

Adobe





Motivation

Recent Generative Compositing Methods require a mask as input, defining the region of generation. This leads to several limitations:

- Drawing an accurate mask can be non-trivial, leading to unnatural composite images.
- It limits the ability to synthesize appropriate **object effects** (long shadows, reflections, ...).
- Background areas around the object tend to be inconsistent with the original background.

We propose:

- Introduce novel task: "Unconstrained Image Compositing"
- Diffusion model for unconstrained image compositing, trained on synthesized paired data



Data Generation



- 1. Segment foreground objects and filter out those too large/ small.
- 2. Detect shadows using instance shadow detector.
- 3. Use heuristics for approximating reflection masks.
- 4. Define inpainting mask as union of object, shadow and reflection masks.
- 5. Apply GAN-based inpainting model followed by Diffusion-based inpainting model for obtaining a clean background image.

Thinking Outside the BBox: Unconstrained Generative Object Compositing

Gemma Canet Tarrés¹, Zhe Lin², Zhifei Zhang², Jianming Zhang², Yizhi Song³, Dan Ruta¹, Andrew Gilbert¹, John Collomosse^{1,2}, Soo Ye Kim²

¹ University of Surrey, ² Adobe Research, ³ Purdue University







a) Long Shadows

b) Long Reflections

Applications

c) Obj-Obj Interaction

d) Multi-Object



Comparison to SoTA Object Placement Prediction

SimOPA \uparrow **LPIPS** \uparrow Ours (w/o bbox) 0.382

 $\vec{\sigma}$ Table 2: Quantitative evaluation of predicted location and scale of our model compared to state-of-the-art object placement prediction models. LPIPS is $\times 10^{-3}$.

Comparison to SoTA Generative Object Compositing

Benefits of unconstrained compositing approach:

(i) Better background preservation (rows 3-6).

(ii) More natural object effects (i.e. shadows and reflections) beyond the bounding box (rows 3-4).

(iii) Can adjust any **misaligned bounding box** (rows 1-2).

| line (Compositing Quality) line (Identity Preservation) | |
|--|--|
| Ours | |
| | |
| Ours | |
| Ours | |
| Ours 53.1 54.0 | |
| Ours | |
| | |

| Method | DreamBooth | | | Pixabay-Comp | | | |
|--|--|---|--------------------------------|------------------------------|---------------------|--|-----------|
| | $\overline{	ext{CLIP-Score}^{\uparrow}}$ | $\mathbf{DINO}	extsf{-}\mathbf{Score}^{\uparrow}$ | $\mathbf{DreamSim} \downarrow$ | $\mathbf{FID}\!\!\downarrow$ | CLIP-Score ↑ | $\mathbf{DINO}	ext{-}\mathbf{Score}\uparrow$ | DreamSim↓ |
| $ObjectStitch^{\dagger}$ [50] | 78.018 | 85.247 | 0.342 | 70.111 | 74.964 | 77.506 | 0.488 |
| $PaintByExample^{\dagger}$ [62] | 77.782 | 79.887 | 0.438 | 82.923 | 76.604 | 75.707 | 0.515 |
| TF-ICON* [36] | 79.094 | 81.781 | 0.341 | 77.368 | 75.694 | 77.810 | 0.485 |
| Any Door^{\ddagger} [9] | 80.619 | 83.632 | 0.272 | 72.996 | 80.284 | 80.829 | 0.399 |
| ControlCom $^{\diamond}$ [68] | 74.312 | 70.497 | 0.424 | 66.071 | 72.006 | 67.476 | 0.614 |
| $\rm Ours~(w/~bbox)$ | 80.946 | 85.646 | 0.285 | 62.406 | 77.129 | 80.896 | 0.395 |
| Table 1: Quantitative comparison of composition quality and identity preservation | | | | | | | |

FID is only computed on Pixabay-Comp, which has ground truth images. [†]: Model finetuned on the same data as Ours.[‡]: Paper version, already includes diverse video and multiview data. *: Paper version, inference-based model that does not require training. \diamond : Paper version, no available training code.

| PA | | | Pixabay-Comp | | | |
|----------------------------|---------------|--|-----------------------------------|-----------|----------------|--|
| $\mathbf{oU} > 0$ | $0.5\uparrow$ | $\overline{\mathbf{mean-IoU}}\uparrow$ | $\overline{{ m IoU}>0.5}\uparrow$ | mean-IoU↑ | LPIPS ↑ | |
| 16.8 $^{\circ}_{\prime}$ | 76 | 0.094 | 48.0~% | 0.246 | 1.218 | |
| 12.2 0 | % | 0.189 | 30.2~% | 0.327 | 2.832 | |
| 11.2 $^{\circ}_{\prime}$ | % | 0.194 | 8.6~% | 0.237 | 2.072 | |
| 10.8 % | % | 0.123 | 12.2~% | 0.230 | 0.000 | |
| 31.4 | % | 0.196 | 65.4~% | 0.562 | 3.158 | |

