

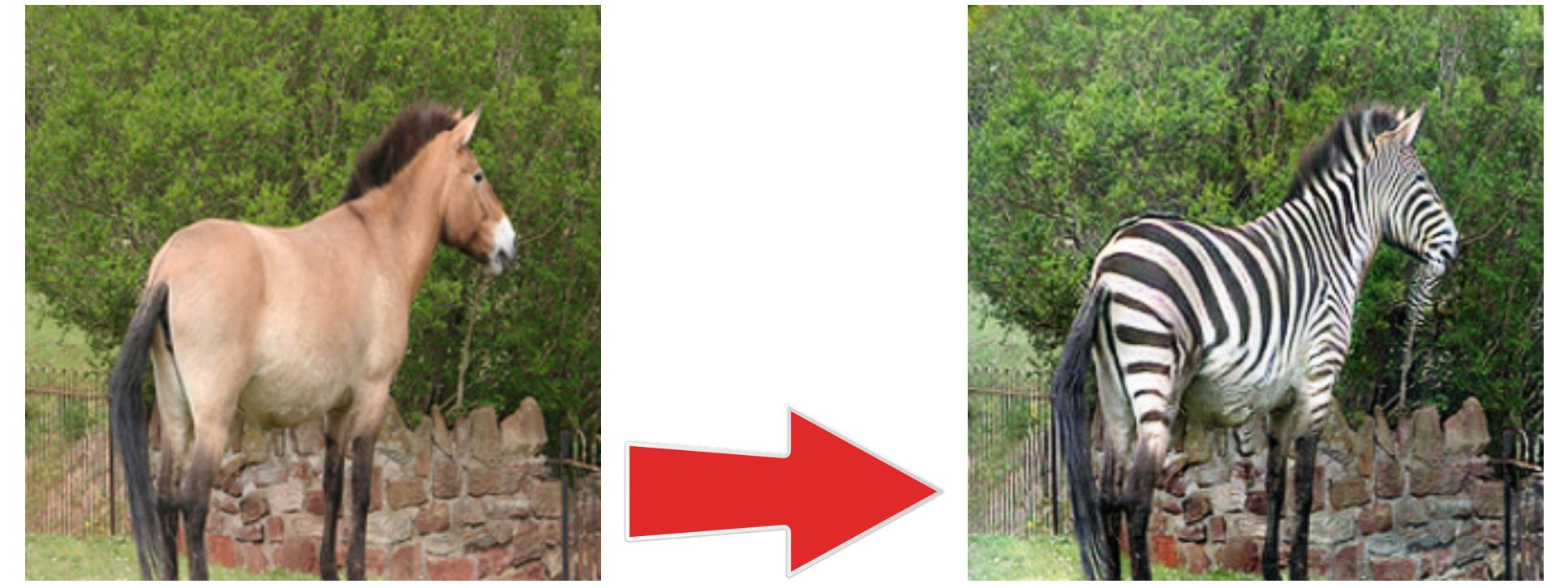
ESTHER: EXTREMELY SIMPLE IMAGE TRANSLATION THROUGH SELF-REGULARIZATION

CHAO "HARRY" YANG*, TAEHWAN KIM†, RUIZHE WANG†, HAO PENG† AND C.-C. JAY KUO*

*UNIVERSITY OF SOUTHERN CALIFORNIA

†OBEN INC.

OVERVIEW



Problem: Unpaired image translation between two domains, where the goal is to learn the mapping from an input image in the source domain to an output image in the target domain.

Challenge: An ill-posed task (no paired data)!

Method: An extremely simple yet effective image translation approach, which consists of a single generator and is trained with a self-regularization term and an adversarial term.

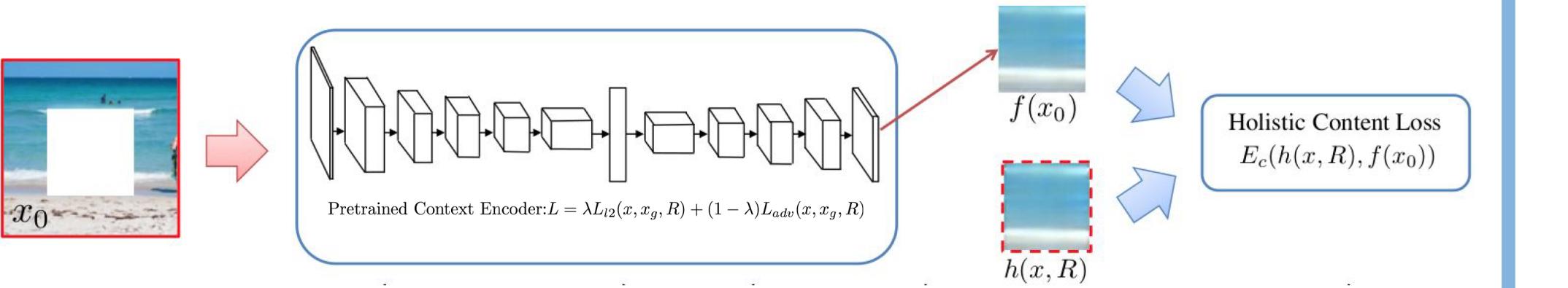
Results: Our model achieves better performance than other methods on a broad range of tasks and applications.

THE 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
- 1: Downsize x to 128×128 .
 - 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 $\tilde{x} = x^s$.
 - 5: Update \tilde{x} that minimizes the joint loss:
 $L = L_{content} + L_{texture} + L_{tv-smoothness}$.
 - 6: Compute x^{s+1} by up-sampling \tilde{x} .
 - 7: **end for**
 - 8: Return \tilde{x}^N .

THE CONTENT NETWORK

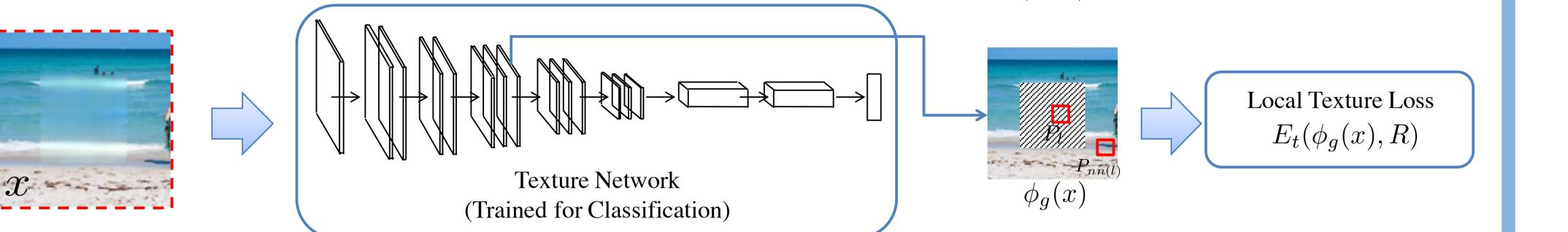


Context Encoder Predicts the Low-Res Content

The content constraint:

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

THE TEXTURE NETWORK



Pre-trained VGG Optimizes the High-Res Texture

The texture constraint:

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

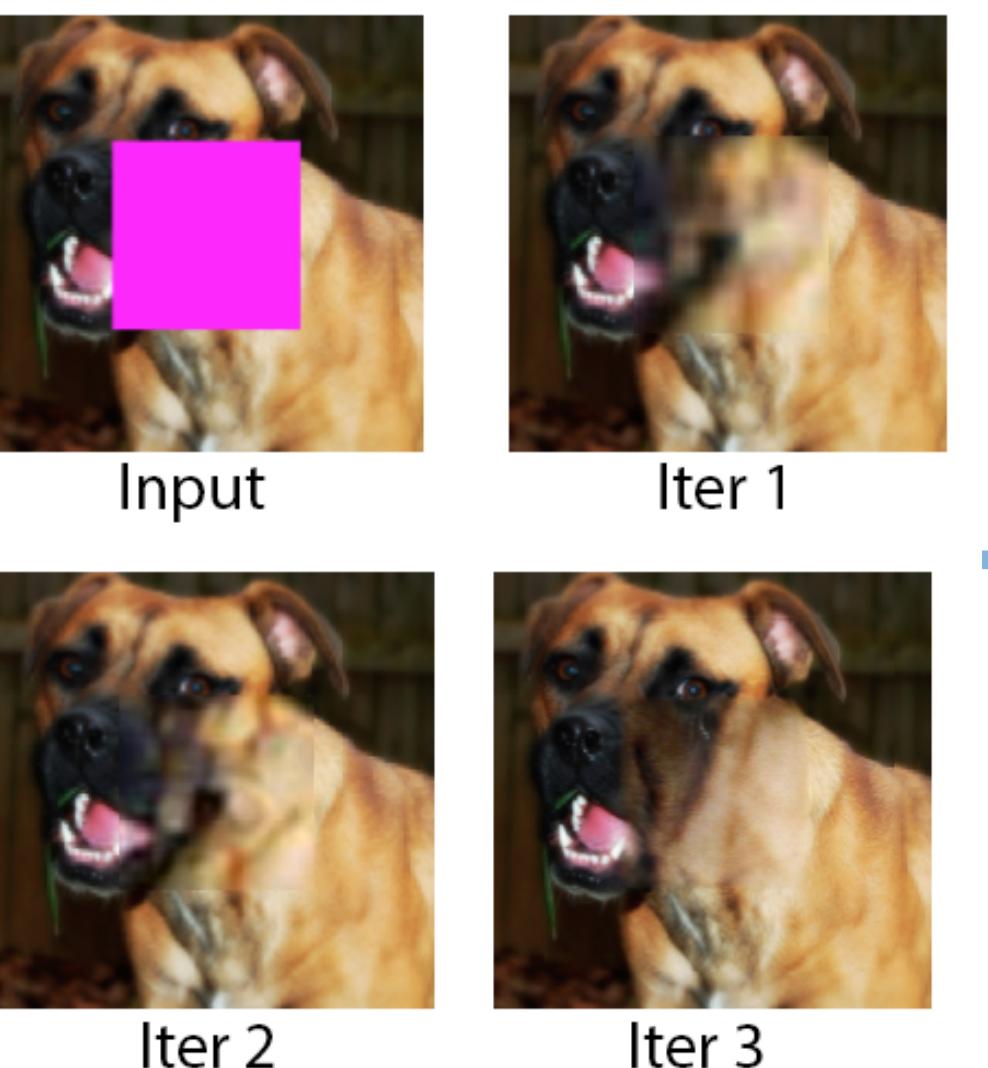
THE JOINT LOSS FUNCTION

At each iteration, we minimize:

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

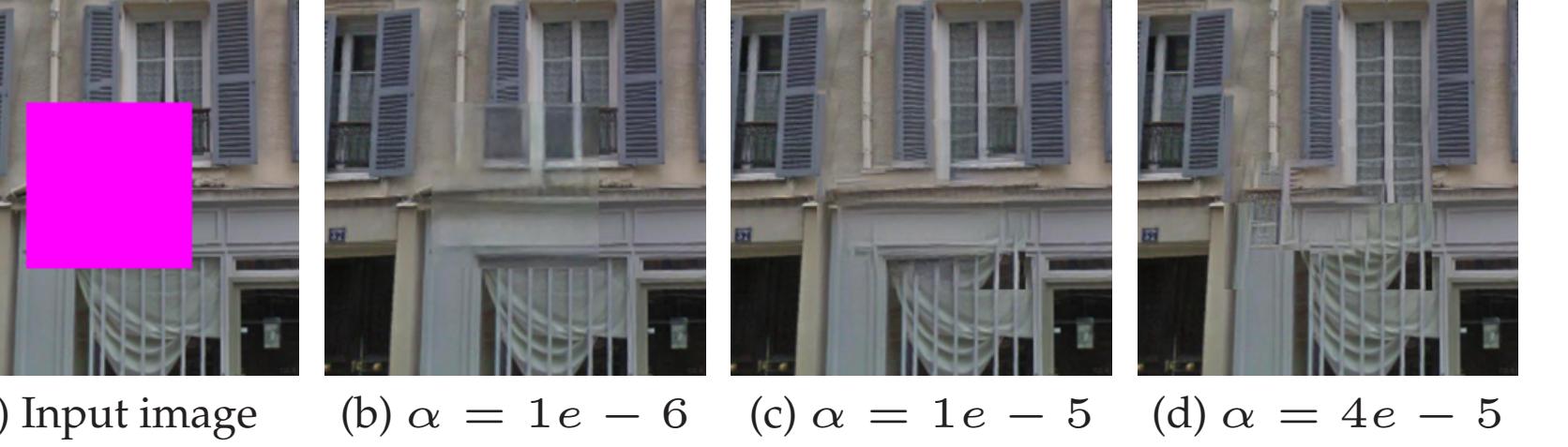
MULTI-SCALE OPTIMIZATION

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



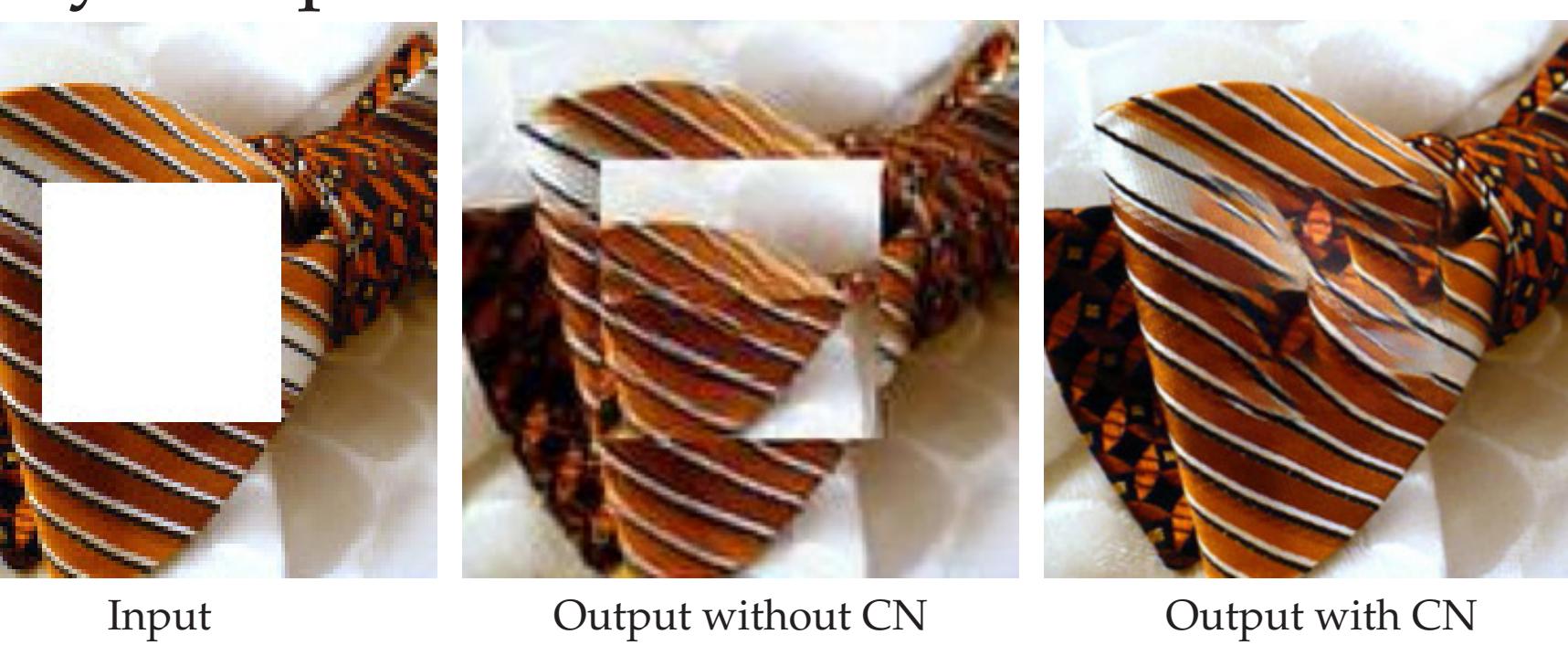
CHANGING THE TEXTURE WEIGHT α

The weight α 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:

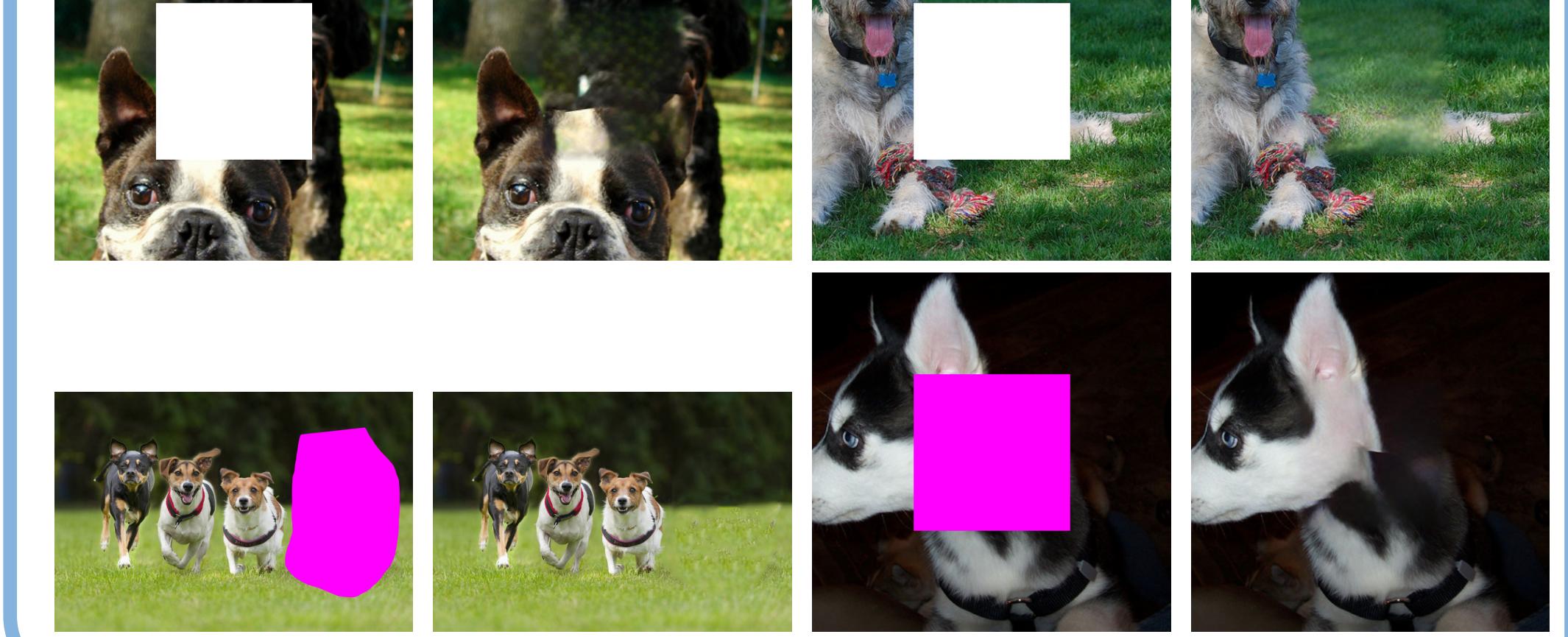
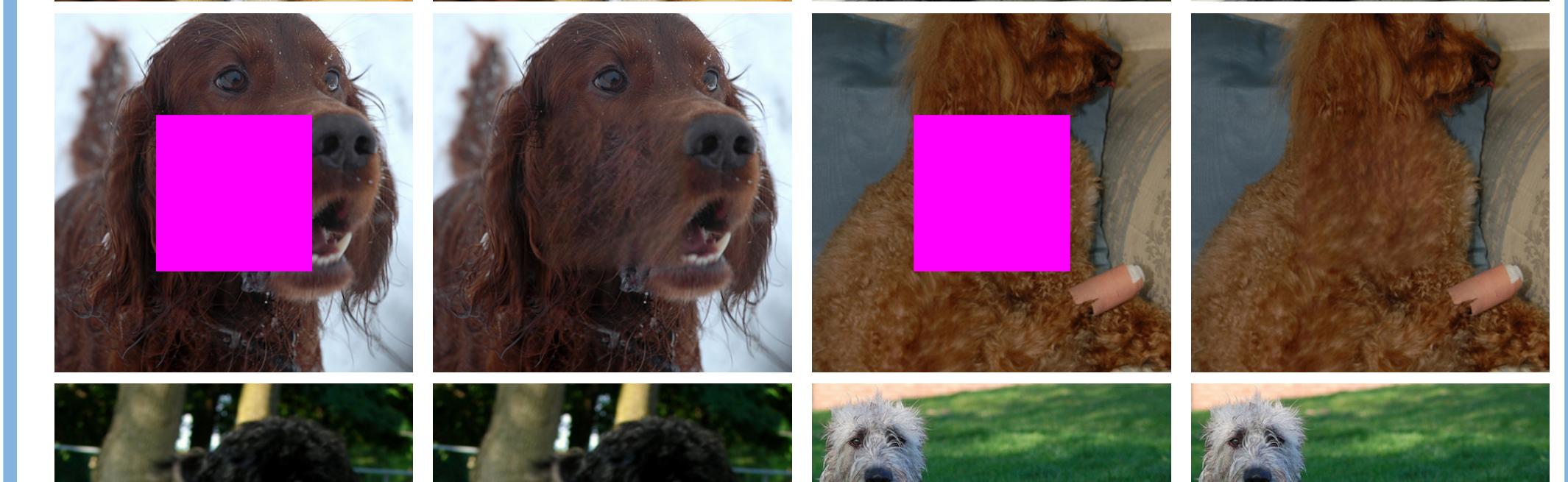
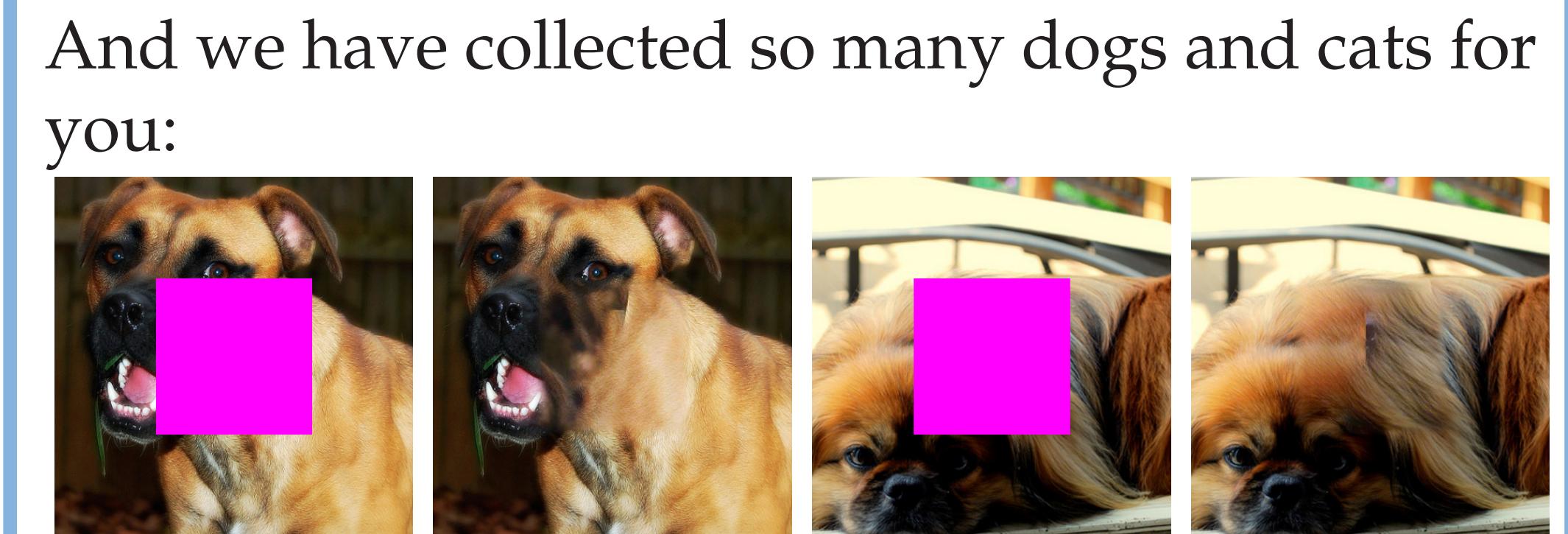


DROPPING THE CONTENT CONSTRAINT

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



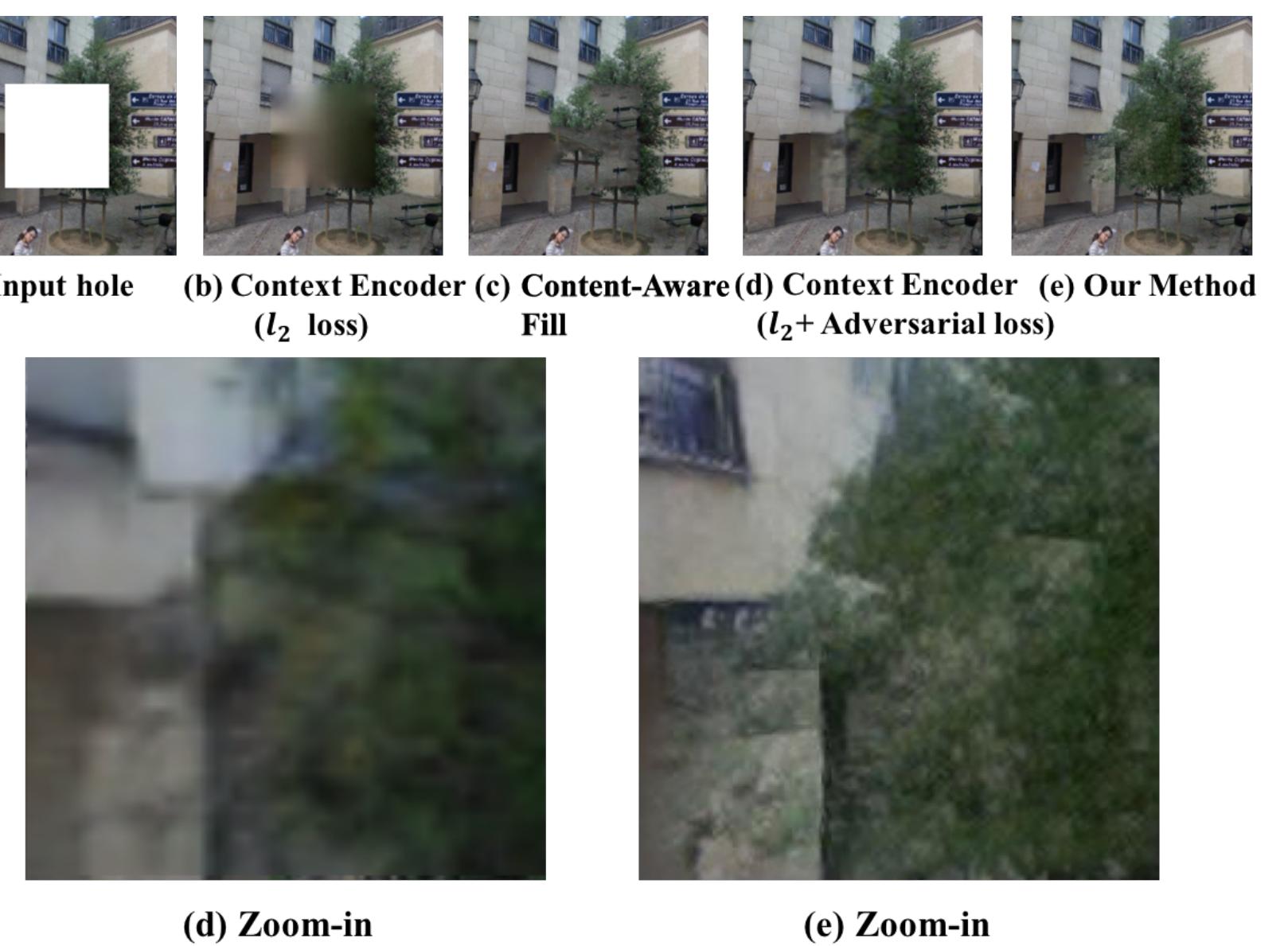
WE LOVE CATS AND DOGS



AND WHAT INPAINTS PARIS LIKE PARIS?



OVERALL COMPARISON WITH OTHER METHODS



FOR PAPER, RESULTS, CODE AND MORE

