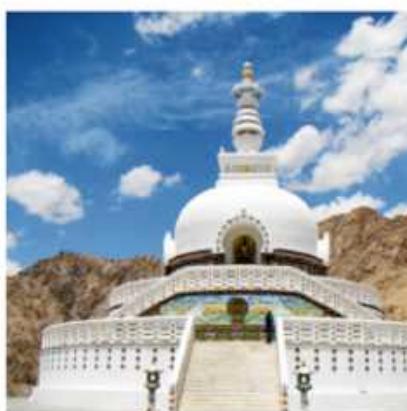
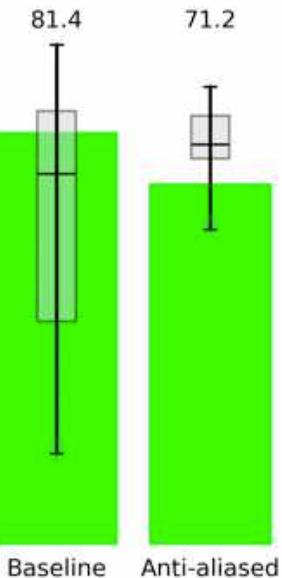
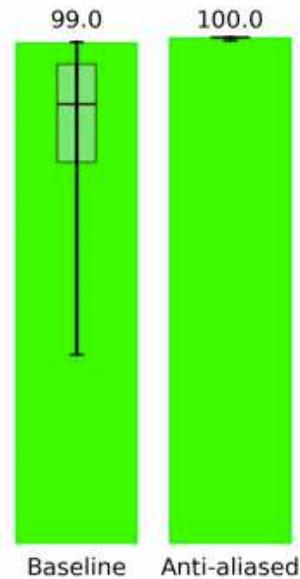
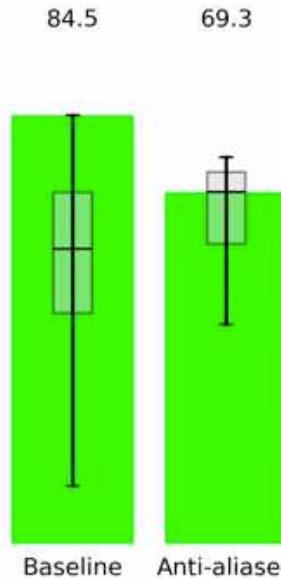
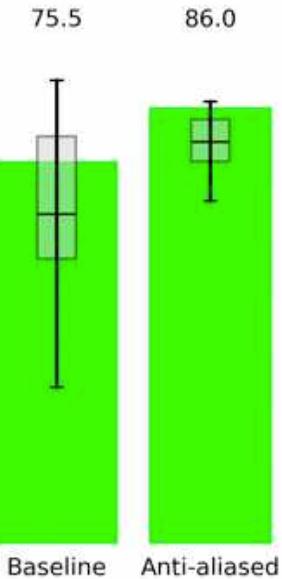
ImageNet



ImageNet



Qualitative examples



Discussion

Striding aliases

Antialiasing improves

- + shift-equivariance, accuracy
- + stability to perturbations
- + robustness to corruptions

Code + models: richzhang.github.io/antialiased-cnns/

Ex: `models.resnet50()` → `models_lpf.resnet50()`

→ Implications on image generation models?

Making fake images is getting easier



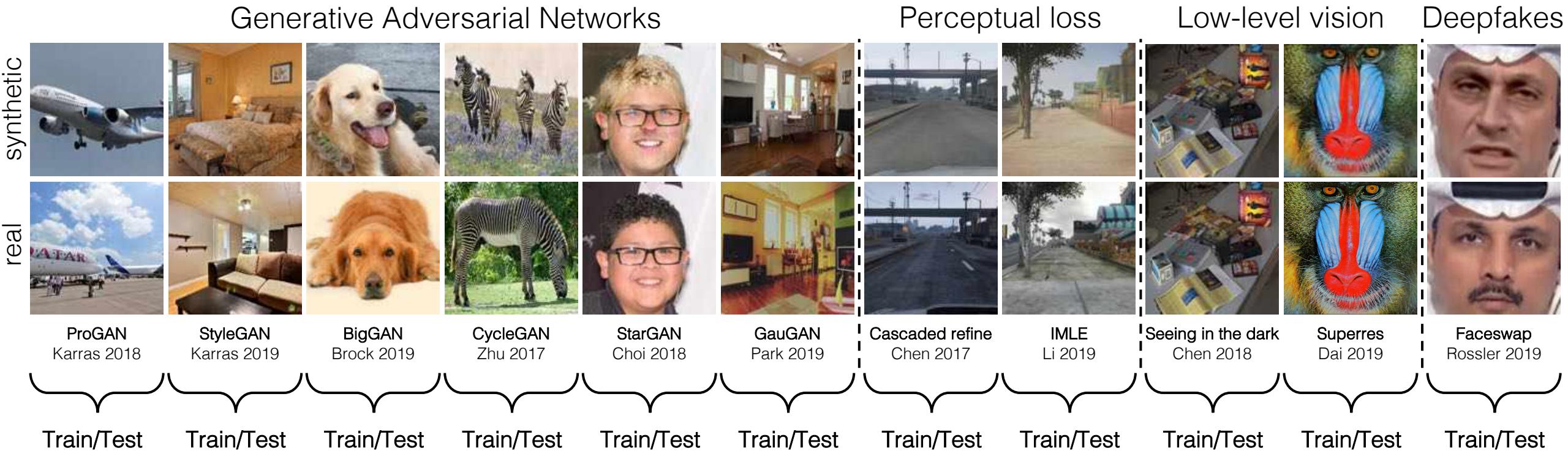
“Deepfakes”



GANs

Can we create a “universal” detector?

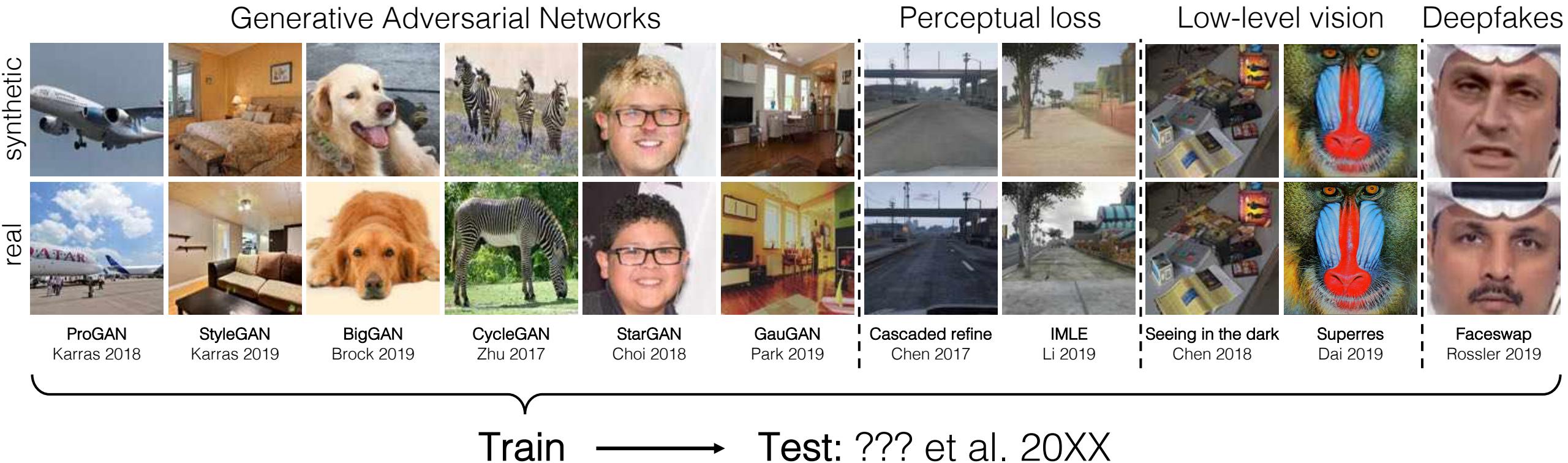
Dataset of CNN-generated fakes



Test: ??? et al. 20XX

Can we create a “universal” detector?

Dataset of CNN-generated fakes



Can we create a “universal” detector?

Dataset of CNN-generated fakes

Generative Adversarial Networks



ProGAN
Karras 2018

Train

eGAN
as 2019

StyleGAN
Karras 2019

BigGAN
Brock 2019

CycleGAN
Zhu 2017

StarGAN
Choi 2018

GauGAN
Park 2019

Perceptual loss



Cascaded refine Chen 2017

IMLE
Li 2019

Low-level vision



Seeing in the dark

Chen 2018

Superres
Dai 2019

Deepfakes



Faceswap
Rossler 2019

Many differences (architecture, dataset, objective)

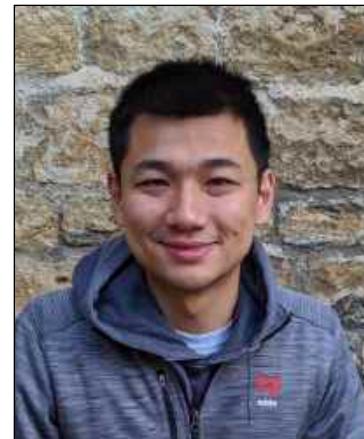
CNN-generated images are surprisingly easy to spot...for now



Sheng-Yu Wang



Oliver Wang



Richard Zhang

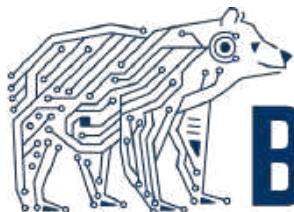


Andrew Owens



Alexei A. Efros

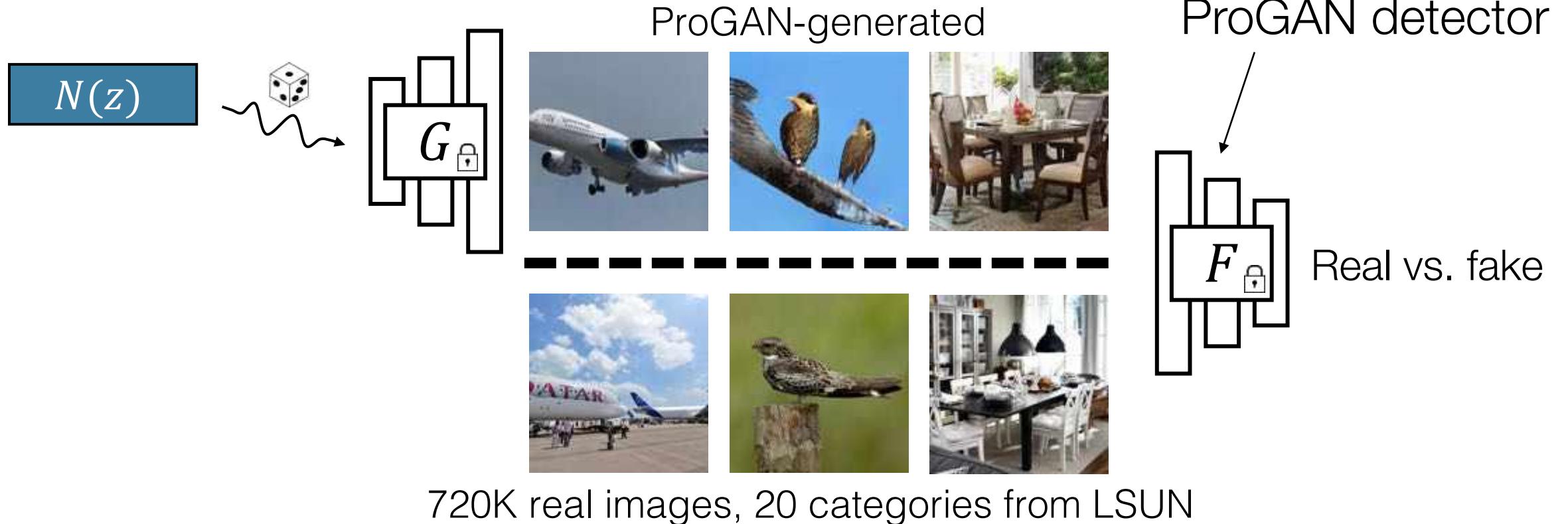
In CVPR, 2020 (oral).



BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH



Training on ProGAN



Testing across architectures

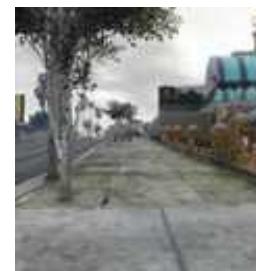
Synthesized images from other CNNs



...

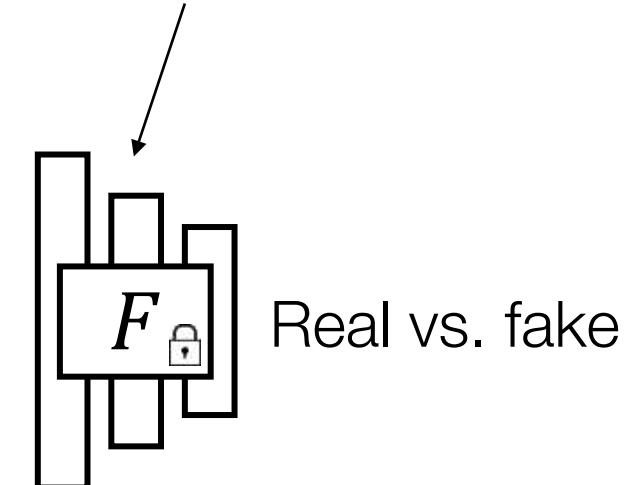


...

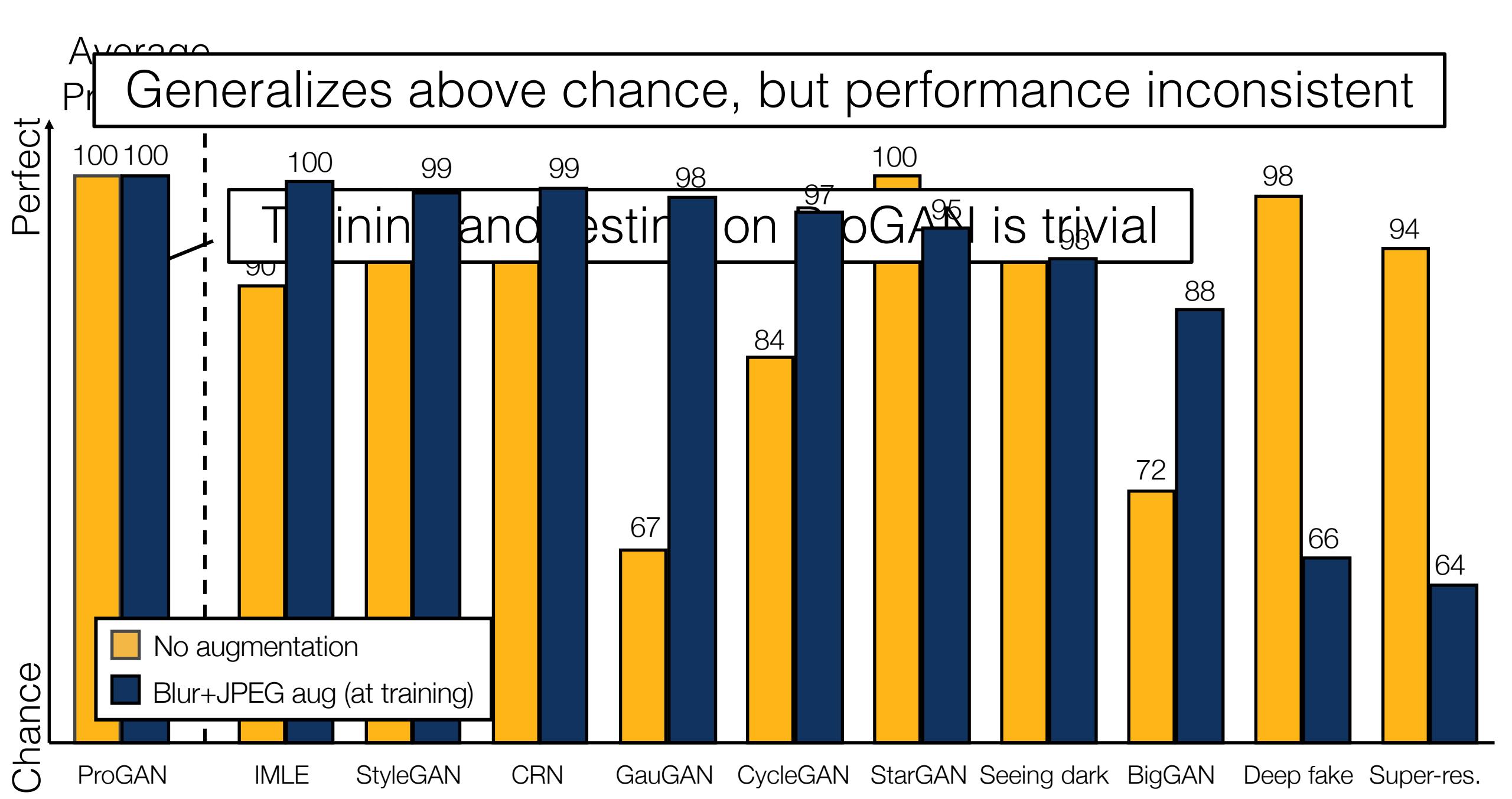


Real images

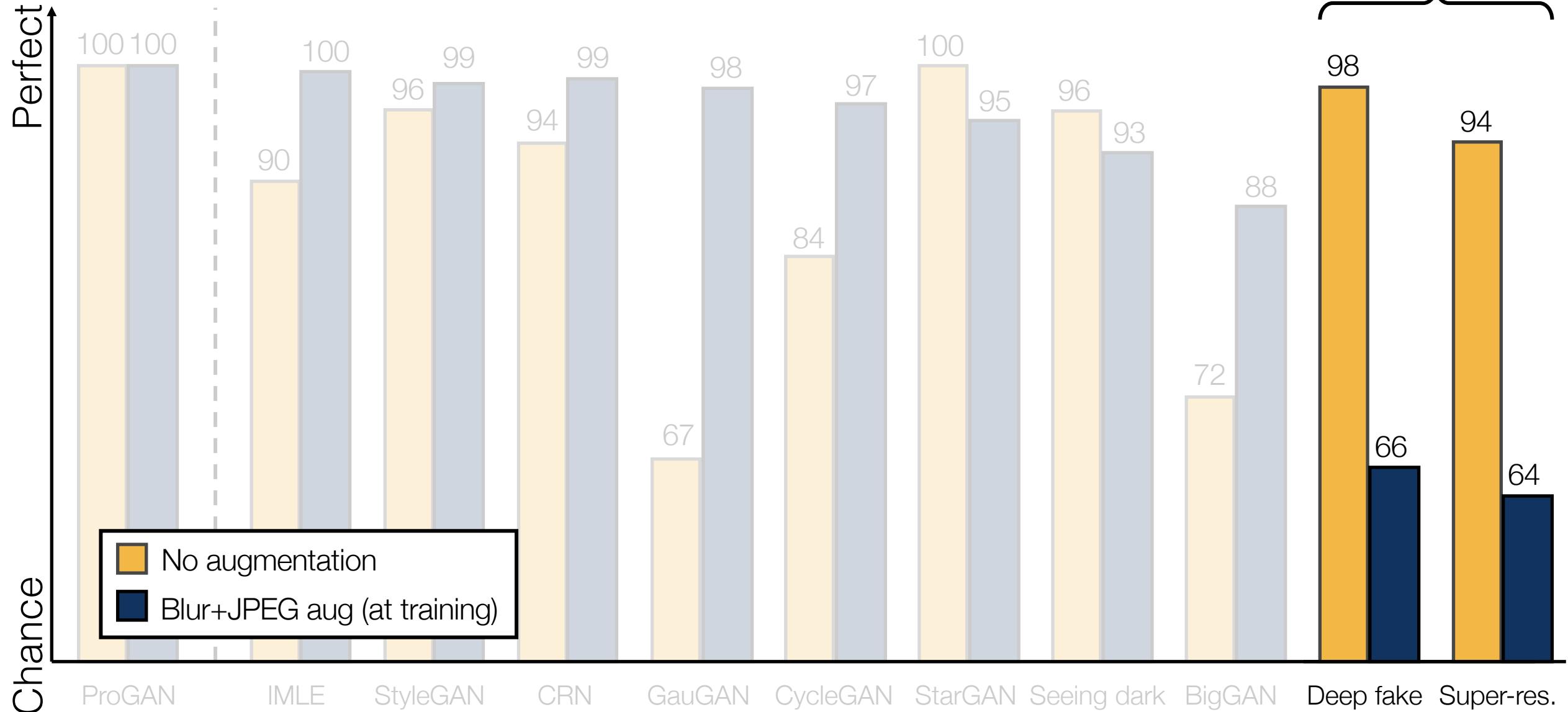
ProGAN detector



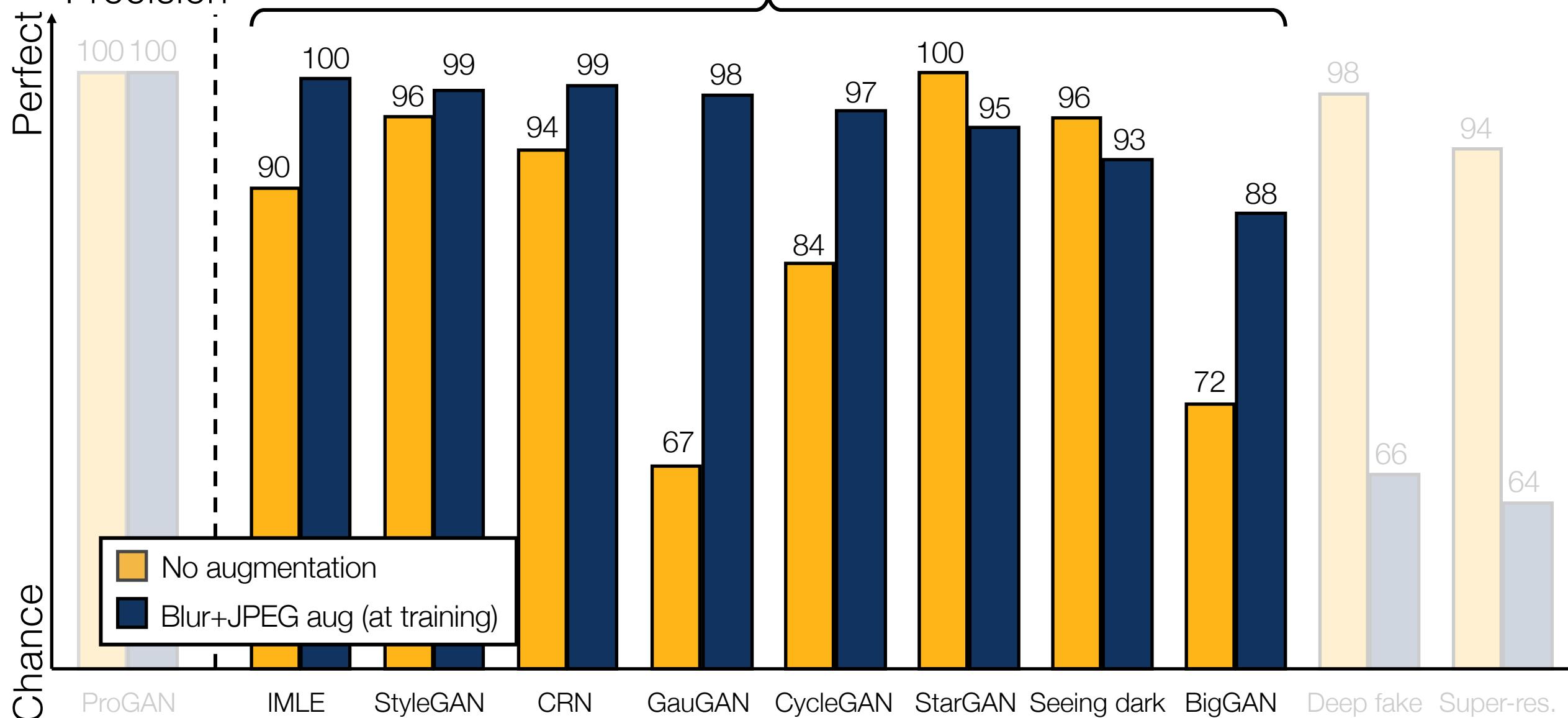
Real vs. fake

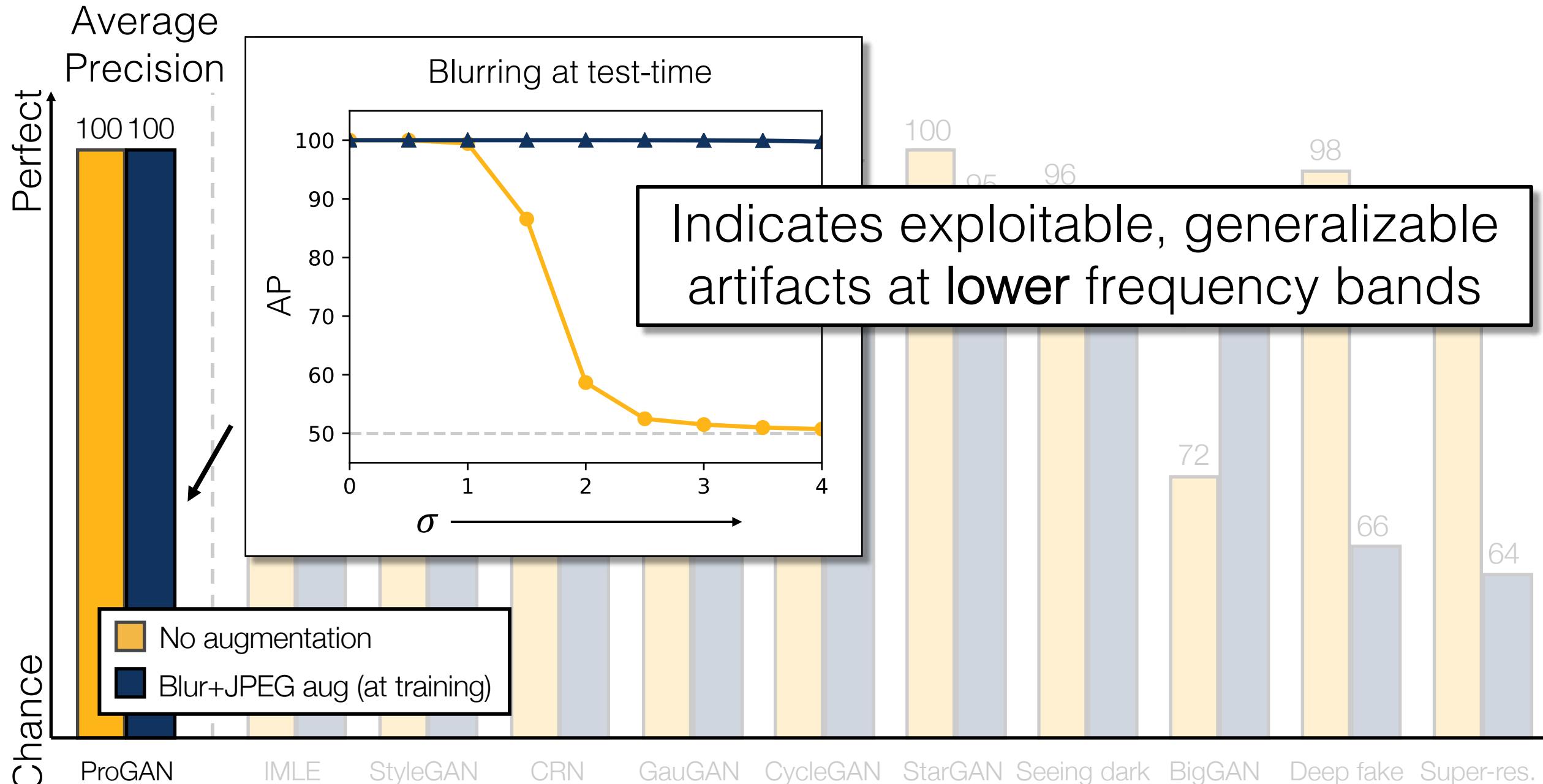


Augmentation is not always appropriate



Aggressive augmentation adds surprising generalization



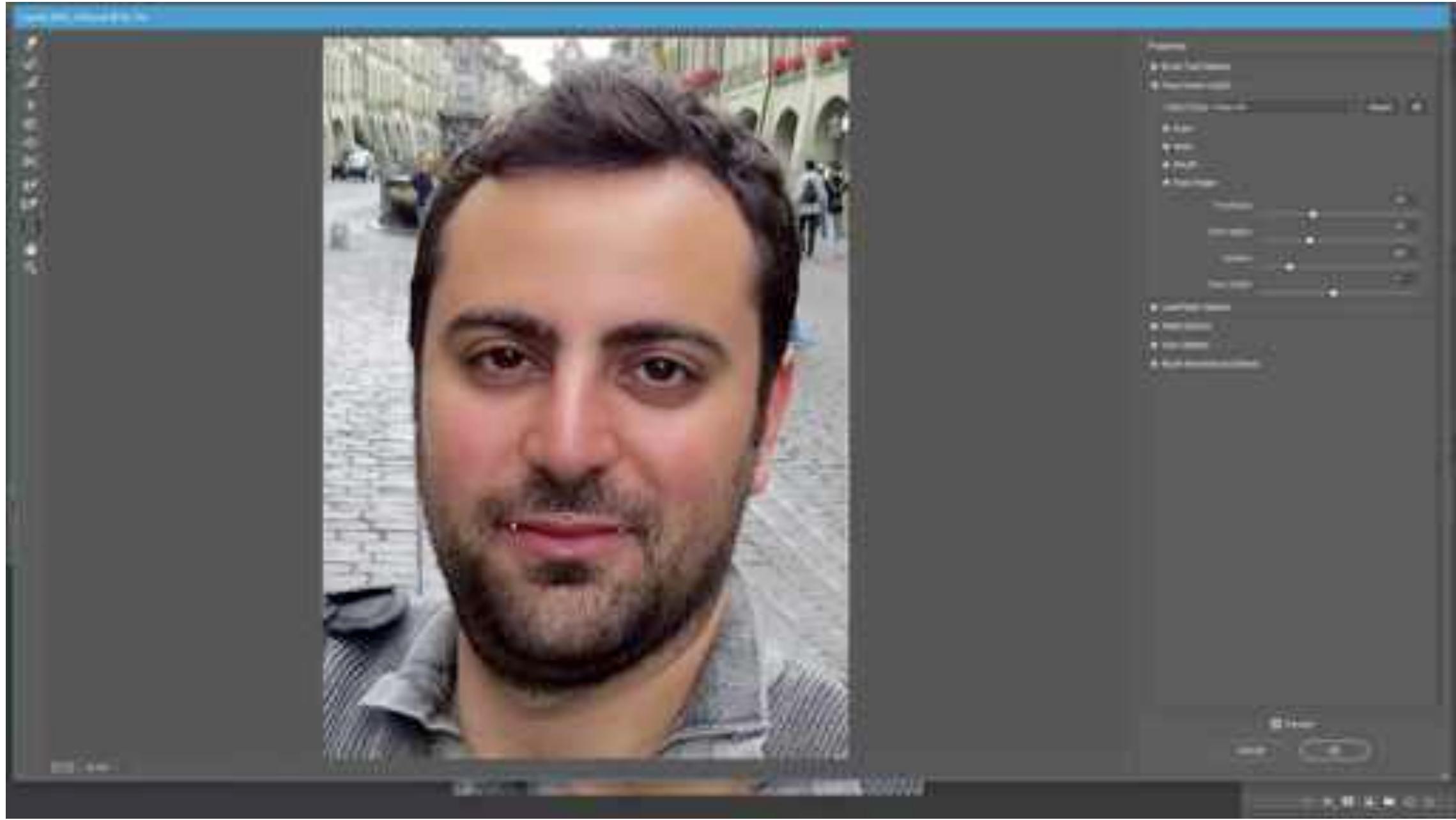


Discussion

- Suggests CNN-generated images have common artifacts
- Artifacts can be detected by a simple classifier!
 - StyleGAN2 (released **after** our submission): 100% AP on FFHQ
 - Swapping Autoencoder (Park et al.): 95% AP on FFHQ
 - **Note:** AP is computed on a collection of images;
a real/fake decision on a per-image basis is more difficult
- Situation may not persist
 - GANs train with a discriminator
 - Future architecture changes
 - “Shallow” fakes, e.g., Photoshop

Media manipulation example





https://www.youtube.com/watch?v=5Qqv_C6iVvQ

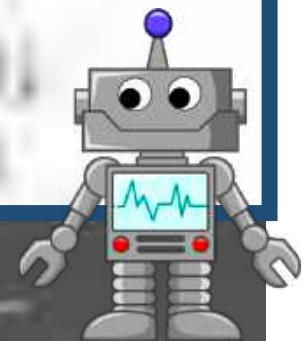


```
var data = new ArrayBuffer(10);
```

```
var dataView = new DataView(data);
```

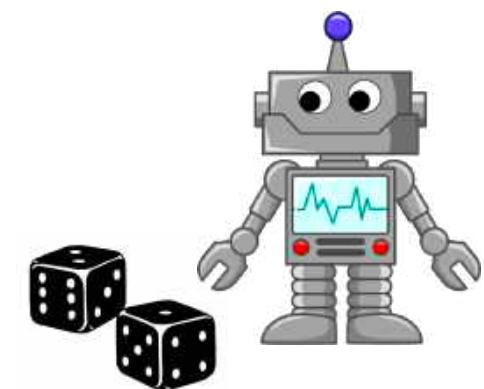
```
data.setFloat32(0, 1234567890.123456789);
```

```
dataView.getFloat32(0);
```





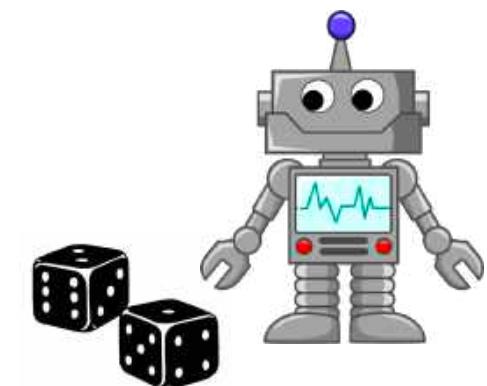
Original

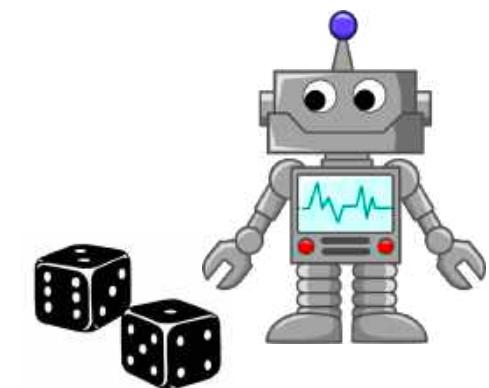


#1 modification

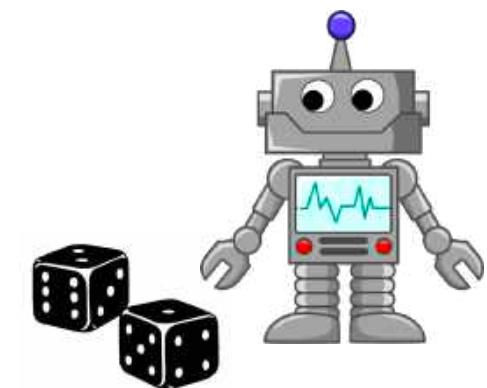


#2 modification



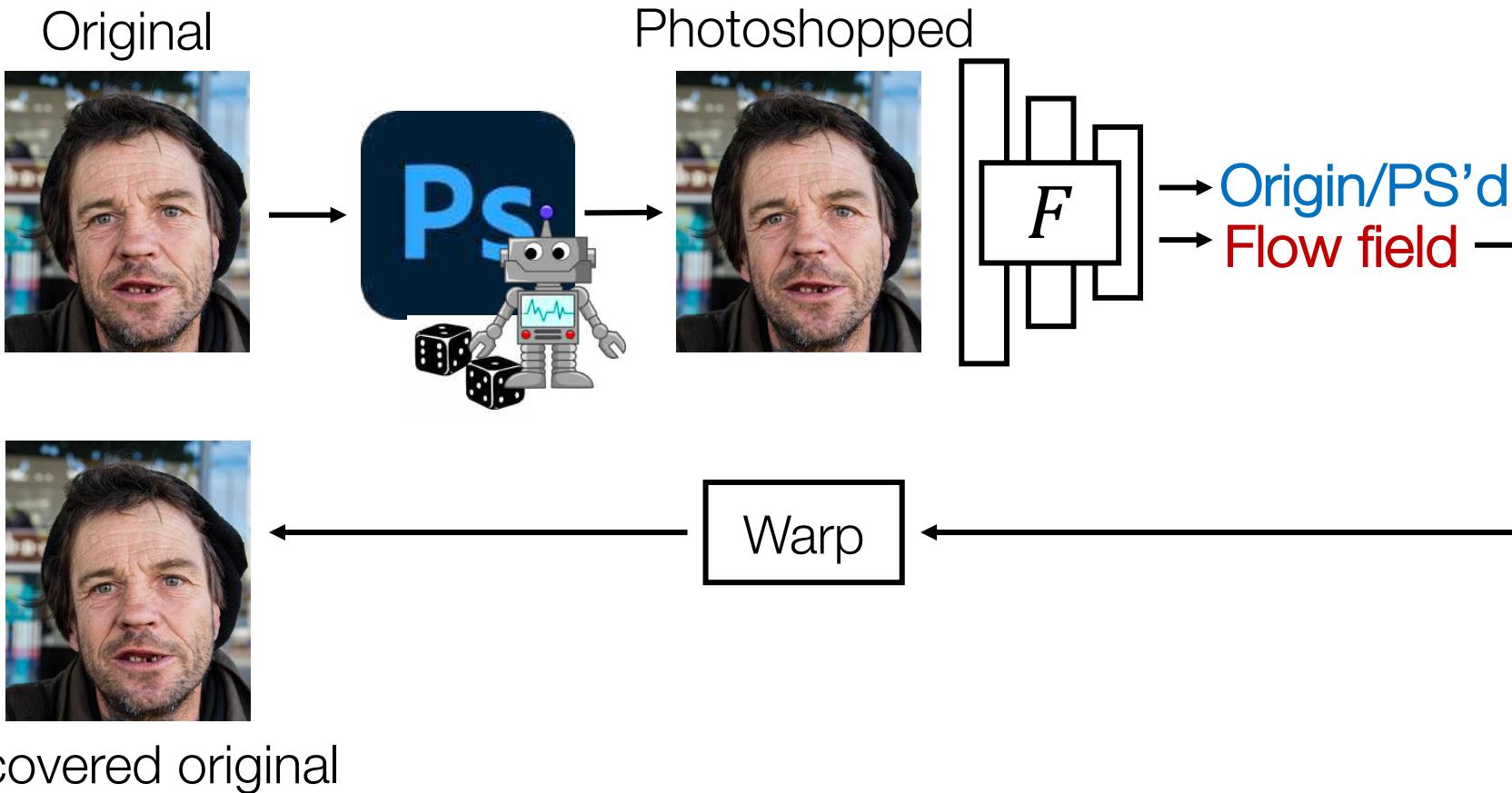


#3 modification

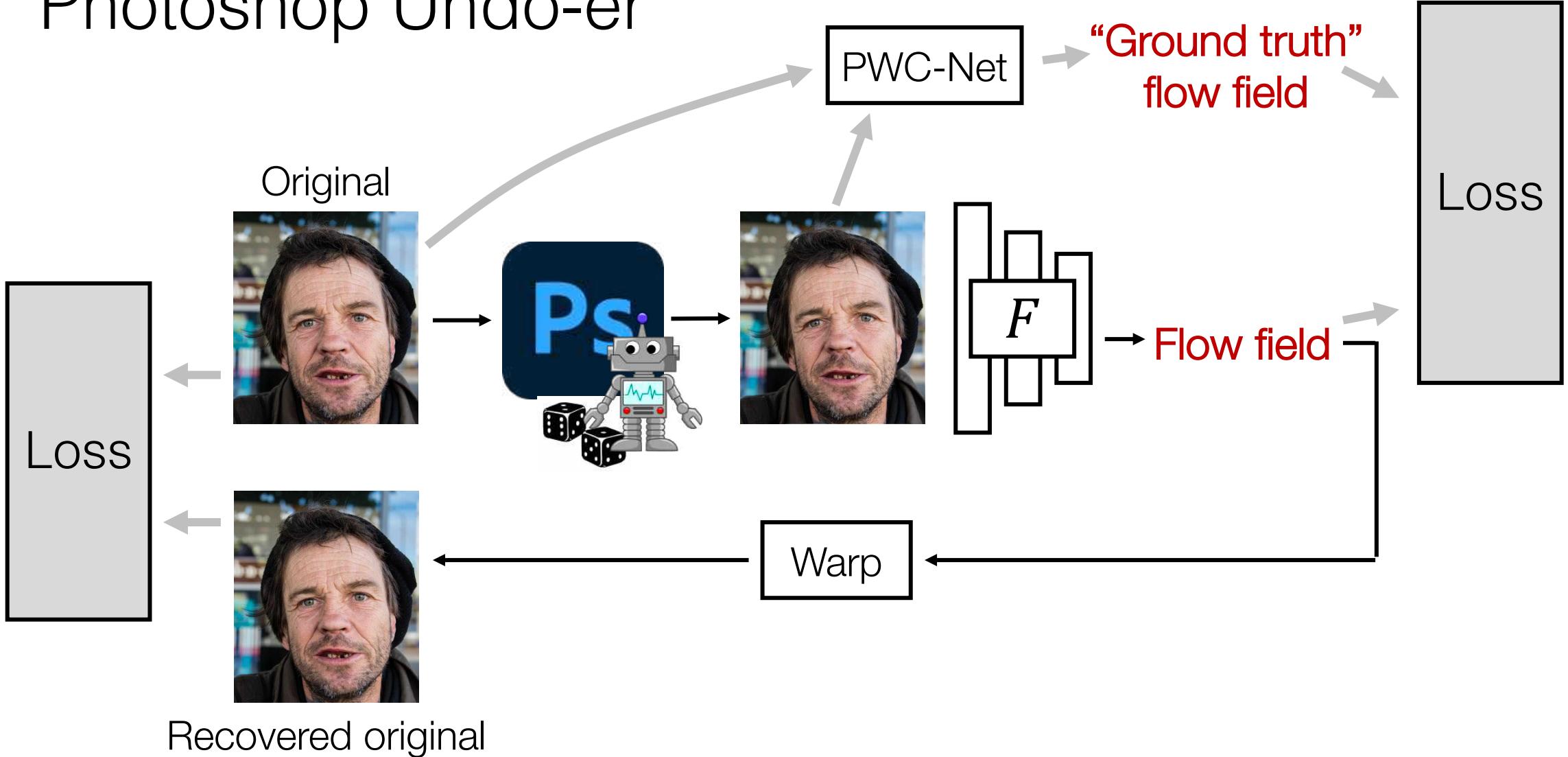


#4 modification

Photoshop Undo-er



Photoshop Undo-er





Manipulated



Flow prediction



Suggested “undo”



Original



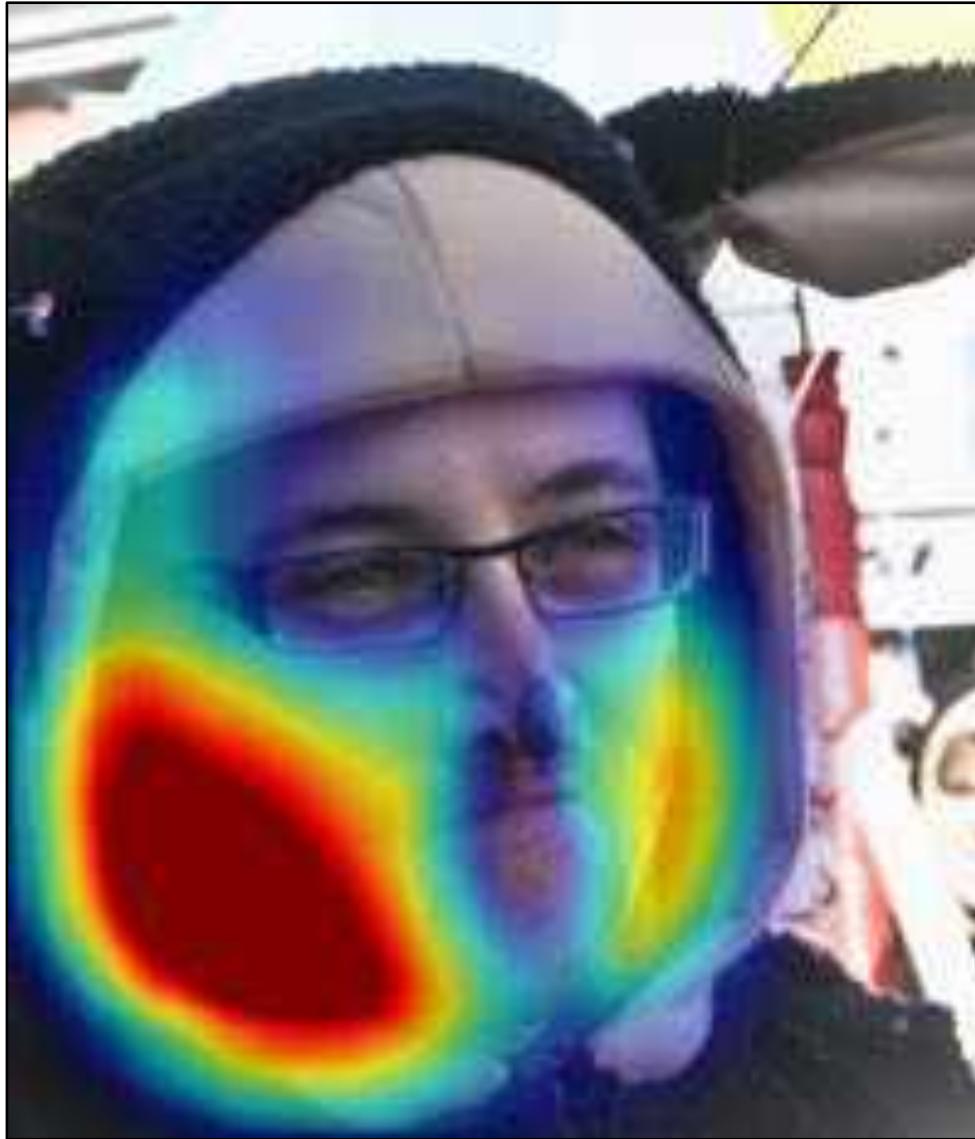
Manipulated vs. Original



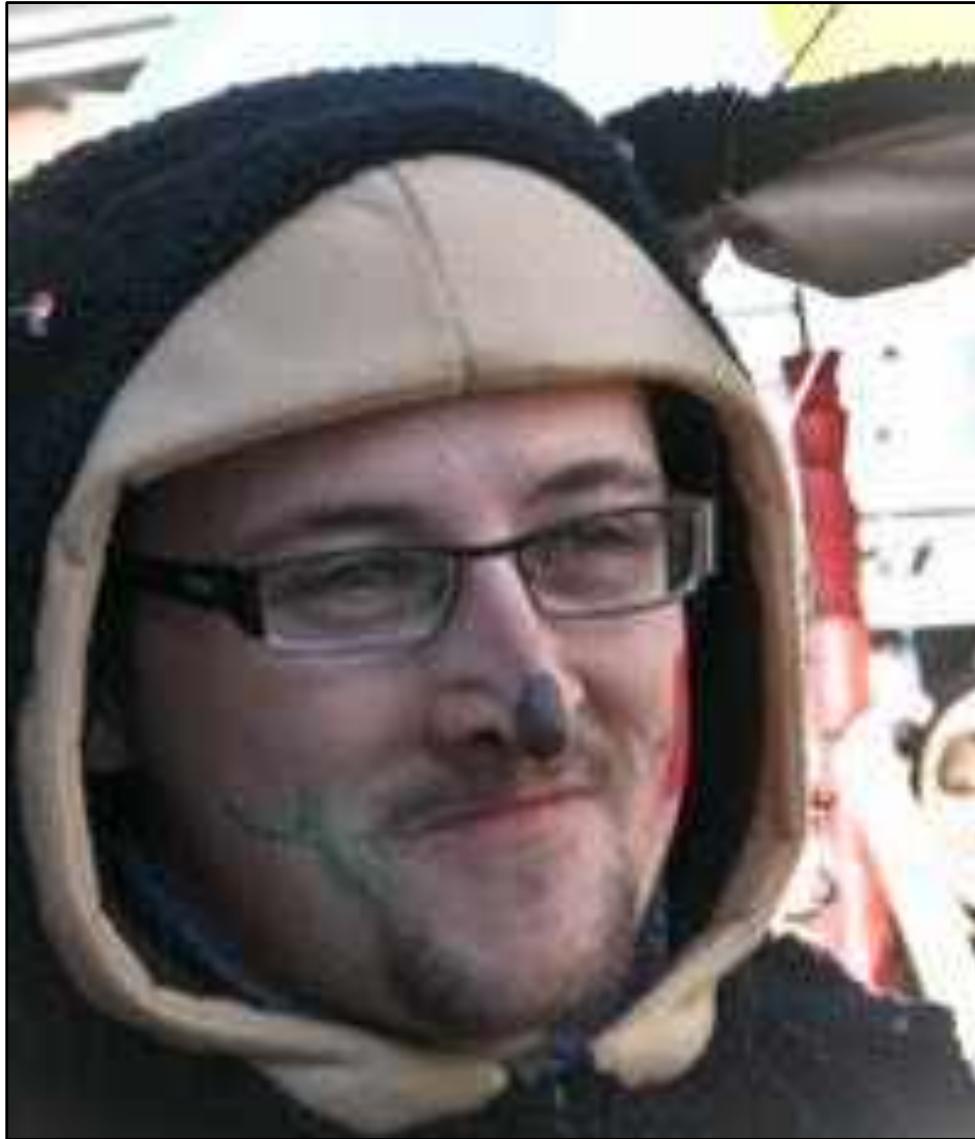
Undo vs. Original



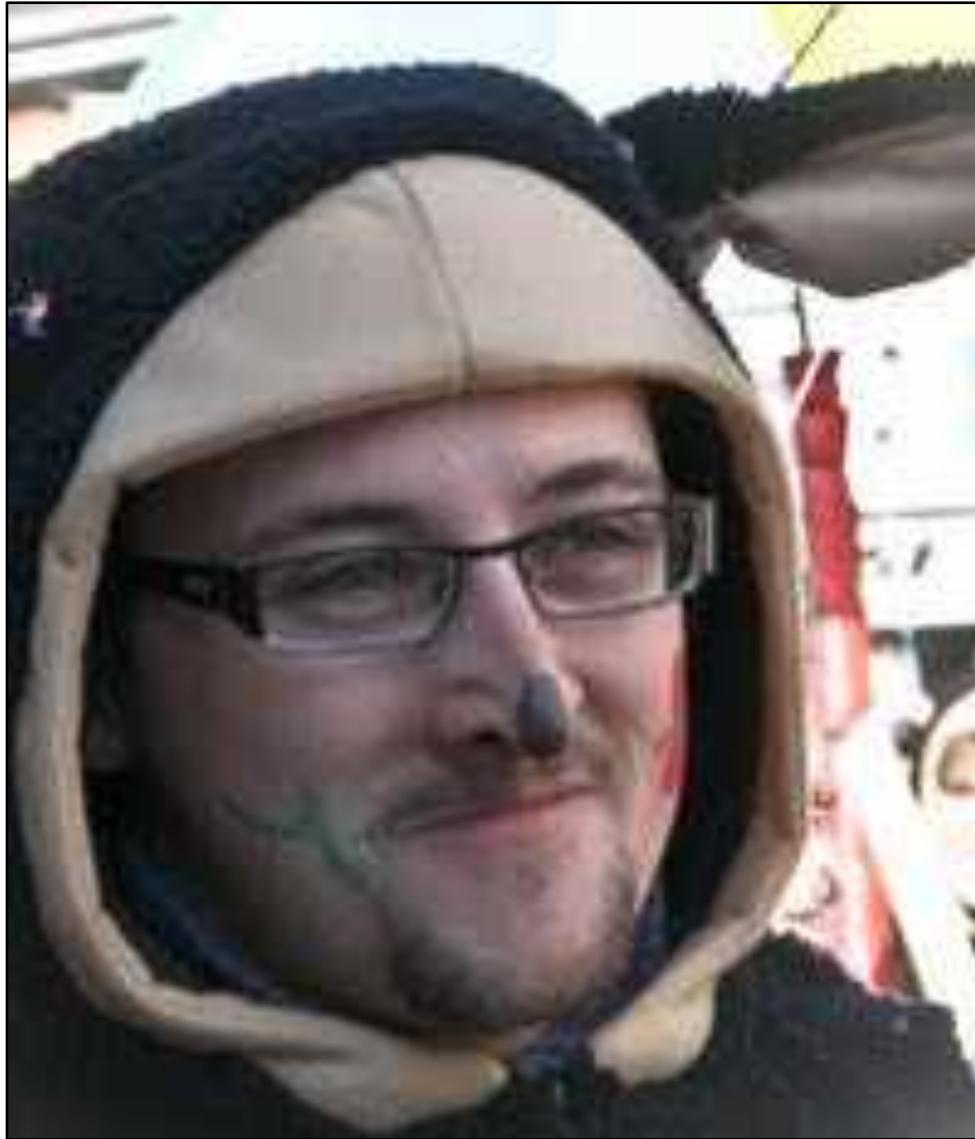
Manipulated



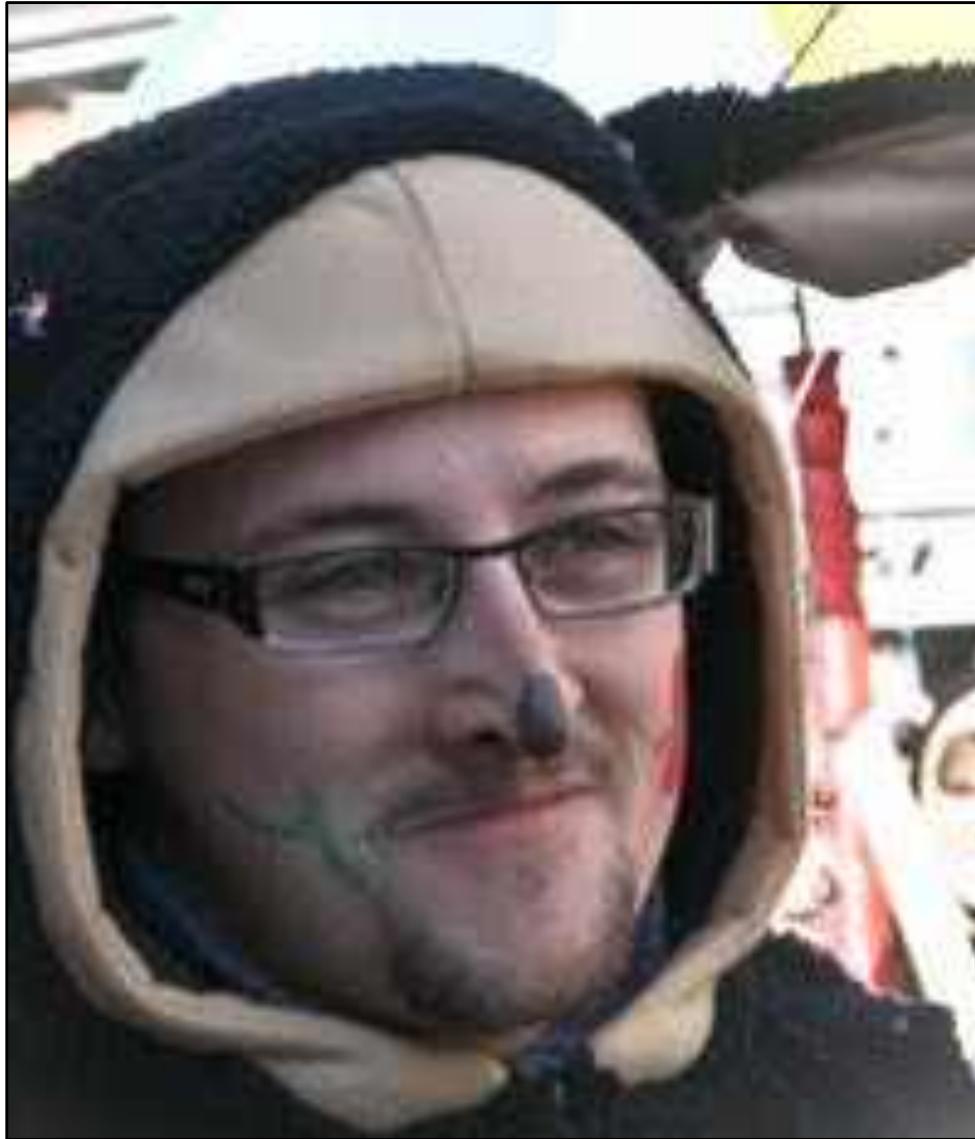
Flow prediction



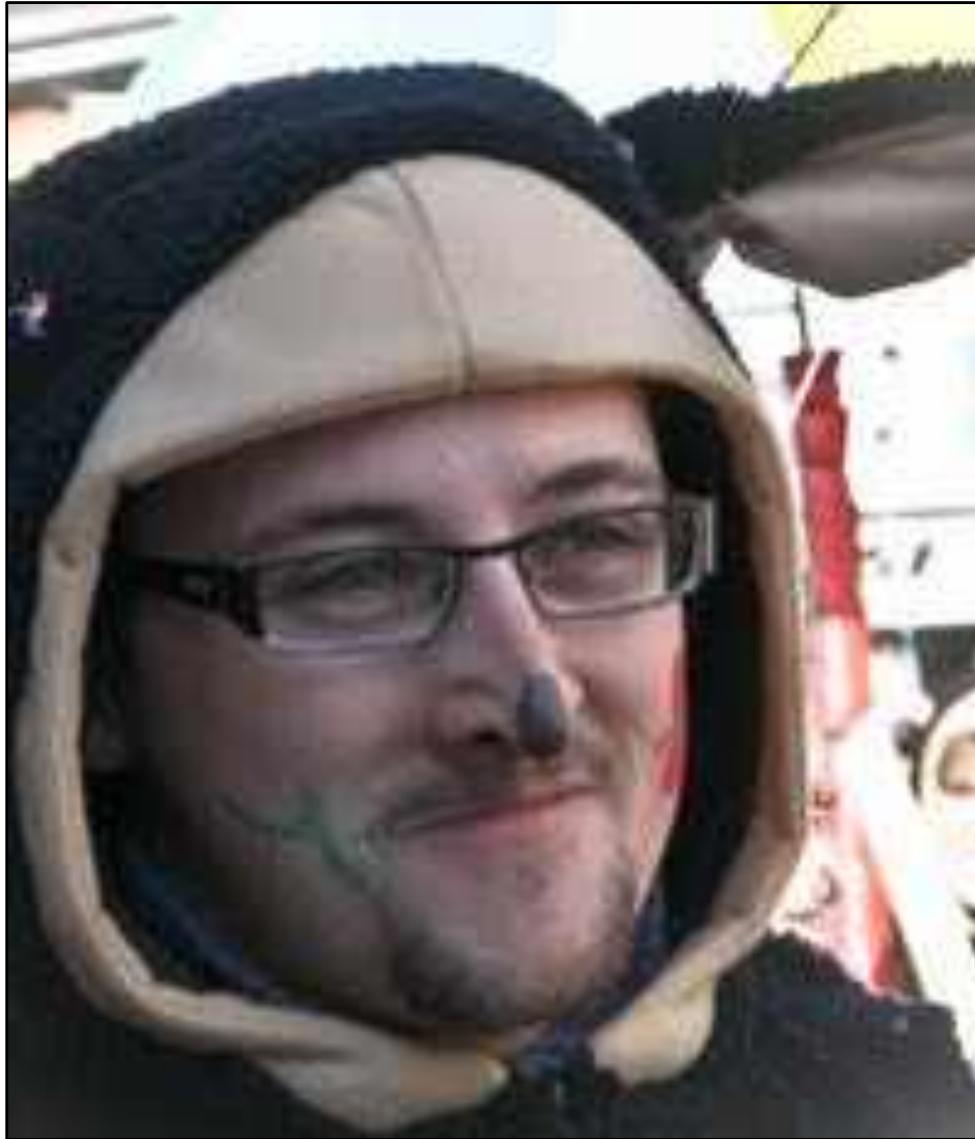
Suggested “undo”



Original



Manipulated vs. Original



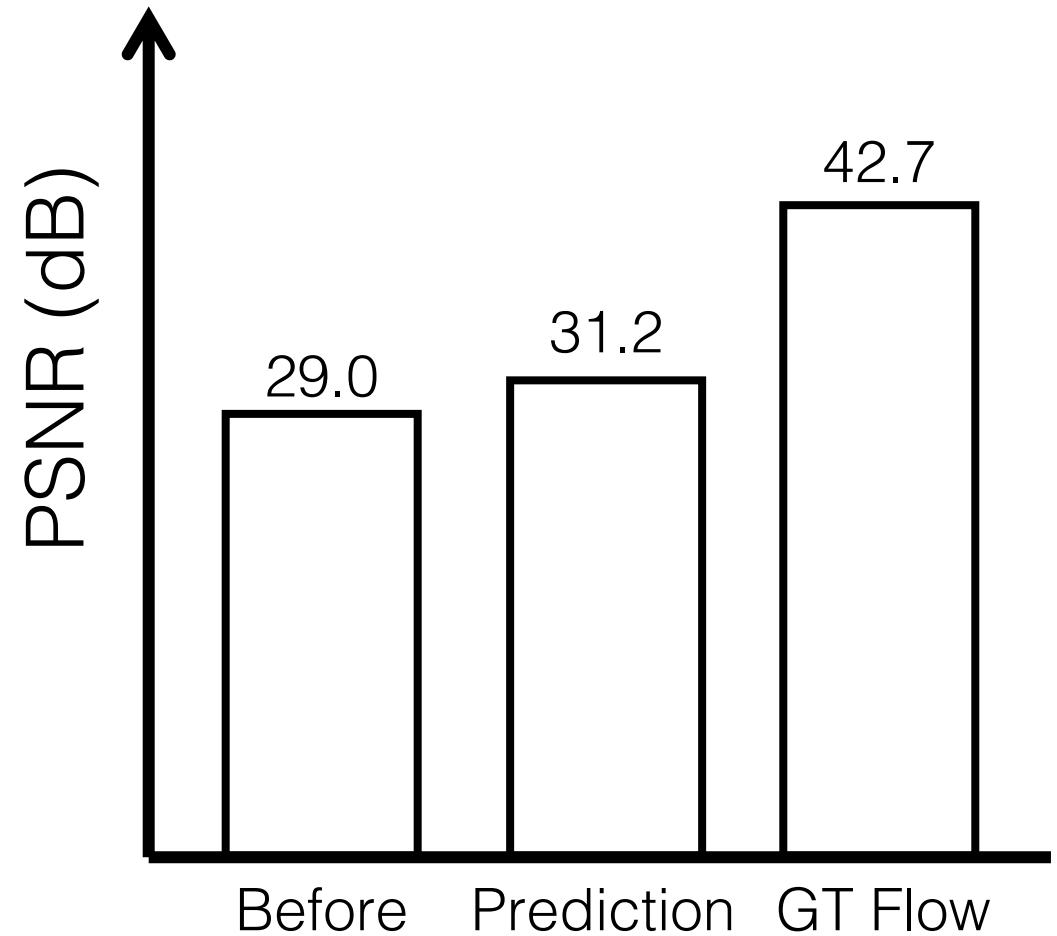
Undo vs. Original

Reversal evaluation

Senses of generalization

- Heldout artist data

Held-out artist generated data

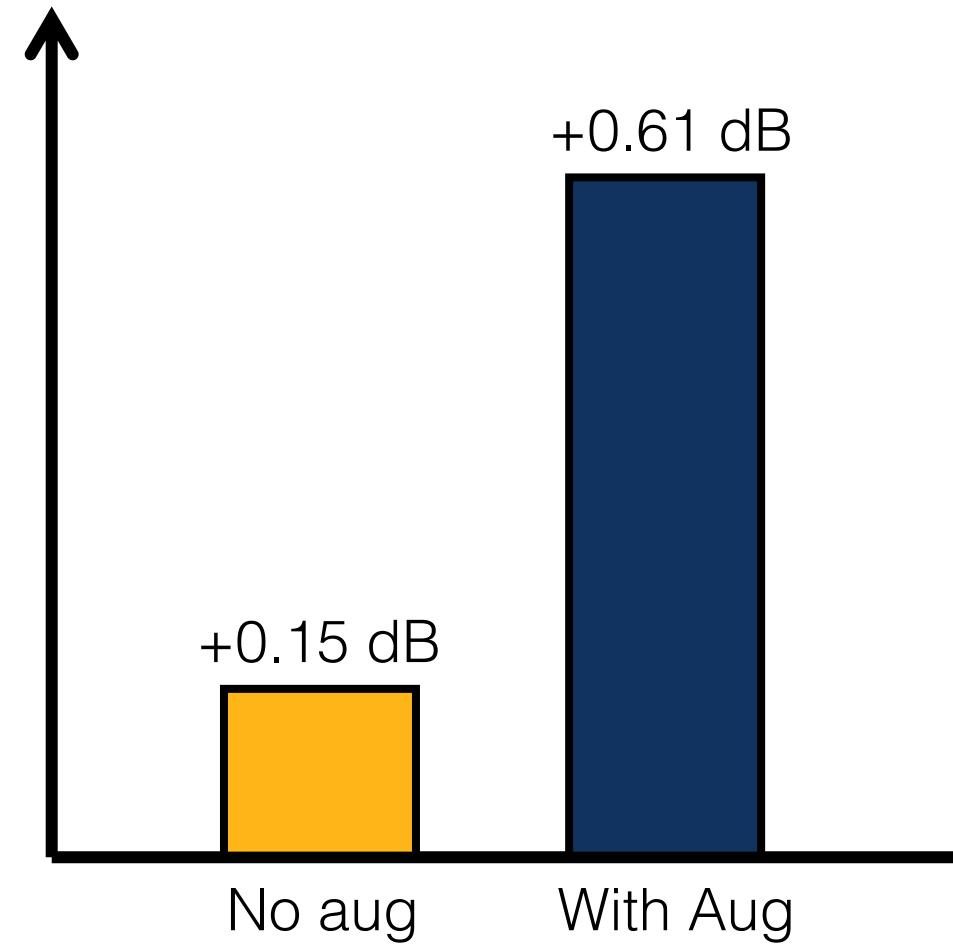


Reversal evaluation

Senses of generalization

- Heldout artist data
- Post-processing

Facebook post-processing



Data augmentation important (again)

Snapchat warps



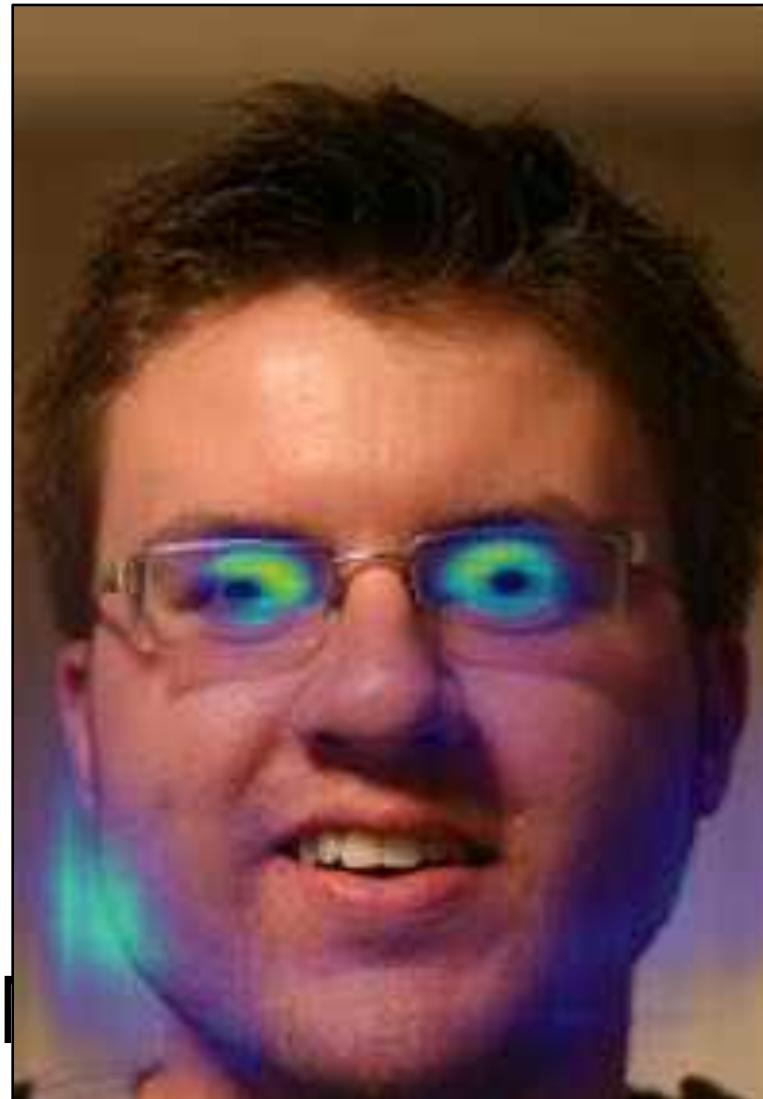
Original Photo

Snapchat warps



Manipulated Photo

Snapchat warps



Flow Prediction

Snapchat warps



Suggested “Undo”

Snapchat warps



Some generalization across warp methods

Original Photo

Different image domain



Different image domain



Different image domain

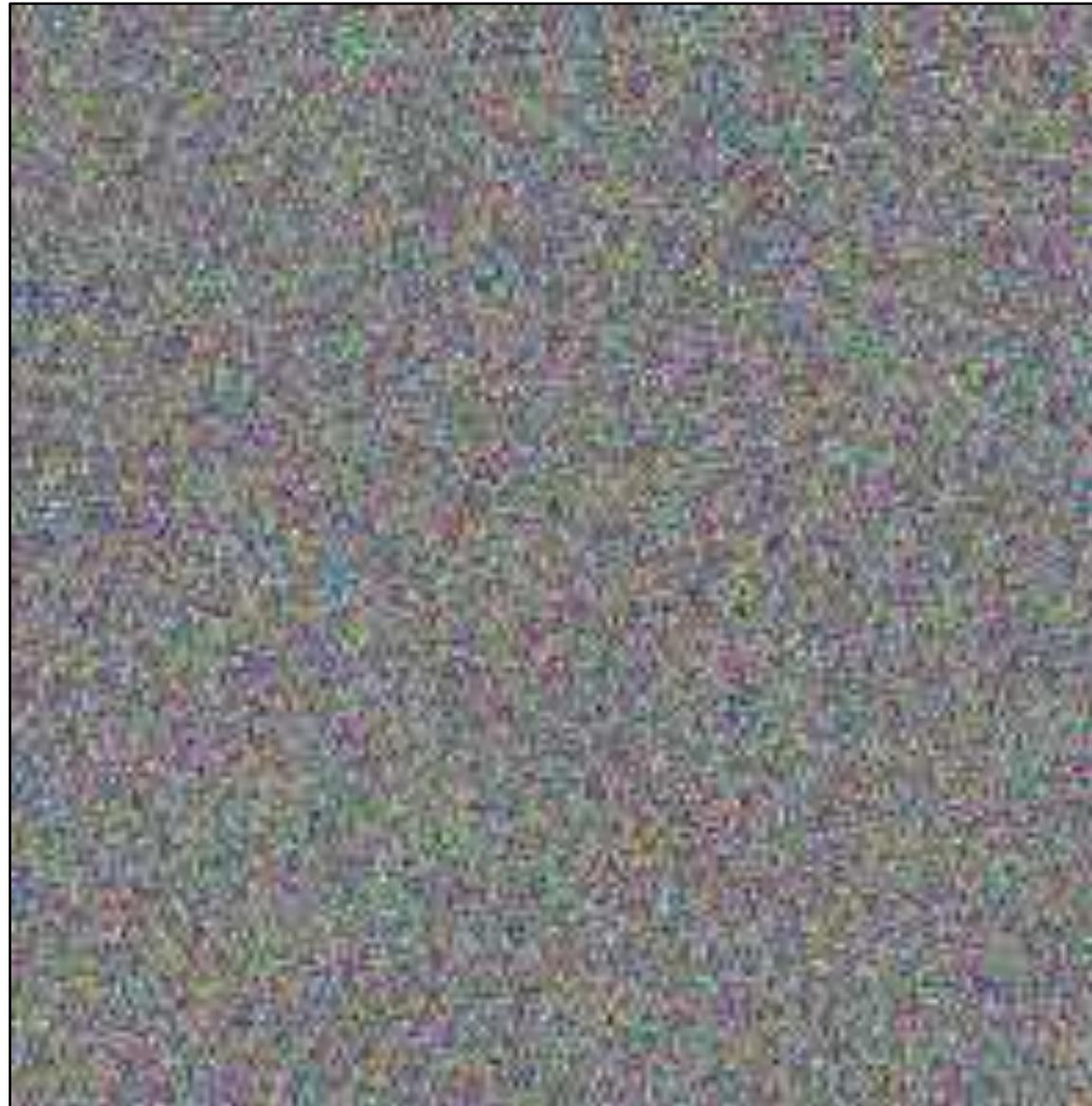


Predicted warp (not successful)

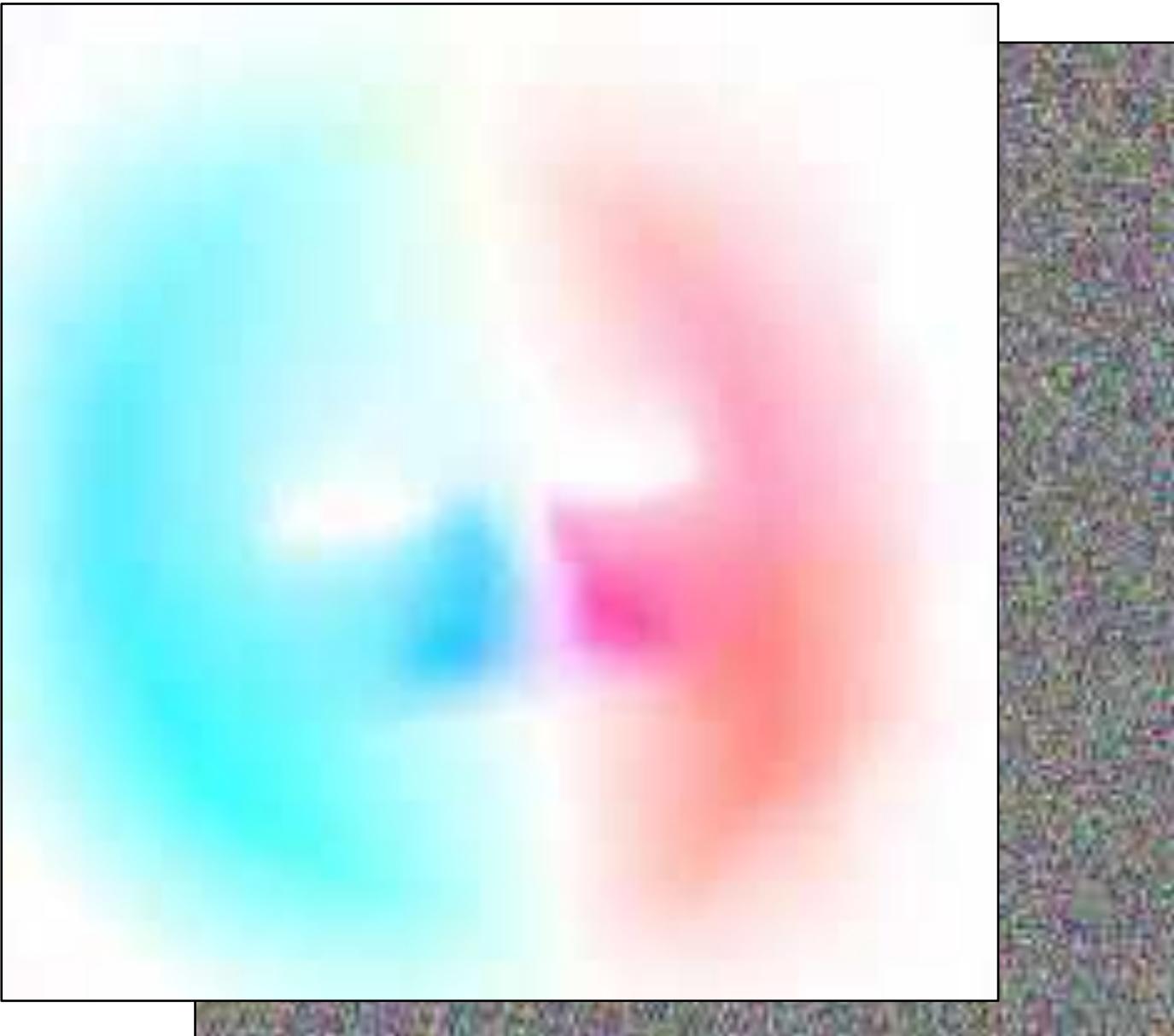


Does not generalize well to arbitrary image

Warped noise



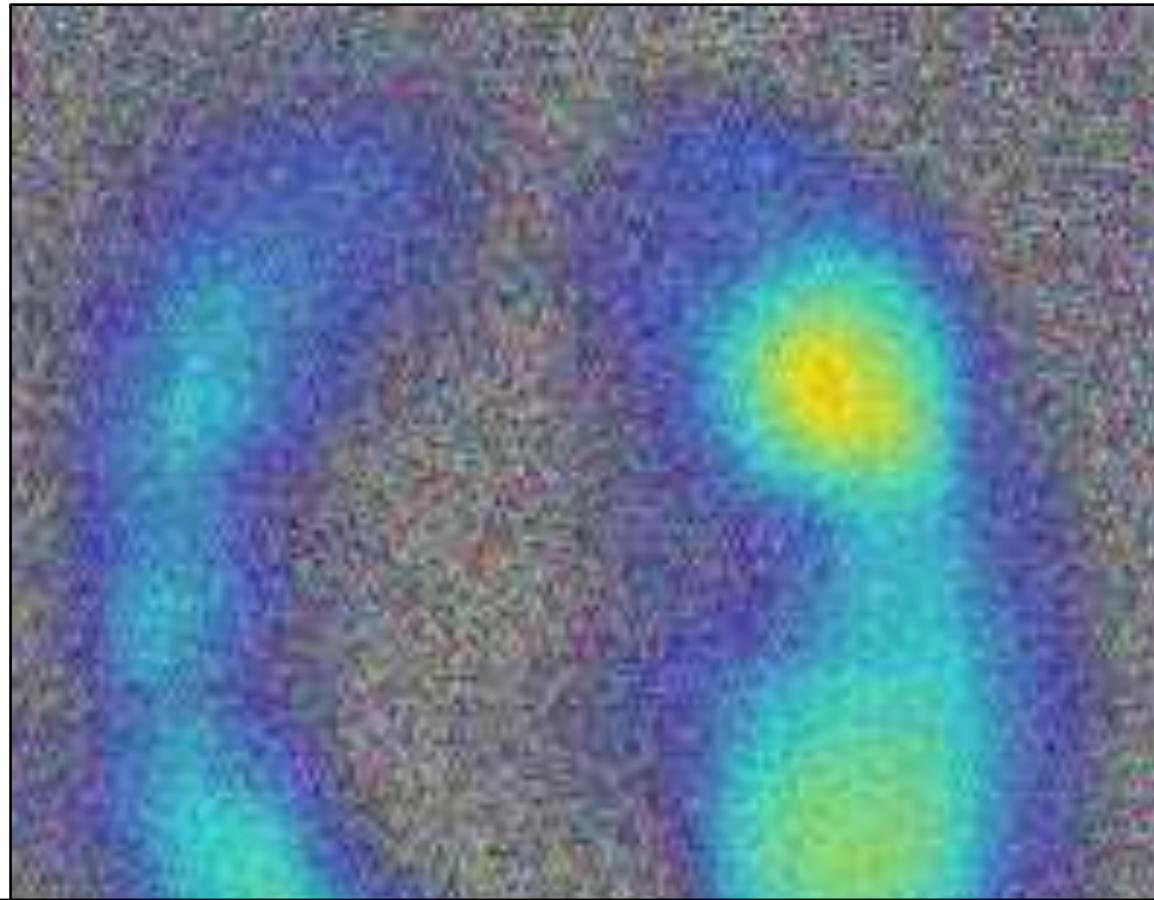
Warped noise



Warped noise



Warped noise

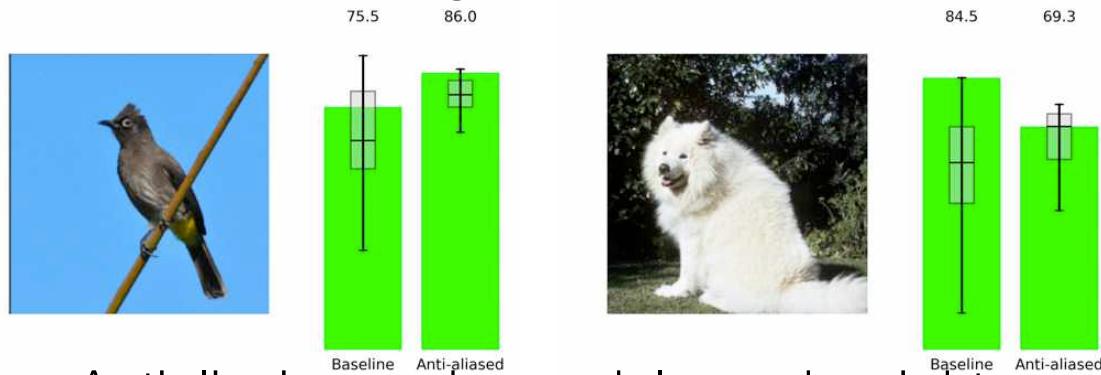


Successfully identifies warp from noise map

Discussion

- Suggests a multi-prong approach
- For rapidly evolving tools, continuous training and generalize
- For a relatively static tool, directly specialize
- Detection is only a piece of the puzzle
 - e.g., Content Authenticity Initiative: <https://contentauthenticity.org/>, collaboration between Adobe, New York Times, and Twitter

Making convolutional networks
shift-invariant again.
R. Zhang. In ICML, 2019.



Antialiasing code, models, and weights
richzhang.github.io/antialiased-cnns/

CNN-generated images are
surprisingly easy to spot...for now.
Wang, Wang, Zhang, Owens, Efros. In CVPR, 2020.

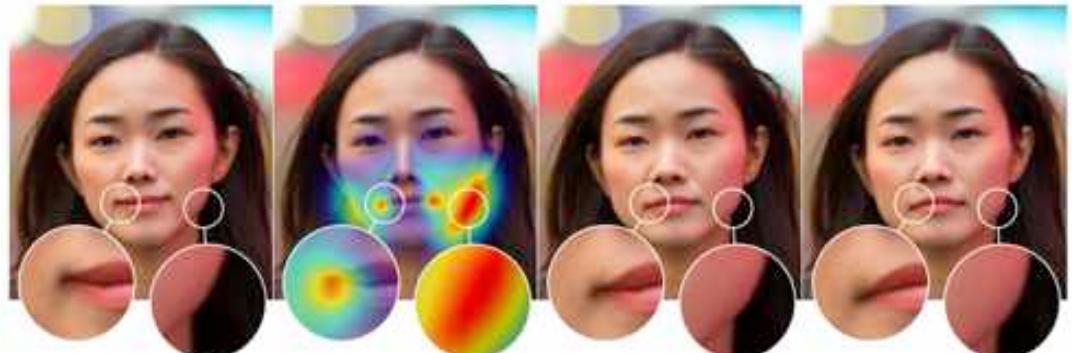


peterwang512.github.io/CNNDetection/



Detecting Photoshopped Images by
Scripting Photoshop.

Wang, Wang, Owens, Zhang, Efros. In ICCV, 2019.



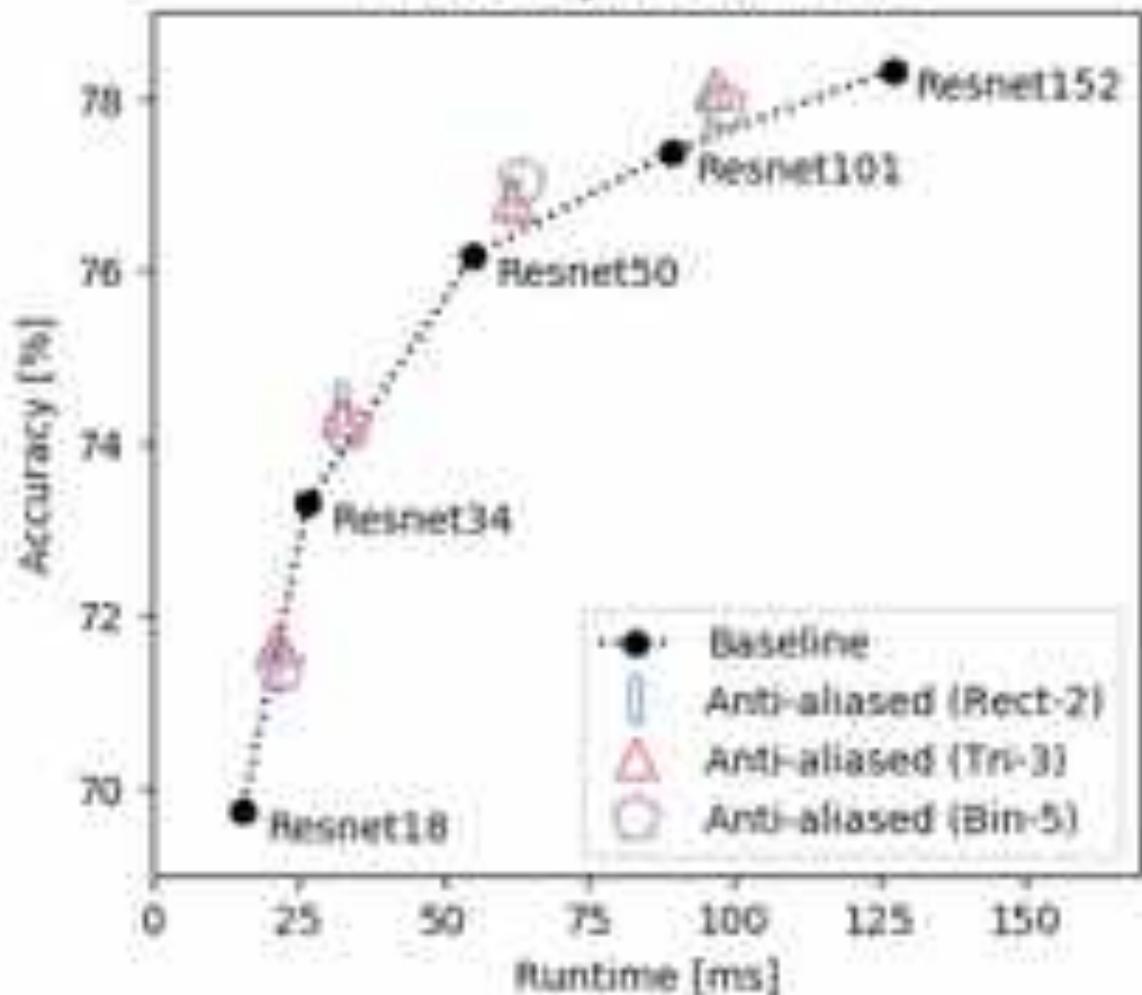
peterwang512.github.io/FALdetector/

Thank You!

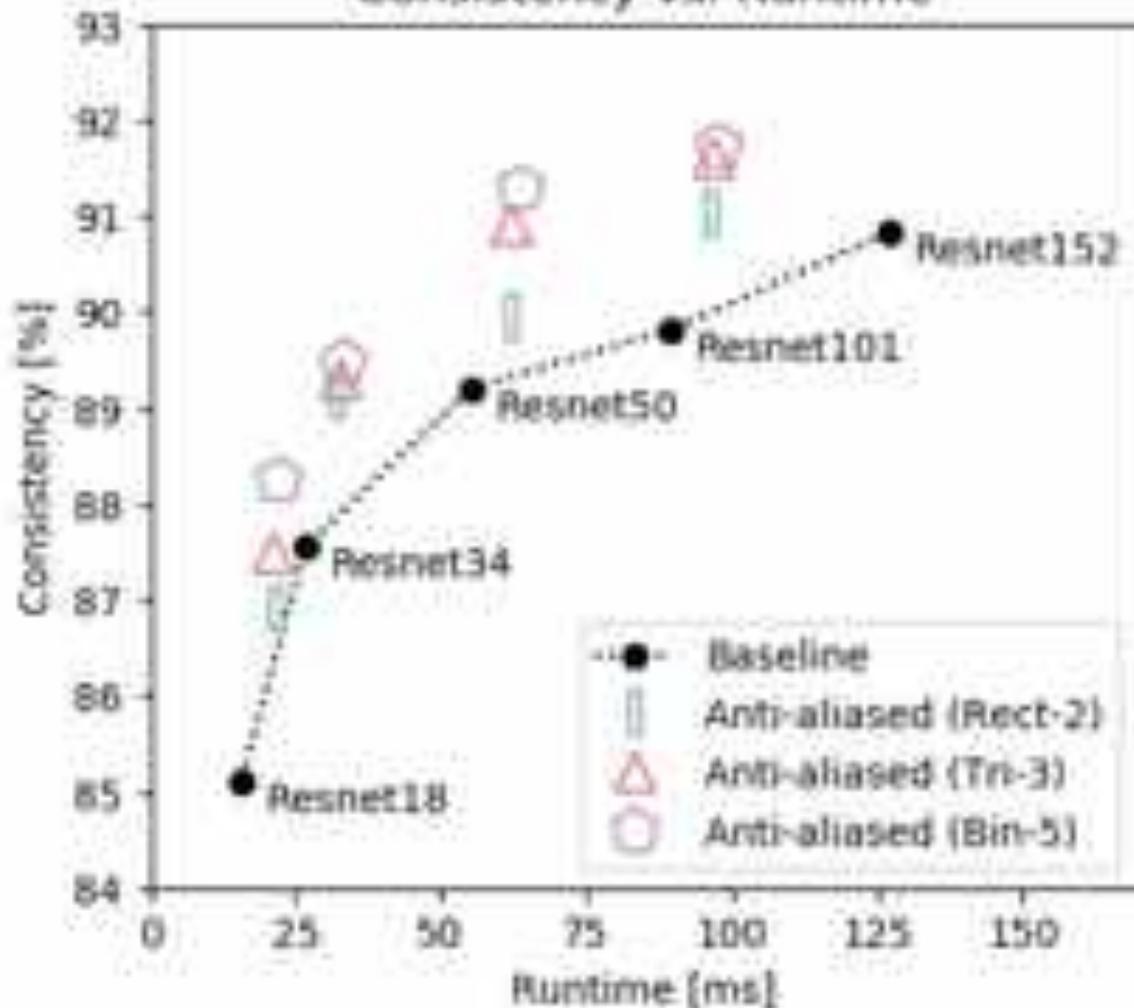


Backup

Accuracy vs. Runtime



Consistency vs. Runtime



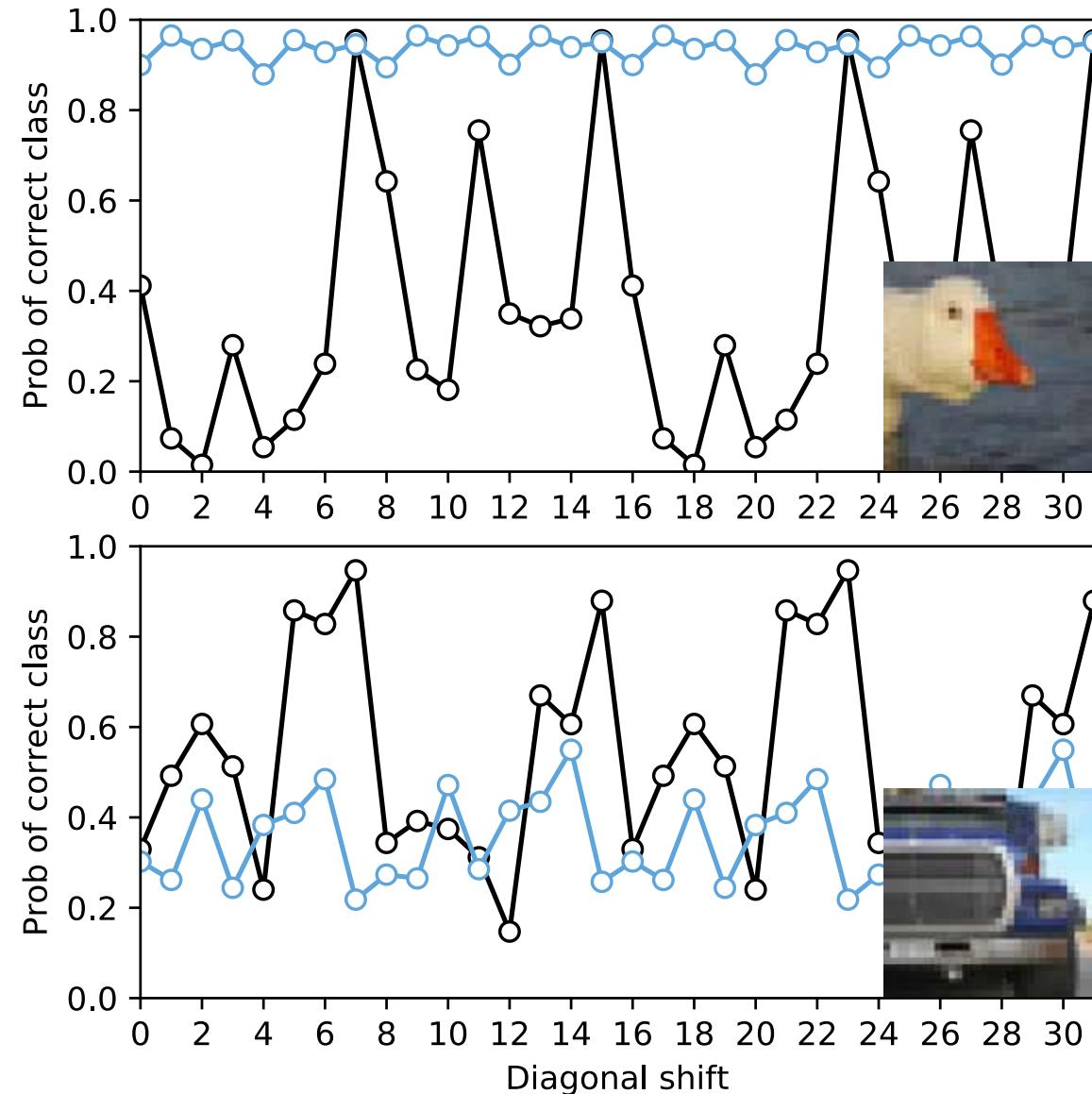
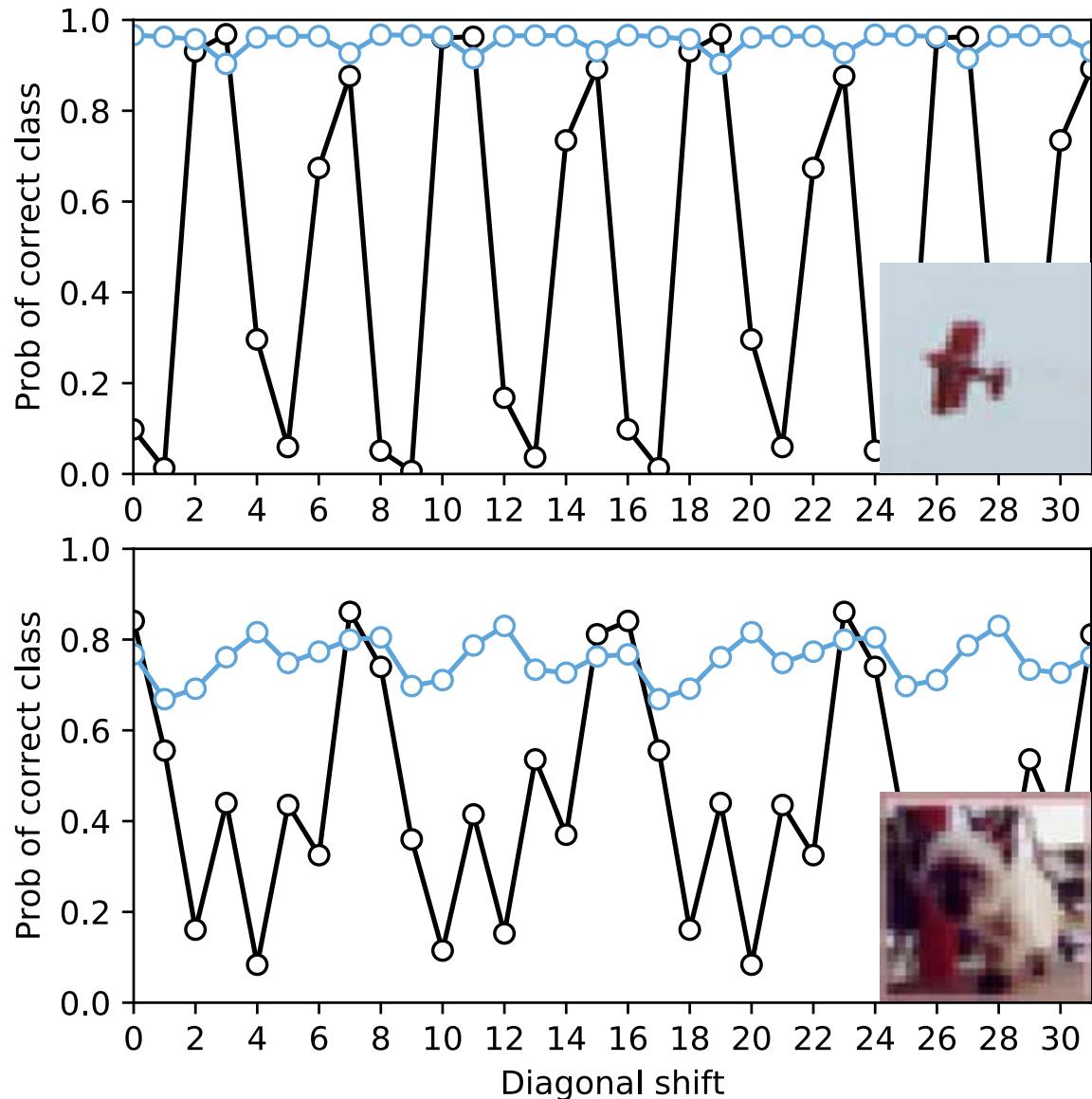
ResNet50 on ImageNet-P (Hendrycks et al., 2019)

	Flip Rate (FR) (lower is better)											
	Noise		Blur		Weather		Geometric				Mean	
	Gauss	Shot	Motion	Zoom	Snow	Bright	Translate	Rotate	Tilt	Scale	Unnorm	Norm
Baseline	14.04	17.38	6.00	4.29	7.54	3.03	4.86	6.79	4.01	11.32	7.92	57.99
Rect-2	14.08	17.16	5.98	4.21	7.34	3.20	4.42	6.43	3.80	10.61	7.72	56.70
Tri-3	12.59	15.57	5.39	3.79	6.98	3.01	3.95	5.80	3.53	9.90	7.05	51.91
Bin-5	12.39	15.22	5.44	3.72	6.76	3.15	3.78	5.67	3.44	9.45	6.90	51.18

ResNet50 on ImageNet-C (Hendrycks et al., 2019)

	Corruption Error (CE) (lower is better)																
	Noise			Blur				Weather				Digital				Mean	
	Gauss	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	Jpeg	Unnorm	Norm
Baseline	68.70	71.10	74.04	61.40	73.39	61.43	63.93	67.76	62.08	54.61	32.04	61.25	55.24	55.24	46.32	60.57	76.43
Rect-2	65.81	68.27	70.49	60.01	72.14	62.19	63.96	68.00	61.83	54.95	32.09	60.25	55.56	53.89	43.62	59.54	75.16
Tri-3	63.86	66.07	69.15	58.36	71.70	60.74	61.58	66.78	60.29	54.40	31.48	58.09	55.26	53.89	43.62	58.35	73.73
Bin-5	64.31	66.39	69.88	60.31	71.37	61.60	61.25	66.82	59.82	51.84	31.51	58.12	55.29	50.81	42.84	58.14	73.41

Baseline vs Anti-aliased



Loss Function

- End point error :

$$\sum_{p \in I} ||U(p) - \hat{U}(p)||_2$$

Loss Function

- Multi-scale gradient error:

$$\sum_{p \in I} \sum_s (\|\nabla_x^s U(p) - \nabla_x^s \hat{U}(p)\|_2 + \|\nabla_y^s U(p) - \nabla_y^s \hat{U}(p)\|_2)$$

$$\nabla_x^s A \equiv A[:, s :] - A[:, : -s]$$

Loss Function

- Reconstruction error:

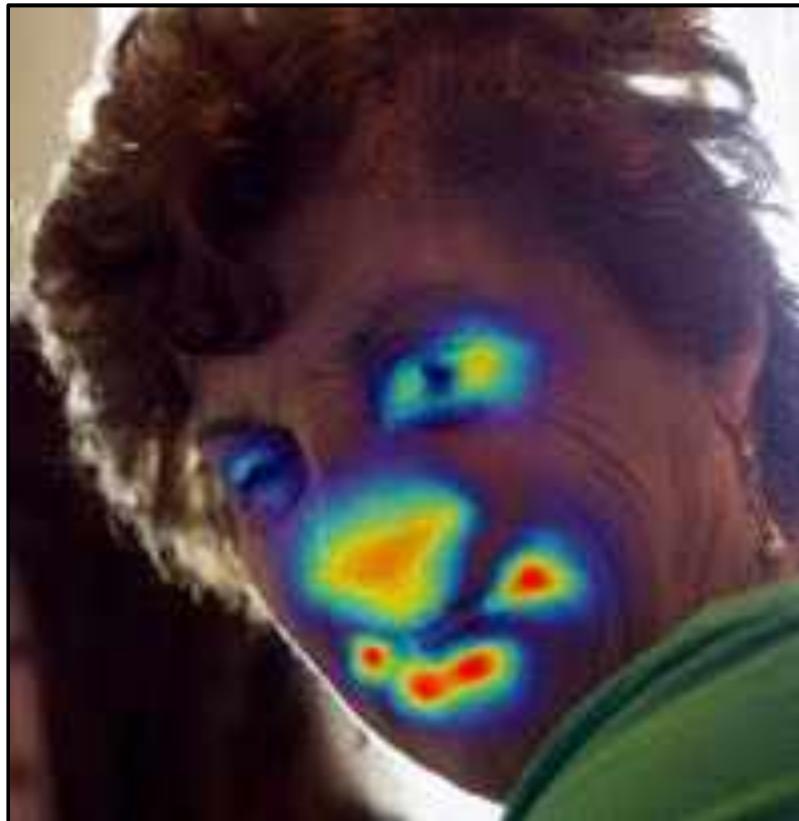
$$\sum_{p \in I} \|X_{orig}(p) - X_{reconst}(p)\|_1$$



Loss Function

- Reconstruction error:

$$\sum_{p \in I} \|X_{orig}(p) - X_{reconst}(p)\|_1$$



Loss Function

- Reconstruction error:

$$\sum_{p \in I} ||X_{orig}(p) - X_{reconst}(p)||_1$$

