

Inpaint Anything: Segment Anything Meets Image Inpainting

Tao Yu¹ Runseng Feng¹ Ruoyu Feng¹ Jinming Liu² Xin Jin² Wenjun Zeng² Zhibo Chen¹

¹University of Science and Technology of China ²Eastern Institute for Advanced Study

{yutao666, fengruns, ustcfry}@mail.ustc.edu.cn

{jmliu, jinxin, wenjunzeng}@eias.ac.cn, chenzhibo@ustc.edu.cn

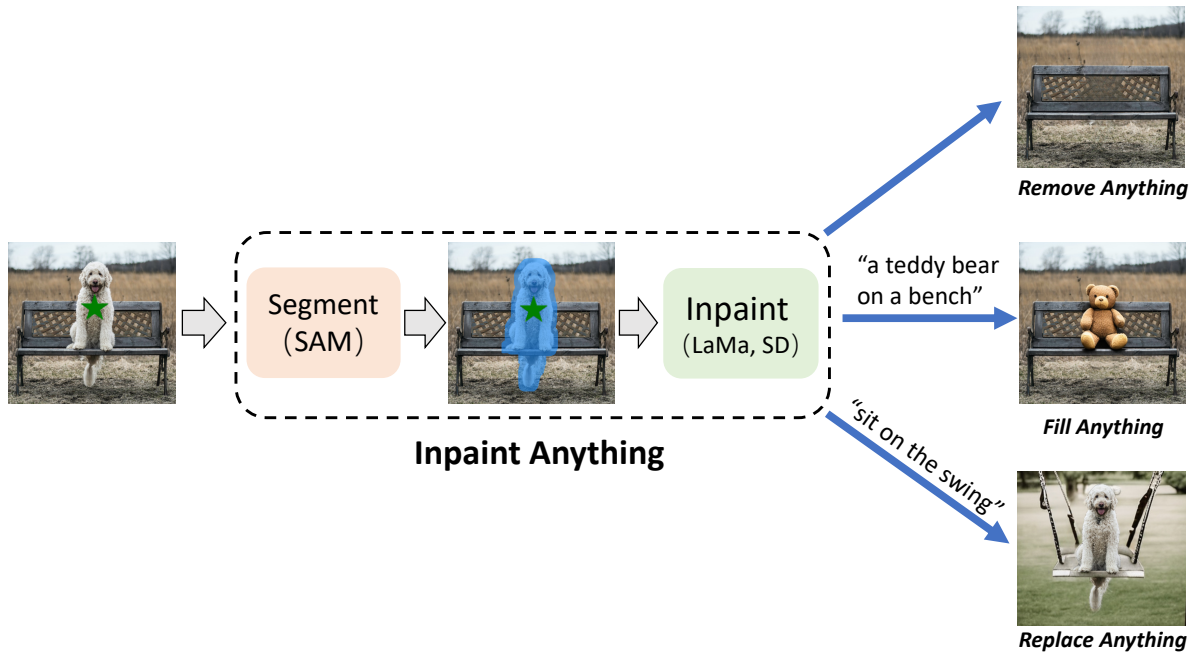


Figure 1: Illustration of our Inpaint Anything. Users can select any object in an image by clicking on it. With powerful vision models, e.g., SAM[7], LaMa [13] and Stable Diffusion (SD) [11], Inpaint Anything is able to remove the object smoothly (i.e., Remove Anything). Further, by inputting text prompts, users can fill the object with any desired content (i.e., Fill Anything) or replace the background of it arbitrarily (i.e., Replace Anything).

Abstract

Modern image inpainting systems, despite the significant progress, often struggle with mask selection and holes filling. Based on Segment-Anything Model (SAM) [7], we make the first attempt to the mask-free image inpainting and propose a new paradigm of “clicking and filling”, which is named as Inpaint Anything (IA). The core idea behind IA is to combine the strengths of different models in order to build a very powerful and user-friendly pipeline for solving inpainting-related problems. IA supports three main features: (i) *Remove Anything*: users could click on an object and IA will remove it and smooth the “hole” with

the context; (ii) *Fill Anything*: after certain objects removal, users could provide text-based prompts to IA, and then it will fill the hole with the corresponding generative content via driving AIGC models like Stable Diffusion [11]; (iii) *Replace Anything*: with IA, users have another option to retain the click-selected object and replace the remaining background with the newly generated scenes. We are also very willing to help everyone share and promote new projects based on our Inpaint Anything (IA). Our codes are available at <https://github.com/geekyutao/Inpaint-Anything>.

1. Motivation and Observation

1.1. Why do we need Inpaint Anything?

- The state-of-the-art (SOTA) image inpainting works, like LaMa [13], Repaint [10], MAT [8], ZITS [4], *etc.*, have achieved great progress. They can successfully inpaint large regions and work well with complex repetitive structures, generalizing well to high-resolution images. However, they typically need fine annotations for each mask, which are essential for training and inference.
- Segment Anything Model (SAM) [7] is a strong segmentation foundation model, producing high quality object masks from input prompts such as points or boxes, and it can be used to generate comprehensive and accurate masks for all objects in an image. However, their mask segmentation predictions have not been fully-explored.
- Besides, the existing inpainting methods can only fill the removed area with the context. AIGC models open up new opportunities for creation, which has the potential to meet massive demand and assists humans to newly generate their wanted content.
- Therefore, by combining the advantages of SAM [7], the SOTA image inpainters [13], and AI generated content (AIGC) models [11], we provide a powerful and user-friendly pipeline for solving more general inpainting-related problems, such as object removal, new content filling, and background replacing.

1.2. What Inpaint Anything can do?

- **SAM + SOTA inpainters for removing anything:** With IA, users can easily remove specific objects from the interface by simply clicking on them. Furthermore, IA provides an option for users to fill the resulting “hole” with contextual data. Oriented at this, we combine the strengths of SAM and some SOTA Inpainters like LaMa. Once manually refined through corrosion and dilation, the mask predictions generated by SAM serve as input for the inpainting models, providing clear indicators for the object areas to be erased and filled.
- **SAM + AIGC models for filling or replacing anything:**
 - (1) After removing objects, IA provides users the option to fill the resulting “hole” either with contextual data or “new content”. Specifically, a strong AI generated content (AIGC) model like Stable Diffusion [11] is utilized to generate new objects via text prompts.

For example, users can use the word of “dog” or a sentence of “a cute dog, sitting on the bench”, to generate a new dog for filling the hole with such newly generated dog.

- (2) In addition, users have another option to take IA to retain the click-selected object and replace the remaining background with the newly generated scene. This scene replacement process of IA supports various ways of prompting AIGC models, such as using a different image as visual prompt or using a short caption as text prompt. For example, users can keep the dog in an image but replace the original indoor background with an outdoor one.

2. Methodology

2.1. Preliminary

Segment Anything Model (SAM). The fundamental CV model of Segment Anything [7] was released last week, which is a large ViT-based model trained on the large visual corpus (SA-1B). SAM has demonstrated promising segmentation capabilities in various scenarios and the great potential of the foundation models for computer vision. This is a ground-breaking step toward visual artificial general intelligence, and SAM was once hailed as “the CV version of ChatGPT”.

SOTA Inpainters. Image inpainting, as an ill-posed inverse problem, is widely explored in the field of computer vision and image processing, which intended to replace missing regions of damaged images with visually plausible structure and texture. The success of deep learning has brought new opportunities [13, 10, 8, 4], and all these SOTA methods can be categorized from multiple perspectives, e.g., inpainting strategies, network structures, and loss functions. For our Inpaint Anything (IA), we investigated the use of a simple, single-stage approach LaMa [13] for mask-based inpainting, which is arguably good in generating repetitive visual structures by combining fast Fourier convolutions (FFCs) [1], perceptual loss [6], and an aggressive training mask generation strategy.

AIGC Models. ChatGPT ¹ and other Generative AI (GAI) techniques all belong to the category of Artificial Intelligence Generated Content (AIGC), which involves the creation of digital content, such as images, music, and natural language, through AI models. It is considered a new type of content creation and has shown to achieve state-of-the-art performance in various content generation [11, 12]. For our work of IA, we directly employ a powerful AIGC model of Stable Diffusion [11] to generate the desired content in the hole based on text-prompting.

¹<https://chat.openai.com/>

2.2. Inpaint Anything

The principle of our proposed Inpaint Anything (IA) is to composite off the shelf foundation models to enable the ability of solving extensive image inpainting problems. By compositing the strengths of various foundation models, IA can generate high-quality inpainted images. Specifically, our IA has three schemes, *i.e.*, *Remove Anything*, *Fill Anything* and *Replace Anything*, which are designed to *remove*, *fill* and *replace* anything, respectively.

Remove Anything. Remove Anything focuses on the object removal problem [2, 3, 5] by allowing users to eliminate any object from an image while ensuring that the resulting image remains visually plausible. Remove Anything consists of three steps: clicking, segmenting, and removing, as shown in Figure 1. In the first step, users select the object they want to remove from an image by clicking on it. Next, a foundation segmentation model, such as Segment Anything [7], is utilized to automatically segment the object based on the click location and create a mask. Finally, a state-of-the-art inpainting model, such as LaMa [13], is used to fill the hole created by the removed object using the mask. Since the object is no longer present in the image, the inpainting model fills the hole with background information. Note that, for the entire process, users only need to click on the object they want to remove from the image.

Fill Anything. Fill Anything allows users to fill any object in an image with any content they want. The tool consists of four steps: clicking, segmenting, text-prompting, and generating. The first two steps of Fill Anything are the same as in Remove Anything. In the third step, users input a text prompt that indicates what they want to fill the object hole with. Finally, a powerful AIGC model, such as Stable Diffusion [11], is adopted to generate the desired content in the hole based on the text-prompt inpainting model.

Replace Anything. Replace Anything is able to replace any object with any background. The process of Replace Anything is similar to that of Fill Anything, but in this case, the AIGC model is prompted to generate visually consistent background that exists outside the specified object.

Practice. Compositing foundation models to solve tasks may encounter problems such as incompatibility or inappropriateness. We should consider intermediate processing for better coordination between models and tasks. In this work, for image inpainting scenario, we summarize several good practice for the composition below.

- **Dilation matters.** We observe that the segmentation result (*i.e.*, object mask) of SAM may contain discontinuous and non-smooth boundaries, or holes inside the

object region. These issues pose challenges for effectively removing or filling objects. Consequently, we employ a dilation operation to refine the mask. Further, for filling objects, large masks give AIGC models more space to create, benefiting the “alignment” to user purpose. Thus, we adopt a large dilation in Fill Anything.

- **Fidelity matters.** Most state-of-the-art AIGC models such as Stable Diffusion require images to be of a fixed resolution, typically 512×512 . Simply resizing images to this resolution may result in a loss of fidelity, which can adversely impact the final inpainted results. Therefore, it is essential to adopt measures that preserve the original image quality, such as utilizing cropping techniques or maintaining the image’s aspect ratio when resizing.
- **Prompt matters.** Our research indicates that text prompts exert a significant influence on AIGC models. However, we observe that simple prompts, such as “a teddy bear on a bench” or “a Picasso painting on the wall”, typically produce satisfactory results in the scenario of text-prompt inpainting. In contrast, longer and more complex prompts may yield impressive outcomes, but they tend to be less user-friendly.

3. Experiment

We evaluate Remove Anything, Fill Anything and Replace Anything in our Inpaint Anything in three cases, *i.e.*, removing objects, filling objects and replacing background, respectively. We collect test images from COCO dataset [9], LaMa test set [13] and photos taken by our phones. The results are in Figure 2, 3 and 4. The experimental results indicate that the proposed Inpaint Anything is both general and robust, effectively inpainting images with diverse content, resolutions, and aspect ratios.

4. Conclusion

Inpaint Anything (IA) is a versatile tool that combines the capabilities of *Remove Anything*, *Fill Anything*, and *Replace Anything*. Based on the vision foundation models of Segmentation Anything, SOTA inpainter and AIGC models, IA enables to realize mask-free image inpainting, and also supports the user-friendly operation of “click for removing, and prompt for filling”. Besides, IA can handle more various and high-quality input images, with any aspect ratio and 2K resolution. We build such interesting project to demonstrate a strong power of fully exploiting the existing LARGE AI models, and reveal the potential of “Composable AI”. We are also very willing to help everyone share and promote new projects based on our Inpaint Anything (IA). In the future, we will further develop our

Inpaint Anything (IA) to support more practical functions, like fine-grained image matting, editing, etc., and apply it to more realistic applications.

References

- [1] Lu Chi, Borui Jiang, and Yadong Mu. Fast fourier convolution. *Advances in Neural Information Processing Systems*, 33:4479–4488, 2020. [2](#)
- [2] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama. Object removal by exemplar-based inpainting. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, volume 2, pages II–II. IEEE, 2003. [3](#)
- [3] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on Image Processing*, 13(9):1200–1212, 2004. [3](#)
- [4] Qiaole Dong, Chenjie Cao, and Yanwei Fu. Incremental transformer structure enhanced image inpainting with masking positional encoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11358–11368, 2022. [2](#)
- [5] Omar Elharrouss, Noor Almaadeed, Somaya Al-Maadeed, and Younes Akbari. Image inpainting: A review. *Neural Processing Letters*, 51:2007–2028, 2020. [3](#)
- [6] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 694–711. Springer, 2016. [2](#)
- [7] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. [1](#), [2](#), [3](#)
- [8] Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. Mat: Mask-aware transformer for large hole image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10758–10768, 2022. [2](#)
- [9] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. [3](#)
- [10] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. [2](#)
- [11] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. [1](#), [2](#), [3](#)
- [12] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. [2](#)
- [13] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2149–2159, 2022. [1](#), [2](#), [3](#)

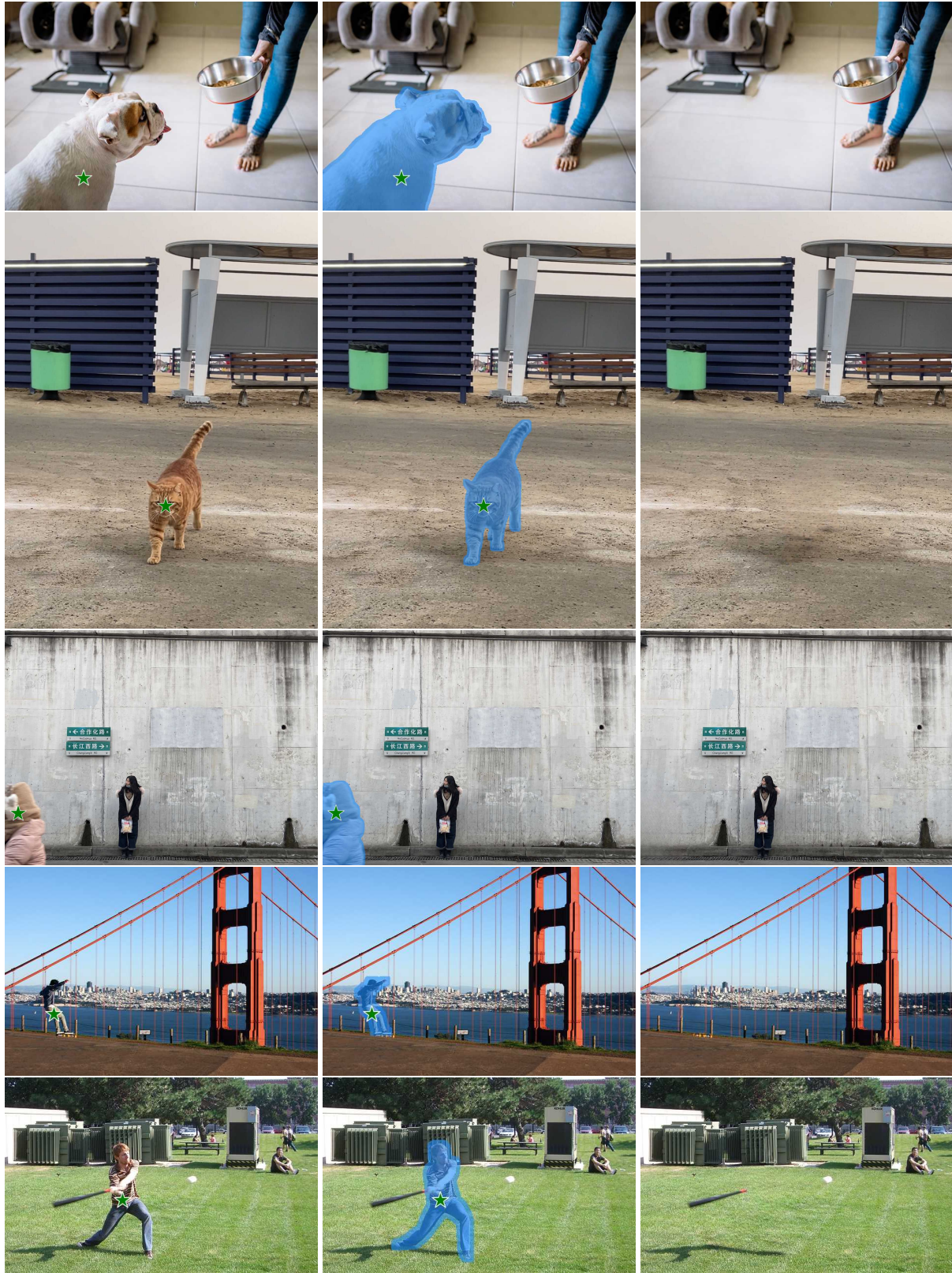


Figure 2: Visualization results of Remove Anything.



(a) Text prompt: a teddy bear on a bench



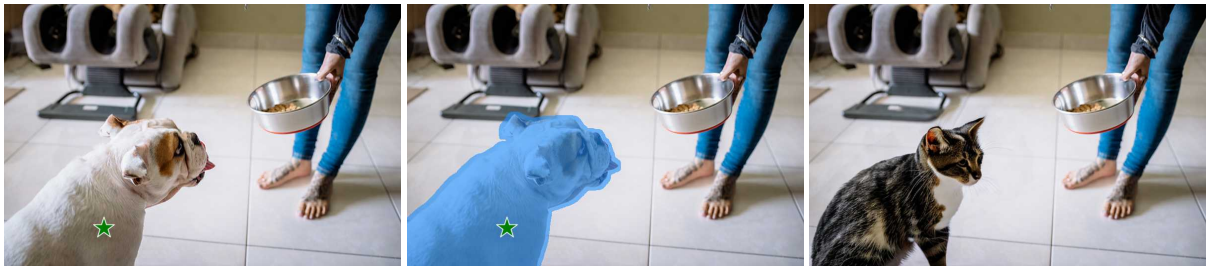
(b) Text prompt: a camera lens in the hand



(c) Text prompt: an aircraft carrier on the sea



(d) Text prompt: a sports car on a road



(e) Text prompt: a cat, waiting for food

Figure 3: Visualization results of Fill Anything.



(a) Text prompt: crossroad in the city



(b) Text prompt: breakfast



(c) Text prompt: a bus, on the center of a country road, summer evening



(d) Text prompt: a man in office

Figure 4: Visualization results of Replace Anything.