# SEMANTIC IMAGE INPAINTING WITH DEEP GENERATIVE MODELS

**I L L I N O I S** — UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

Raymond A. Yeh[*]     Chen Chen[*]     Teck Yian Lim     Alexander G. Schwing     Mark Hasegawa-Johnson     Minh N. Do

Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign

* indicating equal contribution.
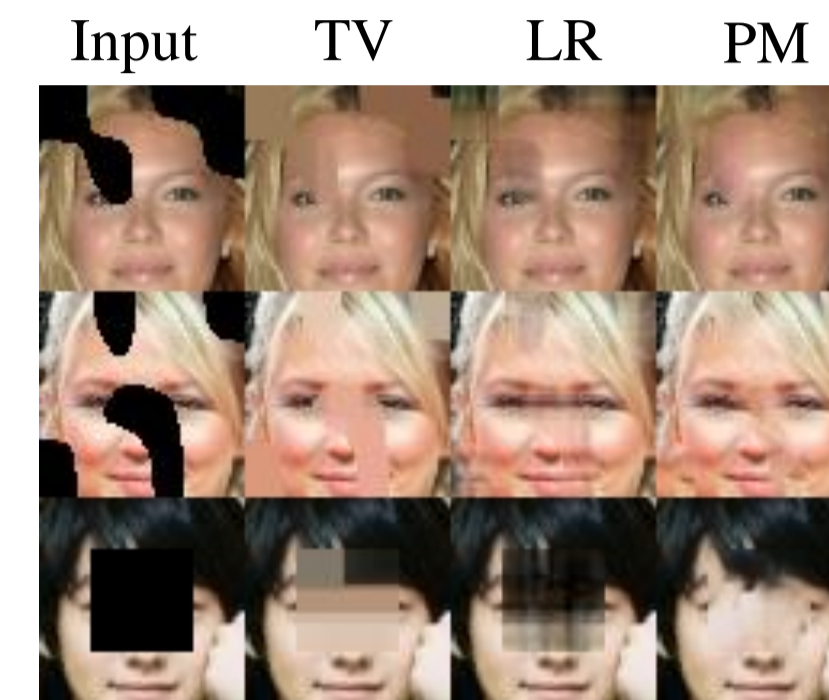
CVPR 2017 — July 21-26 HONOLULU

## 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
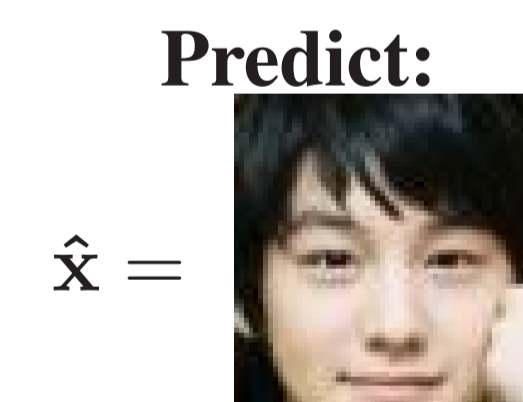- inpainting as constrained optimization problem using **context and prior loss**


Input  TV  LR  PM

## INTRODUCTION

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

**Given:**
$\mathbf{y}$     $\mathbf{M}$



**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:**
$\hat{\mathbf{x}} =$



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

**Generative Adversarial Networks:**
- Generator $G$: deep net mapping perturbation $\mathbf{z}$ to artificial sample
- Discriminator $D$: deep net discriminating between artificial and real sample, $\mathbf{x}$
- Program:

$$\min_G \max_D V(G,D) = \mathbb{E}_{\mathbf{x} \sim p_{data}}[\log(D(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{Z}}(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

## OUR APPROACH

**Intuition of our approach:**
- Hypothesis: image that is not from $p_{data}$ (*e.g.*, corrupted data) should not lie on the learned encoding manifold; use manifold can be used as a prior
- Instead of working in the pixel space, we recover the encoding $\hat{\mathbf{z}}$ "closest" to the corrupted image while constrained to the manifold

**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:

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

**Illustration of the approach:**



$Loss = L_p(z) + L_c(z| $ ,  $)$

Input  $G(\mathbf{z}^{(0)})$  $G(\mathbf{z}^{(1)})$  ...  $G(\hat{\mathbf{z}})$  Blending

Importance weight metric:
$$\mathbf{W}_i = \begin{cases} \sum_{j \in N(i)} \dfrac{(1-\mathbf{M}_j)}{|N(i)|} & \text{if } \mathbf{M}_i \neq 0 \\ 0 & \text{if } \mathbf{M}_i = 0 \end{cases}$$
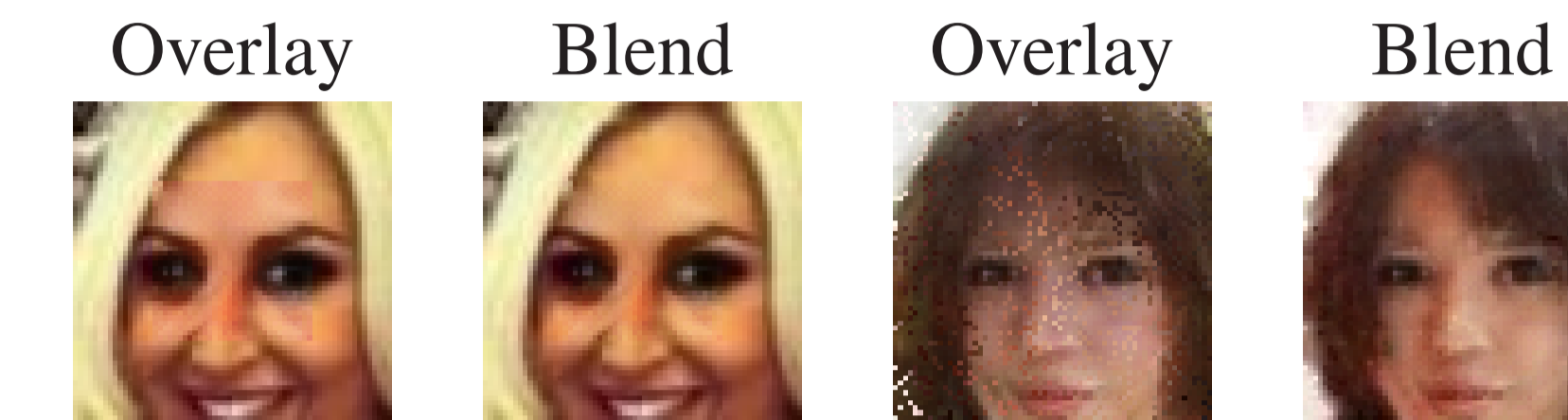
$N(i)$ defines the neighborhood of $i$

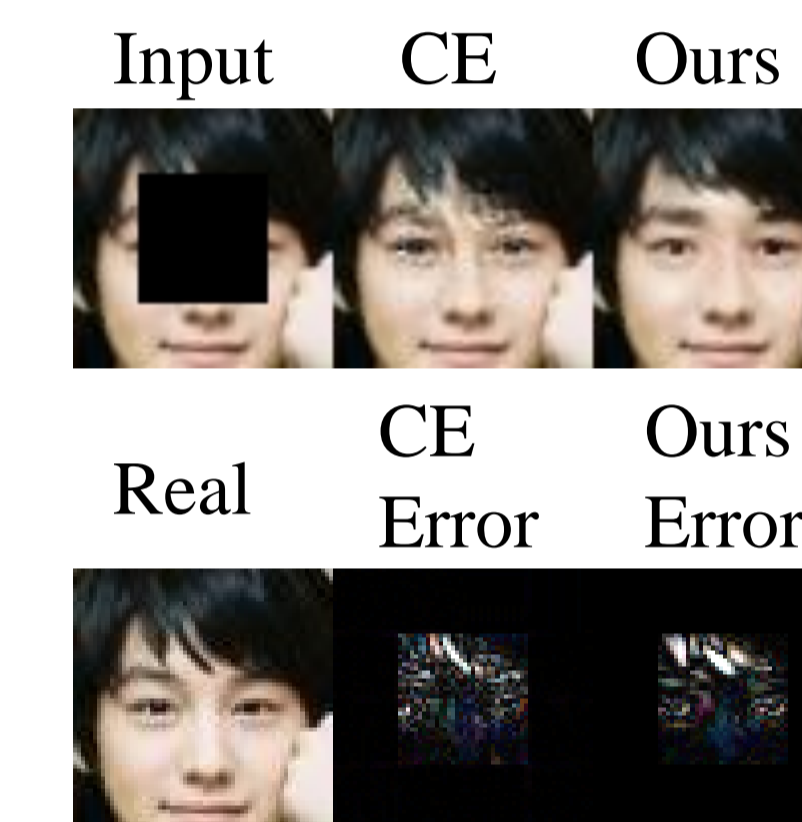**Recovering prediction $\hat{\mathbf{x}}$ via poisson blending rather than simple overlay:**

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

## RESULTS

**Comparison: Poisson Blending vs. Overlay:**


Overlay  Blend  Overlay  Blend

**Quantitative Results:**


Input  CE  Ours
Real  CE Error  Ours Error

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

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

**Qualitative Results:**


Real  Input  CE  Ours     Real  Input  CE  Ours     Real  Input  CE  Ours