



# Image Inpainting

Yan Zhao

May 21, 2018

# Contents:

## 1 Introduction to Image Inpainting

- What is Image Inpainting
- Introduction to Image Inpainting Methods
- Applications of Image Inpainting

# Contents:

## 2 Image Quality Assessment

- PSNR
- SSIM

# Contents:

## 3 Classical Methods to Image Inpainting

- Examplar-based Inpainting
- OpenCV

# Contents:

## 4 Deep Learning to Image Inpainting

- Datasets
- Methods in Recent Papers

# Methods in Recent Papers:

- Globally and Locally Completion
- Semantic Image Inpainting
- High-Resolution Image Inpainting
- Content Encoder

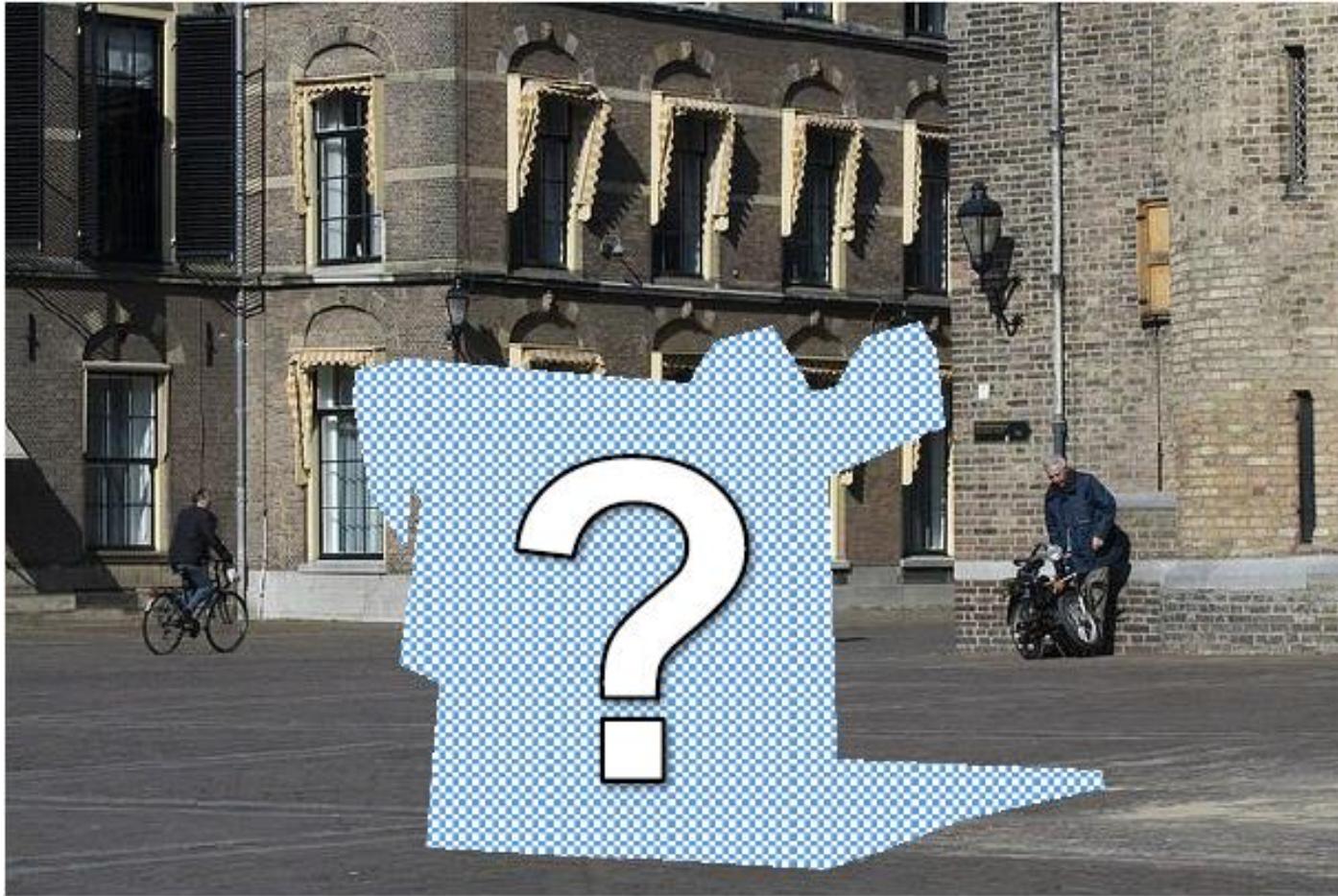


If your photo is damaged, what would you do?

# The image completion problem

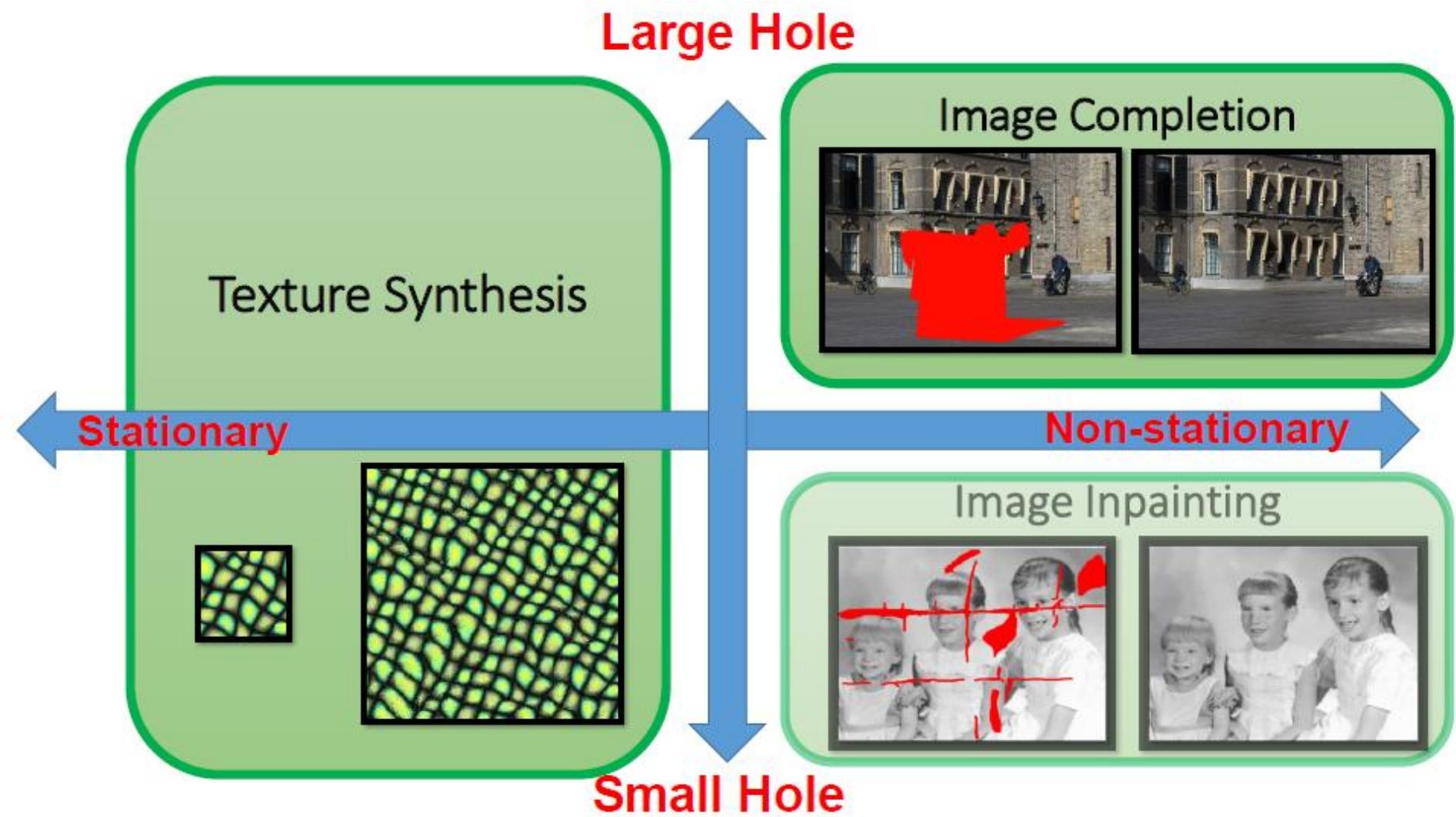


# The image completion problem



# The image completion problem





# What is Image Inpainting ?



Original and restored image

Inpainting is the process of reconstructing lost or deteriorated parts of **images** and **videos**.

# What is Image Inpainting ?



# Objective: Make attentive inpainter



IF BACKGROUND OF TARGET REMOVING OBJECT IS SIMPLE, EXISTING METHOD WORKS FINE



HOWEVER, IF BACKGROUND OF TARGET REMOVING  
OBJECT IS COMPLEX, BETTER NEED ANOTHER METHOD

# What is Image Inpainting ?

- In the **museum** world: in the case of a valuable painting, this task would be carried out by a skilled art conservator or art restorer.
- In the **digital** world: inpainting (also known as *image interpolation* or *video interpolation*) refers to the application of sophisticated **algorithms** to replace lost or corrupted parts of the image data (mainly small regions or to remove small defects).

In the museum word:



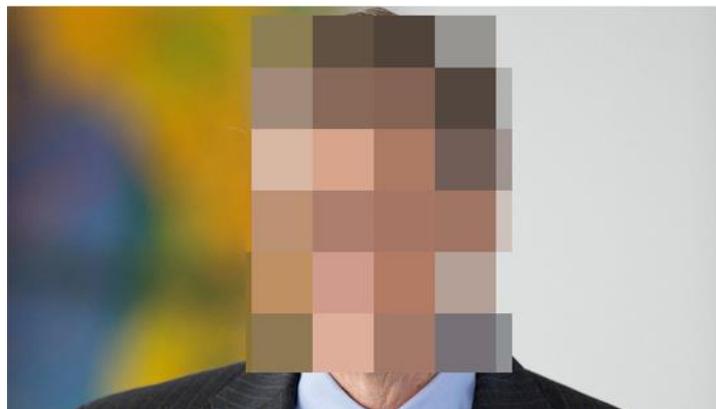
# In the digital word:



**Photo repair**



**De-watermarking**



**De-masking**



**Target removal**

# Image Inpainting Methods

Inpainting is rooted in the restoration of images. Traditionally, inpainting has been done by professional restorers. The underlying methodology of their work is as follows:

- The global picture determines how to fill in the gap. The purpose of inpainting is to restore the unity of the work.
- The structure of the gap surroundings is supposed to be continued into the gap. Contour lines that arrive at the gap boundary are prolonged into the gap.
- The different regions inside a gap, as defined by the contour lines, are filled with colors matching for those of its boundary.
- The small details are painted, i.e. “**texture**” is added.

# Image Inpainting Methods

- Structural inpainting
- Textural inpainting
- Combined Structural and textural inpainting

# Structural inpainting

- Use **geometric** approaches for filling in the missing information in the region which should be inpainted.
- These algorithms focus on the consistency of the geometric structure.

# Textural inpainting

- All the structural inpainting methods are not able to restore texture.
- **Texture** has a repetitive pattern which means that a missing portion cannot be restored by continuing the level lines into the gap.

# Combined Structural and textural inpainting

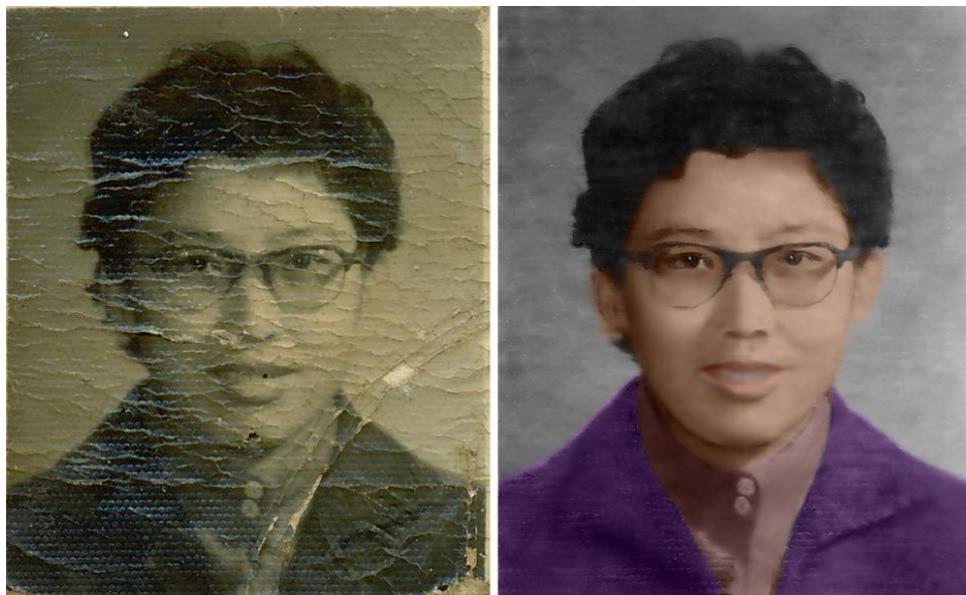
- Perform **texture and structure filling** in regions of missing image information.
- The boundaries between image regions accumulate structural information which is a complex phenomenon. This is the result when blending different textures together.
- The state of the art inpainting method attempts to combine structural and textural inpainting.

# Applications of Image Inpainting

## ➤ **Photo repair**

- De-watermarking
- De-mosaicing
- Object removal
- De-blurring

# Photo repair



# Photo repair - Adobe Photoshop



Resized  
→  
Imitation stamp



Original image

Repaired image

# Applications of Image Inpainting

- Photo repair
- **De-watermarking**
- De-mosaicing
- Object removal
- De-blurring

# De-watermarking – Adobe Photoshop



# De-watermarking – OpenCV (*Open Source Computer Vision*)

- A library of programming functions mainly aimed at real-time computer vision
- OpenCV supports the deep learning frameworks TensorFlow, Torch/PyTorch and Caffe.

OpenCV 's application areas include:

2D and 3D feature toolkits	Egomotion estimation
Facial recognition system	Gesture recognition
Human–computer interaction (HCI)	Augmented reality
Mobile robotics	Motion understanding
Object identification	Segmentation and recognition
Structure from motion (SFM)	Motion tracking

# De-watermarking – OpenCV



Original image

SIFT  
FLANN



Inpaint



Inpainted image



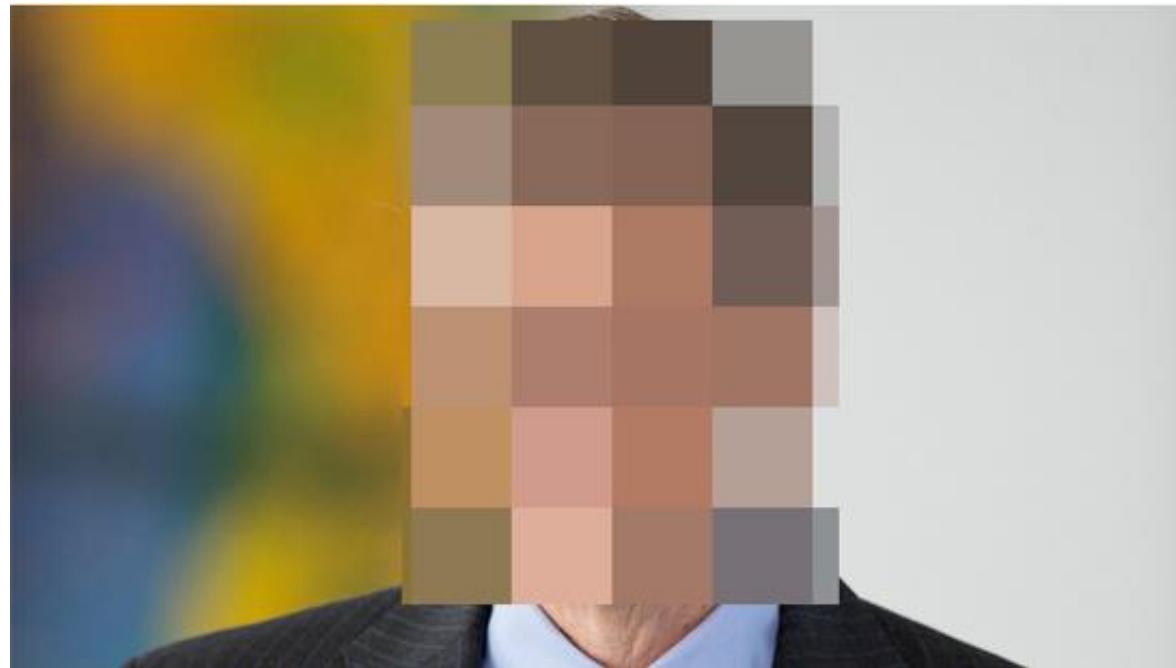
Masked image

# Applications of Image Inpainting

- Photo repair
- De-watermarking
- **De-mosaicing**
- Object removal
- De-blurring

# De-mosaicing

In the Internet era of increasingly sensitive personal information, Internet users are generally used to cover sensitive information on photos with “fuzzy” or “mosaic”.



# De-mosaicing VS Auto-mosaicing

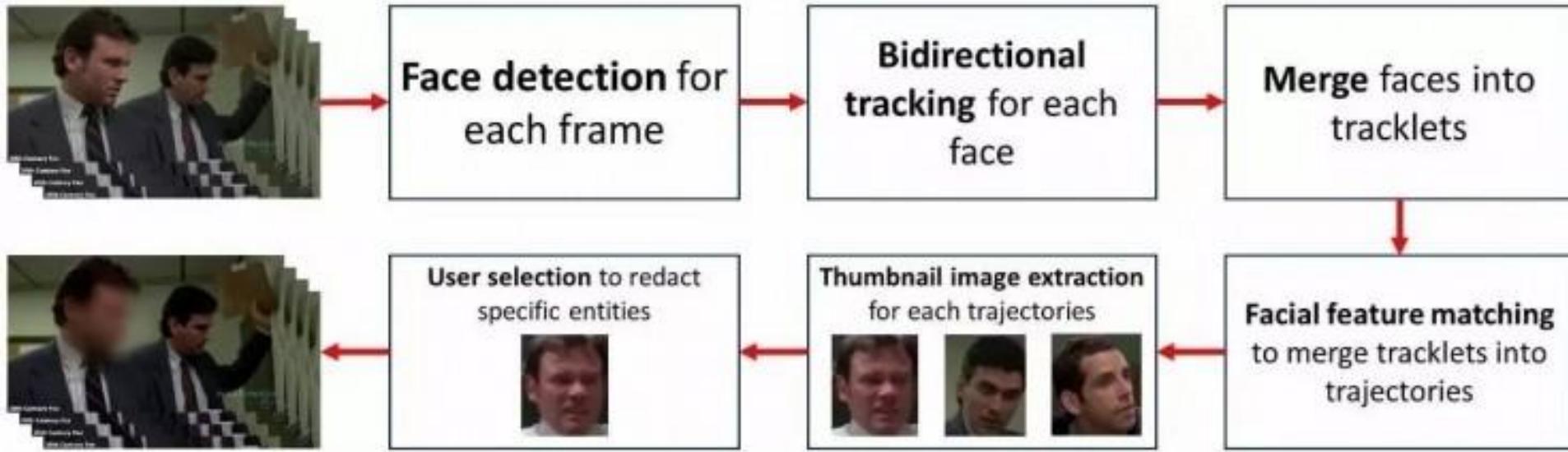


# Microsoft: Auto-mosaicing



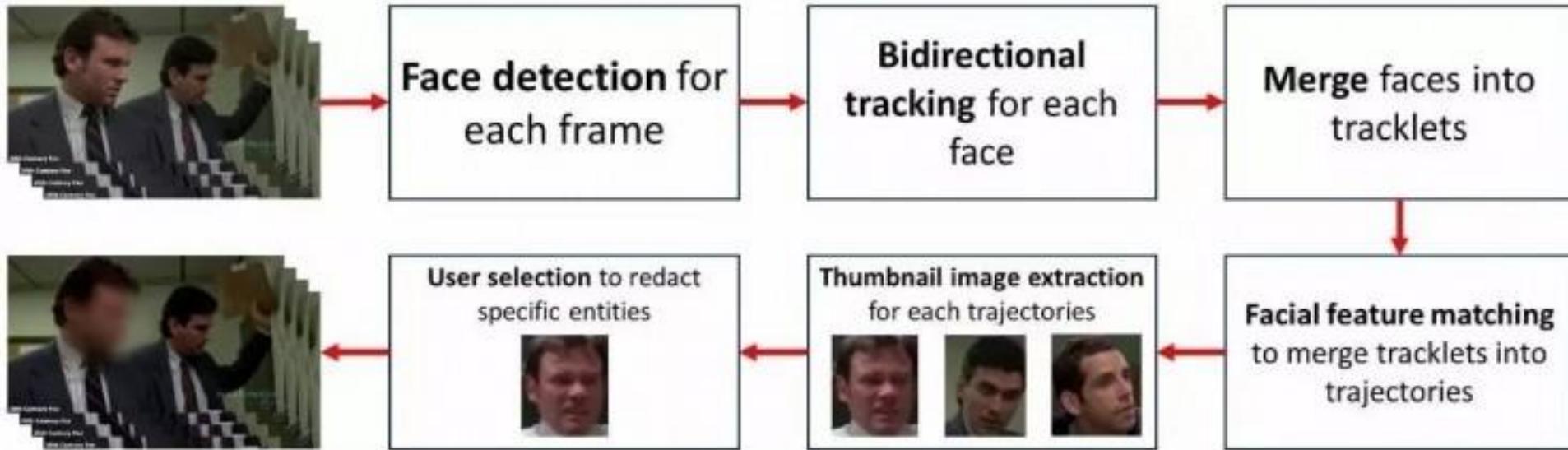
- Original video frame
- When the character appears, the user selects a person to mosaic
- Masked video frame

# Microsoft: Auto-mosaicing



- Find out where all the faces in the video are
- Connect all the faces of the same person.

# Microsoft: Auto-mosaicing

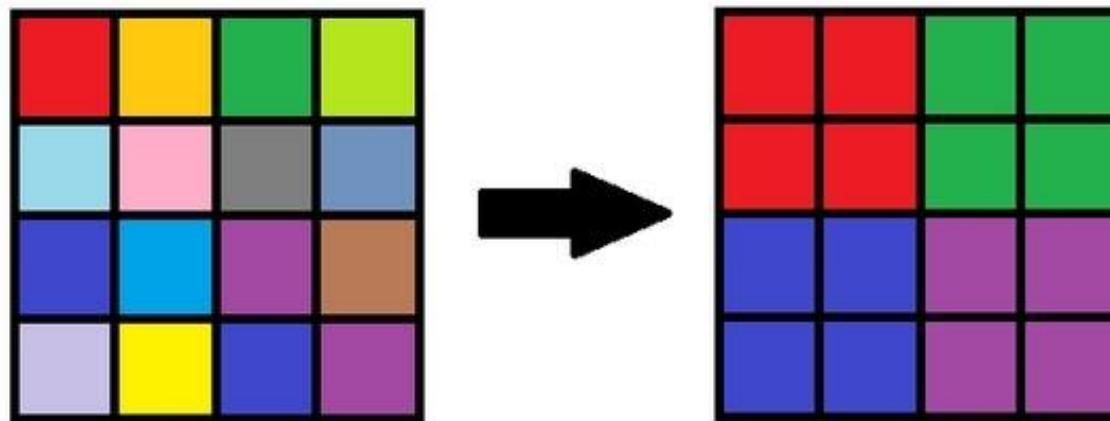


- Face detection
- Face tracking
- Face recognition

# De-mosaicing

In general, mosaic is irreversible loss information, it is difficult to restore.

Mosaicing is a process of reducing the amount of information that has been lost. It is equivalent to sampling the image signal at a lower frequency than the original data.



# De-mosaicing

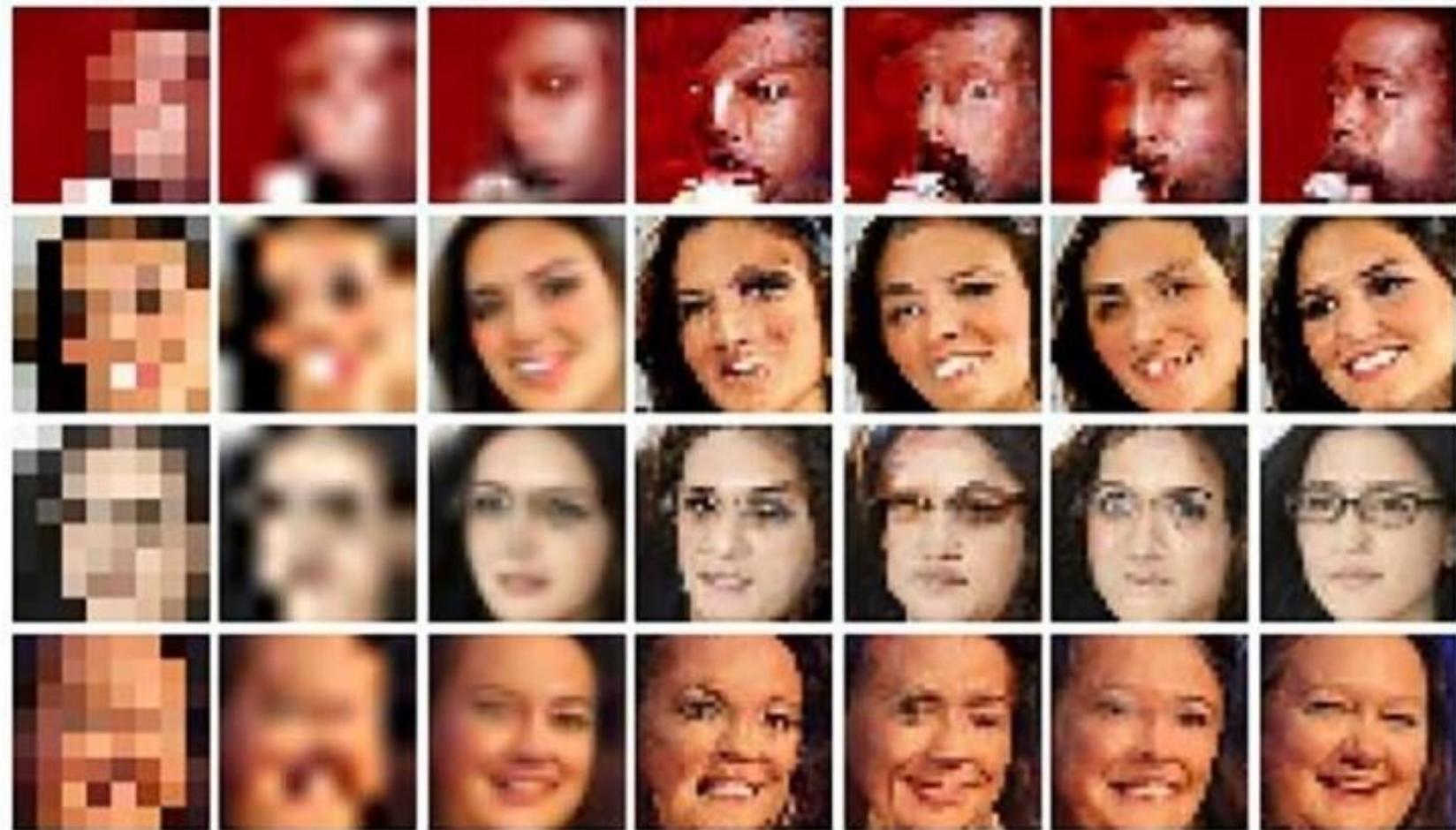
- Generalized mosaics can be of two broad categories: large mosaic blurs and small mosaic blurs.
- Bulk mosaic blurring is often more difficult to handle. And small mosaic blur is easier than it is.



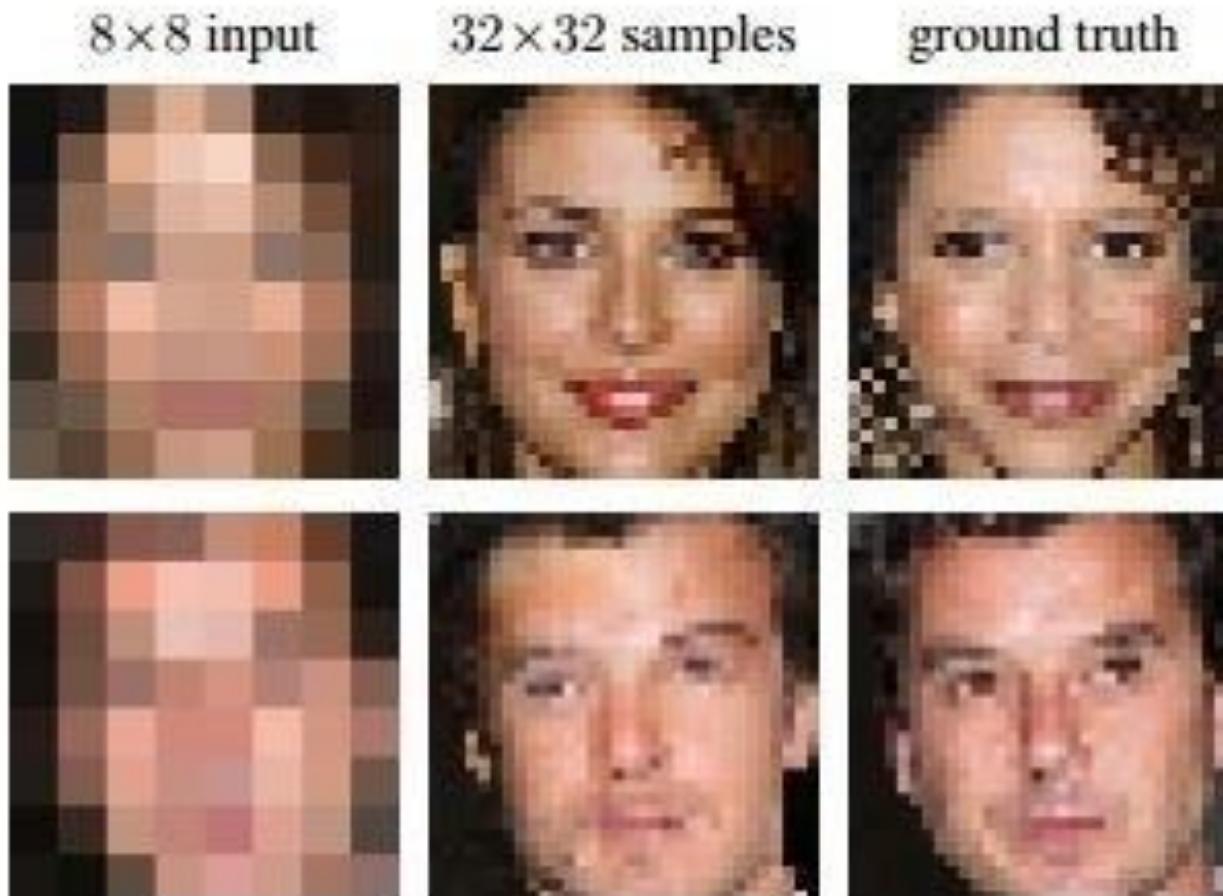
# De-mosaicing- Torch

Dataset	Original	Mosaic					20	P3 10	1
		2 × 2	4 × 4	8 × 8	16 × 16				
MNIST									
CIFAR-10									
AT&T									
FaceScrub									

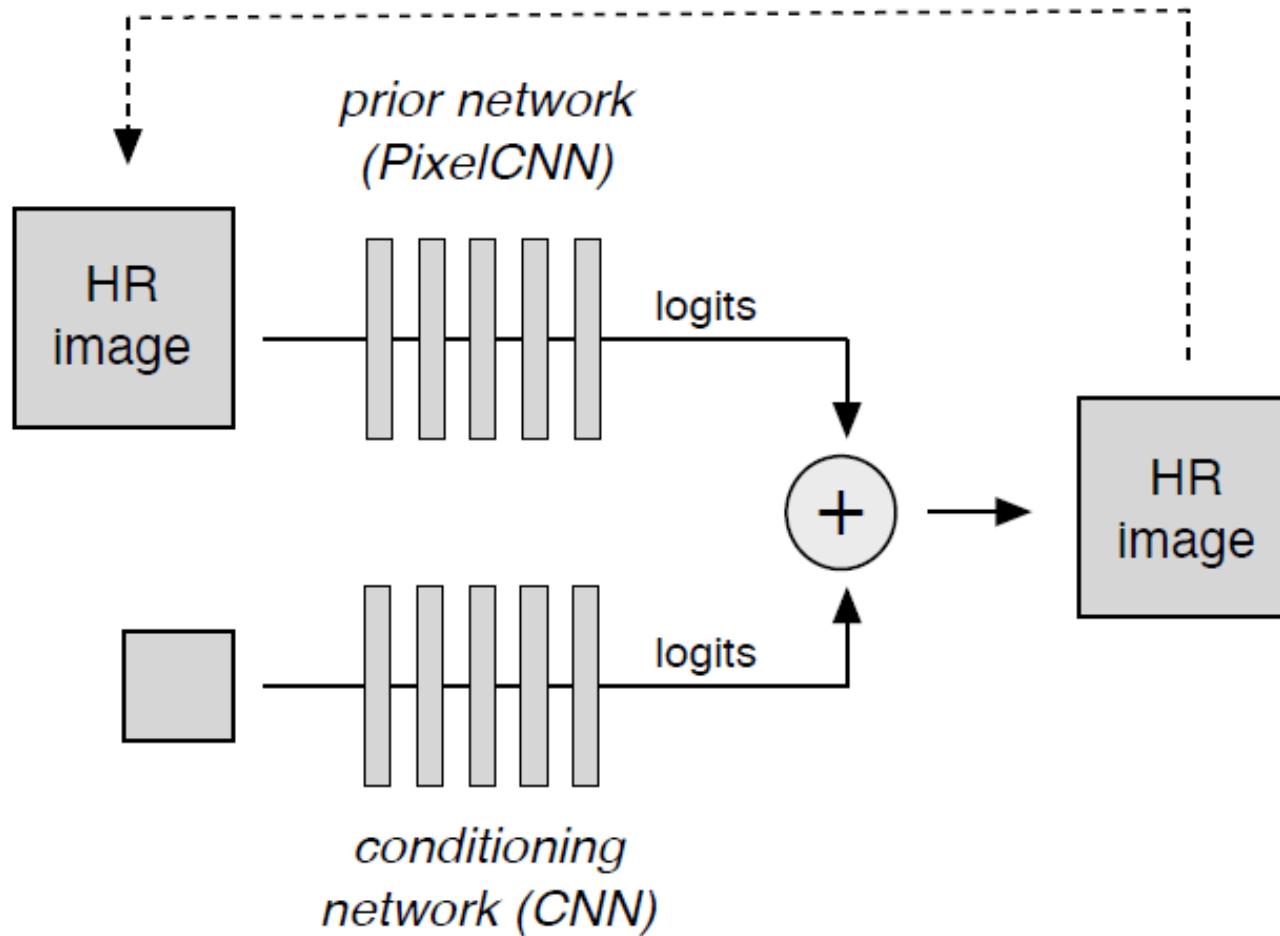
# Google: De-mosaicing



# Google: De-mosaicing



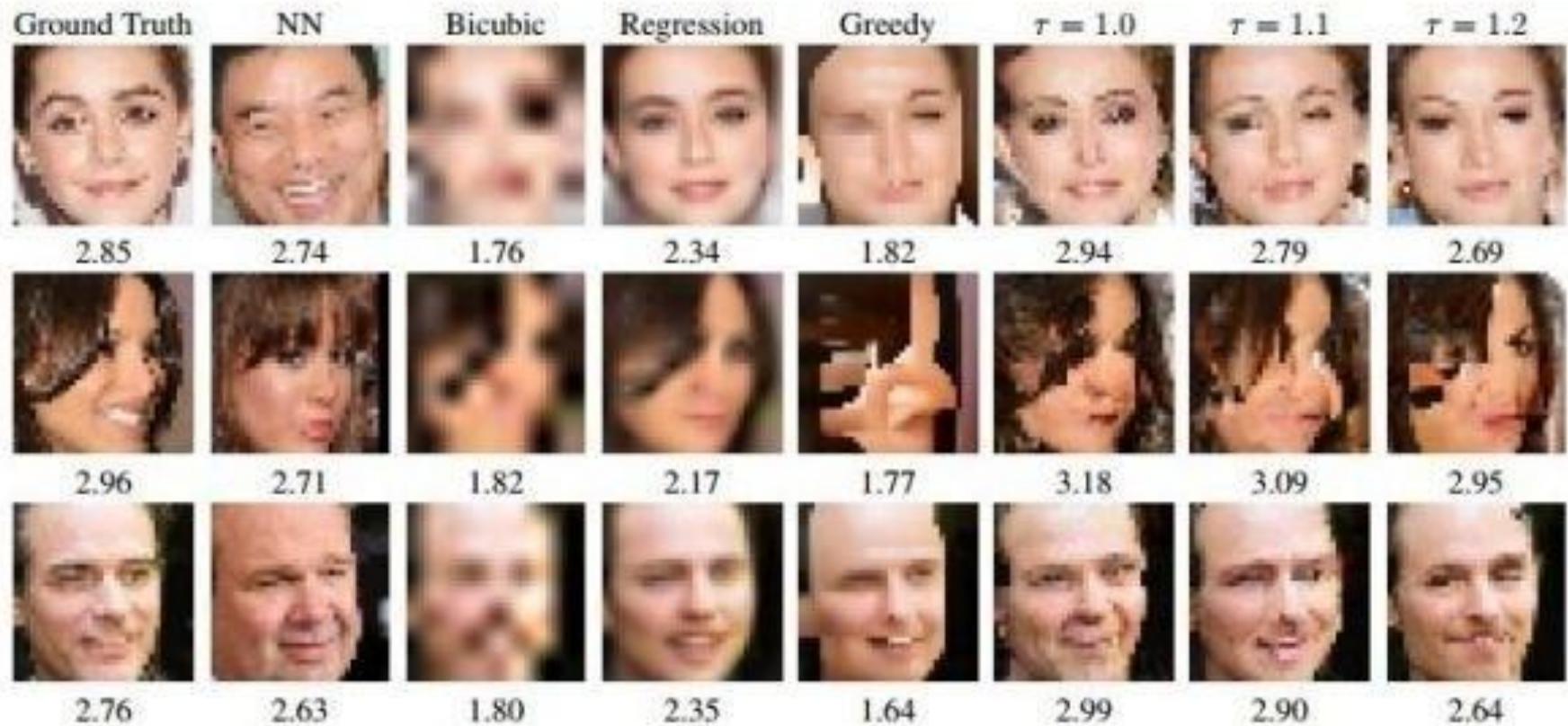
# Network



# Google: De-mosaicing



# Google: De-mosaicing

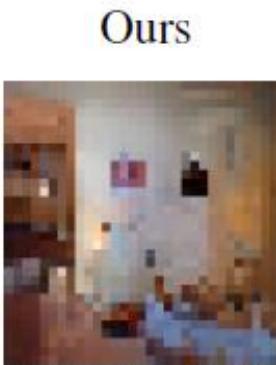


# Samples images that performed best and worst in human ratings:

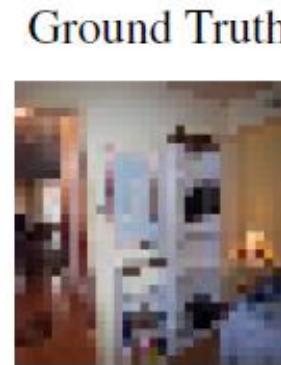
- The best and worst rated images in the human study.  
The fractions below the images denote how many times a person choose that image over the ground truth.



$23/40 = 57\%$



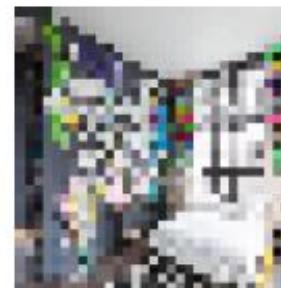
$34/40 = 85\%$



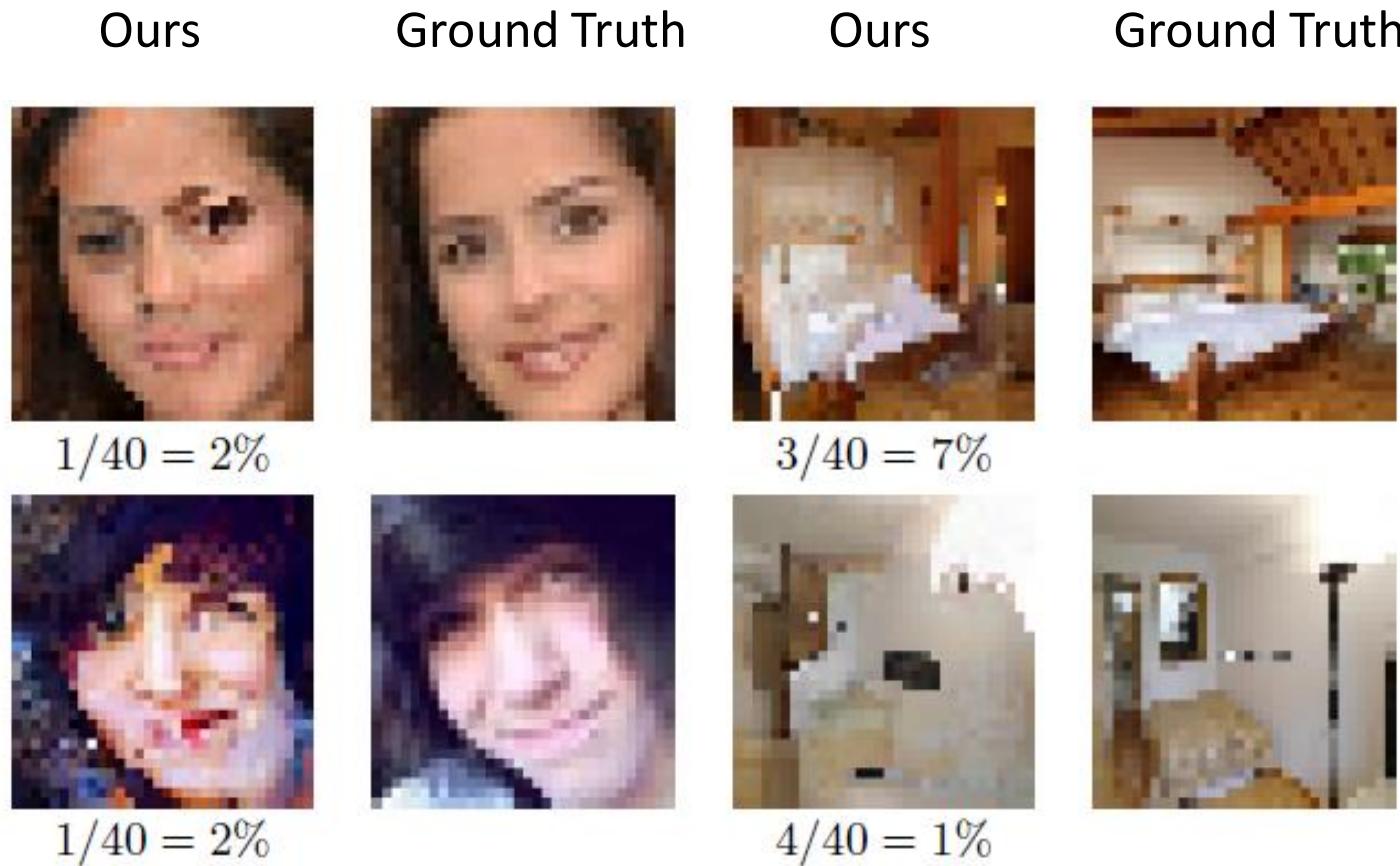
$17/40 = 42\%$



$30/40 = 75\%$



# Samples images that performed best and worst in human ratings:



# Applications of Image Inpainting

- Photo repair
- De-watermarking
- De-mosaicing
- **Object removal**
- De-blurring

# Object removal



Remove unwanted targets

# Object removal

- Idea: Remove object(s) from digital photographs, and then fill the hole with information extracted from the surrounding area.
- Result: Filled region should look “reasonable” to the human eyes.

# Example

Before



After



# Approach

## Texture Synthesis

- Idea:
  - Sample color values of the surrounding area
  - Generate textures with sampling result to fill the hole
- Advantage:
  - Cheap and effective
  - No blur or other degradation
- Disadvantage:
  - May lose linear structure and composite textures

# Object removal



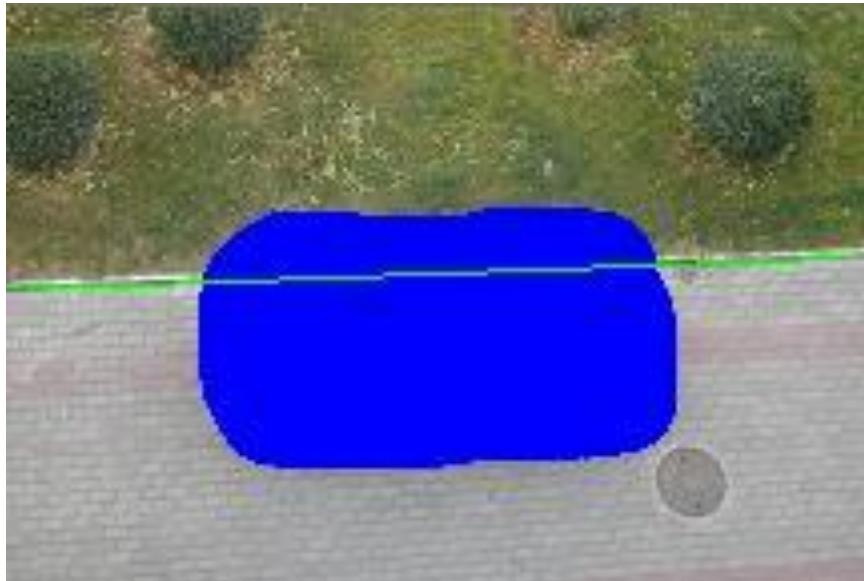
From the left (including objects, such as cars and other objects)  
to the right (after removing the target)

# Object removal



- Erase the area in the image that you want to remove or repair

# Object removal

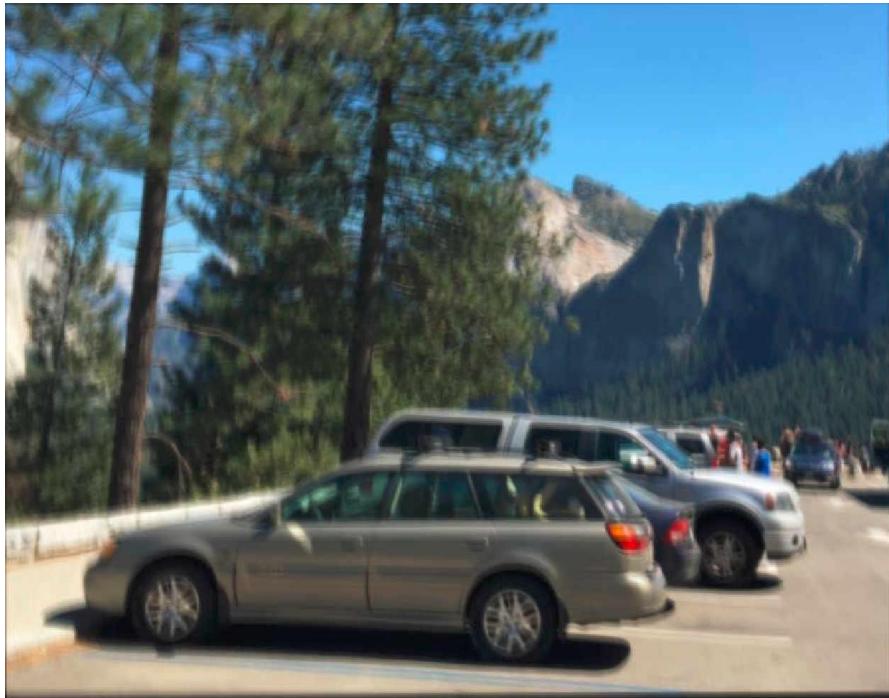


- Draw contour guides to force the position of the contour to be repaired or automatically generated by the system.
- Smartly fills the area to be repaired, making the picture look intact and naturally flowing.

# Applications of Image Inpainting

- Photo repair
  - De-watermarking
  - De-mosaicing
  - Object removal
- **De-blurring**

# Deblurring

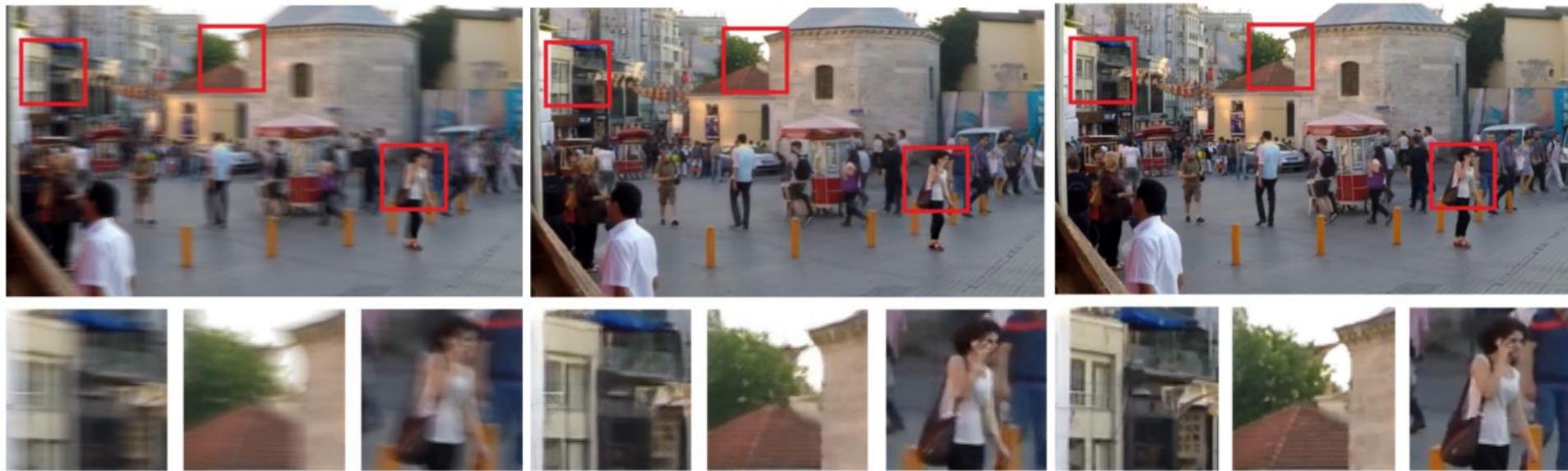


# Motion blur



- Motion blur is a static scene or a series of pictures such as a movie or an object moving rapidly in an animation causing a noticeable blurry dragging trace

# DeblurGAN



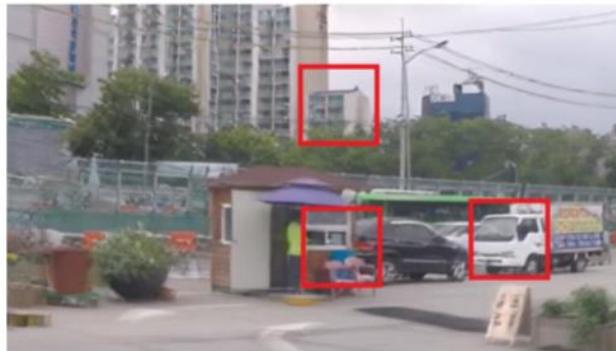
- An end-to-end learned method for motion deblurring. The learning is based on a **conditional GAN** and the content loss
- Achieves state-of-the-art performance both in the structural similarity measure and visual appearance.

# DeblurGAN

Blurred



DeblurGAN

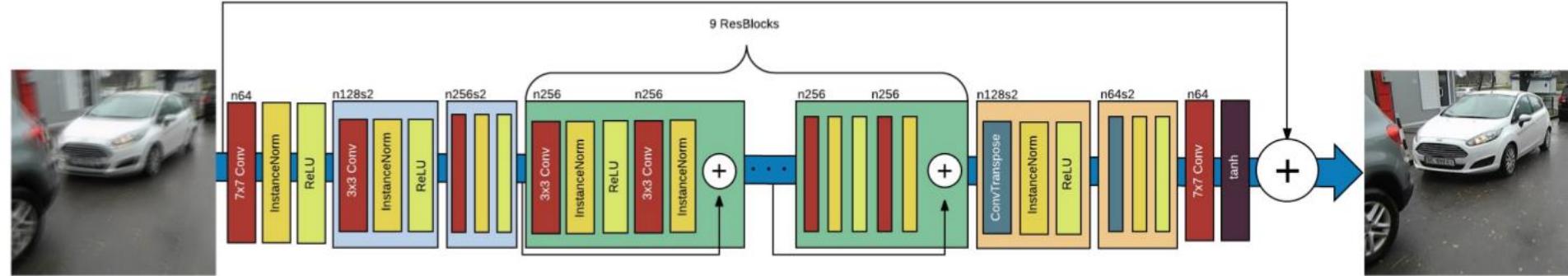


Ground Truth



GoPro images processed by DeblurGAN

# DeblurGAN -generator architecture



- Two strided convolution blocks with stride  $\frac{1}{2}$
- Nine residual blocks
- Two transposed convolution blocks
- Each ResBlock consists of a convolution layer, instance normalization layer, and ReLU activation

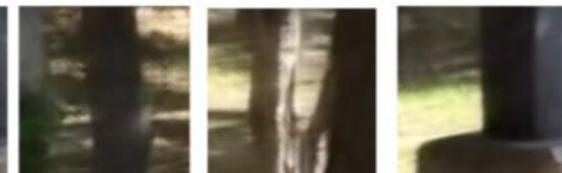
Blurred



Nah et al. [1]



DeblurGAN



S. Nah, T. Hyun, K. Kyoung, and M. Lee. Deep Multi-scale Convolutional Neural Network for Dynamic Scene Deblurring. 2016

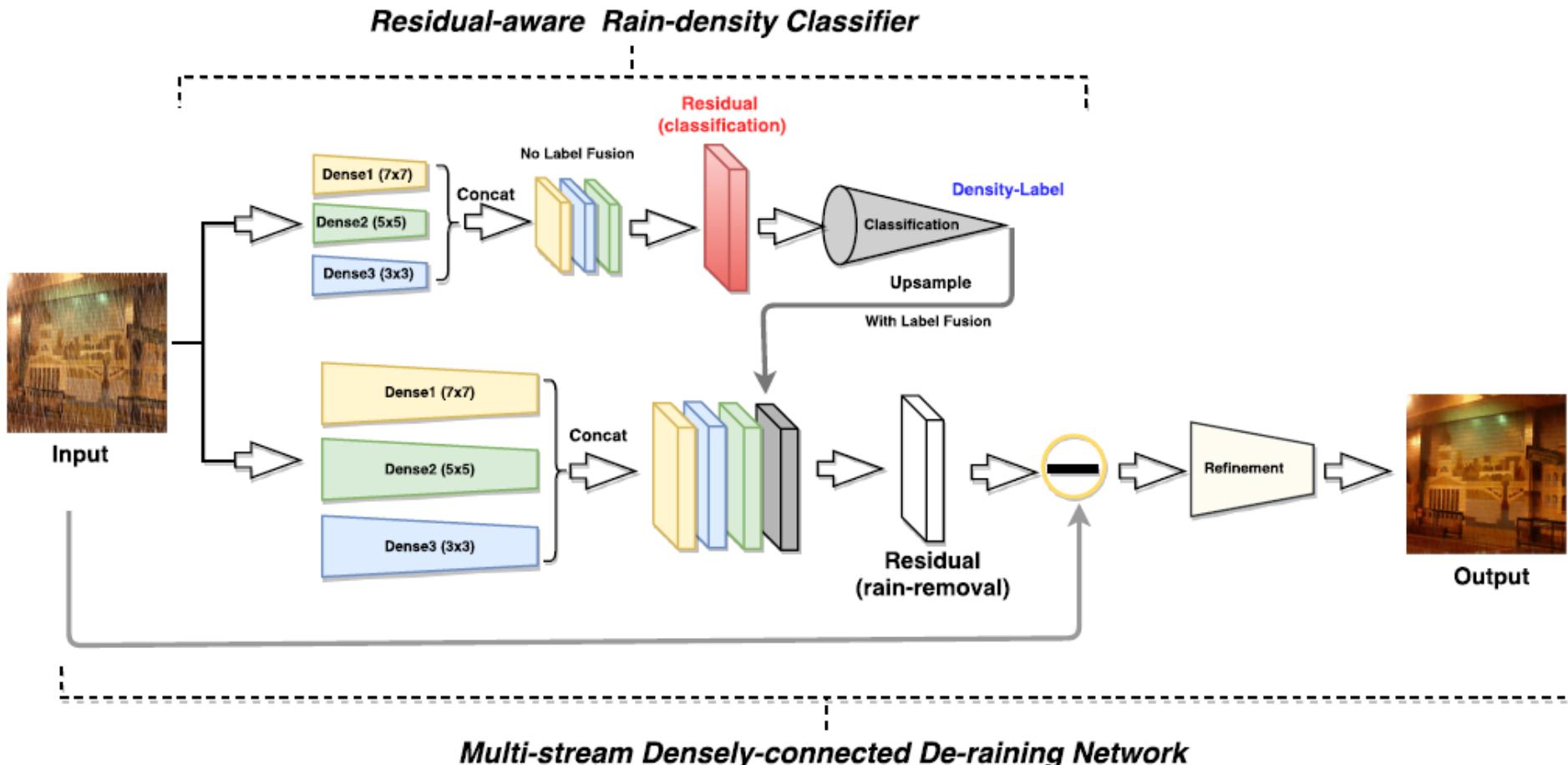
# De-raining



- Single image de-raining
- Density-aware Image De-raining method using a Multi stream Dense Network (DID-MDN)
- Rain-density classification
- Rain streak removal

He Zhang Vishal M. Patel. Density-aware Single Image De-raining using a Multi-stream Dense Network. In CVPR, 2018.

# DID-MDN



He Zhang Vishal M. Patel. Density-aware Single Image De-raining using a Multi-stream Dense Network. In CVPR, 2018.

# Result



Input

De-raining

He Zhang Vishal M. Patel. Density-aware Single Image De-raining using a Multi-stream Dense Network. In CVPR, 2018.



## Image Quality Assessment

# Image Quality Assessment

**Image quality** (often **Image Quality Assessment, IQA**) is a characteristic of an image that measures the perceived image degradation (typically, compared to an ideal or perfect image).

- Full-reference (FR) methods
  - Reduced-reference (RR) methods
  - No-reference (NR) methods

# Image quality metrics

## ➤ Full-reference metrics

- Peak Signal to Noise Ratio (PSNR)
- Structural Similarity (SSIM)
- Visual Information Fidelity (VIF)

# Peak Signal to Noise Ratio (PSNR)

- **PSNR**, is an engineering term for the ratio between the maximum possible power of a signal and the power of corrupting noise that affects the fidelity of its representation.
- Many signals have a very wide dynamic range, PSNR is usually expressed in terms of the logarithmic decibel scale.

PSNR is most easily defined via the **mean squared error (MSE)**  
 Given a noise-free  $m \times n$  monochrome image  $I$  and its noisy approximation  $K$ , MSE is defined as:

$$MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [I(i, j) - K(i, j)]^2$$

The PSNR (in dB) is defined as:

$$\begin{aligned} PSNR &= 10 \cdot \log_{10} \left( \frac{MAX_I^2}{MSE} \right) \\ &= 20 \cdot \log_{10} \left( \frac{MAX_I}{\sqrt{MSE}} \right) \\ &= 20 \cdot \log_{10}(MAX_I) - 10 \cdot \log_{10}(MSE) \end{aligned}$$

$MAX_I$  ( $=255$ ) is the maximum possible pixel value.

# Structural Similarity (SSIM)

- SSIM is a method for predicting the perceived quality of digital television and cinematic pictures, as well as other kinds of digital images and videos.
- SSIM is used for measuring the similarity between two images.
- SSIM is designed to improve on traditional methods such as peak signal-to-noise ratio (PSNR) and mean squared error (MSE).

The SSIM index is calculated on various windows of an image. The measure between two windows  $x$  and  $y$  of common size  $N \times N$  is:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

with:

- $\mu_x$  the average of  $x$ ;  $\mu_y$  the average of  $y$ ;
- $\sigma_x^2$  the variance of  $x$ ;  $\sigma_y^2$  the variance of  $y$ ;
- $\sigma_{xy}$  the covariance of  $x$  and  $y$ ;
- $c_1 = (k_1 L)^2, c_2 = (k_2 L)^2$  two variables to stabilize the division with weak denominator;
- $L$  the dynamic range of the pixel-values;
- $k_1 = 0.01$  and  $k_2 = 0.03$  by default

The SSIM index satisfies the condition of symmetry:  $SSIM(x, y) = SSIM(y, x)$

# PSNR & SSIM

## PSNR:

- Mostly used, but its value does not reflect people's subjective feeling well.
- Value range :  $20 \sim 40$  dB
- The larger the value, the better the video quality.

## SSIM :

- Complicated calculation, and its value can better reflect the subjective feelings of the human eye.
- Value range :  $0 \sim 1$
- The larger the value, the better the video quality.



# Classical Methods to Image Inpainting

# Inpainting: context and issues

Inpainting corresponds to filling holes (i.e. missing areas) in images

Let be an image  $I$  defined as

$$I : \Omega \subset \mathbb{R}^n \mapsto \mathbb{R}^m$$

Let be a degradation operator  $M$

$$M : \Omega \mapsto \{0, 1\}$$

$$M(x) = \begin{cases} 0, & \text{if } x \in U \\ 1, & \text{otherwise} \end{cases}$$

$n = 2$  for a 2D image

$m = 3$  for a (R,G,B) image

$$\Omega = S \cup U$$

S being the known part of  $I$

U the unknown part of  $I$

Let  $F$  the observed image:

$$F = M \circ I$$

- is the Hadamard product

# Inpainting: context and issues

Different configurations according to the definition of  $M$ :



Original image



80% of the pixels  
have been  
removed.

Sparsity and  
low-rank methods



damaged portions  
in black, scratches

Diffusion-based  
methods



Object removal

Examplar-based  
methods

# Exemplar-based Inpainting:

## Texture synthesis

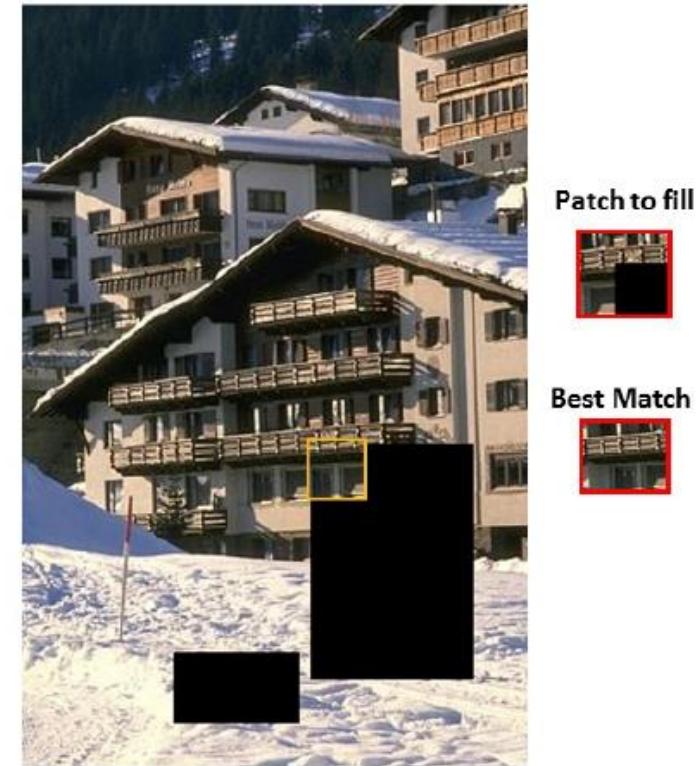
Exemplar-based inpainting methods rely on the assumption that the known part of the image provides a good dictionary which could be used efficiently to restore the unknown part (Efros and Leung, 1999).

The recovered texture is therefore learned from similar regions.

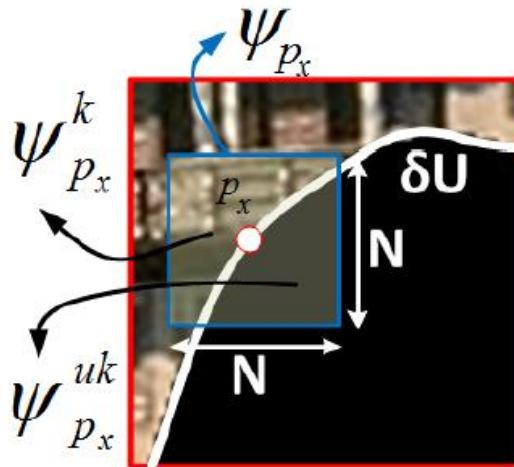
- This can be done simply by sampling, copying or combining patches from the known part of the image;

### Template Matching

- Patches are then stitched together to fill in the missing area.



# Notations:



$$\psi_{p_x} = \begin{bmatrix} I_1(0,0) \\ I_1(0,1) \\ \vdots \\ I_M(N-2, N-2) \\ I_M(N-1, N-1) \end{bmatrix}$$

- a patch  $\psi_{p_x}$  is a discretized  $N \times N$  neighborhood centered on the pixel  $p_x$ . This patch can be vectorized in a raster-scan order as a  $mN^2$ -dimensional vector;
- $\psi^{uk}_{p_x}$  denotes the unknown pixels of the patch;
- $\psi^k_{p_x}$  denotes its known pixels;
- $\psi_{p_x(i)}$  denotes the  $i^{th}$  nearest neighbour of  $\psi_{p_x}$ ;
- $\delta U$  is the front line;

# Exemplar-based Inpainting:

## Criminisi et al.'s algorithm

Criminisi et al. ([Criminisi et al., 2004](#)) has brought a new momentum to inpainting applications and methods. They proposed a new method based on two sequential stages:

- Filling order computation
- Texture synthesis

# Filling order computation:

$$P(p_x) = C(p_x) \times D(p_x)$$

Confidence term

$$C(p_x) = \frac{\sum_{q \in \Psi_{p_x}^k} C(q)}{|\Psi_{p_x}|}$$

Where  $|\Psi_{p_x}|$  is the area of  $\Psi_{p_x}$

Data term

$$D(p_x) = \frac{|\nabla I^\perp(p_x) \cdot n_{p_x}|}{\alpha}$$

Where  $\alpha$  is the a normalization constant in order to ensure that  $D(p_x)$  is in the range 0 to 1

# Texture synthesis:

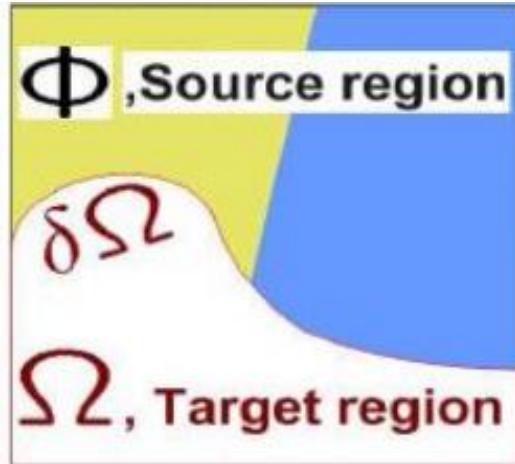
A template matching is performed within a local neighborhood:

$$p_y = \arg \min_{q \in W} d(\Psi_{p_q}^k, \Psi_{p_x^*}^k)$$

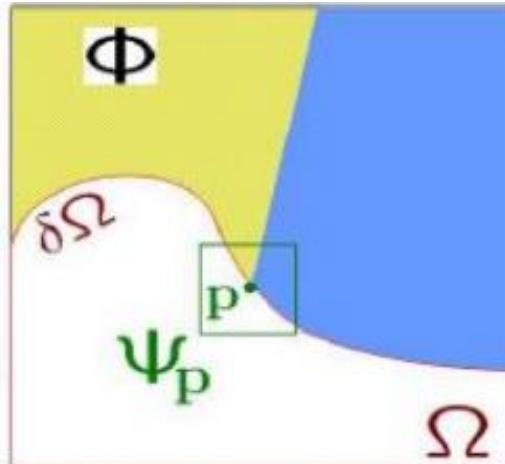
- $W \subseteq S$  is the window search;
- $\Psi_{p_x^*}^k$  are the known pixels of the patch  $\Psi_{p_x^*}$  with the highest priority;
- $\Psi_{p_y}^k$  are the known pixels of the nearest patch neighbor;
- $d(a, b)$  is the sum of squared differences between patches a and b.

The pixels of the patch  $\Psi_{p_y}^{uk}$  are then copied into the unknown pixels of the patch  $\Psi_{p_x^*}$

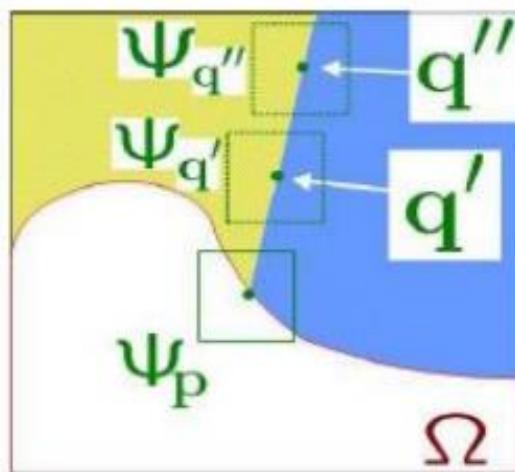
# Texture synthesis



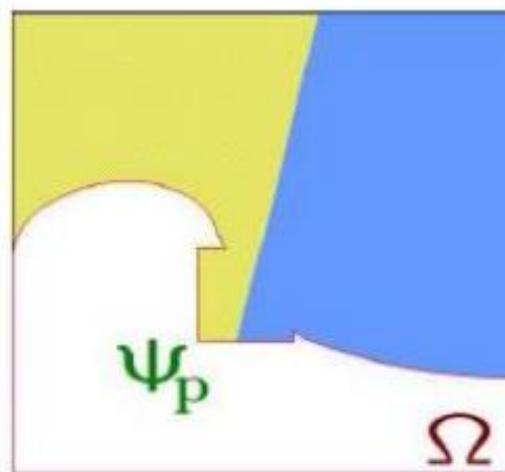
a



b



c



d

- (a) Original image
- (b) Repaired patch
- (c) Match and search
- (d) Copy or synthesis

# Some examples:

Inpainted pictures with (Criminisi et al., 2004)'s method:



# Development

[1] Connelly Barnes, Eli Shechtman, Adam Finkelstein, Dan B Goldman.

**PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing**  
ACM Transactions on Graphics (Proc. SIGGRAPH), 2009.

[2] Connelly Barnes, Eli Shechtman, Dan B Goldman, Adam Finkelstein.

**The generalized PatchMatch correspondence algorithm.**  
In Proceedings of European Conference on Computer Vision (ECCV), 2010.

[3] S Darabi, E Shechtman, C Barnes and et al.

**Image melding: combining inconsistent images using patch-based synthesis.**

[4] K He and J Sun. **Image completion approaches using the statistics of similar patches.**

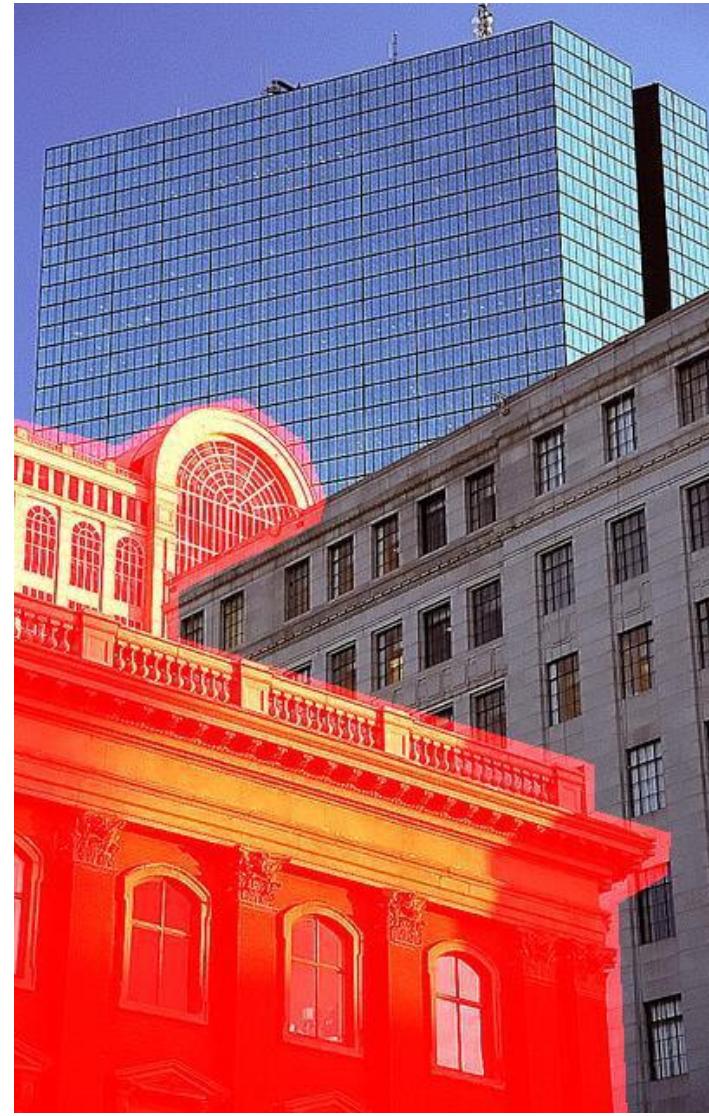
IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 2014.

[5] Jia-Bin Huang, Sing Bing Kang, Narendra Ahuja, Johannes Kopf.

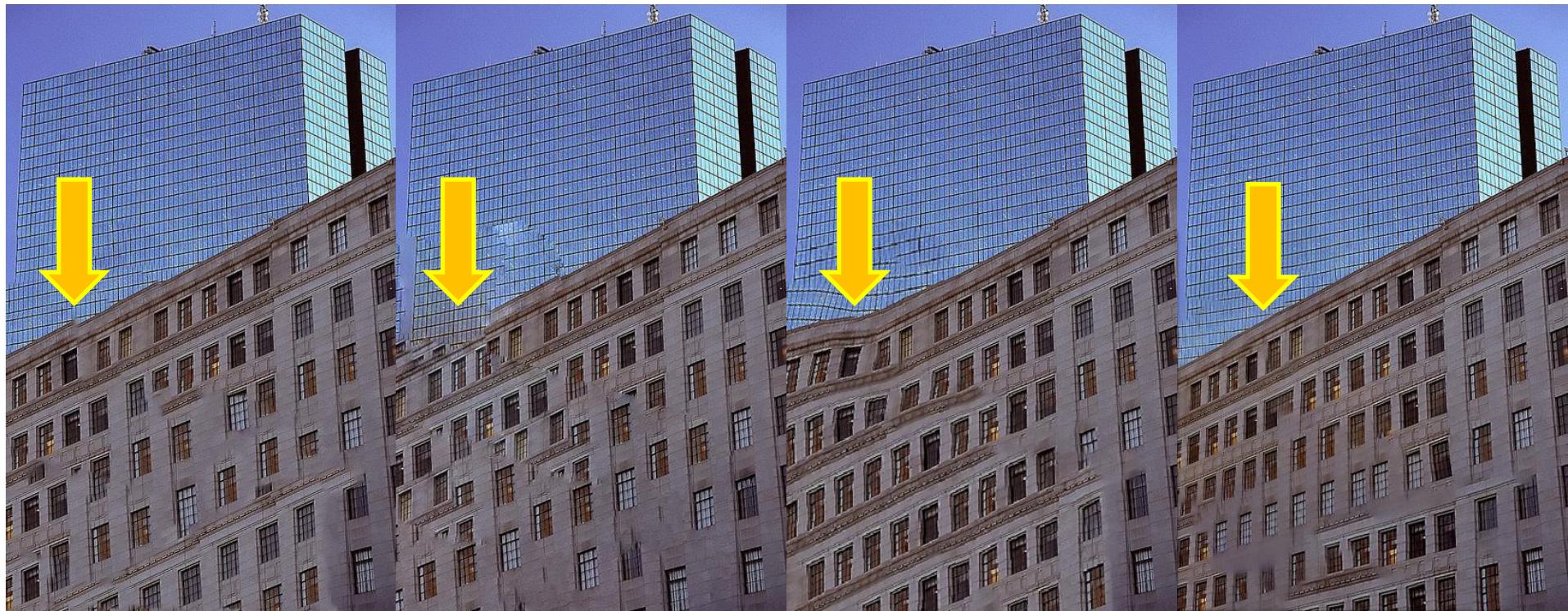
**Image Completion using Planar Structure Guidance.**

In SIGGRAPH, 2014.

# Input and hole



# Result



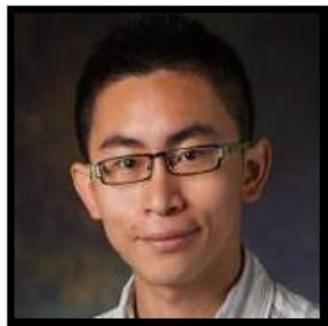
Photoshop CAF  
Patch Match  
[Barnes et al. 2009]

Statistics of Patch Offsets  
[He and Sun 2014]

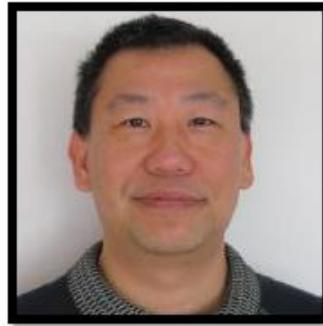
Image Melding  
[Darabi et al. 2012]

Planar Structure  
Guidance  
[Huang et al. 2014]

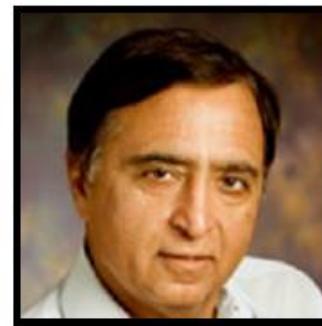
# Image Completion using Planar Structure Guidance



Jia-Bin Huang



Sing Bing Kang



Narendra Ahuja



Johannes Kopf

# Image Completion



Input image

Region to fill



## Photoshop Content Aware Fill [Wexler et al. 2007] [Barnes et al. 2009]



Statistics of Patch Offsets  
[He and Sun ECCV 2012]

Translational patches  
are not sufficient!



Image Melding  
[Darabiet al. 2012]

Searching patch  
transformation space is HARD!



## Planar Structure Guidance

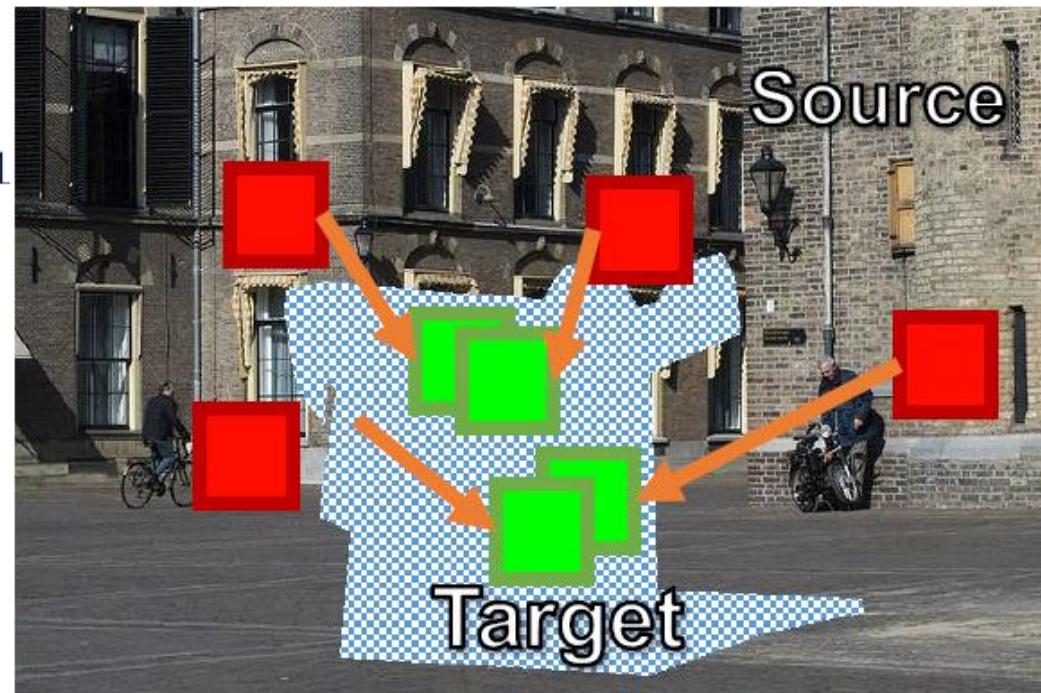
Use mid-level information

# Image Completion using Planar Structure Guidance

- Based on non-parametric framework of [Wexler et al. 2007]
- The basic form  $\min_{\{t_i, s_i\}} \sum E_{color}(t_i, s_i)$
- Patch matching cost

$$E_{color}(t_i, s_i) = || \boxed{\text{green patch}} - \boxed{\text{red patch}} ||_1$$

- Patch Match [Barnes et al. 2009] for efficient nearest neighbor field search.



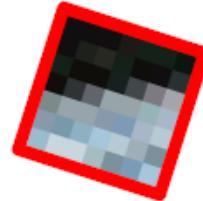
# Image Completion using Planar Structure Guidance

- Low-level processing is not sufficient
- Idea: guide the completion using mid-level information!
- Transformation of source patches

- Translational  
(Photoshop CAF)



- Rotation and scale  
(Image Melding)



- Projective transform  
(Planar Structure  
Guidance)



# Image Completion using Planar Structure Guidance

Mid-Level cues for guidance

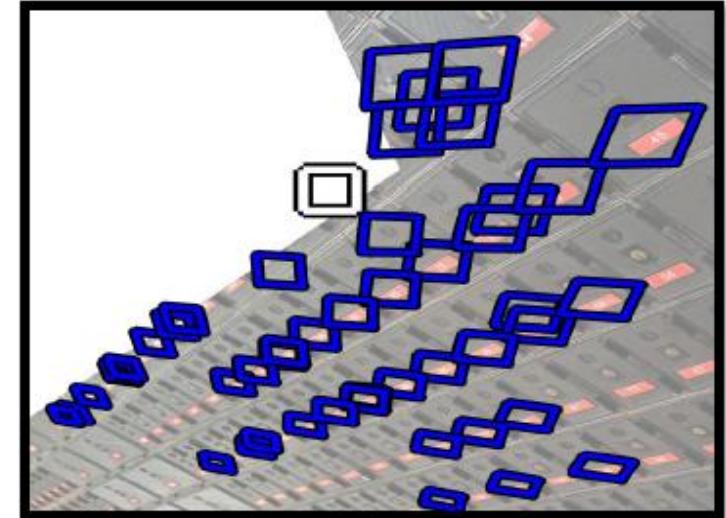
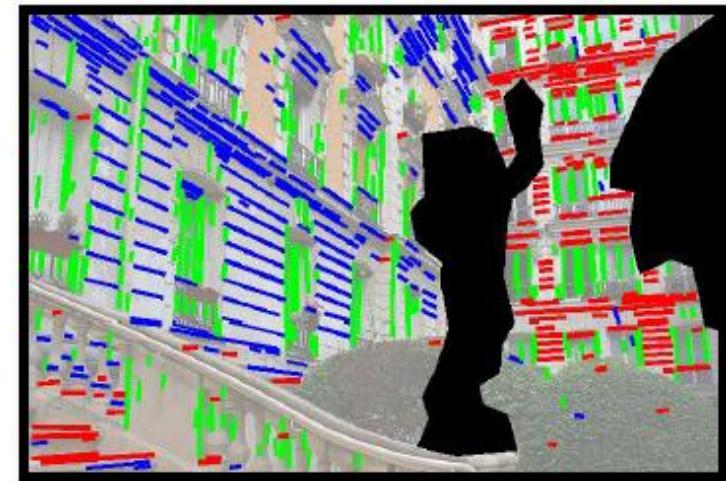


Input and Hole

# Image Completion using Planar Structure Guidance

## Mid-Level cues for guidance

- Planes:  
Cues for how to deform patches  
Plane detection: vanishing point detection and posterior probability
- Regularity:  
Cues for where to sample patches  
Regularity detection:  
Feature detection and matching  
Mode detection



# Image Completion using Planar Structure Guidance

## Result



# Result Comparison



Photoshop Content Aware Fill



Image Melding



Statistics of Patch Offsets



Planar Structure Guidance

# Image Completion using Planar Structure Guidance

## Regularity detection



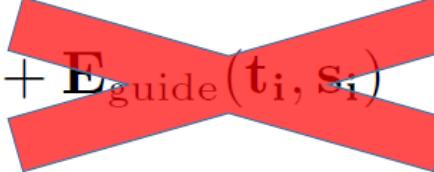
With planar guidance only



Plane and regularity

# What if there are no structures?

- Reduces to baseline image completion algorithms which use translational patches (same as Photoshop content-aware fill)

$$\min_{\{t_i, s_i\}} \sum E_{color}(t_i, s_i) + E_{guide}(t_i, s_i)$$


- Works well with man-made scenes
- Backward-compatible with natural scenes

# What if there are no structures?



Input and Mask



Completion

# What if there are no structures?



Input and Mask

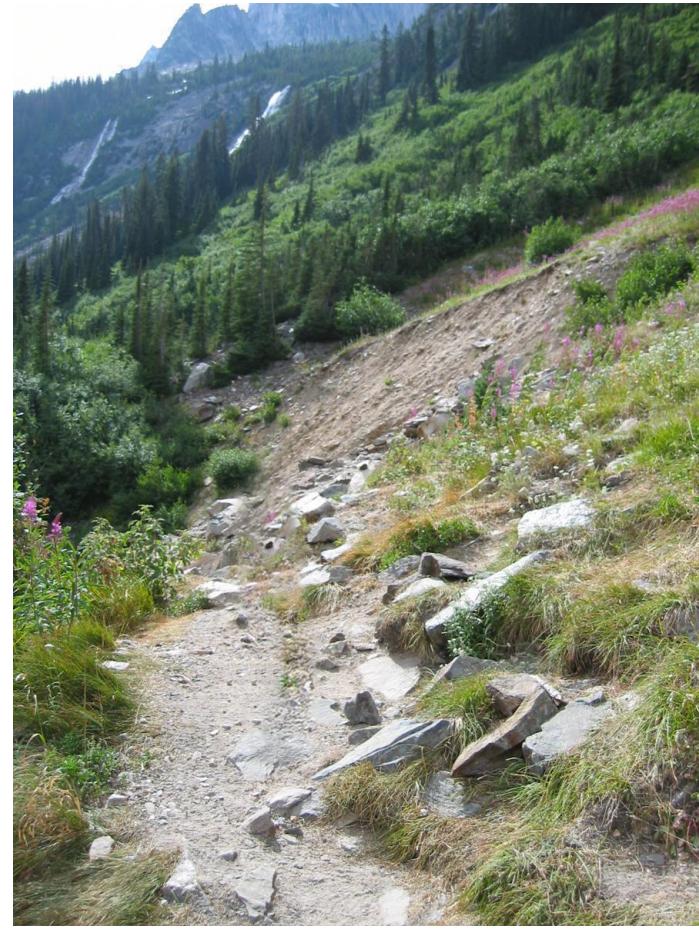


Completion

# What if there are no structures?



Input and Mask



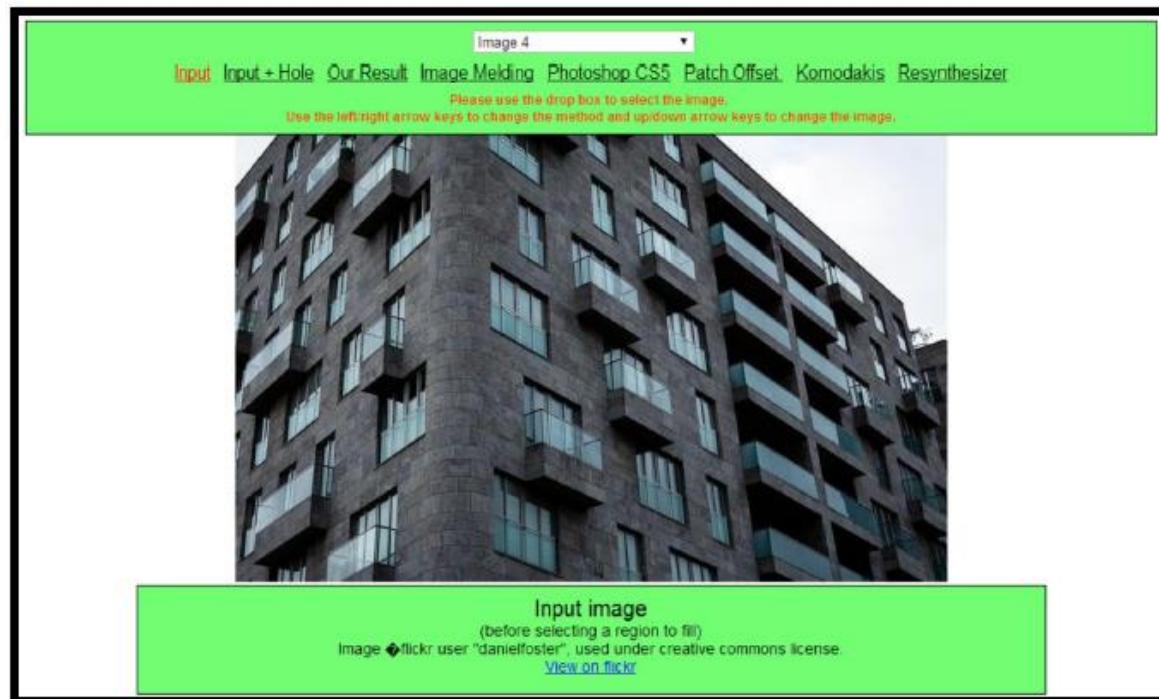
Completion

# Project Website

- <http://bit.ly/planarcompletion>

- Full comparison of 110 images  
( 85 urban scenes + 25 natural scenes )

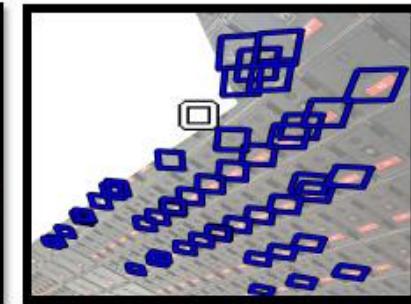
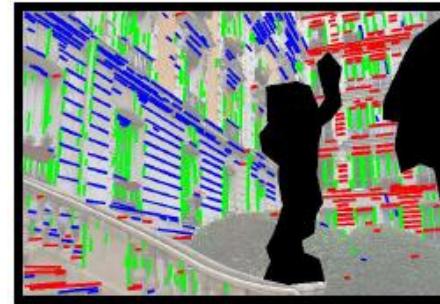
- Code



# Conclusions

- Mid-level constraints for the win!

Our work: planes and regularity



- Great for man-made scenes



- Backward-compatible with natural scenes



# OpenCV-Image Inpainting

## Function: **Inpaint( )**

```
void inpaint( InputArray src, InputArray inpaintMask,  
              OutputArray dst, double inpaintRadius, int flags);
```

- **src**: input single-channel or RGB image
- **inpaintMask**: the mask of image, single-channel image, the size is the same as the original image, the pixels of the mask are all 1, others are all 0.
- **dst**: repaired image
- **inpaintRadius**: the radius of the neighborhood taken by the repair algorithm and is used to calculate the difference of the current pixel.
- **flags**: repair algorithm, **INPAINT\_NS** and **INPAINT\_TELEA**;

# Results



Input and Mask



Inpainted

Other parts have changed color!

# Results



Input and Mask



Inpainted

The repair effect is good!

# Results



Input and Mask



Inpainted

The complex background is not well repaired!

# Results



Input and Mask



Inpainted

In order not to affect images outside the area, artificially limit the areas that need to be repaired

# Results



Input and Mask



Inpainted

The background part of the repair effect is good, but the repair of the edge details is not good!

4

# Deep Learning to Image Inpainting

# Datasets

- CelebA
- MS COCO
- PASCAL VOC
- ImageNet
- Paris Street View
- SVHN
- Car Dataset

# Large-scale CelebFaces Attributes (CelebA) Dataset



# Details

- **CelebFaces Attributes Dataset (CelebA)** is a large-scale face attributes dataset with more than **200K** celebrity images, each with **40** attribute annotations.
- Images cover large pose variations and background clutter.
- CelebA has large diversities, large quantities, and rich annotations, including
  - **10,177** number of **identities**
  - **202,599** number of **face images**
  - **5 landmark locations, 40 binary attributes** annotations per image

# Application

The dataset can be employed as the training and test sets for the following computer vision tasks:

- Face attribute recognition
- Face detection
- Landmark (or facial part) localization

# Sample Images

Eyeglasses



Wearing Hat



Bangs



Wavy Hair



# Sample Images

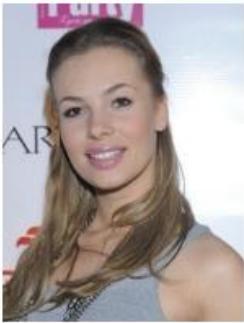
Pointy  
Nose



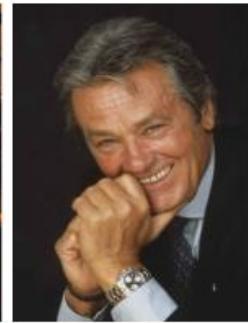
Mustache



Oval Face



Smiling



# MS COCO

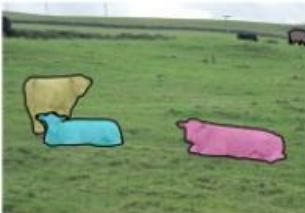
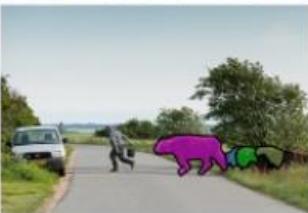
Person



Dog

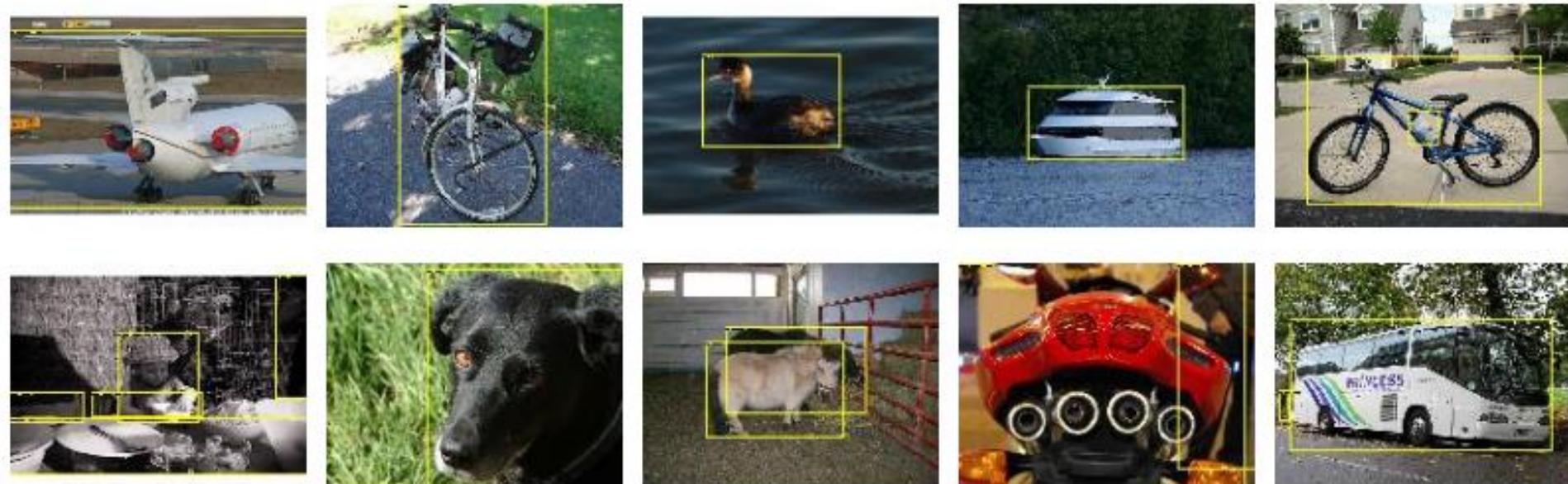


Cow



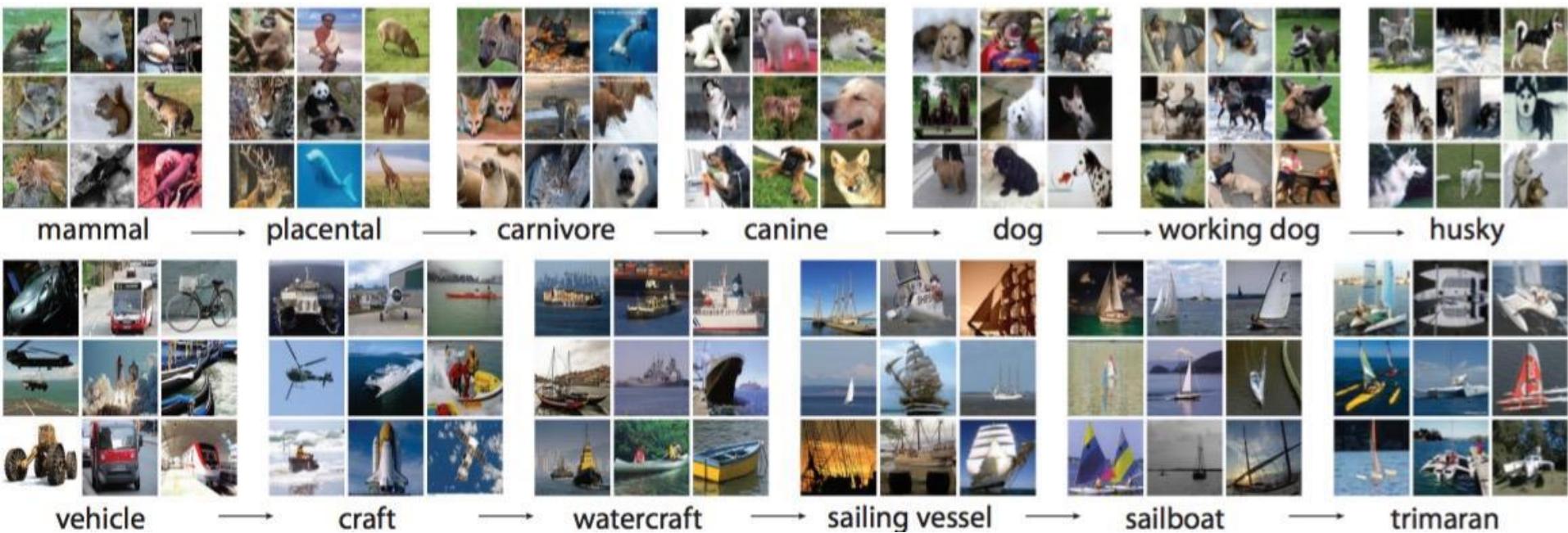
- 102 classes
- 165,482 training images
- 81,208 validation images
- 81,434 test instances

# PASCAL VOC 2012



- 20 classes
- 11,530 train/val images
- 27,450 object instances
- 6,929 segmentations

# ImageNet-1000



- 1000 classes
- 1.3M training images
- 100K testing images

# Paris Street View



- 15,000 images

C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. A. Efros. What makes paris look like paris?  
Communications of the ACM, 58(12):103–110, 2015.

# The Street View House Numbers(SVHN)



- 73,257 training images
- 26,032 test images

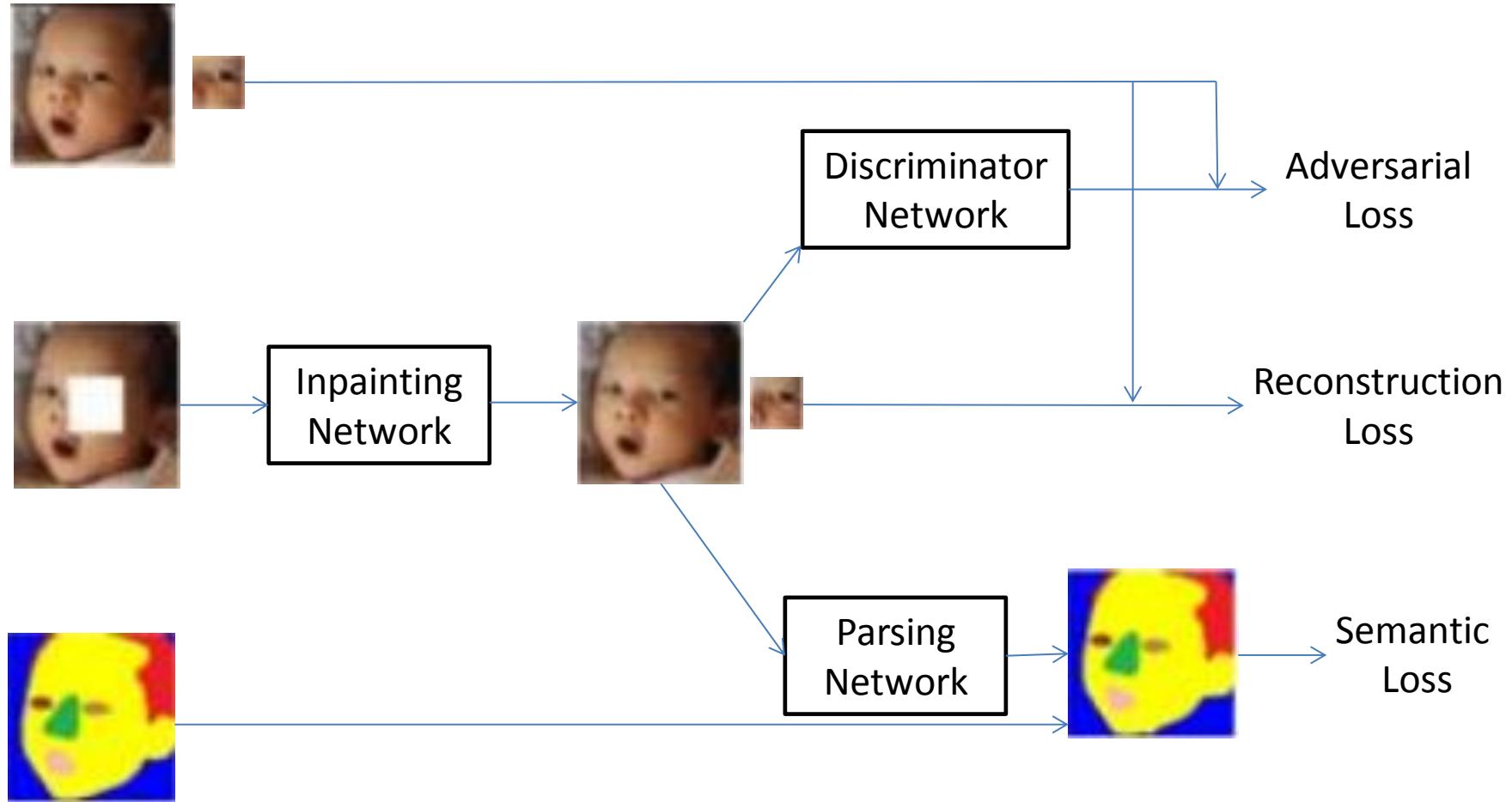
# Cars Dataset



- 196 classes of cars
- 8,144 training images
- 8,041 testing images

## Methods in Recent Papers

# Image Inpainting Framework— Deep Learning



➤ **Globally and Locally Completion**

- Semantic Image Inpainting
- High-Resolution Image Inpainting
- Content Encoder

# Globally and Locally Consistent Image Completion

Satoshi Iizuka, Edgar Simo-Serra, Hiroshi Ishikawa  
Waseda University. In SIGGRAPH ,2017.



[https://github.com/satoshiiiizuka/siggraph2017\\_inpainting](https://github.com/satoshiiiizuka/siggraph2017_inpainting)

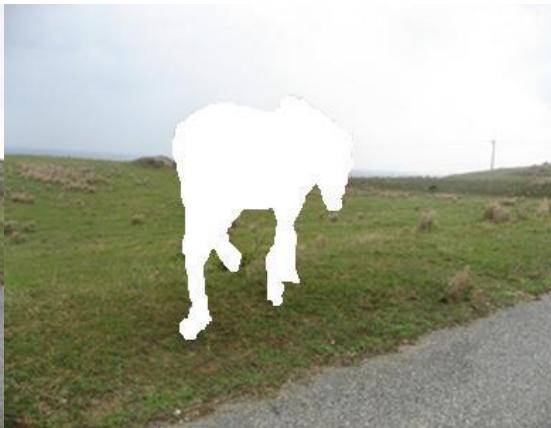
<https://www.slideshare.net/siizuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Image Completion

- Filling in a part of images
  - Object removal
  - Texture generation for occluded regions
  - ...



Input



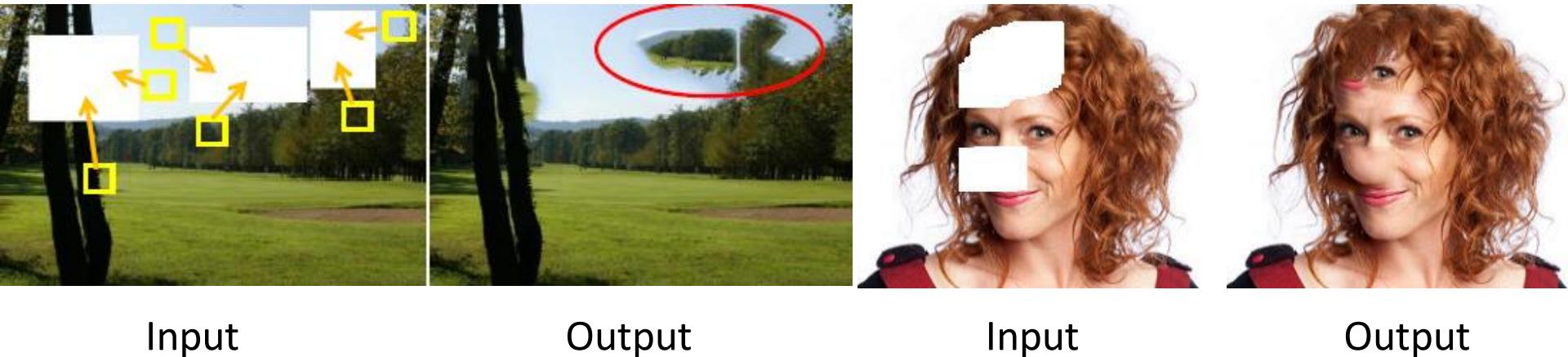
Mask



Image Completion

# Related Work

- Patch-based inpainting
  - Synthesize texture by collecting small image patches
  - Cannot preserve global structures
  - Cannot generate novel objects

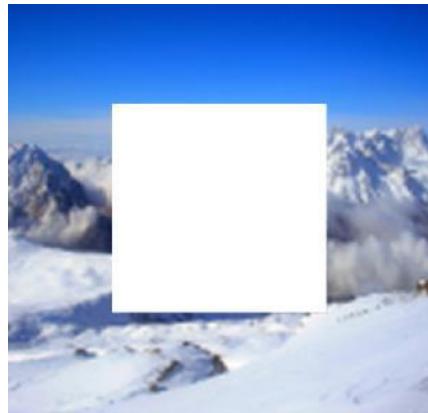


[https://github.com/satoshiiizuka/siggraph2017\\_inpainting](https://github.com/satoshiiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siizuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Related Work

- Learning-based inpainting
  - Learning inpainting with GAN
  - Fixed image resolution( $128 \times 128$  pixels)
  - Fixed mask position and size (center,  $64 \times 64$  pixels)
  - Tends to generate texture that is inconsistent with an input image



Input



Output



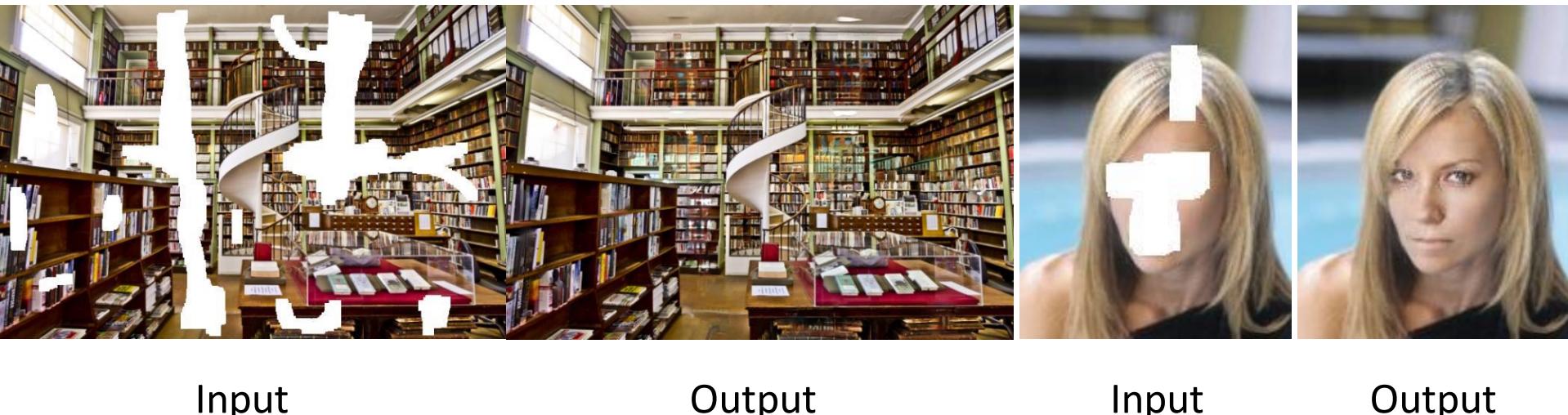
Input



Output

# Our Method

- Novel network for globally and locally consistent image completion
  - Completion network that is able to inpaint arbitrary regions
  - Adversarial training with two auxiliary networks
  - Can generate novel objects



[https://github.com/satoshiizuka/siggraph2017\\_inpainting](https://github.com/satoshiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siiuzuka/siggraph-2017-globally-and-locally-consistent-image-completion>



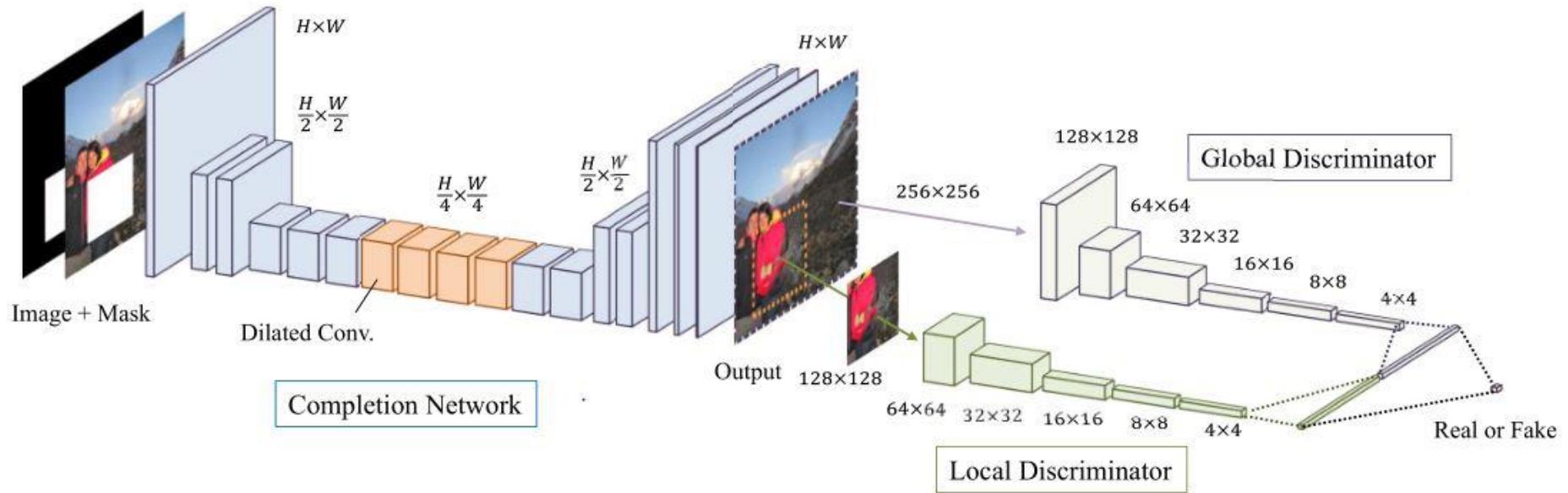
Input

Completion

Satoshi Iizuka, Edgar Simo-Serra, Hiroshi Ishikawa. Globally and Locally Consistent Image Completion.  
In SIGGRAPH ,2017

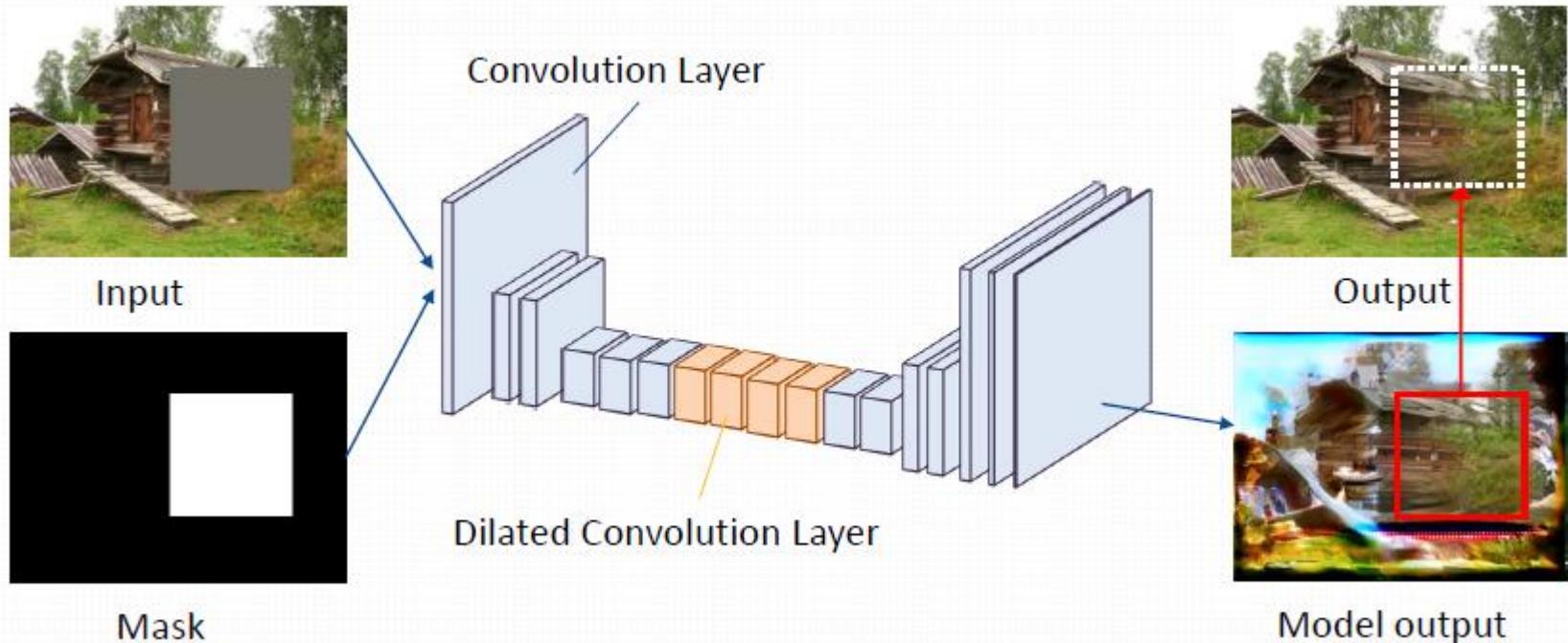
# Demo

# Overview of architecture



# Completion Network

- Composed of 17 convolutional layers
- Output an image filled in holes

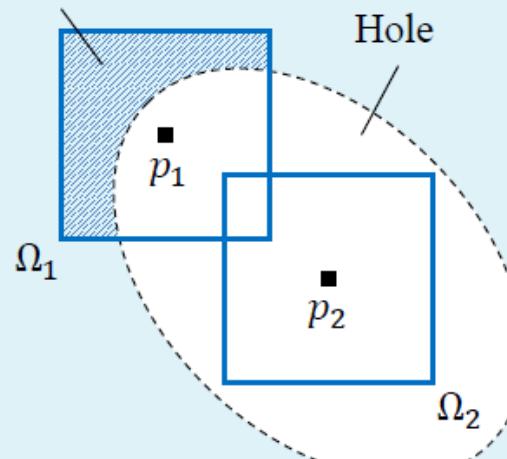


[https://github.com/satoshiizuka/siggraph2017\\_inpainting](https://github.com/satoshiizuka/siggraph2017_inpainting)

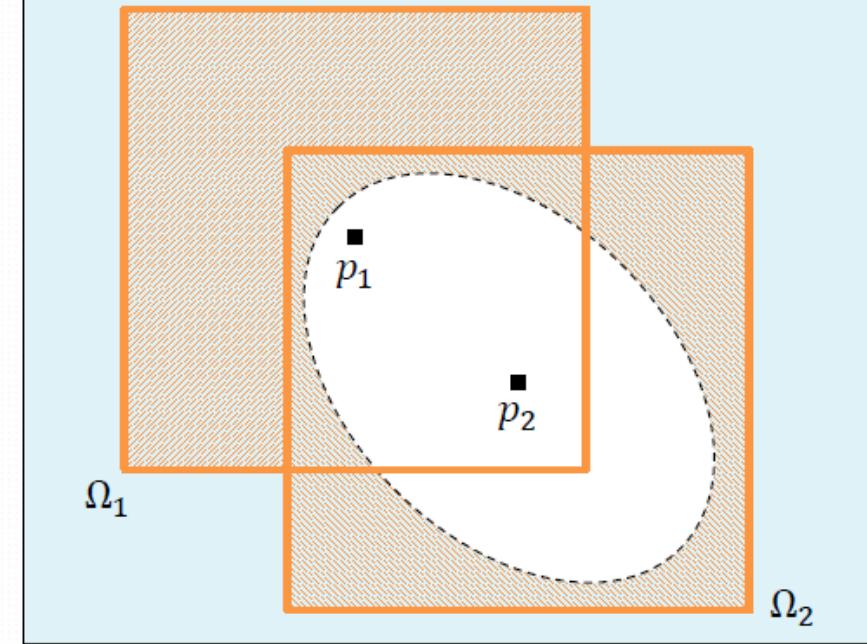
<https://www.slideshare.net/siiuzuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Importance of Spatial Support

Influencing region



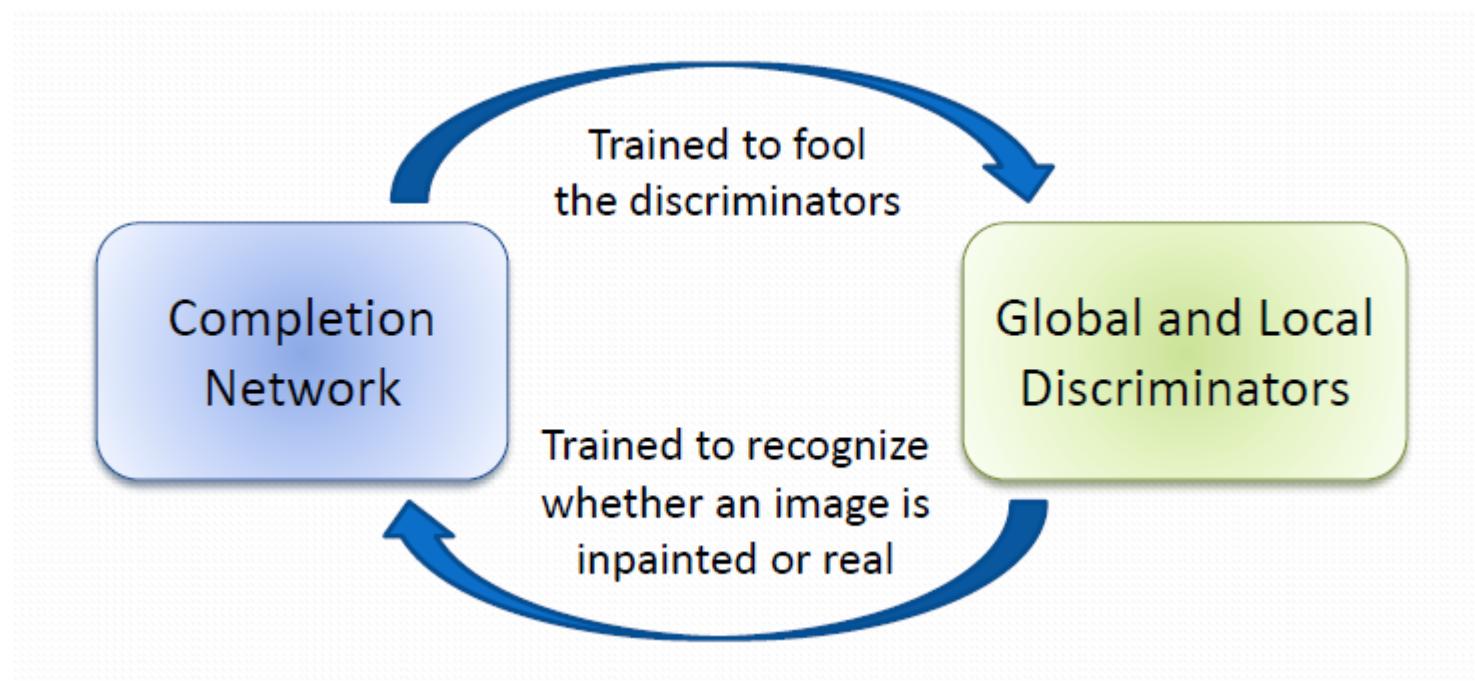
$p_1$  can be computed  
 $p_2$  cannot be computed



Both can be computed

# Training

- Alternately update the completion network and discriminators
  - Based on Generative Adversarial Networks
  - MSE(Mean Squared Error)



[https://github.com/satoshiiiizuka/siggraph2017\\_inpainting](https://github.com/satoshiiiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siizuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Loss function

- MSE loss:

$$L(x, M_c) = \| M_c \odot (C(x, M_c) - x) \|^2$$

- GAN optimization:

$$\min_C \max_D \mathbb{E} [ \log D(x, M_d) + \log(1 - D(C(x, M_c), M_c)) ]$$

- Combined optimization:

$$\begin{aligned} \min_C \max_D \mathbb{E} [ & L(x, M_c) + \alpha \log D(x, M_d) \\ & + \alpha \log(1 - D(C(x, M_c), M_c)) ] \end{aligned}$$

---

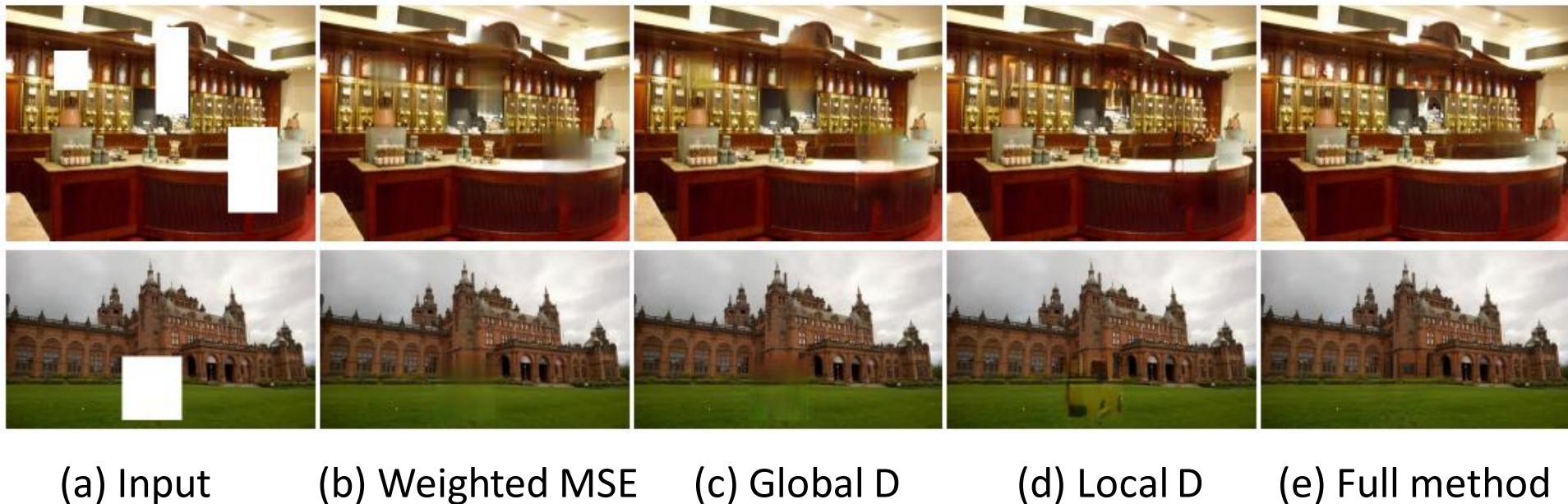
**Algorithm 1** Training procedure of the image completion network.

---

```
1: while iterations  $t < T_{train}$  do
2:   Sample a minibatch of images  $x$  from training data.
3:   Generate masks  $M_c$  with random holes for each image  $x$  in
   the minibatch.
4:   if  $t < T_C$  then
5:     Update the completion network  $C$  with the weighted MSE
     loss (Eq. (2)) using  $(x, M_c)$ .
6:   else
7:     Generate masks  $M_d$  with random holes for each image  $x$ 
     in the minibatch.
8:     Update the discriminators  $D$  with the binary cross entropy
     loss with both  $(C(x, M_c), M_c)$  and  $(x, M_d)$ .
9:     if  $t > T_C + T_D$  then
10:       Update the completion network  $C$  with the joint loss
        gradients (Eq. (5)) using  $(x, M_c)$ , and  $D$ .
11:     end if
12:   end if
13: end while
```

---

# Comparison of training with different discriminator configurations



# Training with Different Discriminator Configurations



Input



Mean Squared Error(MSE)



MSE + Global discriminator



MSE + Local discriminator



Full method

# Post-Processing



(a) Input

(b) Output

(c) Post-processed

Satoshi Iizuka, Edgar Simo-Serra, Hiroshi Ishikawa. Globally and Locally Consistent Image Completion. In SIGGRAPH ,2017

# Dataset

- Places2 dataset
  - About 8 million images with various scenes
  - Randomly generate a hole for training



Input

Ground truth

Places2 dataset



[https://github.com/satoshiizuka/siggraph2017\\_inpainting](https://github.com/satoshiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siiuzuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Result: Image Completion



[https://github.com/satoshiizuka/siggraph2017\\_inpainting](https://github.com/satoshiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siizuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Result: Object Removal



[https://github.com/satoshiizuka/siggraph2017\\_inpainting](https://github.com/satoshiizuka/siggraph2017_inpainting)

<https://www.slideshare.net/siiuzuka/siggraph-2017-globally-and-locally-consistent-image-completion>

# Comparison

Comparison with Photoshop CS6 Content Aware Fill (PatchMatch [Barnes et al. 2009]), [Darabi et al. 2012], [Huang et al. 2014], and [Pathak et al. 2016] on a diverse set of scenes using random masks.

- Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. [PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing..](#) ACM Transactions on Graphics (Proceedings of SIGGRAPH) 28, 3 (2009), 24:1–24:11.
- Soheil Darabi, Eli Shechtman, Connelly Barnes, Dan B Goldman, and Pradeep Sen. [Image Melding: Combining Inconsistent Images using Patch-based Synthesis..](#) ACM Transactions on Graphics (Proceedings of SIGGRAPH) 31, 4 (2012), 82:1–82:10.
- Jia-Bin Huang, Sing Bing Kang, Narendra Ahuja, and Johannes Kopf. [Image Completion Using Planar Structure Guidance..](#) ACM Transactions on Graphics (Proceedings of SIGGRAPH) 33, 4 (2014), 129:1-129:10.
- Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei Efros. [Context Encoders: Feature Learning by Inpainting](#). In IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2016.

# Comparison



Input



Patch Match  
[Barnes et al. 2009]



Image Melding  
[Darabi et al. 2012]



Content Encoder  
[Pathak et al. 2016]



Ours

# Comparison



Input



Patch Match  
[Barnes et al. 2009]



Image Melding  
[Darabi et al. 2012]



Content Encoder  
[Pathak et al. 2016]



Ours

# Application to Specific Dataset

- Fine-tuning the model using a specific dataset
  - Achieves more complicated inpainting

- Face dataset

- 200,000 training images
  - 2499 test images



Large-scale Celeb Faces Attributes Dataset(CelebA)

# Result: Face Completion



# Result: Removing Sunglasses



Original

Input

Output

# Failure Case



Input

Ours

Ground truth

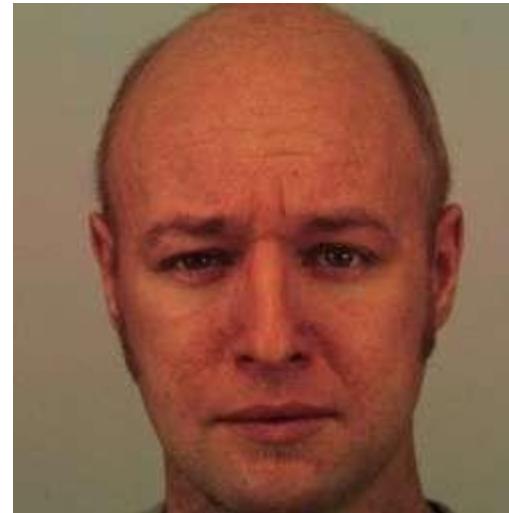
# Conclusion

- Novel network for globally and locally consistent image completion
  - Completion network for arbitrary region inpainting
  - Adversarial training with global and local discriminators

- Globally and Locally Completion
  - **Semantic Image Inpainting**
- High-Resolution Image Inpainting
- Content Encoders

# Semantic Image Inpainting with Deep Generative Models

Raymond A. Yeh\*, Chen Chen\*, Teck Yian Lim, Alexander G. Schwing, Mark Hasegawa-Johnson, Minh N. Do  
In CVPR 2017



[Image from KDEF dataset]



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

# Motivation

**Task:** Semantic image inpainting (filling large missing regions)

- ill-posed task
- requires strong prior knowledge on the data
- extracting information from only a single image produces unsatisfactory results

**Contributions:**

- deep generative models produce missing content by conditioning on available data
- inpainting as constrained optimization problem using context and prior loss



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

# Introduction

## Problem Formulation:

- Corrupted image:  $\mathbf{y}$
- Binary mask:  $\mathbf{M}$
- Task: predict uncorrupted version  $\hat{\mathbf{x}}$

Given:



## Baselines:

- Total Variation and Low Rank assume smoothness in the pixel space
- Context Encoder is a deep model which treats inpainting as a regression problem

Predict:



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

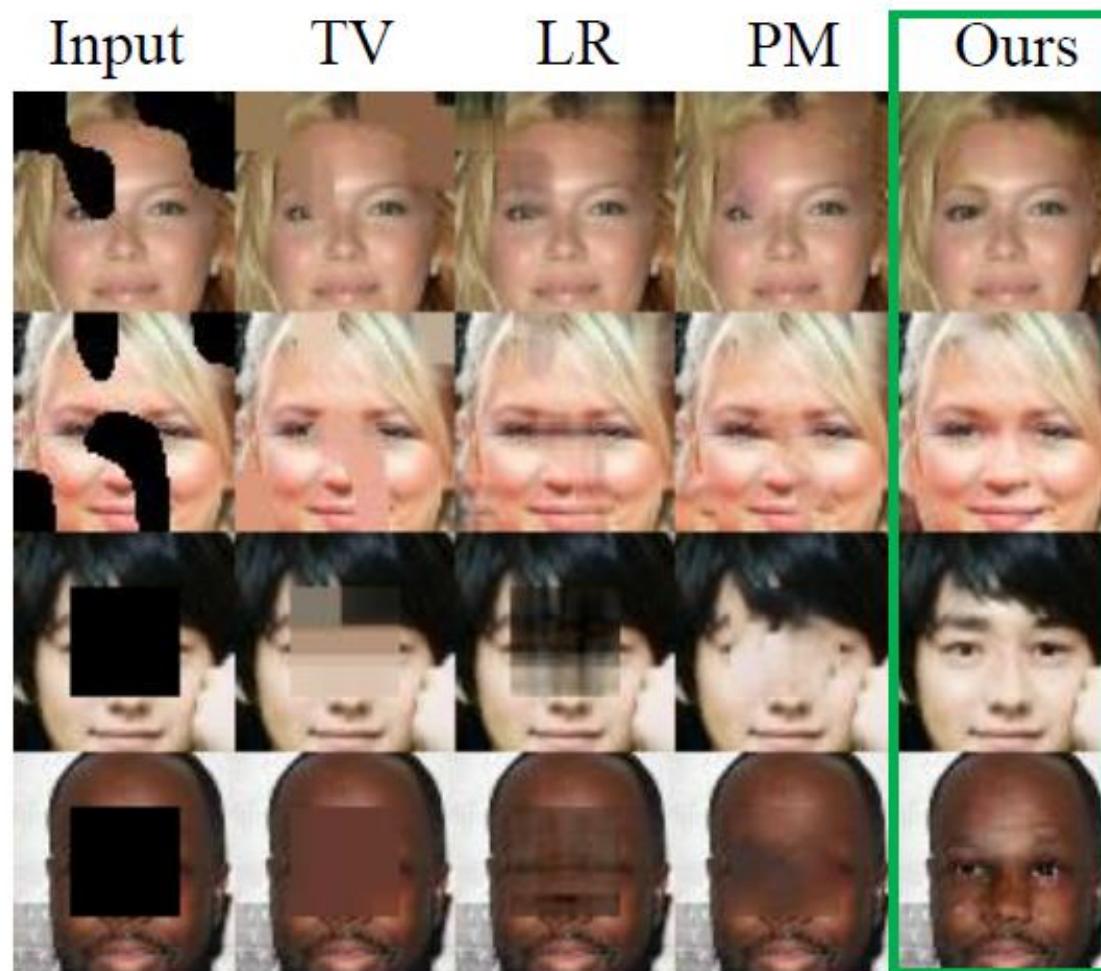


Figure 1. Semantic inpainting results by TV, LR, PM and our method. Holes are marked by black color.

# Introduction

- DCGAN-Based
- Not end-to-end
  - Train generator first (uncorrupted data)
  - Find  $\hat{z}$  for inpainting



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

<https://www.slideshare.net/ssuser73ec8f/semantic-image-inpainting>

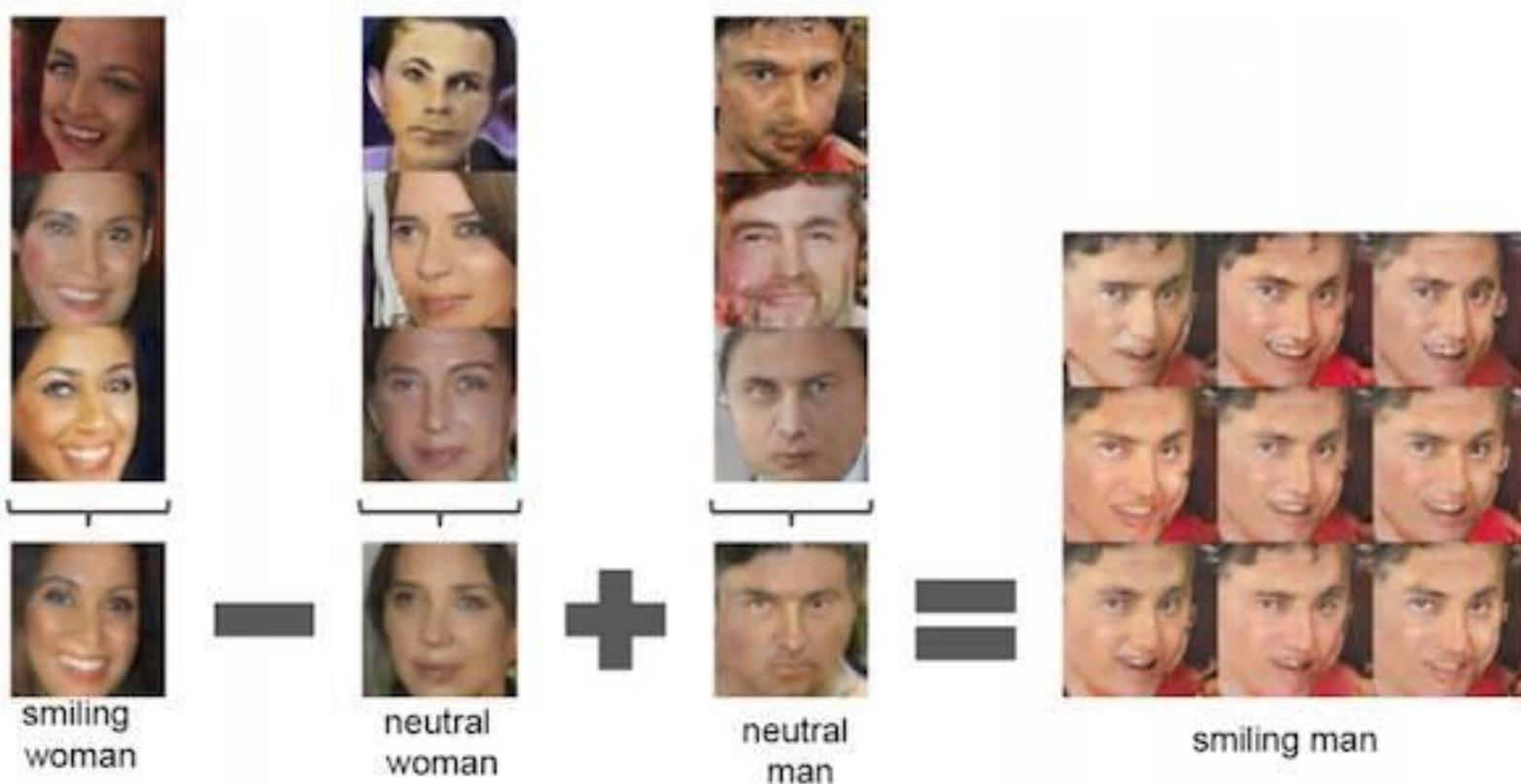
# Introduction

- DCGAN-Based



# Introduction

- DCGAN-Based



# Image Inpainting



Raymond A. Yeh, Chen Chen, Teck Yian Lim, Alexander G. Schwing, Mark Hasegawa-Johnson, Minh N. Semantic Image Inpainting with Deep Generative Models .In CVPR, 2017.

# Introduction

Instead of explicitly defining the prior, we utilize deep generative models to capture prior information

## Generative Adversarial Networks:

- Generator G: deep net mapping perturbation z to artificial sample
- Discriminator D: deep net discriminating between artificial and real sample, x
- Program:

$$\min_G \max_D V(G, D) = E_{x \sim P_{data}} [\log(D(x))] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

# Approach

## Intuition of our approach:

- **Hypothesis:** Trained  $G$  is efficient-image not from  $p_{data}$  (e.g., corrupted data) **should not lie** on the learned encoding manifold; use manifold can be used as a prior
- **Objective:** Find encoding  $\hat{z}$ : “closest” to the corrupted image while constrained to the manifold



[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

# Approach

**Solving for the “closest” encoding  $\hat{\mathbf{z}}$  :**

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \mathcal{L}_c(\mathbf{z} | \mathbf{y}, \mathbf{M}) + \mathcal{L}_p(\mathbf{z})$$

**Context Loss:** importance weighted metric  $\mathbf{W}$  to enforce similarity to the uncorrupted regions:

$$\mathcal{L}_c(\mathbf{z} | \mathbf{y}, \mathbf{M}) = \|\mathbf{W} \odot (G(\mathbf{z}) - \mathbf{y})\|_1$$

**Prior Loss:** prior penalizing unrealistic images based on the discriminator:

$$\mathcal{L}_p(\mathbf{z}) = \lambda \log(1 - D(G(\mathbf{z})))$$

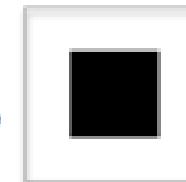
# Approach

CONTEXTUAL LOSS

PRIOR LOSS

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \mathcal{L}_c(\mathbf{z}|\mathbf{y}, \mathbf{M}) + \mathcal{L}_p(\mathbf{z})$$

$$Loss = L_p(\mathbf{z}) + L_c(\mathbf{z}|$$

 $\mathbf{y}$  $\mathbf{M}$ 

$\mathbf{y}$  : corrupted image

$\mathbf{M}$ : binary mask(size equal to the image)

# Approach

**Importance weight metric:** :



$$\mathbf{W}_i = \begin{cases} \sum_{j \in N(i)} \frac{(1 - \mathbf{M}_j)}{|N(i)|} & \text{if } \mathbf{M}_i \neq 0 \\ 0 & \text{if } \mathbf{M}_i = 0 \end{cases}$$

Wi: importance weight at pixel location i

N(i) : set of neighbors of pixel i in a local window

# Approach

## Context Loss:

Importance weighted metric  $\mathbf{W}$  to enforce similarity to the uncorrupted regions:

$$\mathcal{L}_c(\mathbf{z}|\mathbf{y}, \mathbf{M}) = \|\mathbf{W} \odot (G(\mathbf{z}) - \mathbf{y})\|_1$$

$\mathbf{y}$  : corrupted image

$\mathbf{M}$ : binary mask(size equal to the image)

## Prior Loss:

- How realistic the generated image is
- Prior penalizing unrealistic images based on the discriminator

$$\mathcal{L}_p(\mathbf{z}) = \lambda \log(1 - D(G(\mathbf{z})))$$

Without prior loss, the mapping from  $y$  to  $z$  may converge to a perceptually implausible result

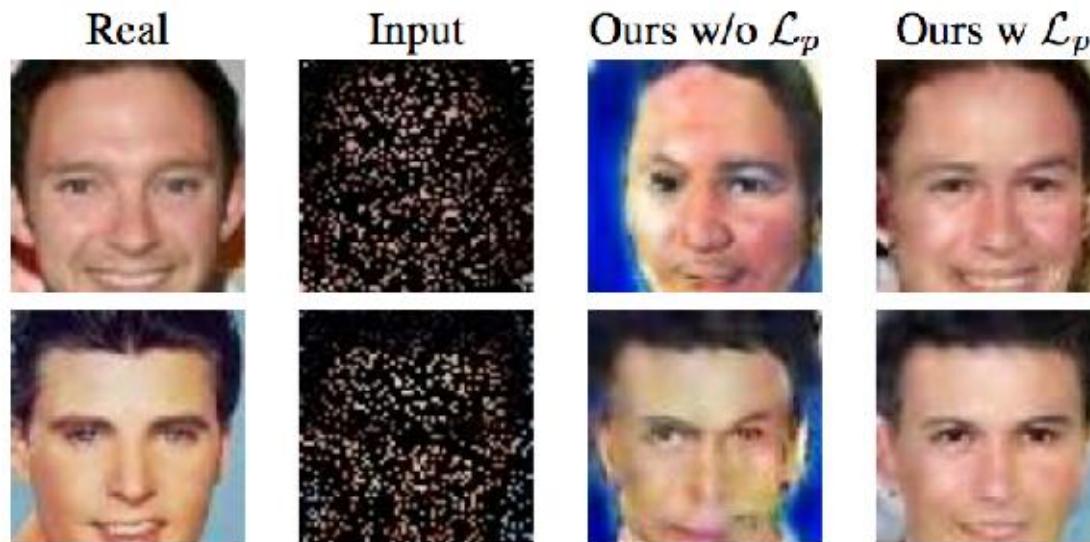
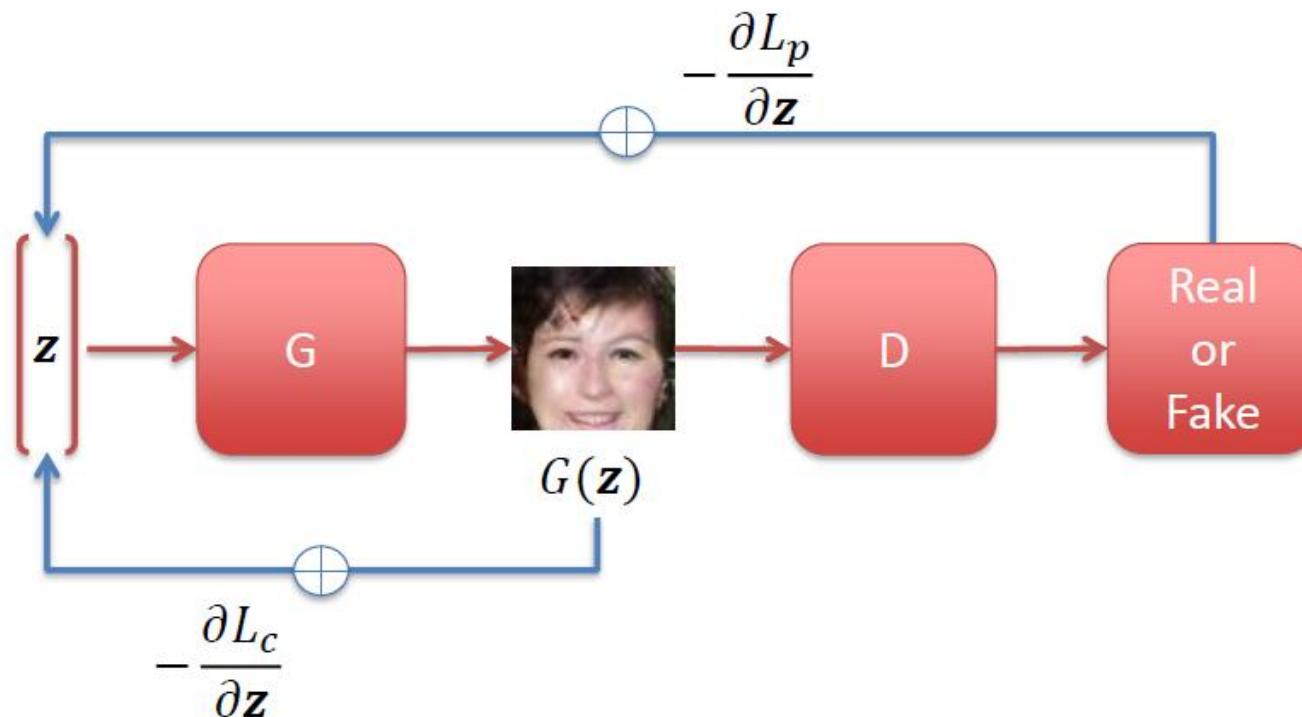
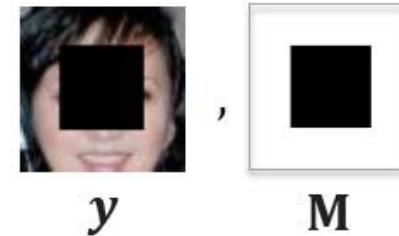


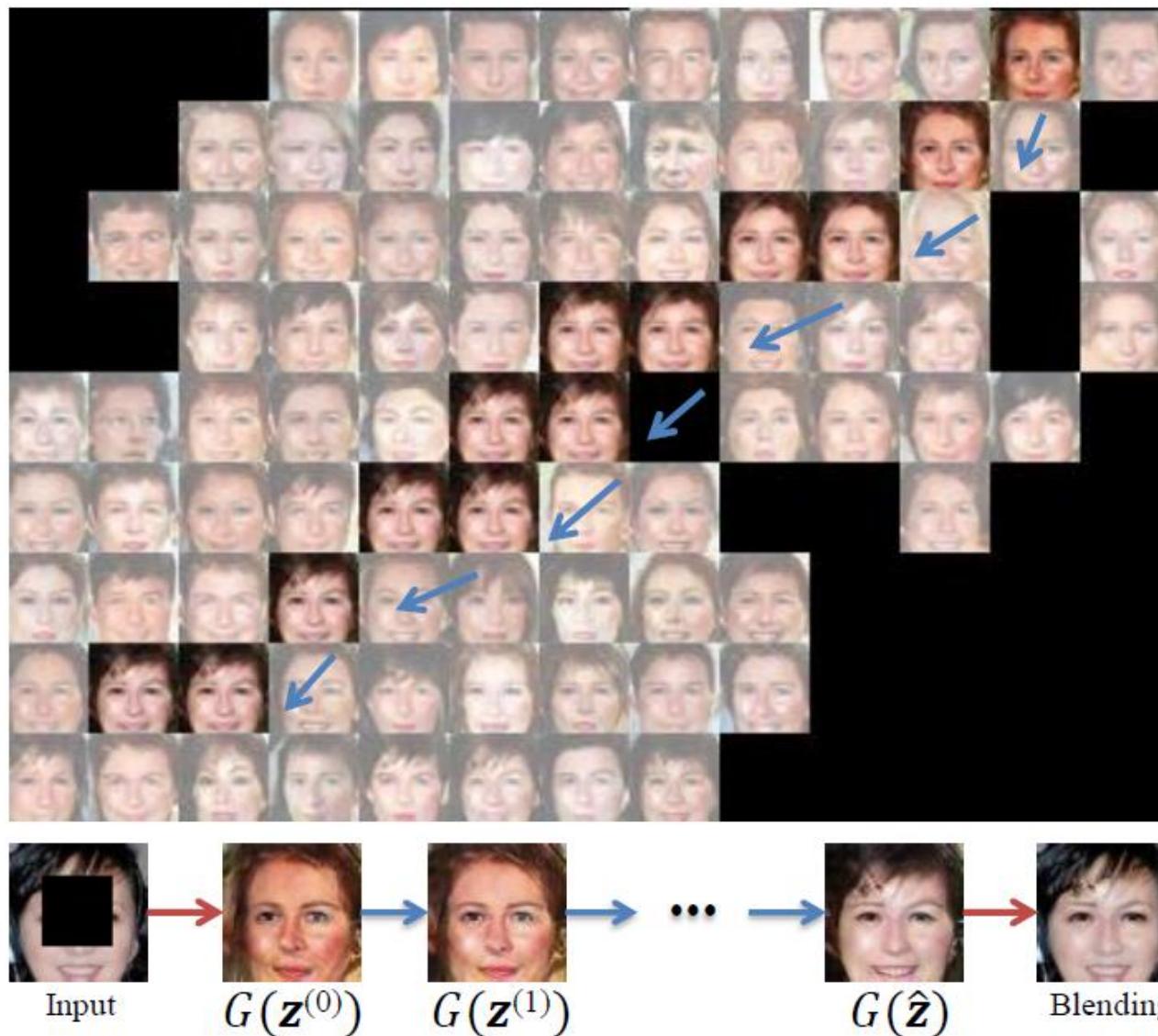
Figure 4. Inpainting with and without the prior loss.

# Illustration of the approach:



$$Loss = L_p(\mathbf{z}) + L_c(\mathbf{z} | \mathbf{y}, \mathbf{M})$$





Raymond A. Yeh, Chen Chen, Teck Yian Lim, Alexander G. Schwing, Mark Hasegawa-Johnson, Minh N. Semantic Image Inpainting with Deep Generative Models .In CVPR, 2017.

# Approach

## Poisson Blending:

Recovering prediction  $\hat{\mathbf{x}}$  via poisson blending rather than simple overlay:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\nabla \mathbf{x} - \nabla G(\hat{\mathbf{z}})\|_2^2 \quad \text{s.t. } \mathbf{x}_i = \mathbf{y}_i \text{ for } \mathbf{M}_i = 1$$

## Comparison: Poisson Blending vs. Overlay:

Overlay



Blend



Overlay

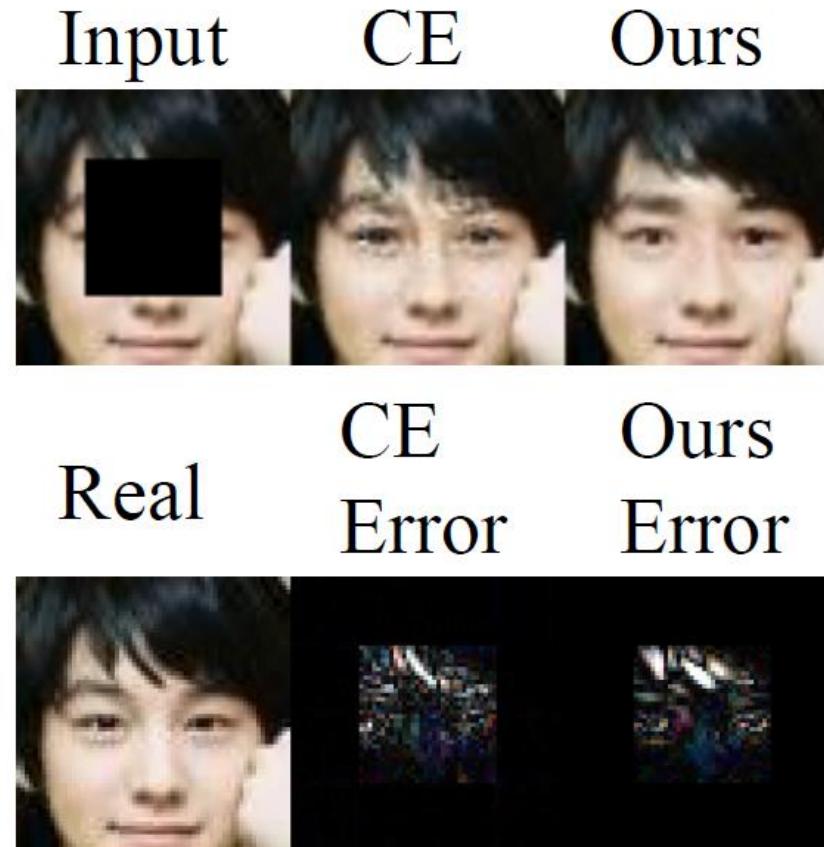


Blend



# Results

## Quantitative Results:



# Quantitative Results:

The PSNR values (dB) on the test sets. Left/right results are by Context Encoder (CE)/ours:

Masks/Dataset	CelebA	SVHN	Cars
Center	<b>21.3</b> /19.4	<b>22.3</b> /19.0	<b>14.1</b> /13.5
pattern	<b>19.2</b> /17.4	<b>22.3</b> /19.8	14.0/ <b>14.1</b>
random	20.6/ <b>22.8</b>	24.1/ <b>33.0</b>	16.1/ <b>18.9</b>
half	<b>15.5</b> /13.7	<b>19.1</b> /14.6	<b>12.6</b> /11.1

- In the figure above, PSNR for CE is 24.71 dB and ours is 22.98 dB
- Higher PSNR does not mean better visual quality
- The solution is not unique, many hallucinations are reasonable

# Datasets:

Images: 64\*64

- CelebA dataset

202,599 face images (2,000 for test)

- Car dataset

16,185 car images (8,041 for test)

- SVHN dataset

99,289 number images (26,032 for test)



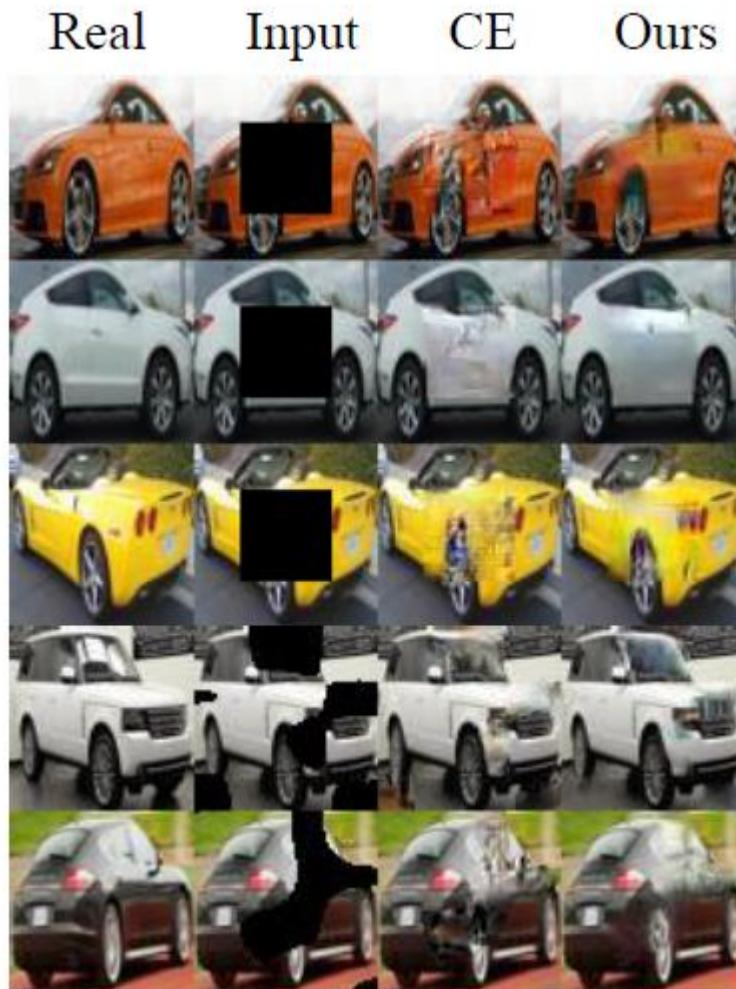
[https://github.com/moodoki/semantic\\_image\\_inpainting](https://github.com/moodoki/semantic_image_inpainting)

[http://www.isle.illinois.edu/~yeh17/projects/semantic\\_inpaint/poster.pdf](http://www.isle.illinois.edu/~yeh17/projects/semantic_inpaint/poster.pdf)

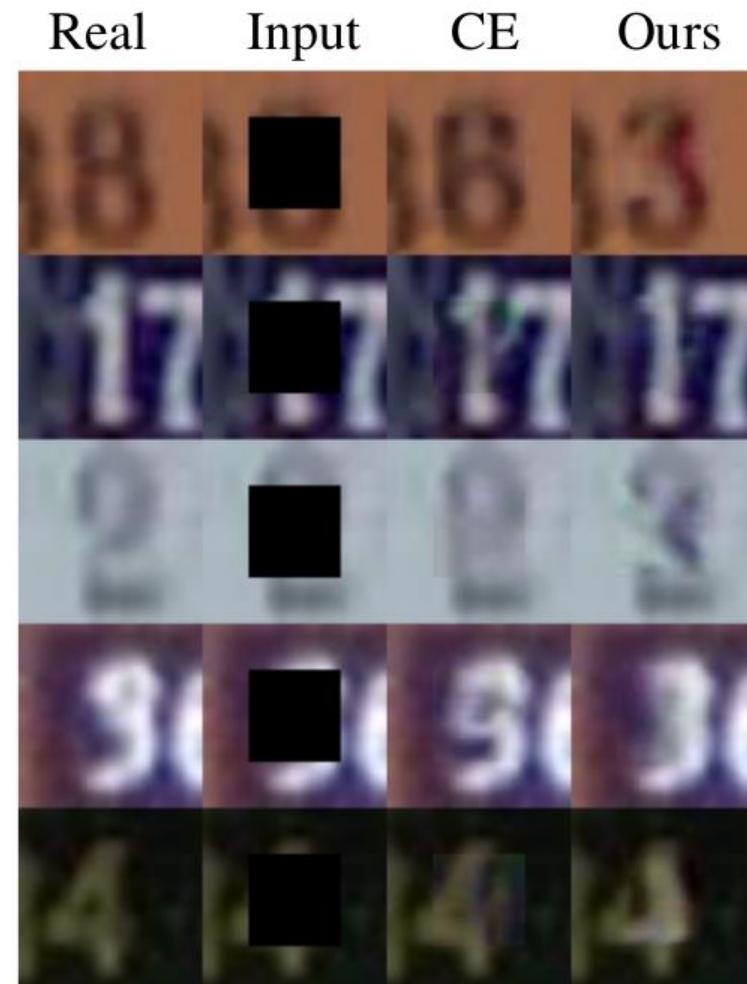
# Results - CelebA:



# Results – The Stanford Car:



# Results – SVHN:



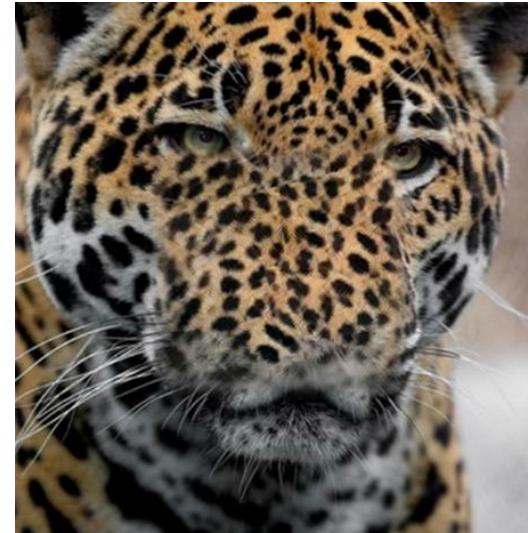
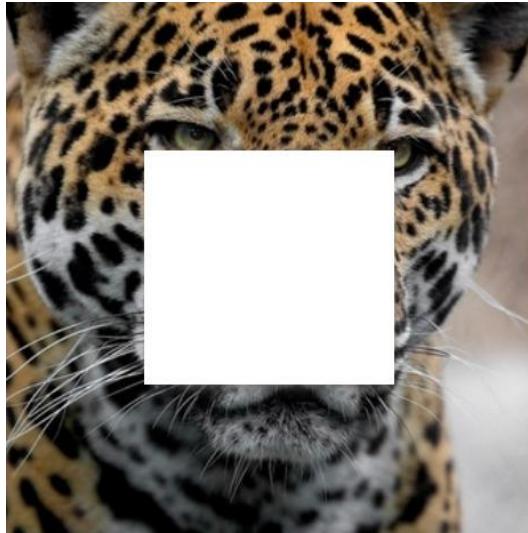
For paper, results, code and more:



- Globally and Locally Completion
- Semantic Image Inpainting
- **High-Resolution Image Inpainting**
- Content Encoders

# High-Resolution Image Inpainting using Multi-Scale Neural Patch Synthesis

Chao Yang, Xin Lu, Zhe Lin, Eli Shechtman, Oliver Wang, and Hao Li  
In CVPR 2017



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Bad things happen



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Introduction

Imaging you would like to edit a photo after breaking up, or restore an old picture from damages, we designed a MULTI-SCALE DEEP LEARNING algorithm to help you! Our approach

- Proposed a joint optimization framework that can hallucinates missing image regions by modeling a global content constraint and local texture constraint with convolutional neural networks.
- Further introduced a multi-scale neural patch synthesis algorithm for high-resolution image inpainting based on the joint optimization framework.



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Algorithm

---

## High-Resolution Hole Filling with Multi-Scale Neural Patch Synthesis

---

**Input:** Image  $x$ ,  
the content network  $f$ ,  
the texture network  $t$ ,  
the number of scales  $N$



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Algorithm

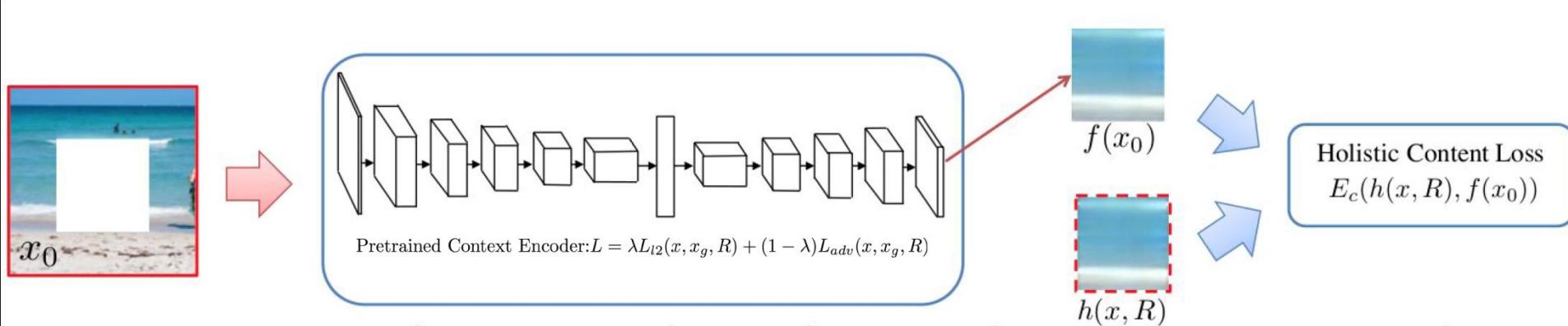
---

## High-Resolution Hole Filling with Multi-Scale Neural Patch Synthesis

---

- 1: Downsize  $x$  to  $128 \times 128$ .
- 2: Compute the initial content reference  $x^1$ .  
by giving  $x$  as input to  $f$ .
- 3: **for**  $s \in [1, 2, \dots, N]$ :
- 4: Initialize  $\tilde{x} = x^s$ .
- 4: Update  $\tilde{x}$  that minimizes the joint loss:  
$$L = L_{content} + L_{texture} + L_{tv-smoothness}.$$
- 5: Compute  $x^{s+1}$  by up-sampling  $\tilde{x}$ .
- 7: **end for**
- 8: Return  $\tilde{x}^N$ .

# The Content Network



## Context Encoder Predicts the Low-Res Content

The content constraint:

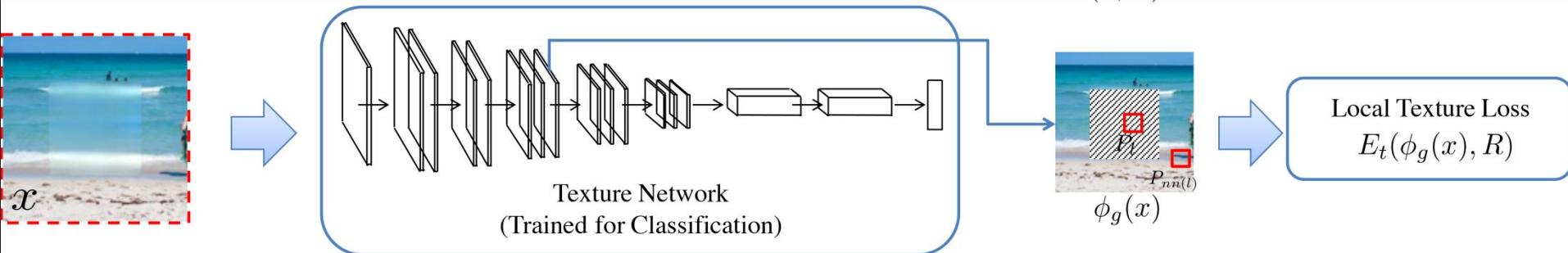
$$E_c(h(x, R), h(x_i, R)) = \| h(x, R) - h(x_i, R) \|_2^2$$



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# The Texture Network



## Pre-trained VGG Optimizes the High-Res Texture

The texture constraint:

$$E_t(\phi_t(x), R) = \frac{1}{|R^\phi|} \sum_{i \in R^\phi} \| h(\phi_t(x), P_i) - h(\phi_t(x), P_{nn(i)}) \|_2^2$$

 GitHub <https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# The Joint Loss Function

At each iteration, we minimize:

$$\begin{aligned} \tilde{x}_{i+1} = & \arg \min_x E_c(h(x, R), h(x_i, R)) \\ & + \alpha E_t(\phi_t(x), R^\phi) + \beta \gamma(x) \end{aligned}$$

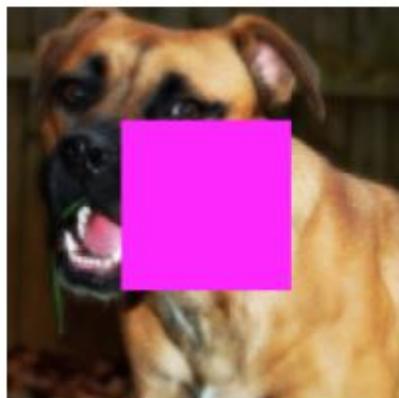


<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Multi-Scale Optimization

We optimize at three scales: **128, 256 and 512**:



Input



Iter 1



Iter 2



Iter 3



[GitHub](https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting) <https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# Changing the texture weight $\alpha$

The weight  $\alpha$  measures the contribution of the texture constraint relative to the content constraint.

It is a trade off between the sharpness of the texture and coherence of the structure:



(a) Input image

(b)  $\alpha = 1e - 6$

(c)  $\alpha = 1e - 5$

(d)  $\alpha = 4e - 5$

# Dropping the content constraint

Without using the content term to guide the optimization, the structure of the inpainting results is completely incorrect, although they are visually sharp:



Input



Output without CN



Output with CN

# Dropping the adversarial loss

When the initial prediction is blurry (using L2 loss only), the final result becomes blurrier as well comparing with using the content network trained with both L2 and adversarial loss.



(a) CN with  $\ell_2$  loss

(b) (a)'s final output

(c) CN with  $\ell_2$  & Adv loss

(d) (c)'s final output.

# Overall comparison with other methods

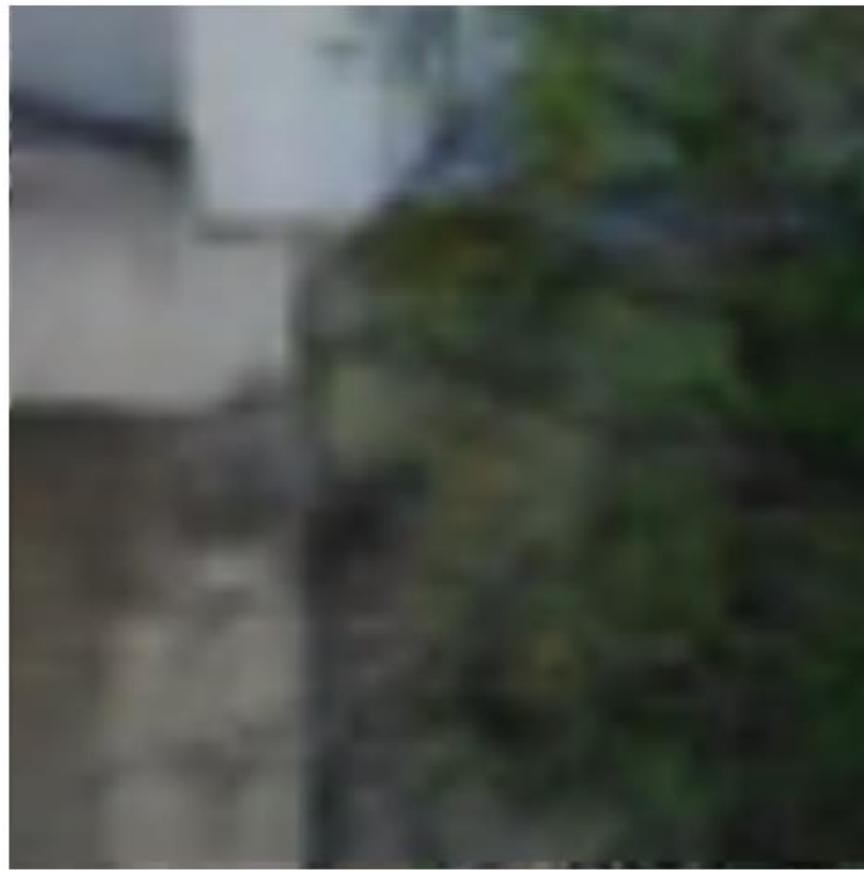


**(a) Input hole    (b) Context Encoder ( $l_2$  loss)    (c) Content-Aware Fill    (d) Context Encoder ( $l_2 +$  Adversarial loss)    (e) Our Method ( $l_2 +$  Adversarial loss)**



<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)



**(d) Zoom-in**



**(e) Zoom-in**

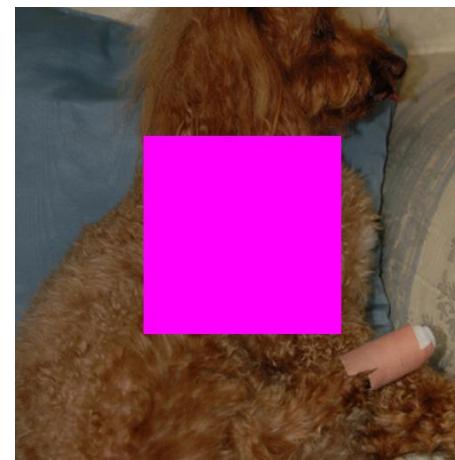
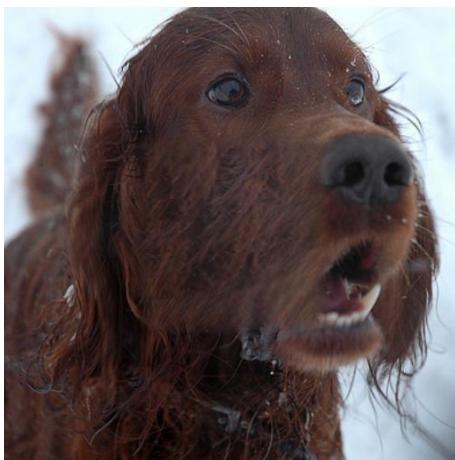
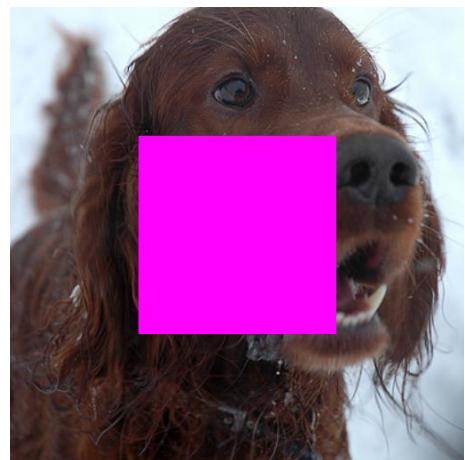
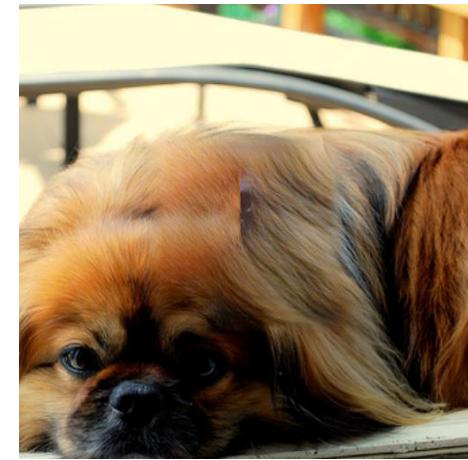
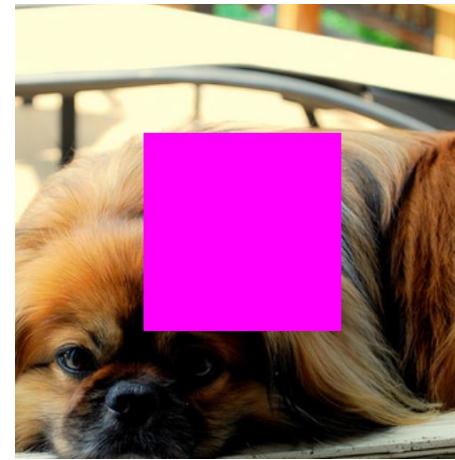


<https://github.com/leehomyc/Faster-High-Res-Neural-Inpainting>

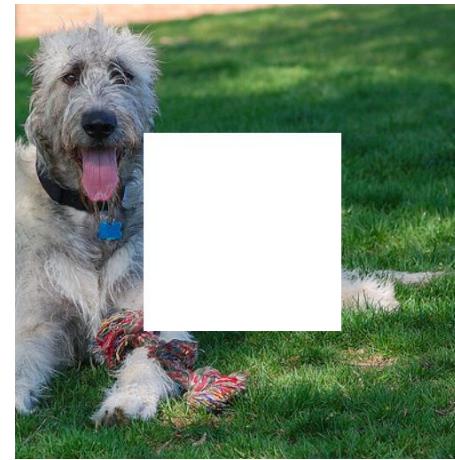
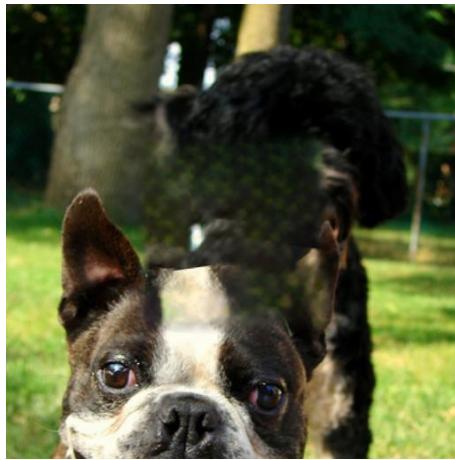
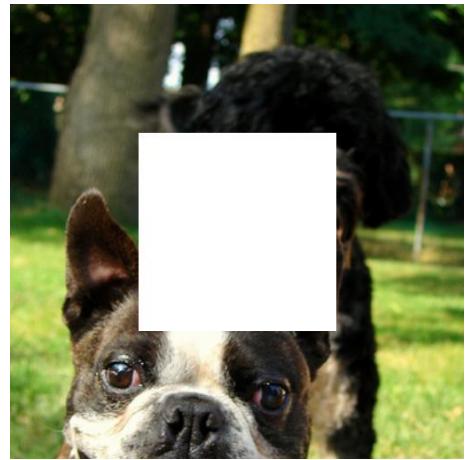
[https://github.com/leehomyc/CVPR\\_2017\\_Poster/](https://github.com/leehomyc/CVPR_2017_Poster/)

# More Results

And we have collected so many dogs and cats for you:



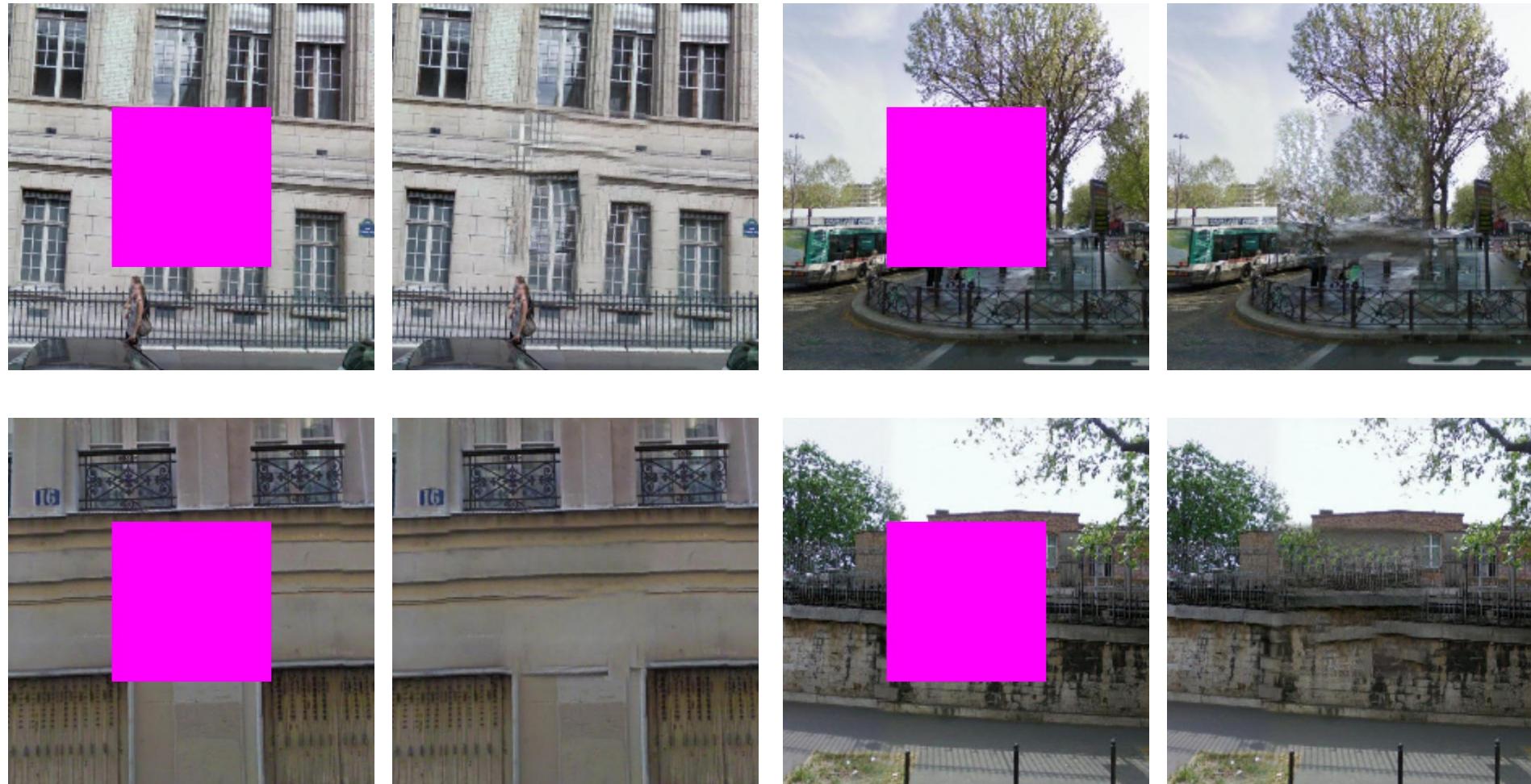
And we have collected so many dogs and cats for you:



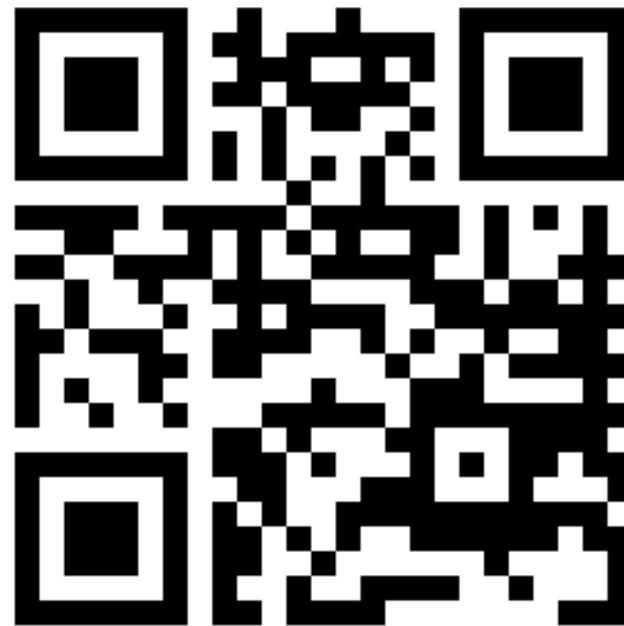
# And what inpaints Paris like Paris ?



# And what inpaints Paris like Paris ?



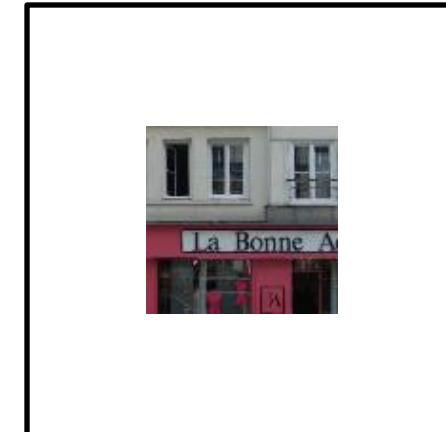
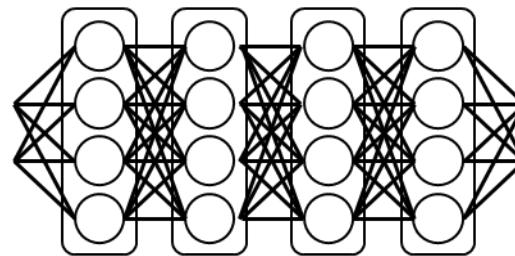
For paper, results, code and more:



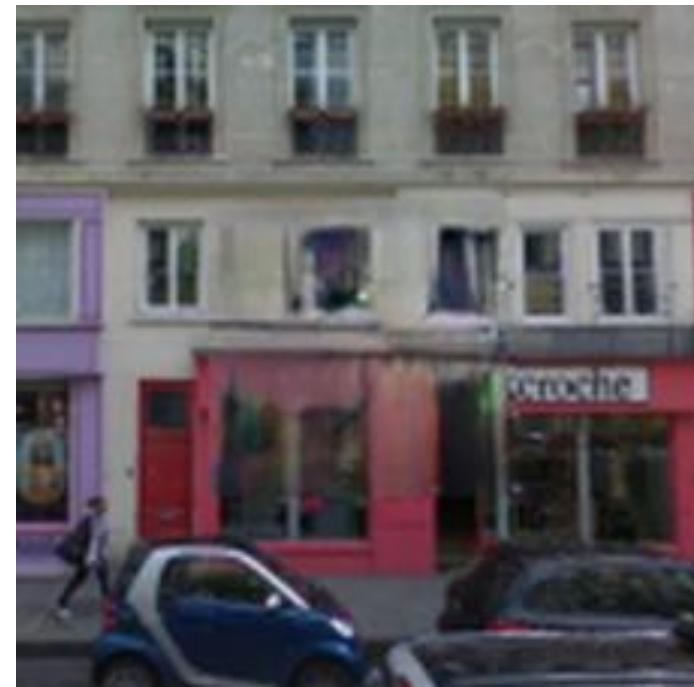
- Globally and Locally Completion
  - Semantic Image Inpainting
  - High-Resolution Image Inpainting
- **Context Encoders**

# Context Encoders: Feature Learning by Inpainting

Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor  
Darrell, Alexei Efros  
In CVPR 2016



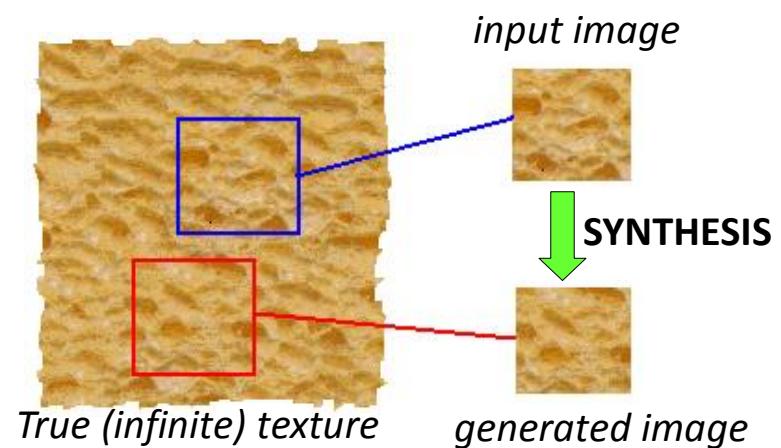
# Context Encoders: Feature Learning by Inpainting



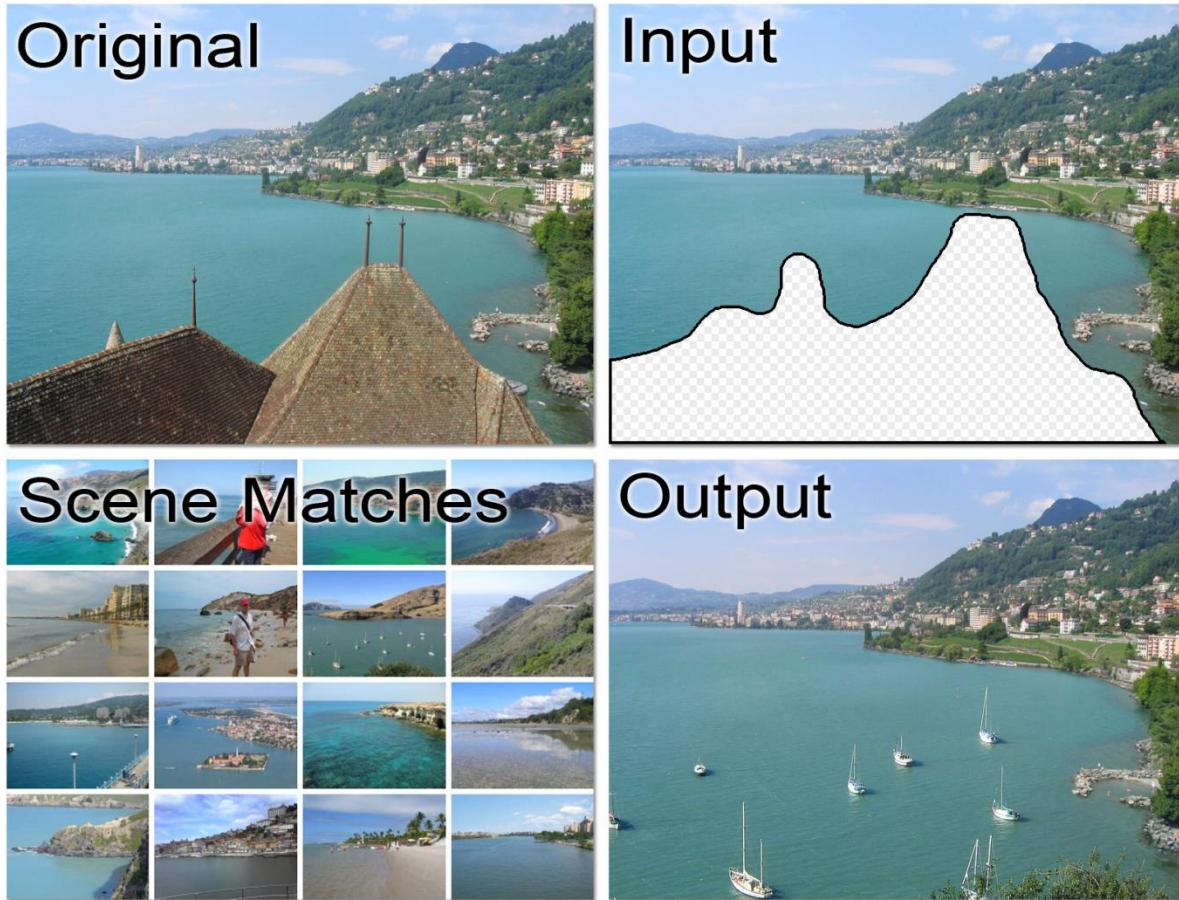
# Image Inpainting

- Inpaint large missing regions in image: not possible for local texture synthesis methods  
[Efros & Leung 1999]

- Nearest neighbor based scene completion  
[Hays & Efros 2007]



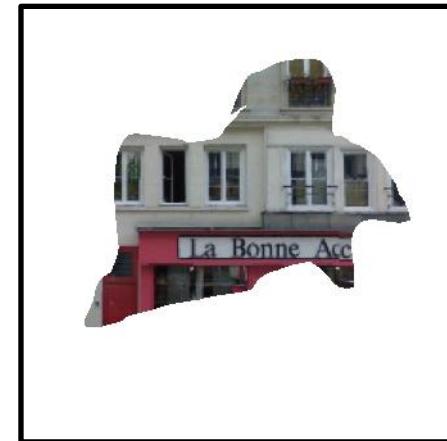
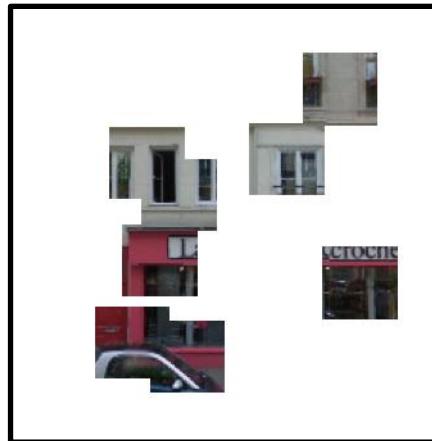
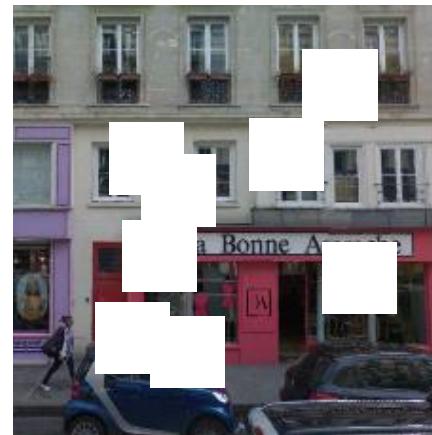
# Image Inpainting



<https://github.com/pathak22/context-encoder>

[https://people.eecs.berkeley.edu/~pathak/context\\_encoder/slides/](https://people.eecs.berkeley.edu/~pathak/context_encoder/slides/)

# Schemes for Region Dropout

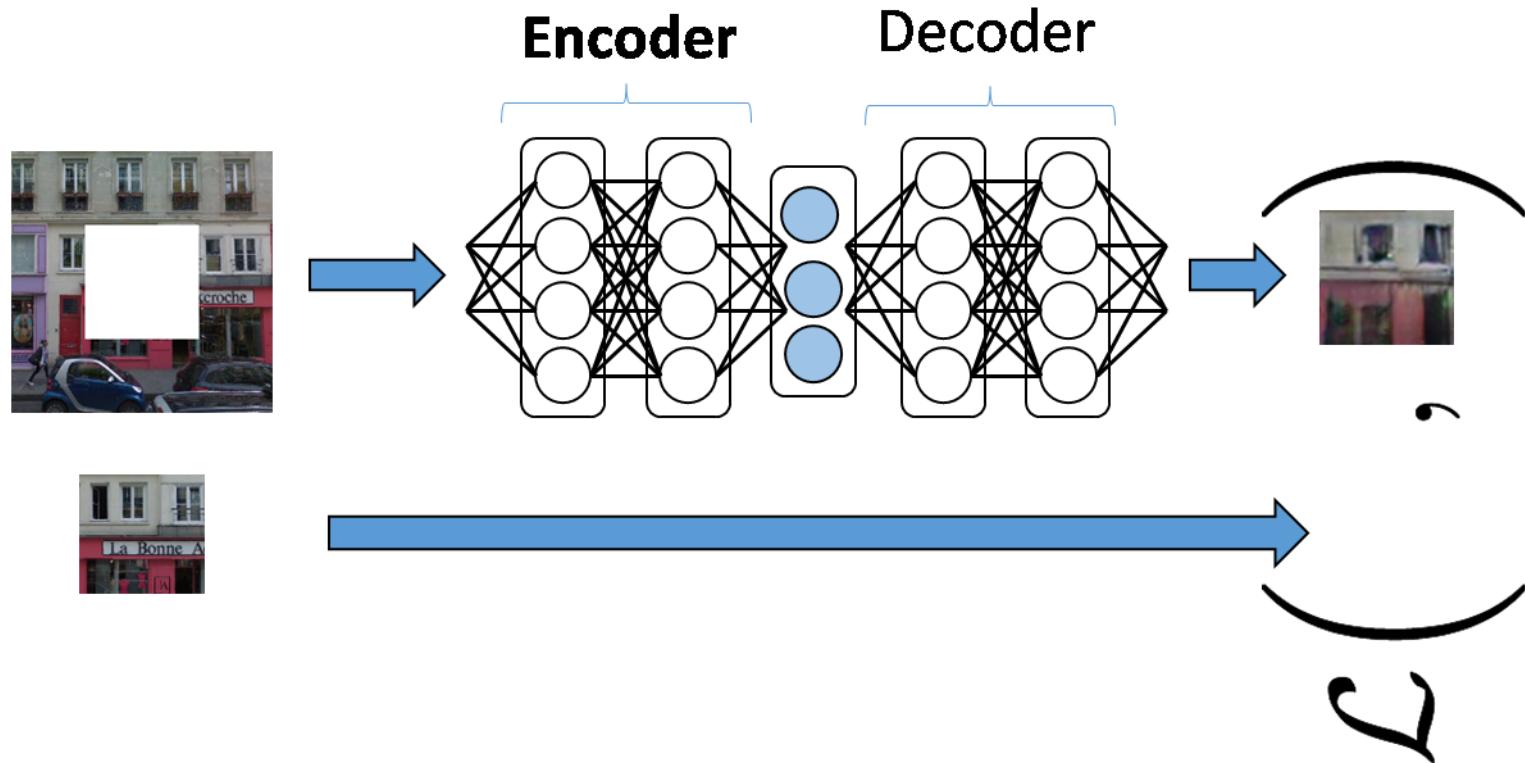


(a) Center Region

(b) Random Blocks

(c) Random Shapes

# Context Encoders



- Encoder can be substituted with any network architecture like AlexNet etc.
- Decoder is a set of UpConv/deconv/frac-strided-conv layers

# Nearest Neighbors with L2 Loss

- Nearest Neighbors are good in the representation space
- But generated output is blurry => Artifact of L2 !!



# Need Mode Picking => GAN

arXiv.org > stat > arXiv:1406.2661

Search or /

(Help | Advanced)

Statistics > Machine Learning

## Generative Adversarial Networks

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio

(Submitted on 10 Jun 2014)

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions  $G$  and  $D$ , a unique solution exists, with  $G$  recovering the training data distribution and  $D$  equal to  $1/2$  everywhere. In the case where  $G$  and  $D$  are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

Subjects: Machine Learning (stat.ML); Learning (cs.LG)

Cite as: arXiv:1406.2661 [stat.ML]

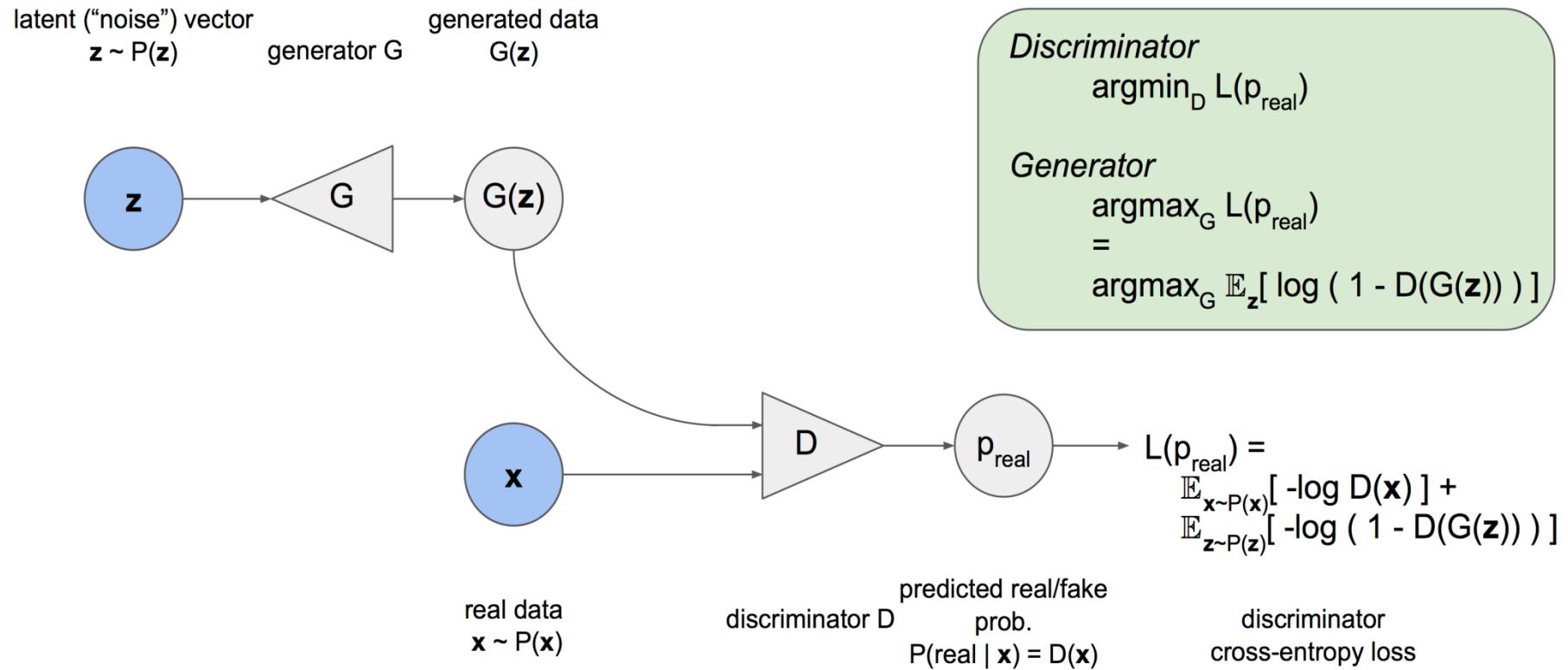
(or arXiv:1406.2661v1 [stat.ML] for this version)

### Submission history

From: Ian Goodfellow [view email]

[v1] Tue, 10 Jun 2014 18:58:17 GMT (1257kb,D)

# Generative Adversarial Networks



[Slide adapted from Jeff Donahue]

# LOSS:

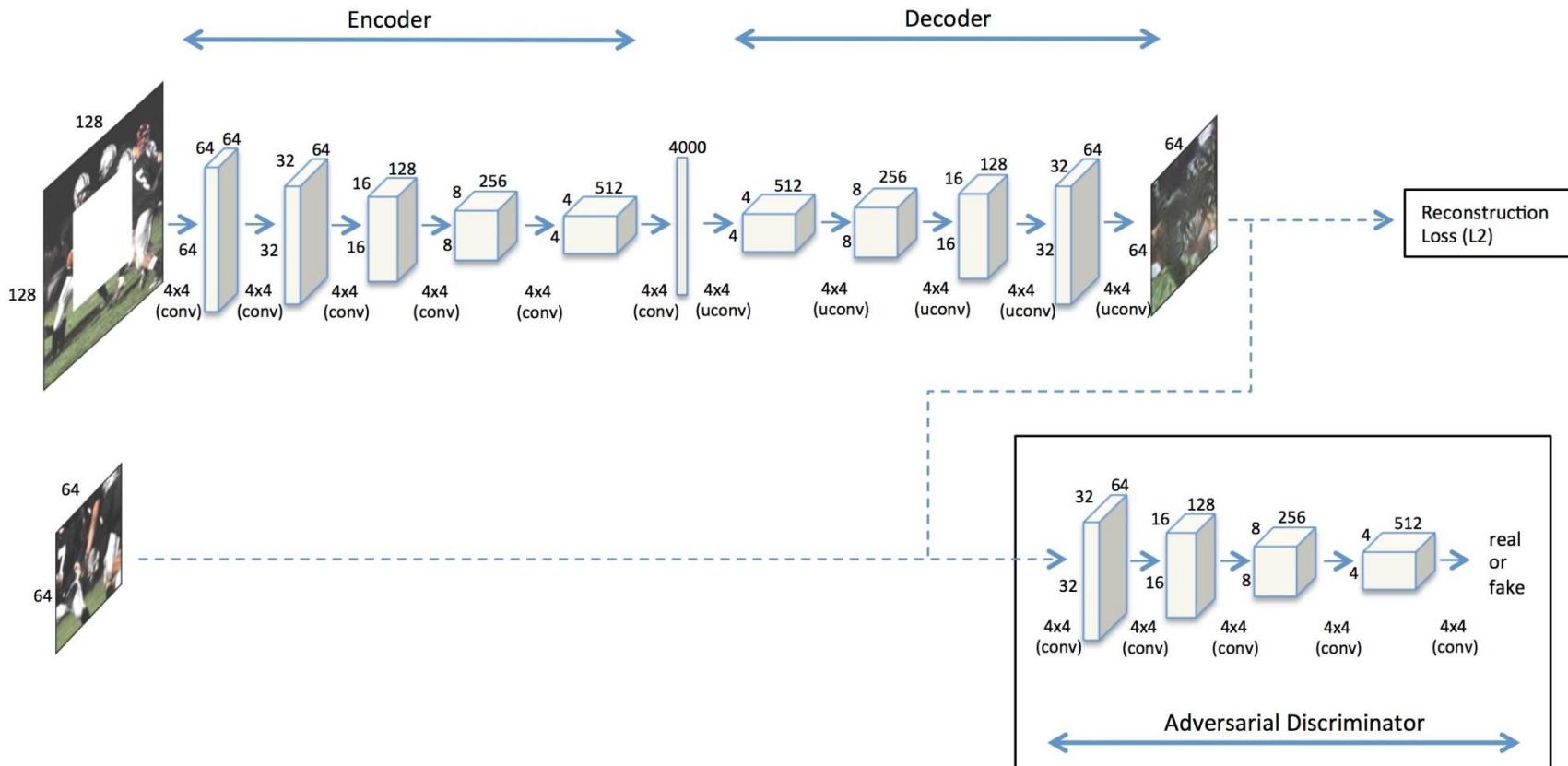
$$\mathcal{L}_{rec}(x) = \|\hat{M} \odot (x - F((1 - \hat{M}) \odot x))\|_2^2$$

$$\mathcal{L}_{adv}(x) = \max_D \mathbb{E}_{x \in \mathcal{X}} [\log(D(x)) + \log(1 - D(F((1 - \hat{M}) \odot x)))]$$

$$\mathcal{L}_{joint}(x) = \lambda_{rec}\mathcal{L}_{rec}(x) + \lambda_{adv}\mathcal{L}_{adv}(x)$$

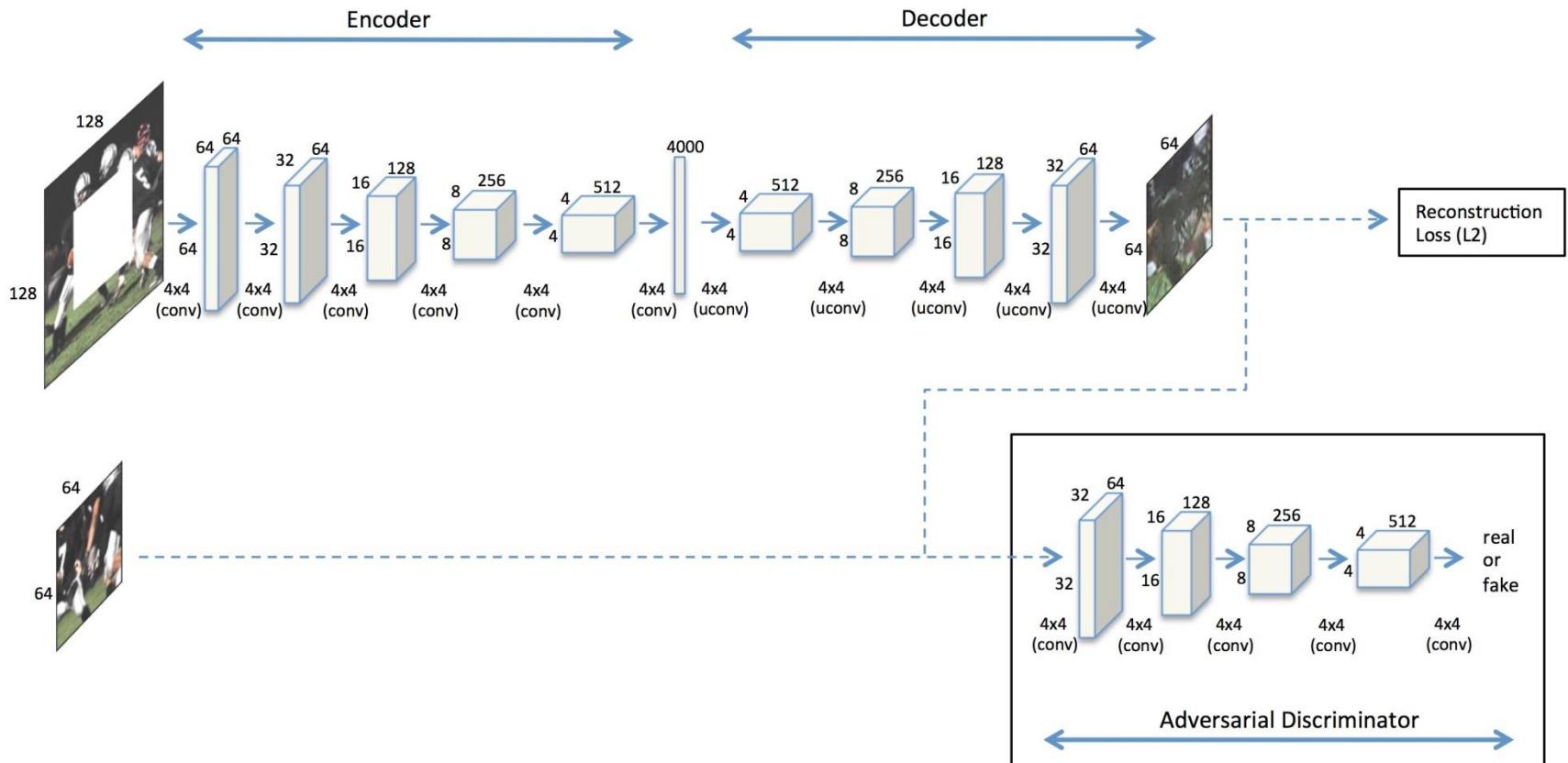
# Key Details

- Encoder is not conditioned on noise
- Discriminator is not conditioned on input

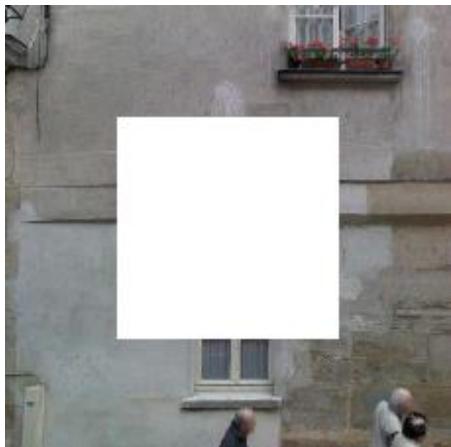


# Decoupled Loss

- Reconstruction  $L_2$  loss ensures “**correctness**”
- Adversarial Loss ensures “**realness**”



# Results: Decoupled Loss



Input Image



L2 Loss



Adversarial Loss



Joint Loss

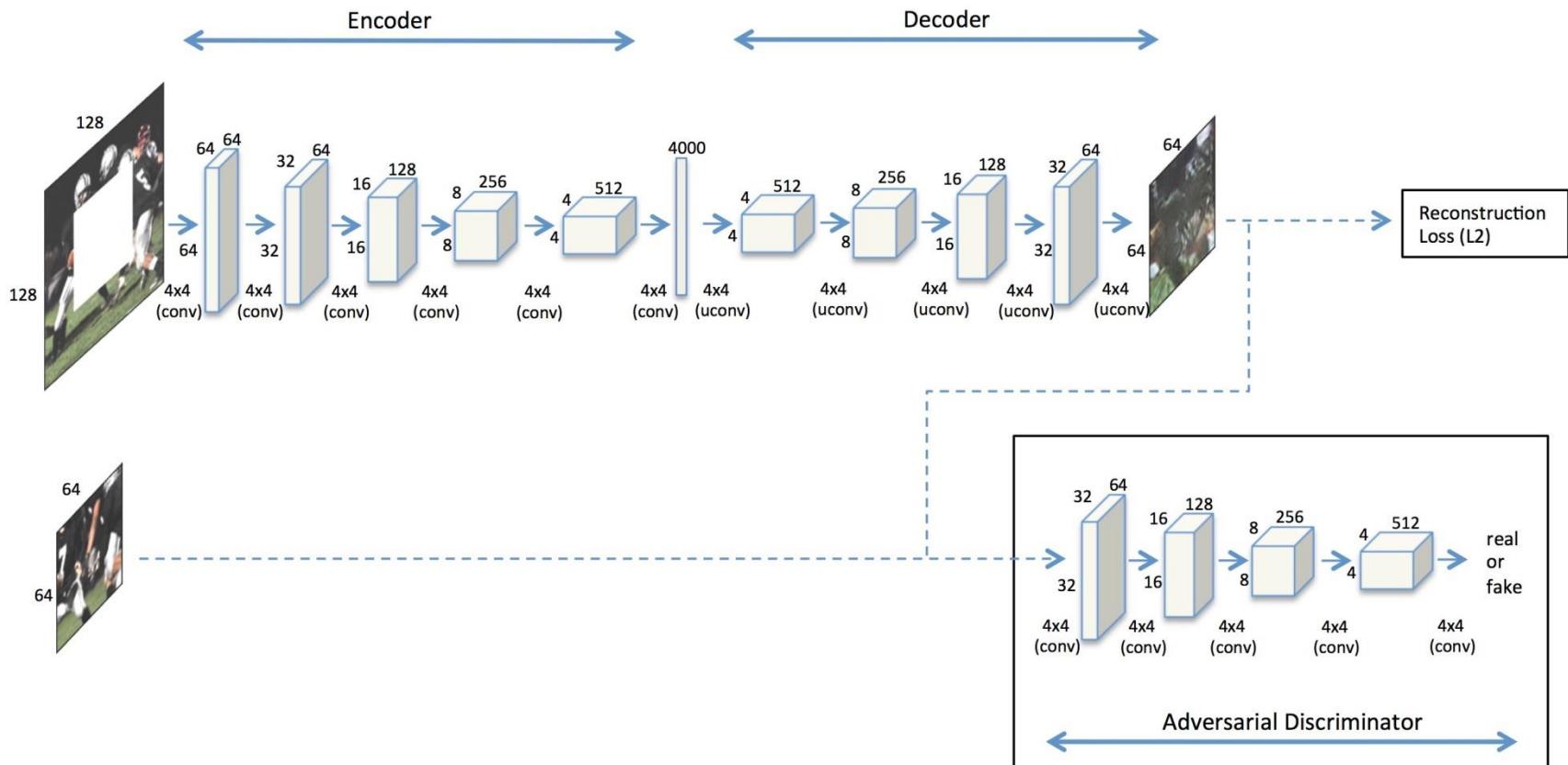
# ImageNet Results



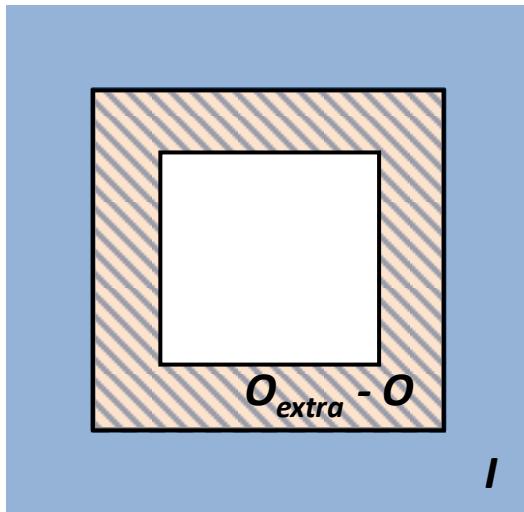
Has artifacts near the boundary !

# Boundary Artifacts

- Discriminator is not conditioned on input !



# Indirect Conditioning



Input Area =  $I$

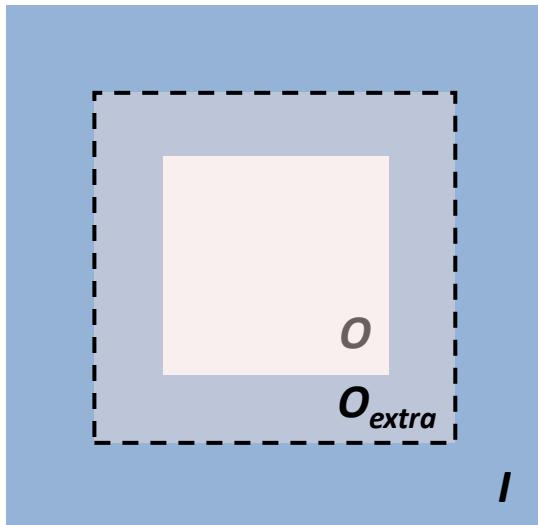
Target Area =  $O$

Predict Area =  $O_{extra}$

Overlap =  $O_{extra} - O$

Penalize L2 loss heavily in overlapping area  
=> so that network memorizes input in the overlap

# Indirect Conditioning

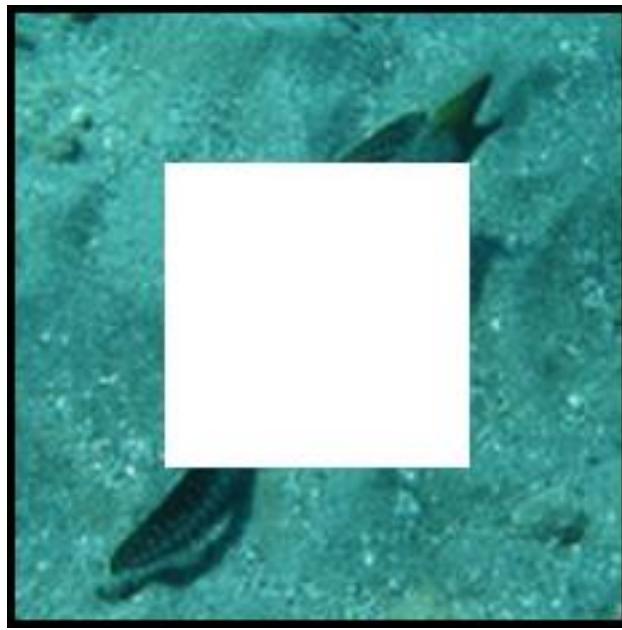


## INDIRECT CONSISTENCY

- Adversary ensures  $O_{extra}$  is coherent (to be realistic)
  - L2 Penalty ensures predicted  $O - O_{extra}$  is same as overlapping area with input  $I$
- Predicted  $O$  is coherent with  $I$

Penalize L2 loss heavily in overlapping area  
=> so that network memorizes input in the overlap

# Indirect Conditioning Results



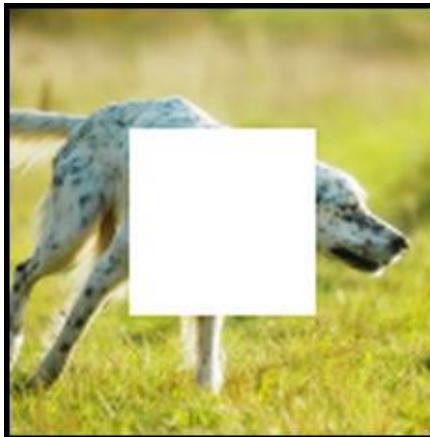
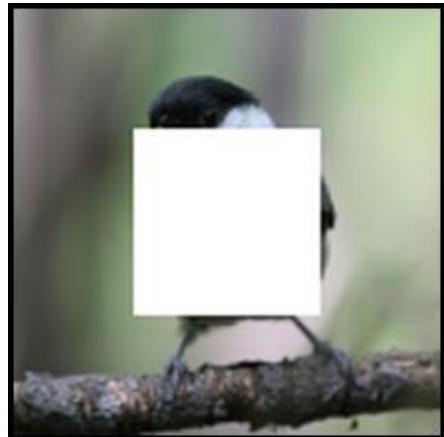
# Indirect Conditioning Results



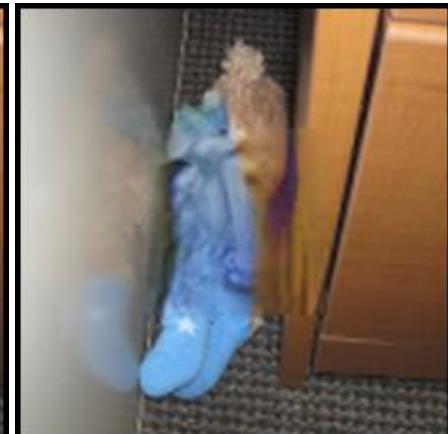
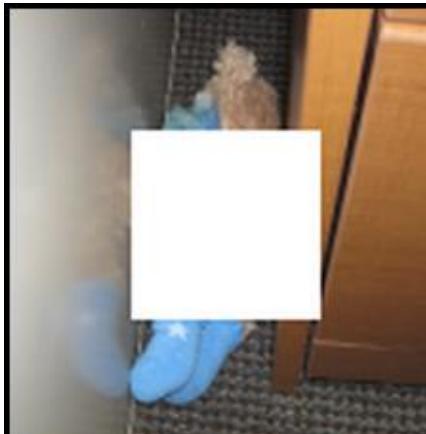
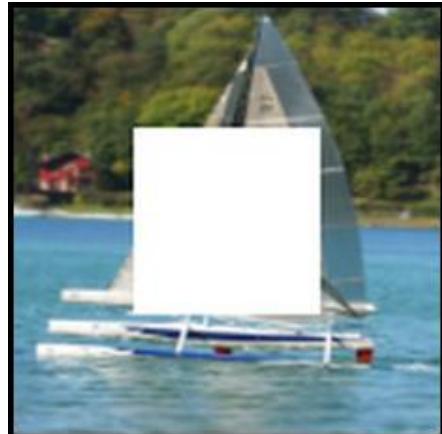
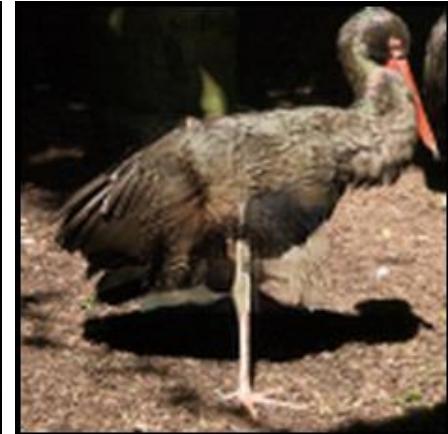
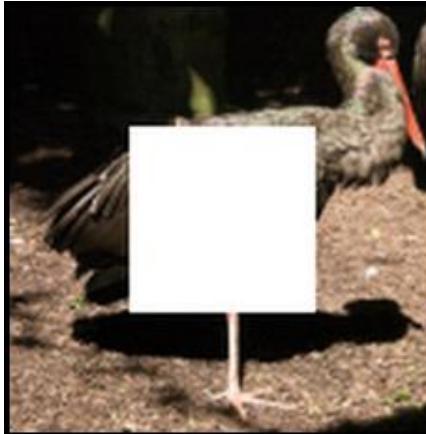
# Indirect Conditioning Results



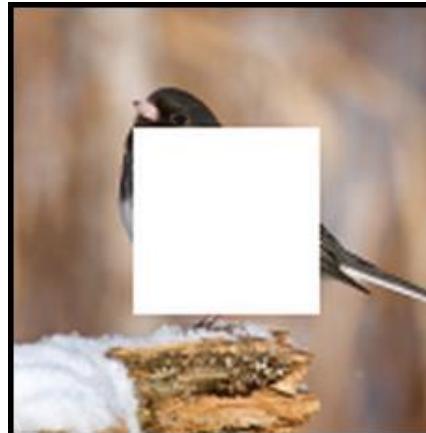
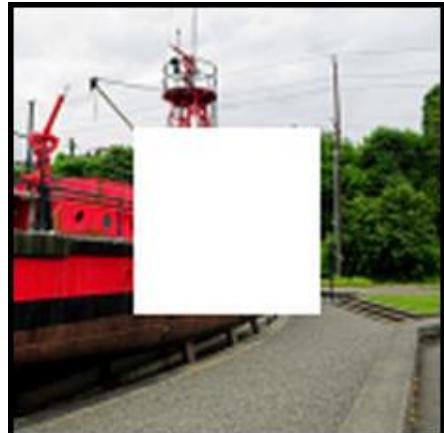
# More Results



# More Results



# More Results



# More Results



# Compare (algorithm in Adobe Photoshop)



Input



Ours



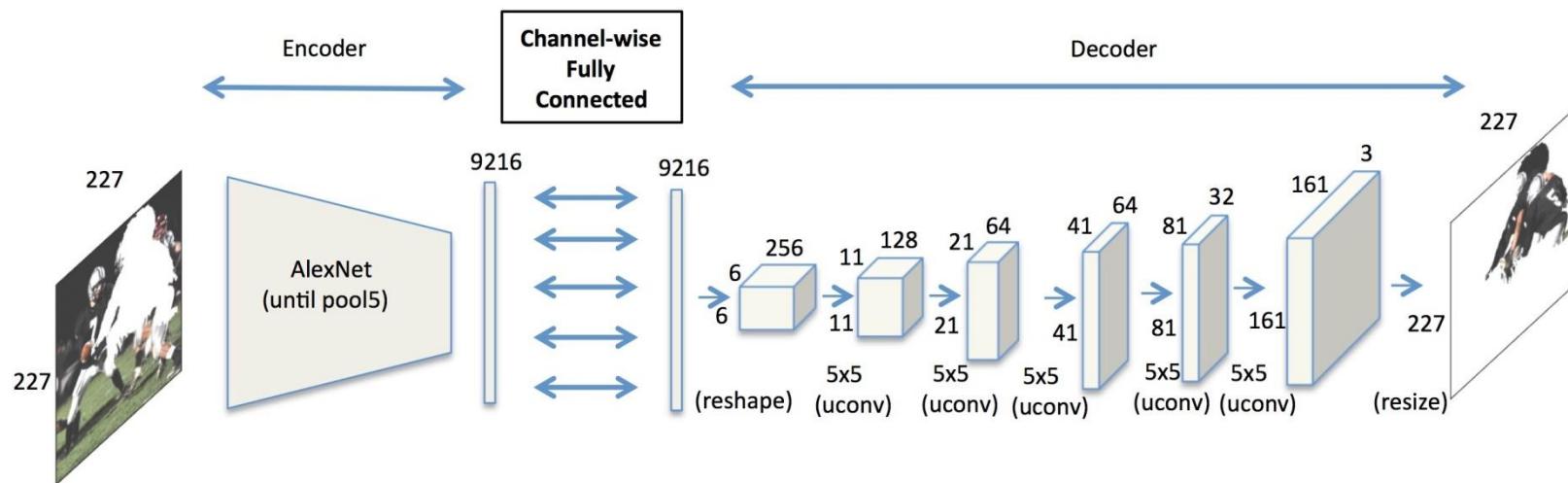
Content-Aware

# Nearest Neighbor Inpainting



# Pre-trained Feature Evaluation

- Replace encoder with VGG or AlexNet etc.



- Channel-wise FC reduces  $O(n^2)$  to  $O(k^2n)$  parameters where  $k \ll n$

# Pre-trained Feature Evaluation

Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	-	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al.</i> [1]	egomotion	10 hours	52.9%	41.8%	-
Wang <i>et al.</i> [39]	motion	1 week	58.7%	47.4%	-
Doersch <i>et al.</i> [7]	relative context	4 weeks	55.3%	46.6%	-
Ours	context	14 hours	56.5%	44.5%	30.0%

# Open Source

 [pathak22 / context-encoder](https://github.com/pathak22/context-encoder) Watch ▾ 32 Unstar 375 Fork 115

Code Issues 10 Pull requests 0 Projects 0 Wiki Insights

[CVPR 2016] Unsupervised Feature Learning by Image Inpainting using GANs <https://people.eecs.berkeley.edu/~pat...>

image-inpainting context-encoders unsupervised-learning machine-learning generative-adversarial-network deep-learning computer-vision gan  
dcgan computer-graphics

7 commits 1 branch 0 releases 1 contributor

Branch: master ▾ New pull request Create new file Upload files Find file Clone or download ▾

 pathak22 add note about paris street-view dataset Latest commit 7332006 on 17 Feb

 GitHub <https://github.com/pathak22/context-encoder>

[https://people.eecs.berkeley.edu/~pathak/context\\_encoder/slides/](https://people.eecs.berkeley.edu/~pathak/context_encoder/slides/)

*Thank you !*