

Week 2: Thursday Paper 1: [Segment Anything](#)

explained by [Segment Anything Model \(SAM\) - The Complete 2024 Guide - viso.ai](#), [Introducing Segment Anything: Working toward the first foundation model for image segmentation](#)

Supplemental: [How to Use the Segment Anything Model \(SAM\)](#)

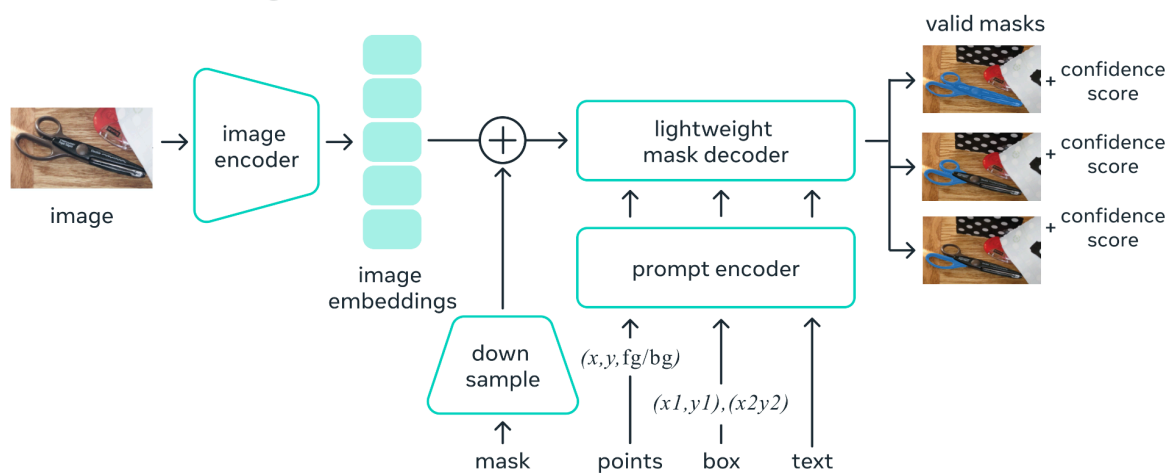
1. What are the different inputs and outputs in “SAM”? How is the output of the model evaluated? Be precise.

Ans:

Inputs: SAM takes segmentation prompts such as points, boxes, masks, or text along with images.

Outputs: It outputs segmentation masks corresponding to the prompts.

Universal segmentation model



2. What are the different types of layers used? Give an intuitive description of what role each of them plays.

Ans:

Image Encoder: This is used to extract embeddings from images

Prompt Encoder: This embeds prompts using positional encodings and text encoders.

Lightweight Mask Decoder: This combines image embeddings and prompt embeddings to predict segmentation masks.

3. What are the architectural differences between U-Net and SAM? Which is better for medical image segmentation tasks and why?

Ans: U-Net is a classic segmentation model with a symmetric encoder-decoder structure.

SAM uses separate encoders for images and prompts and a decoder for mask generation. This setup allows SAM to handle flexible prompts and achieve real-time performance.

SAM's flexibility and prompt-based approach might offer advantages in customized segmentation tasks, but U-Net has a proven track record in medical imaging. So, it would really depend on the specific case.

4. How are ambiguous prompts handled? What alternative approaches would you suggest?

Ans: SAM addresses ambiguous prompts by predicting multiple valid masks for a single prompt, allowing it to handle scenarios where a prompt could refer to multiple objects. I feel this is the best and intuitive way to handle ambiguous prompts and I am not able to think of alternative approaches.

5. What do you think SAM's limitations are? How would you address these limitations?

Ans: SAM might struggle with highly ambiguous prompts and its performance might be limited by the quality and diversity of training data and enhancing the diversity of training data, refining the model's ability to interpret prompts, and integrating feedback mechanisms for continuous learning could be potential ways to address these limitations.

6. What questions do you have about the paper?

Week 2: Thursday Paper 2: [ConvNets Match Vision Transformers at Scale](#)

explained by [ConvNets Match Vision Transformers at Scale — Paper Summary](#) | by Anuj Dutt | Medium

Supplemental: [An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

1. What are the central arguments presented in the paper?

Ans: The paper argues against the belief that Vision Transformers outperform ConvNets when trained with large-scale datasets. It demonstrates that ConvNets, specifically NFNet, can match or exceed the performance of ViTs when both are trained on a similarly extensive dataset and with comparable computational resources.

2. Why (or when) are transformers better than traditional ConvNets for vision?

Ans: Transformers are considered better than traditional ConvNets mainly due to their ability to capture long-range dependencies in the data, which is particularly beneficial for large-scale datasets where these relationships are more complex and nuanced.

3. What different evaluation metrics are used in the paper? Do the metrics used do a good job of supporting the arguments presented? Why or why not?

Ans: The paper uses metrics like Top-1 accuracy on ImageNet and validation loss on JFT-4B for evaluation. These metrics effectively support the arguments since they are standard benchmarks for assessing model performance in computer vision tasks, allowing for a fair comparison between ConvNets and Vision Transformers.

4. What alternative approaches could be used to support the arguments?

Ans: Alternative approaches could involve comparative studies with diverse datasets or experiments focusing on different aspects of model performance like speed, efficiency, or ability to generalize, could further support the arguments.

5. What is your key takeaway after reading this paper?

Ans: My key takeaway is that the superiority of Vision Transformers over ConvNets is not absolute. When we have similar training conditions, ConvNets can perform on par with or even surpass Vision Transformers, challenging the prevailing trend of favoring transformer-based models in computer vision tasks.