



COMPRESSING HEAVY VISION TRANSFORMERS FOR EDGE DEVICES

MSc. Computer Vision

July, 2024

Presented by:

Atanda Abdullahi Adewale

Supervised by:

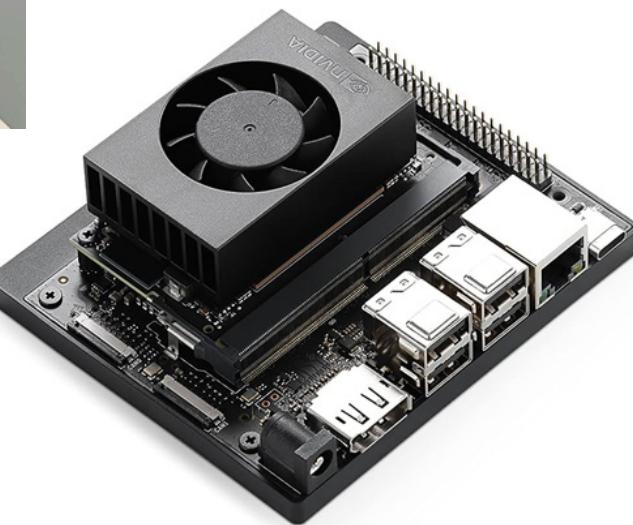
Walid Brahim

CONTENTS



- **Introduction**
 - Wasoria
 - Problem Statement
 - Why Vision Transformers
 - MHSA
 - Segment Anything Model (SAM)
- **Compression Techniques**
 - Model Pruning and Quantization
 - Knowledge Distillation
- **Methods**
 - Dataset Design and Preprocessing
 - Teacher Model
 - MobileSAM (TinyViT)
- **Student Models**
 - ResNet-18
 - EfficientNet-b0
 - MobileViT-s
- **Model Design**
 - Huber Loss
 - Predictor
 - Training Configuration
- **Results**
 - Performance Metrics (mIOU, Params, Flops)
 - Visualizer and Demo on the Edge
- **Conclusion**
- **References**

WASORIA, FRANCE



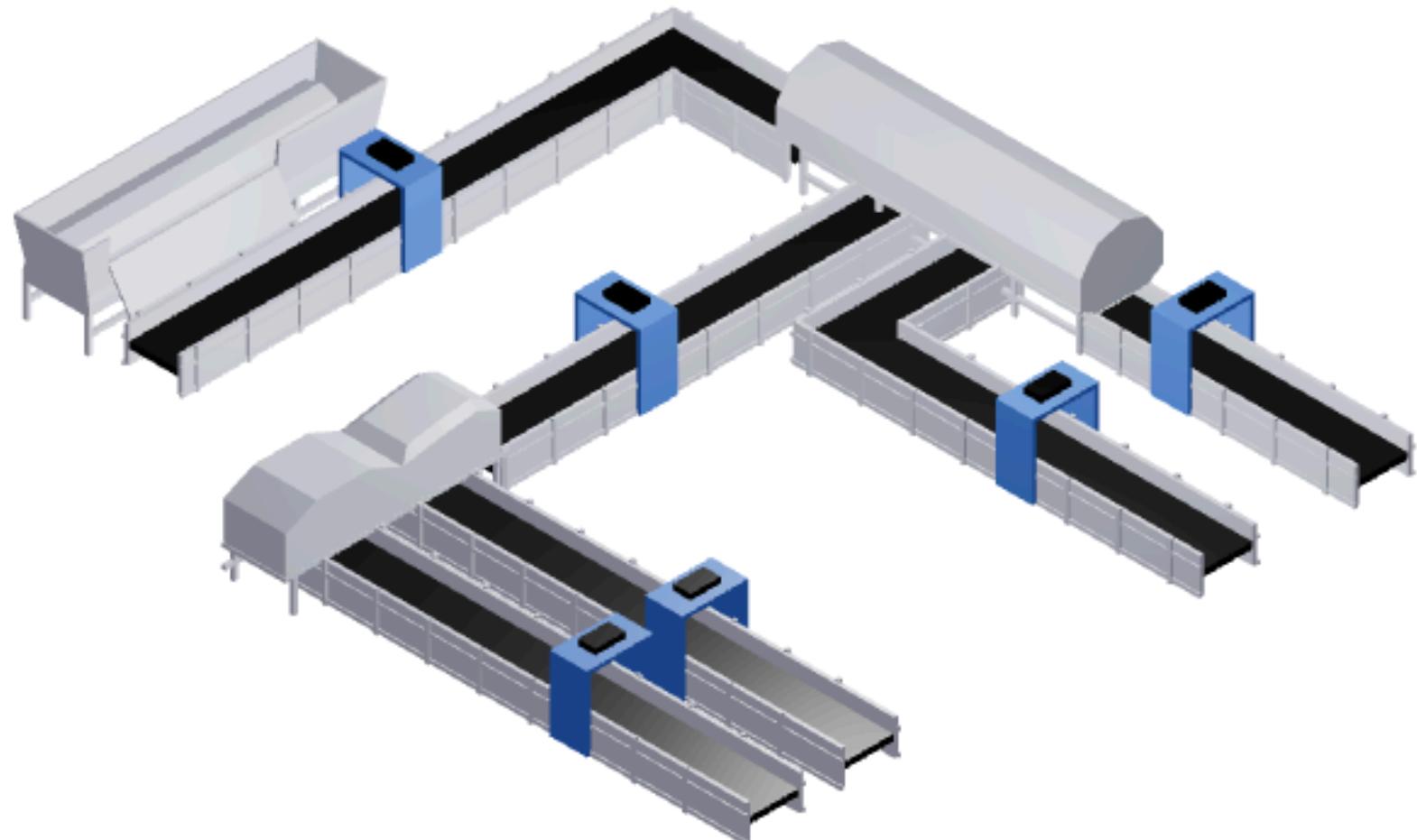
Overview

- Located in Le Creusot, France.
- Waste management and recycling through AI integration.
- Collaborates with local research labs for continuous innovation

Innovative Solutions

- Develops AI Box (CARAC) with controlled light intensity.
- Mounted on conveyor belts for real-time waste sorting
- Instant alerts for sorting issues and quality data gathering.

WASORIA, FRANCE



Overview

- Located in Le Creusot, France.
- Waste management and recycling through AI integration.
- Collaborates with local research labs for continuous innovation

Innovative Solutions

- Develops CARAC AI Box with controlled light intensity.
- Mounted on conveyor belts for real-time waste sorting
- Instant alerts for sorting issues and quality data gathering.

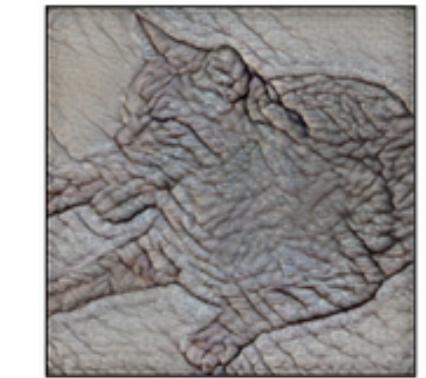
Problem



(a) Texture image
81.4% **Indian elephant**
10.3% indri
8.2% black swan



(b) Content image
71.1% **tabby cat**
17.3% grey fox
3.3% Siamese cat



(c) Texture-shape cue conflict
63.9% **Indian elephant**
26.4% indri
9.6% black swan

Problem:

- Current methods struggle with distinguishing objects of similar textures.
- Convolutional Neural Networks (CNNs) exhibit texture bias, leading to misclassification.
- Vision Transformers (ViTs) are too heavy for edge devices in low-latency sorting.

Research Objectives:

- Develop ViT/CNN detectors for accurate waste detection.
- Integrate encoder into compressed SAM.
- Distill knowledge from SAM to smaller model.
- Optimize model for edge deployment.

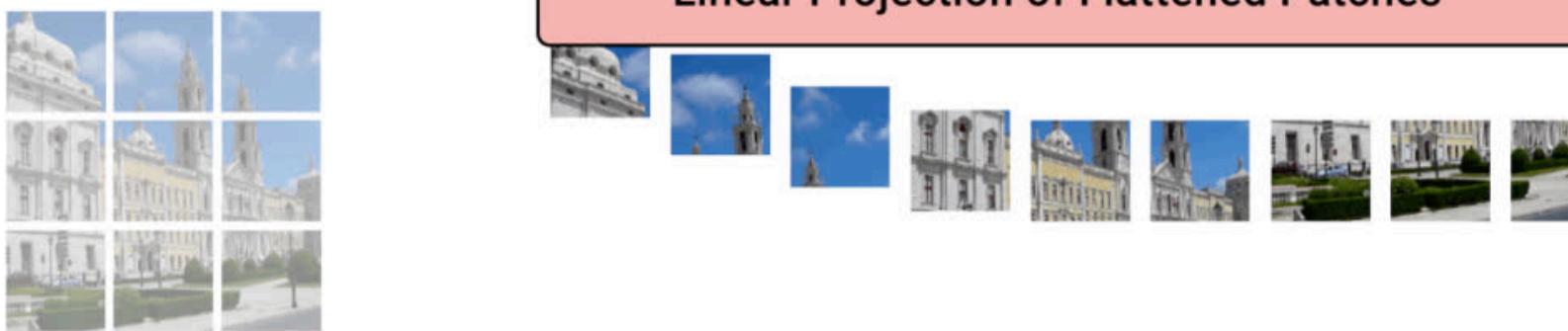
Vision Transformers (ViT)



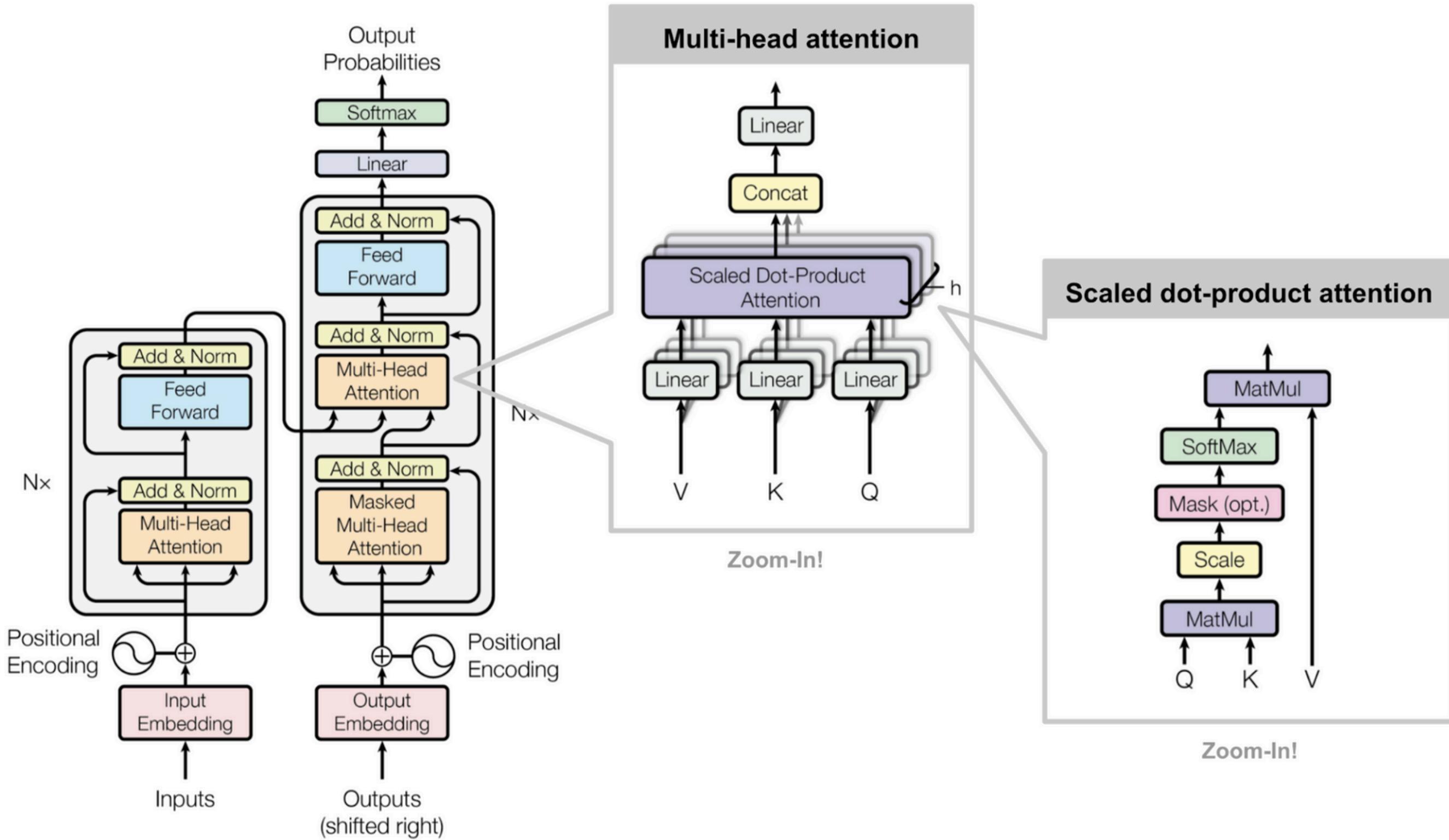
Vision Transformer ViT Architecture

- Split an image into fixed-size patches.
- Flatten the image patches.
- Create lower-dimensional linear embeddings from the patches.
- Include positional embeddings.
- Feed the sequence as input to a state-of-the-art (SOTA) transformer encoder.
- Pre-train the ViT model with image labels on a large dataset.
- Fine-tune on a downstream dataset for image classification

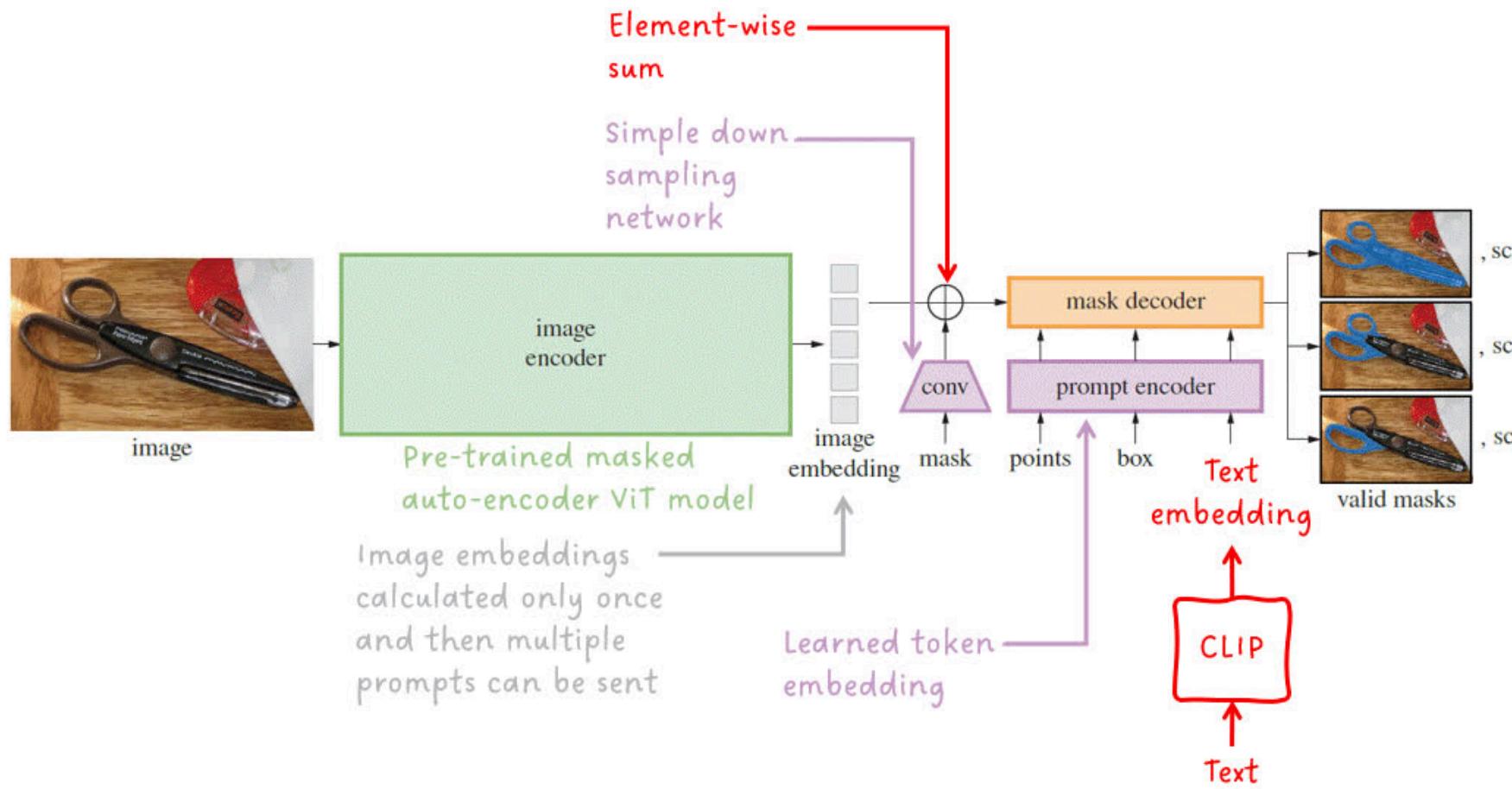
Linear Projection of Flattened Patches



Multi-head Self Attention



Segment Anything Model (SAM)

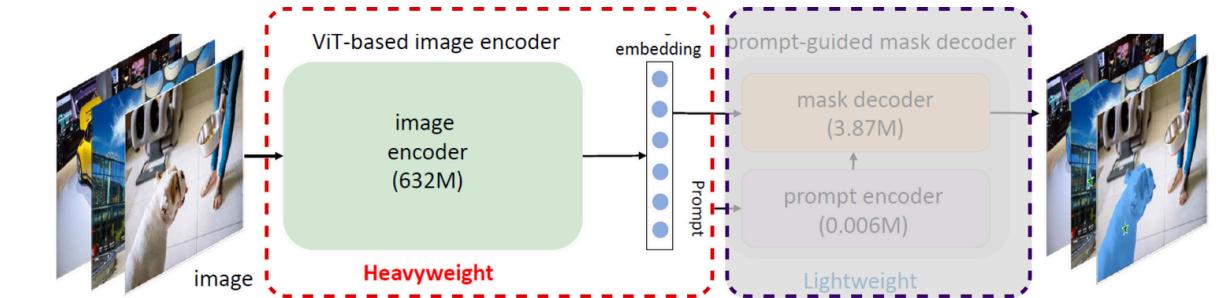


Vision foundation model for semantic image segmentation

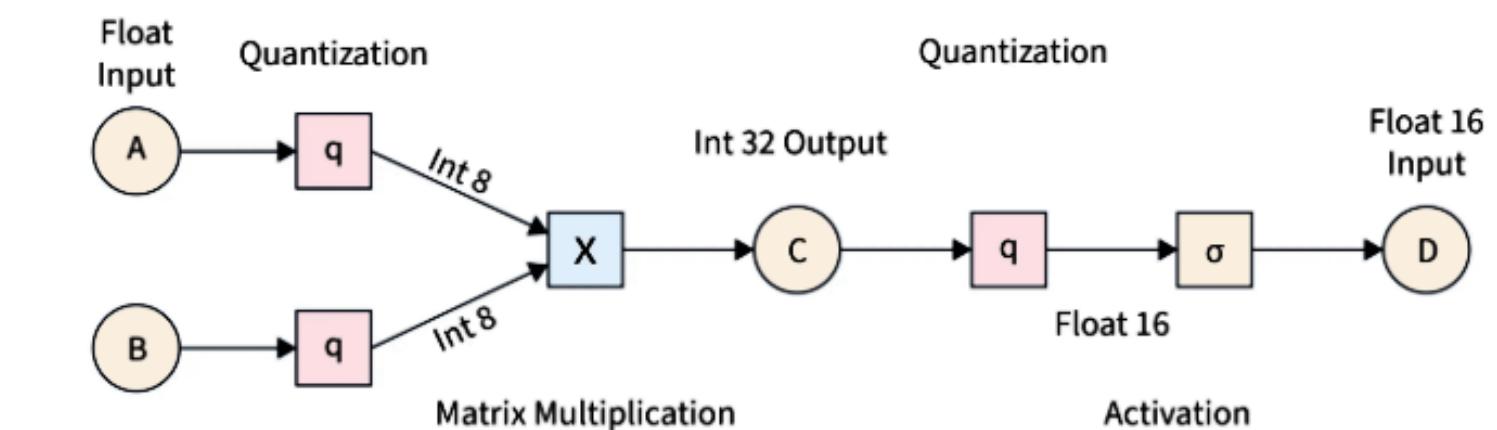
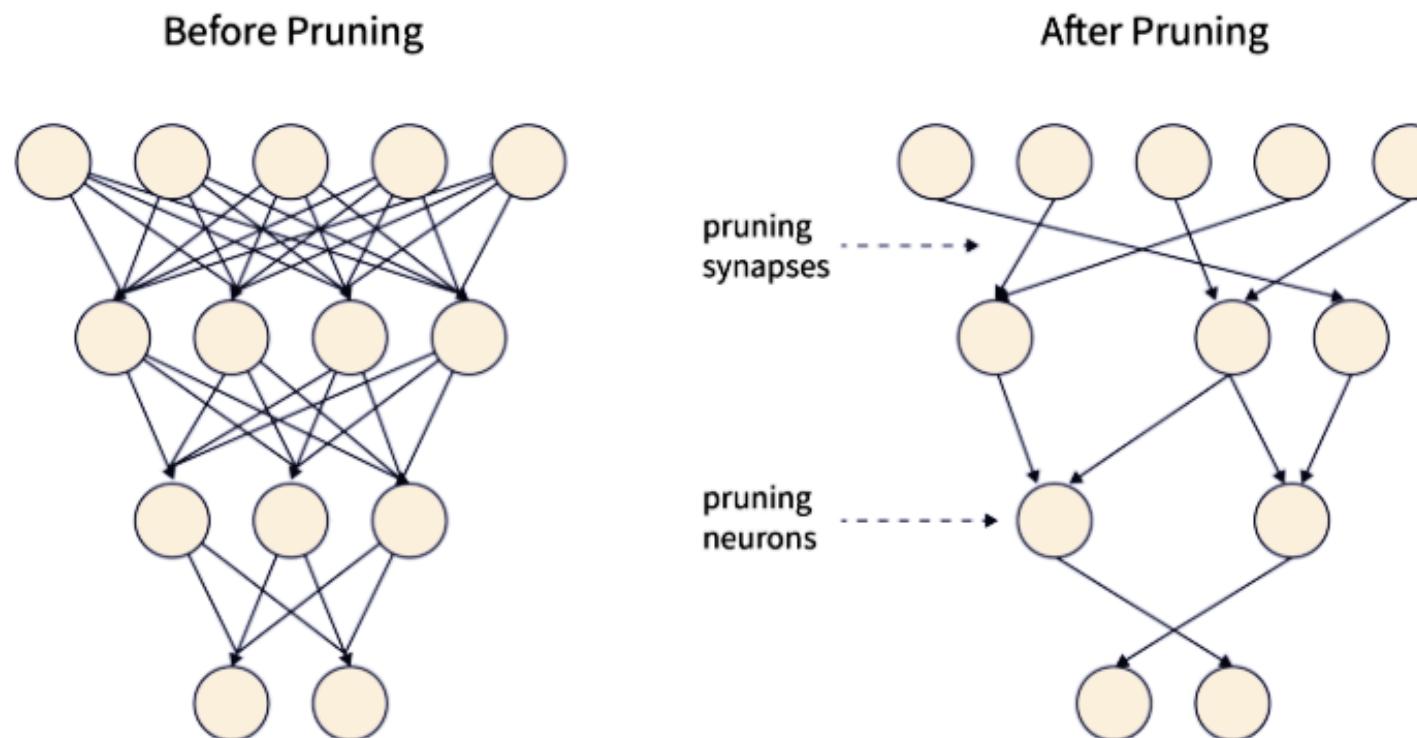
SA1B dataset : 1.1 billion masks from approximately 11 million images

SAM knows a lot about a lot

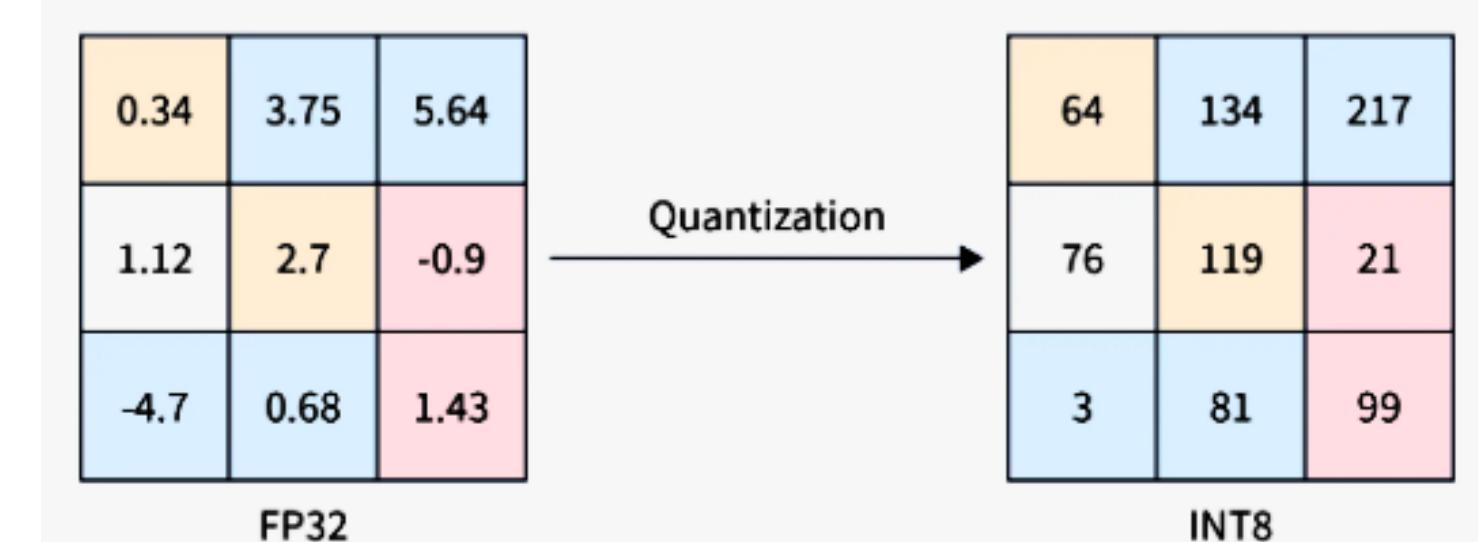
Decoupled approach:
We want to know a lot about a little



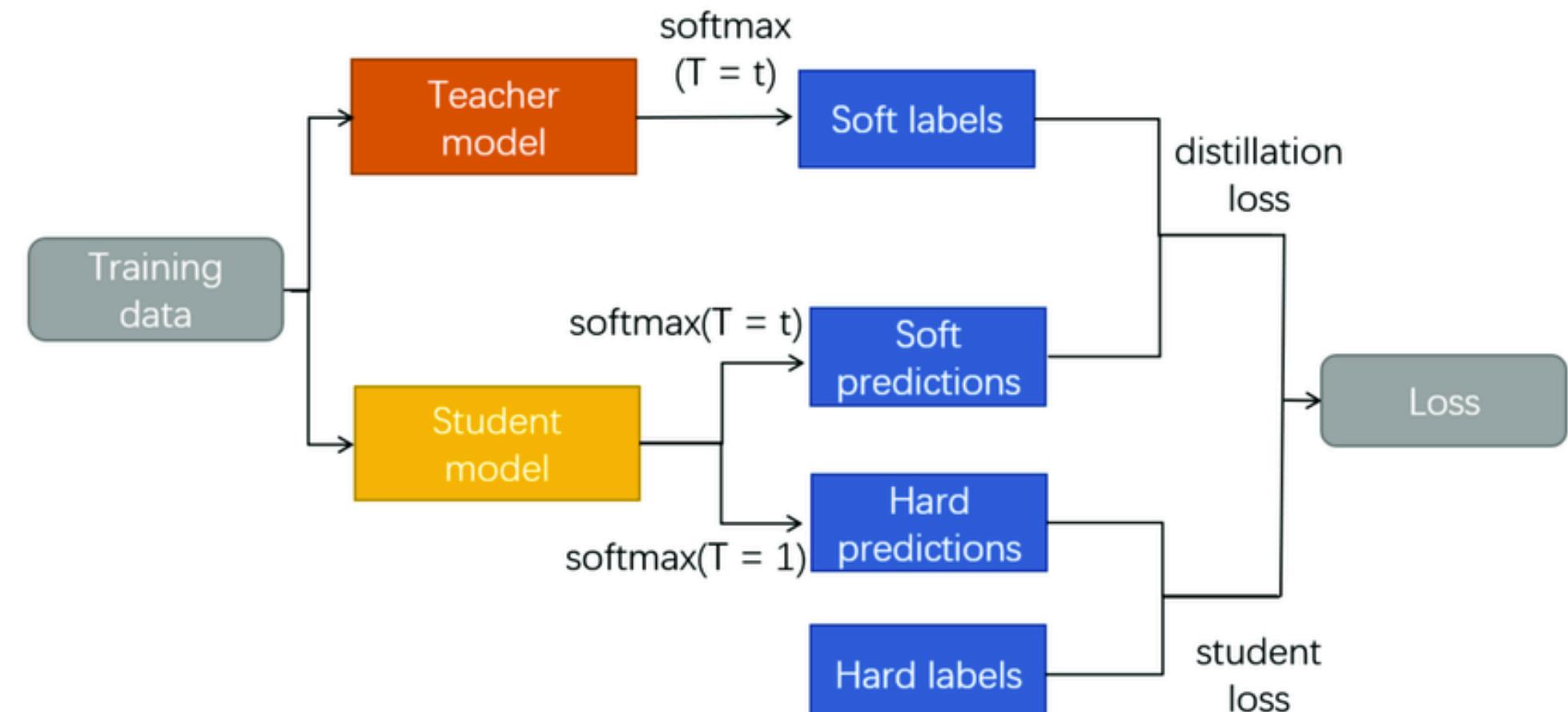
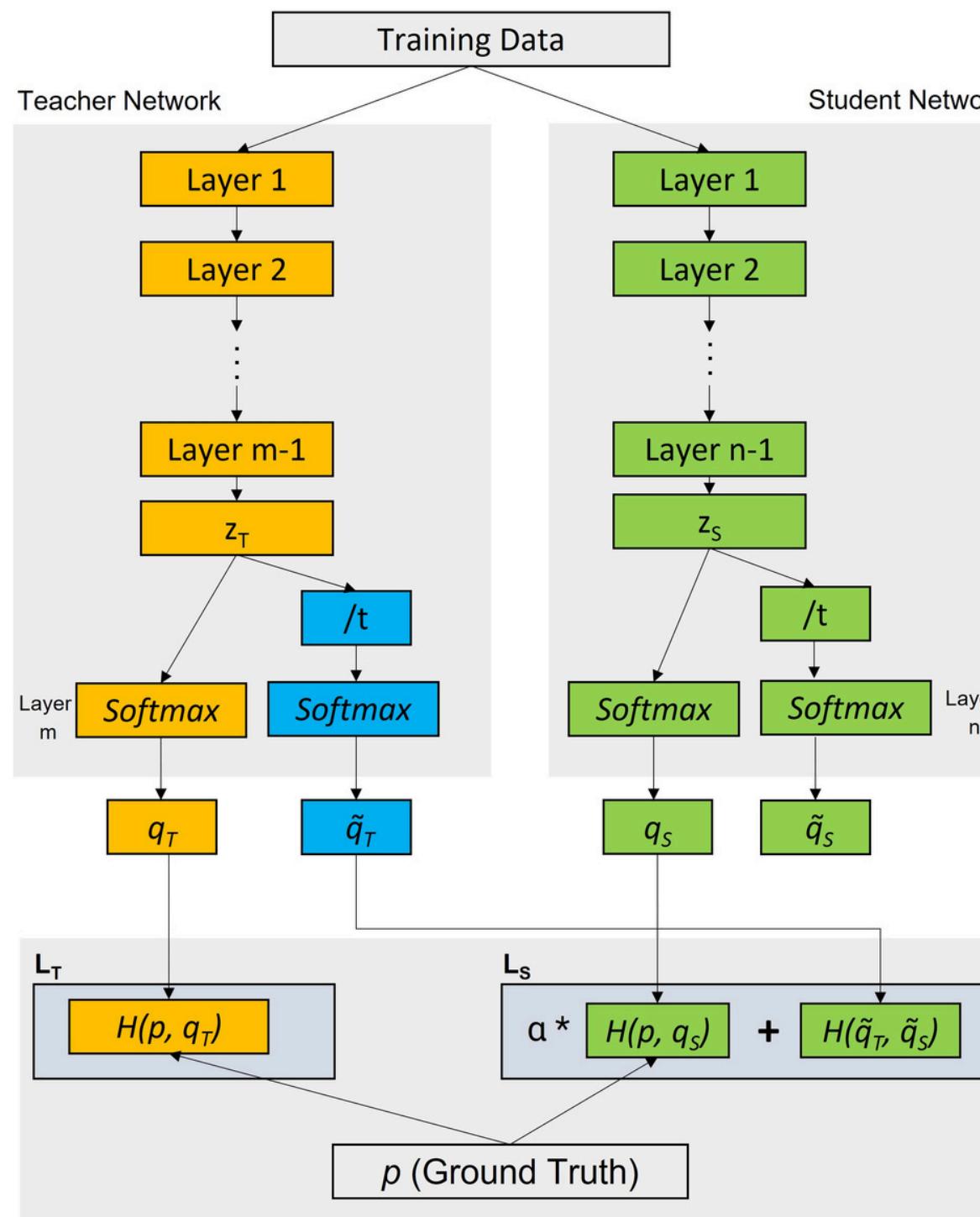
Model Pruning and Quantization



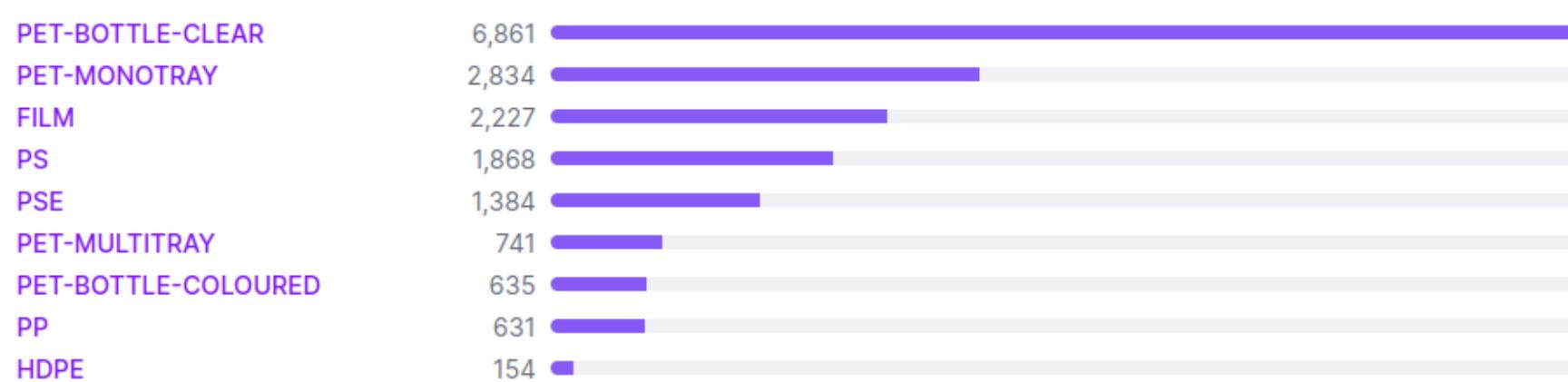
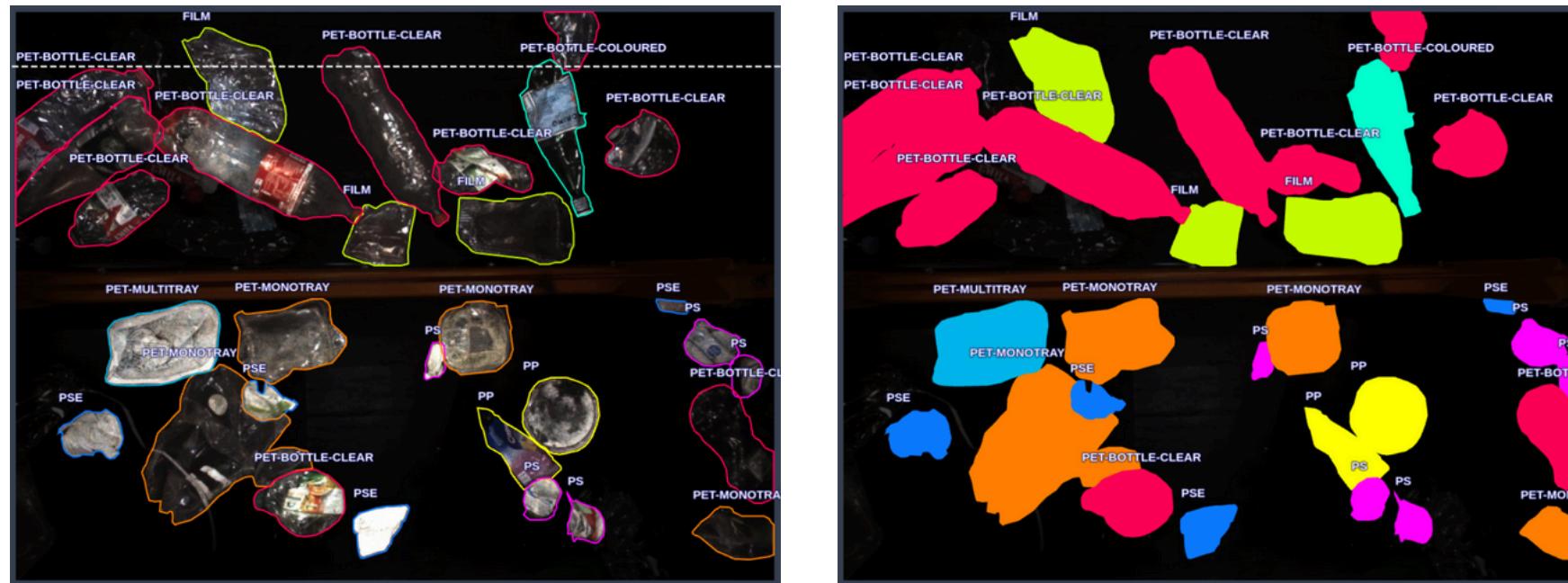
Method	Description	Benefits	Trade-offs
Post-Training Quantization	Lower precision after training	Simple, immediate reduction	Accuracy drop, esp. small models
Quantization-Aware Training	Train with quantization	Better accuracy	Complex, longer training
Unstructured Pruning	Removes individual weights	Large size reduction	Irregular memory access
Structured Pruning	Removes neurons, filters, heads	Hardware-friendly	Higher accuracy loss



Knowledge Distillation



WASORIA 5K Dataset



- Obtained from a waste sorting center in France; designed for instance segmentation tasks.
- COCO Format: (image metadata, categories, bounding boxes, segmentations).
- 5080 images with 18,810 annotated instances across eight waste material classes.
- Data Split: Training set: 4523 images; validation set: 558 images.
- Annotated with Roboflow and Annolid, with detailed polygon annotations.
- Images auto-oriented and resized to 640x640 pixels.
- Augmentations: horizontal flipping, cropping, and rotation.

MobileSAM

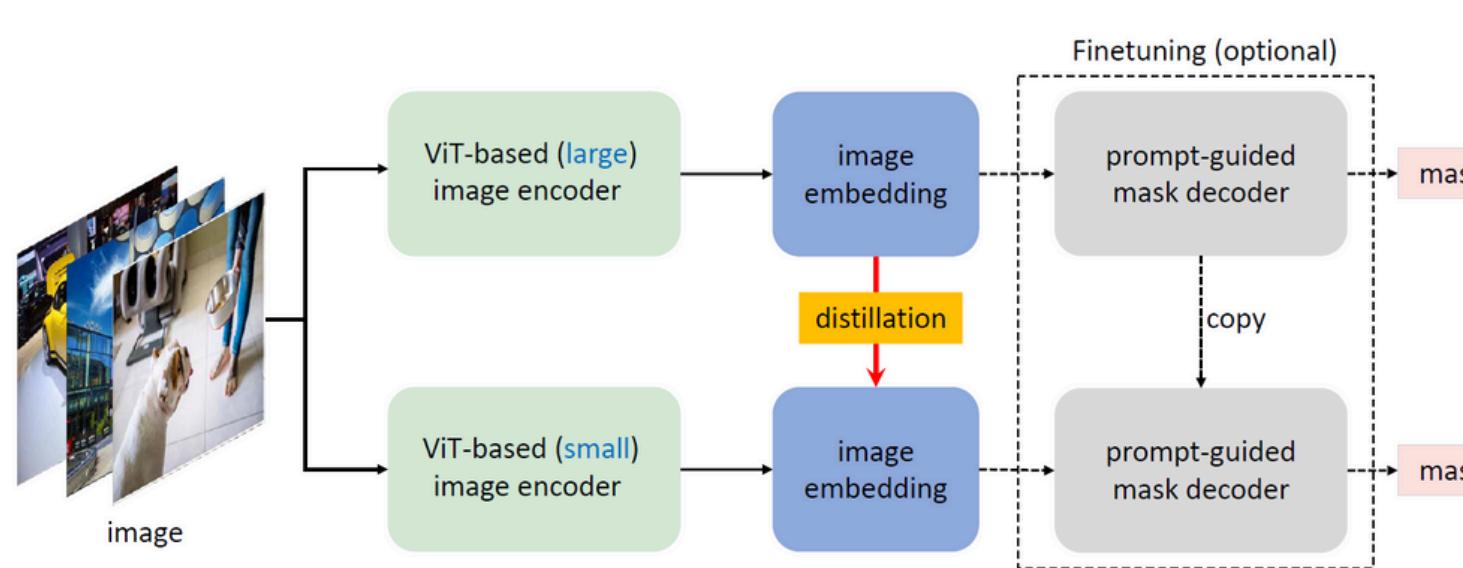
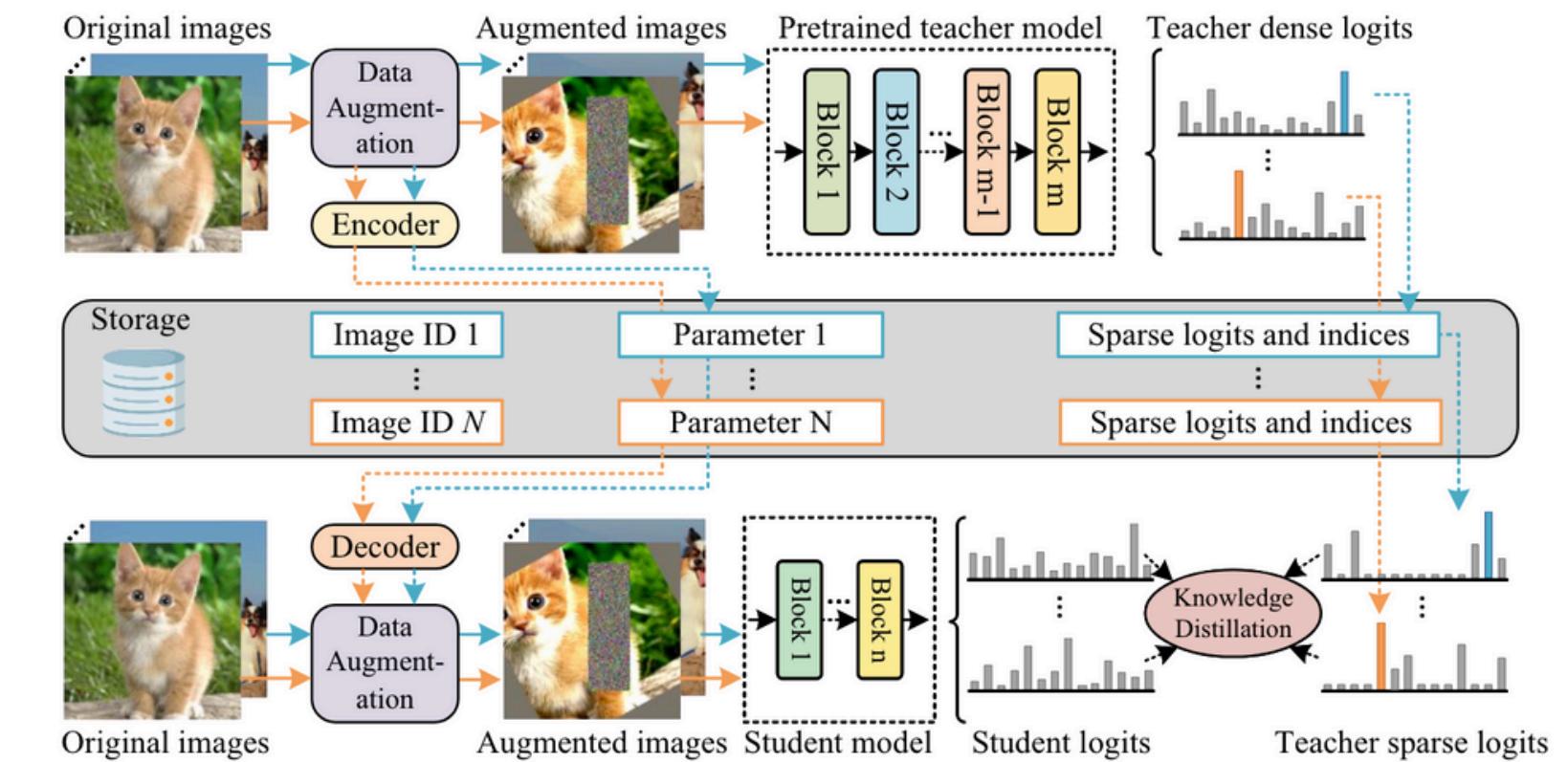


Figure 3: Decoupled distillation for SAM.

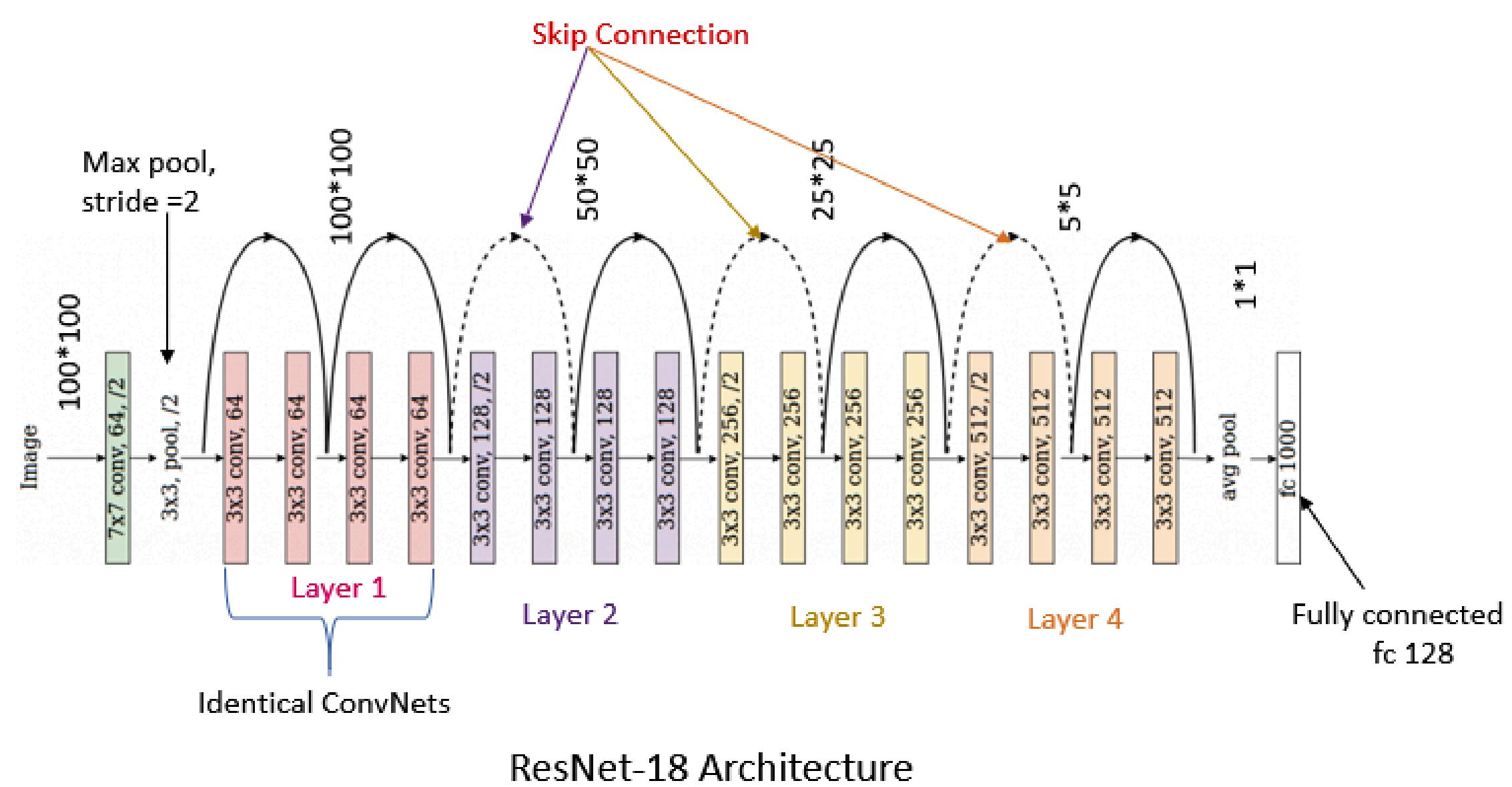
- Replace the heavy (ViT-H) image encoder in SAM with a lightweight TinyViT encoder
- Maintains the same processing pipeline as the original SAM
- Total parameter count of 9.66 million
- Inference speed is 60 times faster than SAM

TinyViT

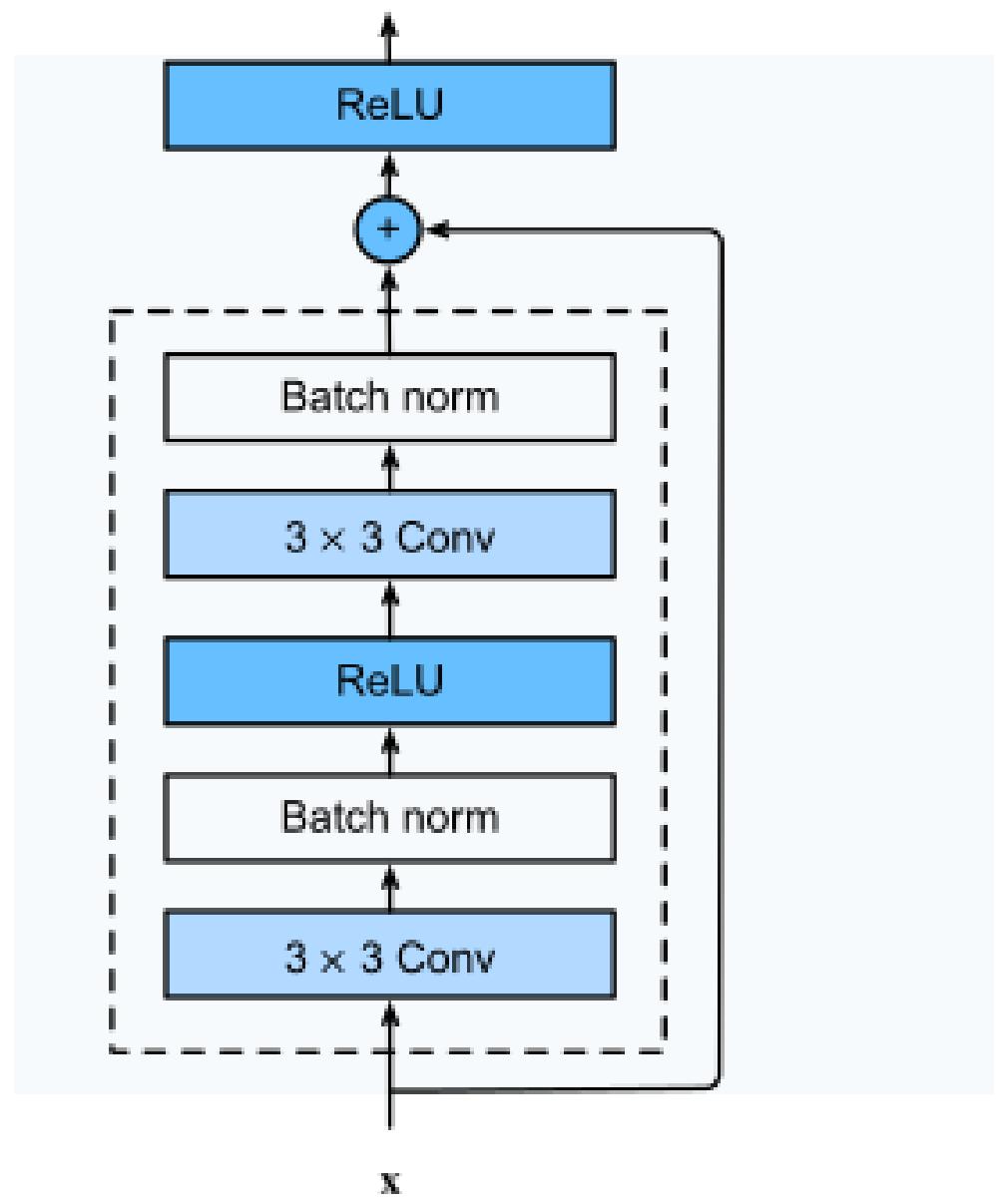


- Efficient training with precomputed data and soft labels.
- Stores sparse soft labels.
- Encodes augmentation parameters into a single parameter.
- Uses MBConvs in early stages, transformer blocks with window attention in later stages.
- Includes attention biases, depthwise convolutions, residual connections, and GELU activation functions.

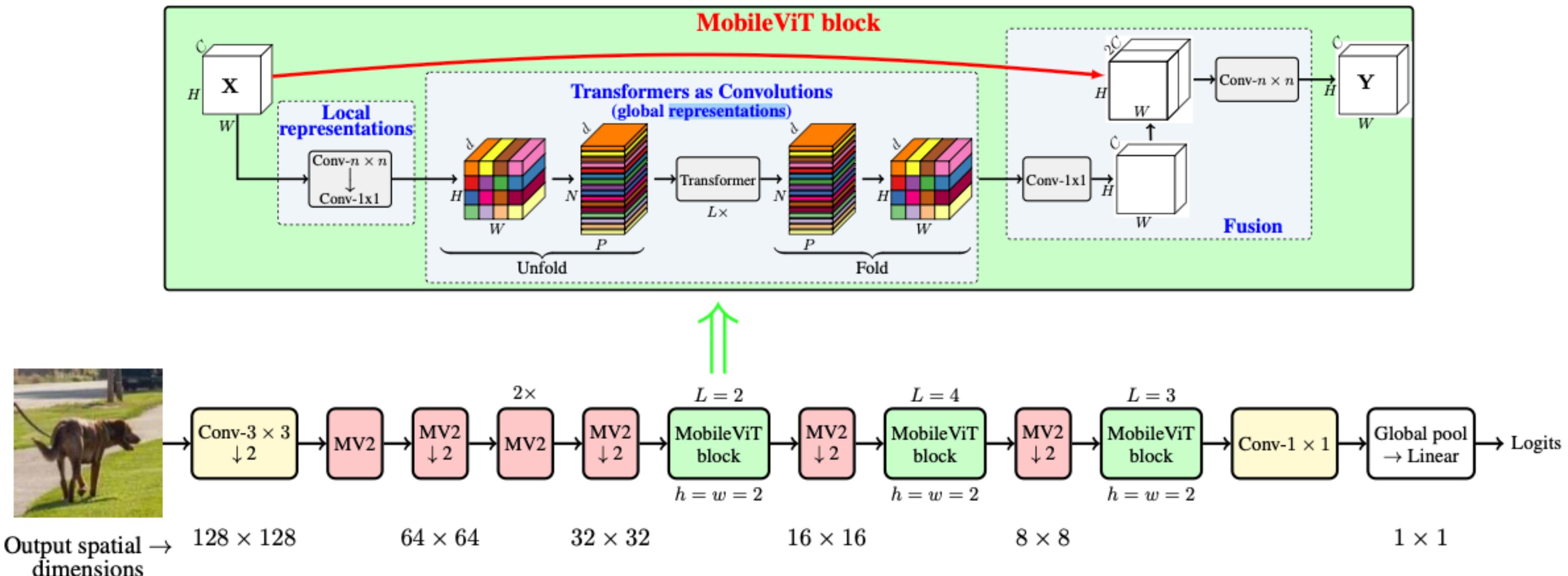
ResNet-18



Fruit 360 Input Image size= 100×100 px

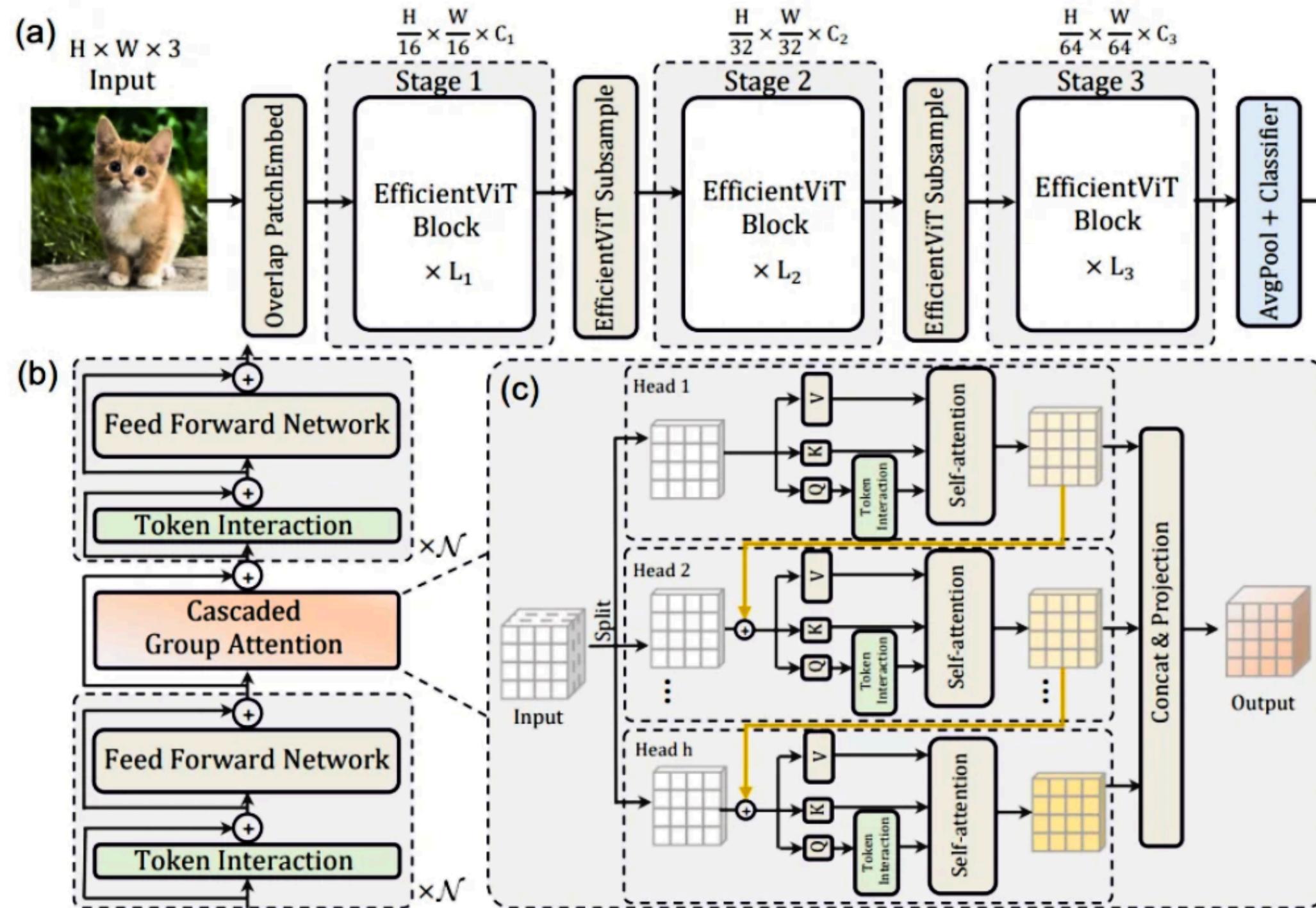


MobileViT-s



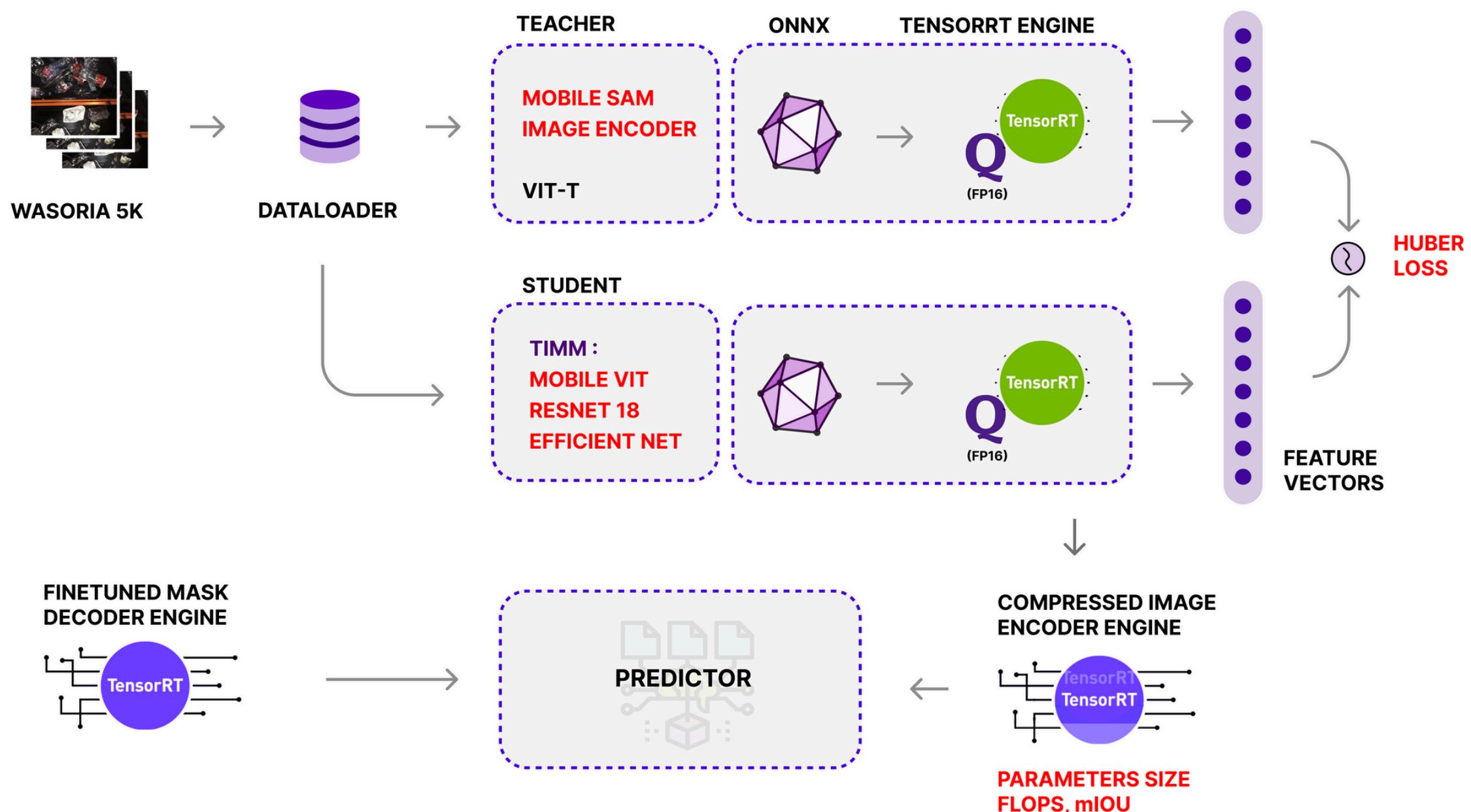
- input X feature maps, unfolded into non-overlapping patches of size P .for transformers, fold patches back to their original spatial dimensions
- $\mathbf{XY}=\mathbf{Conv-1x1}(\mathbf{X}+\mathbf{G})$:) The addition $X+G$ combines the original and global features before applying the final 1×1 convolution.

EfficientViT



- Implements group cascade attention for efficient spatial dependency capture.
- Utilizes a sandwich layout with interleaved attention and convolution layers.
- Offers various size variants for different deployment scenarios.
- Optimizes local and global feature extraction for improved performance

Framework outline



- NVIDIA TensorRT optimizes and accelerates deep learning inference on NVIDIA GPUs for reduced latency and increased throughput.
- ONNX provides an open-source format for model interoperability between different deep learning frameworks.

Huber Loss

Mean Square Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Mean Absolute Error (MAE)

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

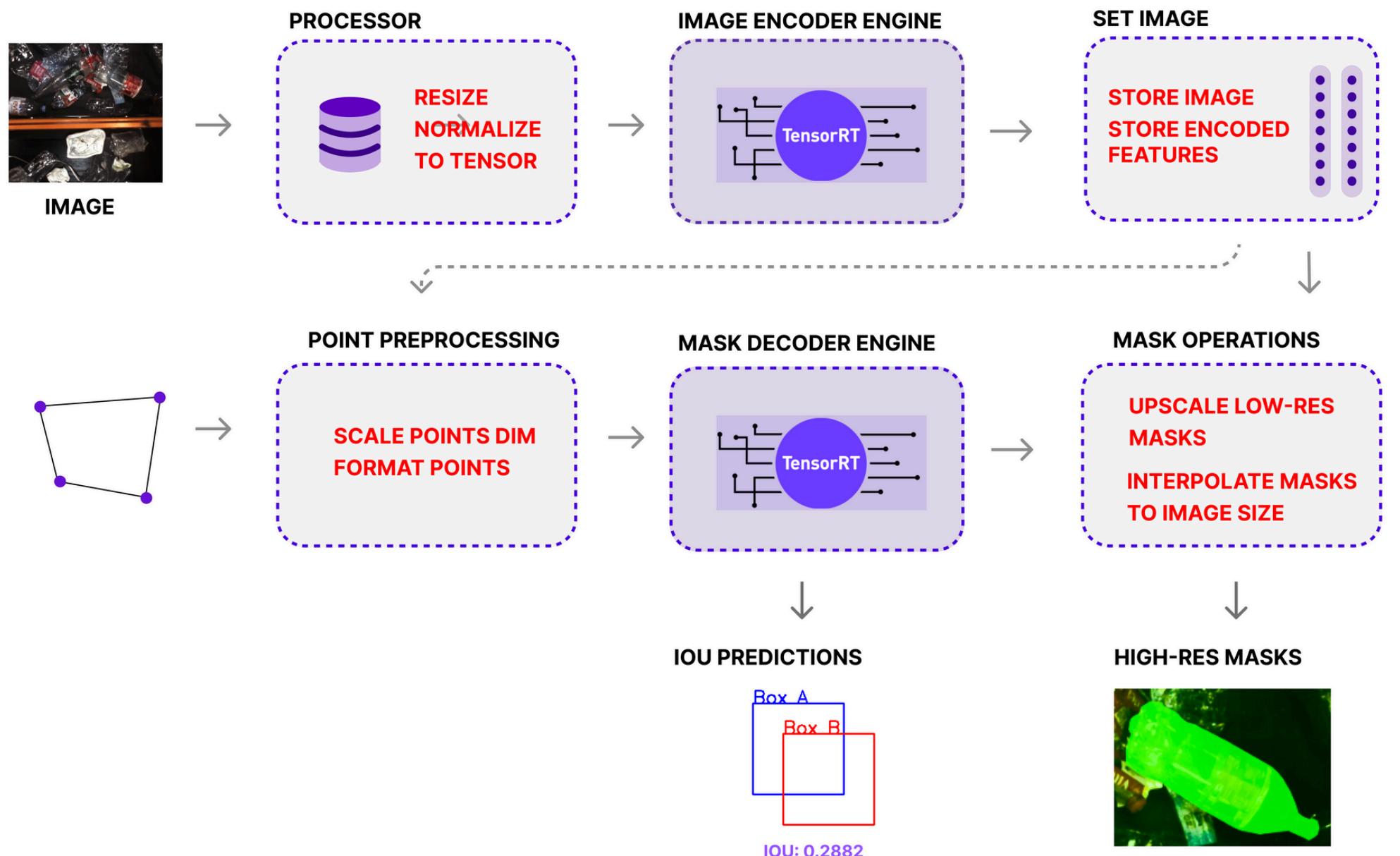
$$L_\delta = \begin{cases} \frac{1}{2}(y - f(x))^2, & \text{if } |y - f(x)| \leq \delta \\ \delta|y - f(x)| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$$

Quadratic

Linear

- **δ :** The threshold parameter that controls the point where Huber Loss transitions from quadratic to linear.
- loss function that is less sensitive to outliers than Mean Squared Error (MSE).
- **Quadratic for Small Errors:** When the error is small ($|y-f(x)| \leq \delta$), the loss behaves like MSE.
- **Linear for Large Errors:** When the error is large ($|y-f(x)| > \delta$), the loss behaves like MAE
- Improves features alignment between complex teacher and student models by penalizing large errors less severely

Predictor



Training configuration

Configuration Parameter	Description	Value
images	Path to the image dataset for distillation	data/coco/train2017
output_dir	Directory to store checkpoints and visualizations	data/wasoria_5k/train
model_name	Name of the student model architecture	resnet18, mobile_vit ...
student_size	Size of the image to feed to the student model	1024
num_images	Limit on the number of images per epoch	None
num_epochs	Number of training epochs	200
batch_size	Number of images per batch	16
num_workers	Number of data loader workers	8
learning_rate	Learning rate for the optimizer	3e-4
distillation_loss	Loss function used for distillation	huber
weight_decay	Weight decay (L2 regularization)	1e-5
data_augmentation	Data augmentation techniques in training	random crop, flip
checkpoint_path	Path to save and load model checkpoints	data/model/checkpoint.pth
log_file	Path to save training logs	data/model/log.txt

```
[wasoria-abdi@wasoria-:~/Desktop/ML_STUDY/nanosam$ python ./nanosam/tools/train.py --images /home/wasoria-abdi/Desktop/ML_STUDY/nanosam/data/train --output_dir /home/wasoria-abdi/Desktop/ML_STUDY/nanosam/data/mobilevit_s --model_name mobilevit_s --num_epochs 200 --batch_size 1 --teacher_image_encoder_engine data/mobile_sam_image_encoder_bs1.engine
]
[17/2024-16:06:10] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchronization. Please non-default stream instead.
[|] 891/891 [04:35<00:00, 3.23it/s]
[|] 891/891 [04:37<00:00, 3.21it/s]
[|] 891/891 [04:35<00:00, 3.24it/s]
[|] 891/891 [04:37<00:00, 3.21it/s]
[|] 891/891 [04:35<00:00, 3.24it/s]
[|] 891/891 [04:38<00:00, 3.20it/s]
[|] 891/891 [04:44<00:00, 3.13it/s]
[|] 891/891 [04:45<00:00, 3.13it/s]
[|] 891/891 [04:43<00:00, 3.15it/s]
[|] 891/891 [04:47<00:00, 3.10it/s]
[|] 891/891 [04:46<00:00, 3.11it/s]
[|] 891/891 [04:42<00:00, 3.15it/s]
[|] 891/891 [04:41<00:00, 3.17it/s]
[|] 891/891 [04:40<00:00, 3.18it/s]
[|] 891/891 [04:41<00:00, 3.16it/s]
[|] 891/891 [04:42<00:00, 3.16it/s]
[|] 891/891 [04:47<00:00, 3.10it/s]
[|] 891/891 [04:51<00:00, 3.05it/s]
[|] 891/891 [04:53<00:00, 3.04it/s]
[|] 891/891 [04:50<00:00, 3.07it/s]
[|] 891/891 [04:39<00:00, 3.19it/s]
[|] 891/891 [04:37<00:00, 3.21it/s]
[|] 891/891 [04:39<00:00, 3.19it/s]
[|] 891/891 [04:36<00:00, 3.23it/s]
```

```
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.00025185859333185477
| 891/891 [04:36<00:00, 3.23it/s]
|- 0.0002528433368689342
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.00025089412498942185
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.0002488935135565458
| 891/891 [04:36<00:00, 3.22it/s]
|- 0.0002569314453603771
| 891/891 [04:36<00:00, 3.23it/s]
|- 0.0002516271961743418
| 891/891 [04:37<00:00, 3.21it/s]
|- 0.000252431181308319
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.0002486367678677813
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.0002468732477778017
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.0002467635430937307
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.00025173573724415434
| 891/891 [04:35<00:00, 3.23it/s]
|- 0.00025061086529599844
| 891/891 [04:35<00:00, 3.23it/s]

no sam wasoria-abdi@wasoria:~/Desktop/ML_STUDY/nanosatS
```



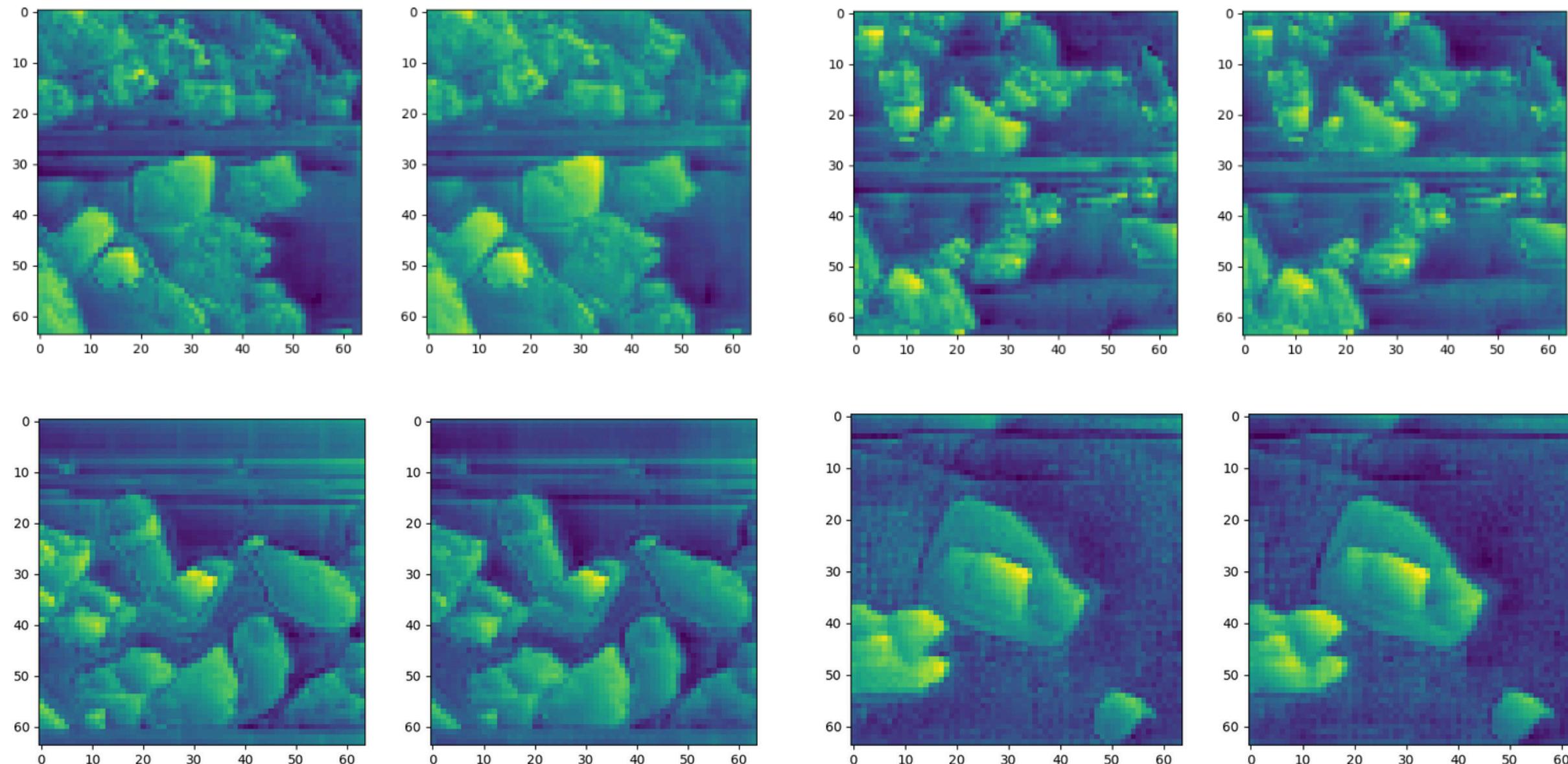
Performance metrics

Model	Default		After Distillation		
	Parameter Size	FLOPs	mIOU all	Parameter Size	FLOPs
MobileSAM	>5M	>10G	-	-	-
ResNet18	11.6M	1.8G	0.73	15.83 M	0.08G
MobileViT-s	9.88M	2.0G	0.75	5.6M	0.10G
EfficientViT-b0	5.3M	0.8G	0.68	4.45M	0.11G

MobileViT-s emerged as the best all-around student model with substantial reductions in parameter size (5.6M) and FLOPs (100 Mega FLOPs), maintaining a high mIOU(). Ideal balance between performance and efficiency.

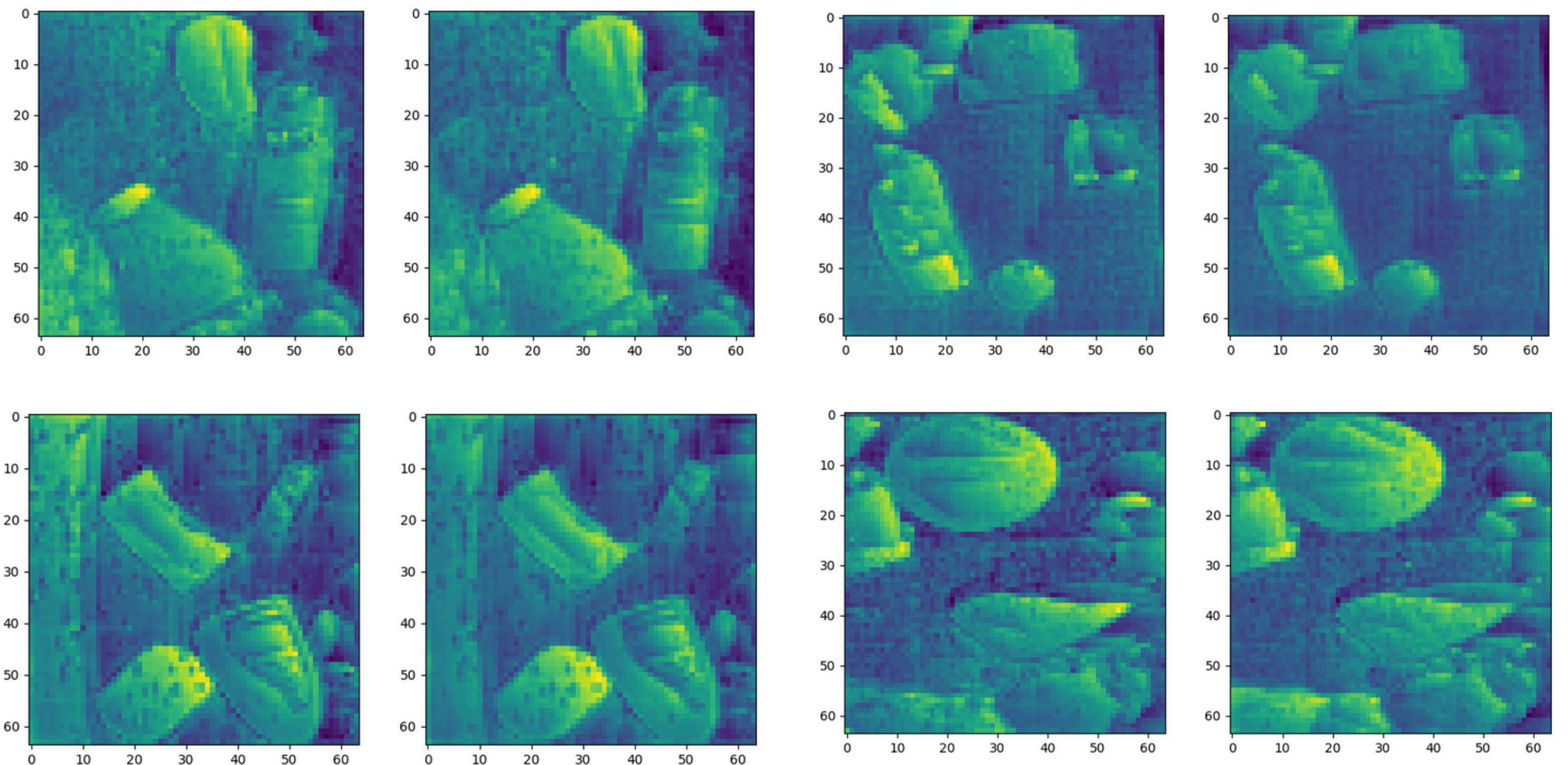
```
(nano_sam) wasoria-abdi@wasoria:~/Desktop/ML_STUDY/nano
/home/wasoria-abdi/Desktop/ML_STUDY/nanosam/data/resne
ONNX Model Analysis
-----
Resnet18
Total Parameters: 15.83M
FLOPs: 0.08G
(nano_sam) wasoria-abdi@wasoria:~/Desktop/ML_STUDY/nano
/home/wasoria-abdi/Desktop/ML_STUDY/nanosam/data/model
onnx
ONNX Model Analysis
-----
Mobilevit S Image Encoder
Total Parameters: 9.88M
FLOPs: 0.10G
(nano_sam) wasoria-abdi@wasoria:~/Desktop/ML_STUDY/nano
/home/wasoria-abdi/Desktop/ML_STUDY/nanosam/data/model
onnx Model Analysis
-----
Efficientvit B0
Total Parameters: 4.45M
FLOPs: 0.11G
```

MobileViT-s



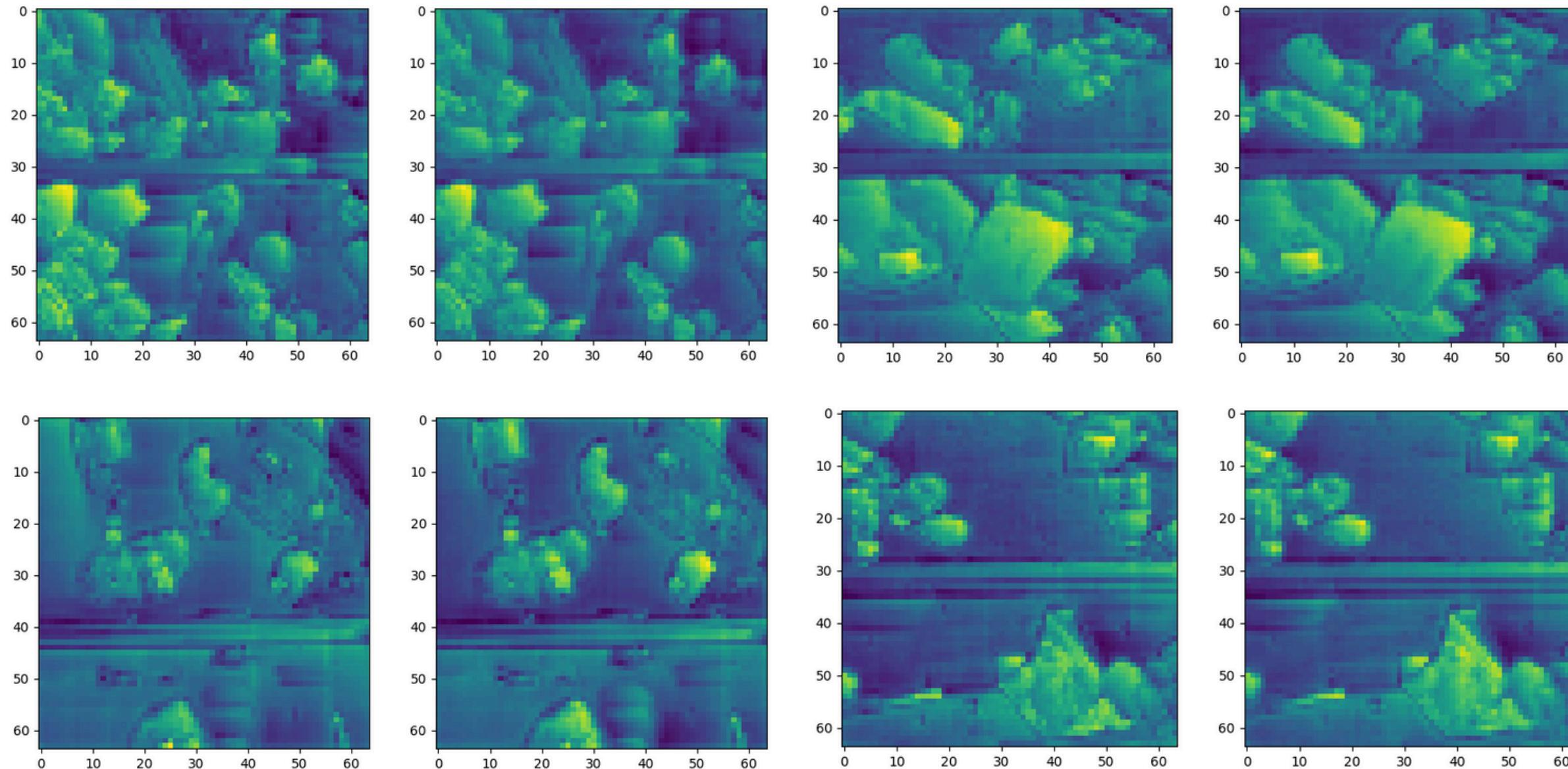
Left: Teacher (MobileSAM), Right: Student (MobileViT-s) showing successful knowledge transfer

ResNet-18



Left: Teacher (MobileSAM), Right: Student (Resnet-18) showing successful knowledge transfer

EfficientViT-b0



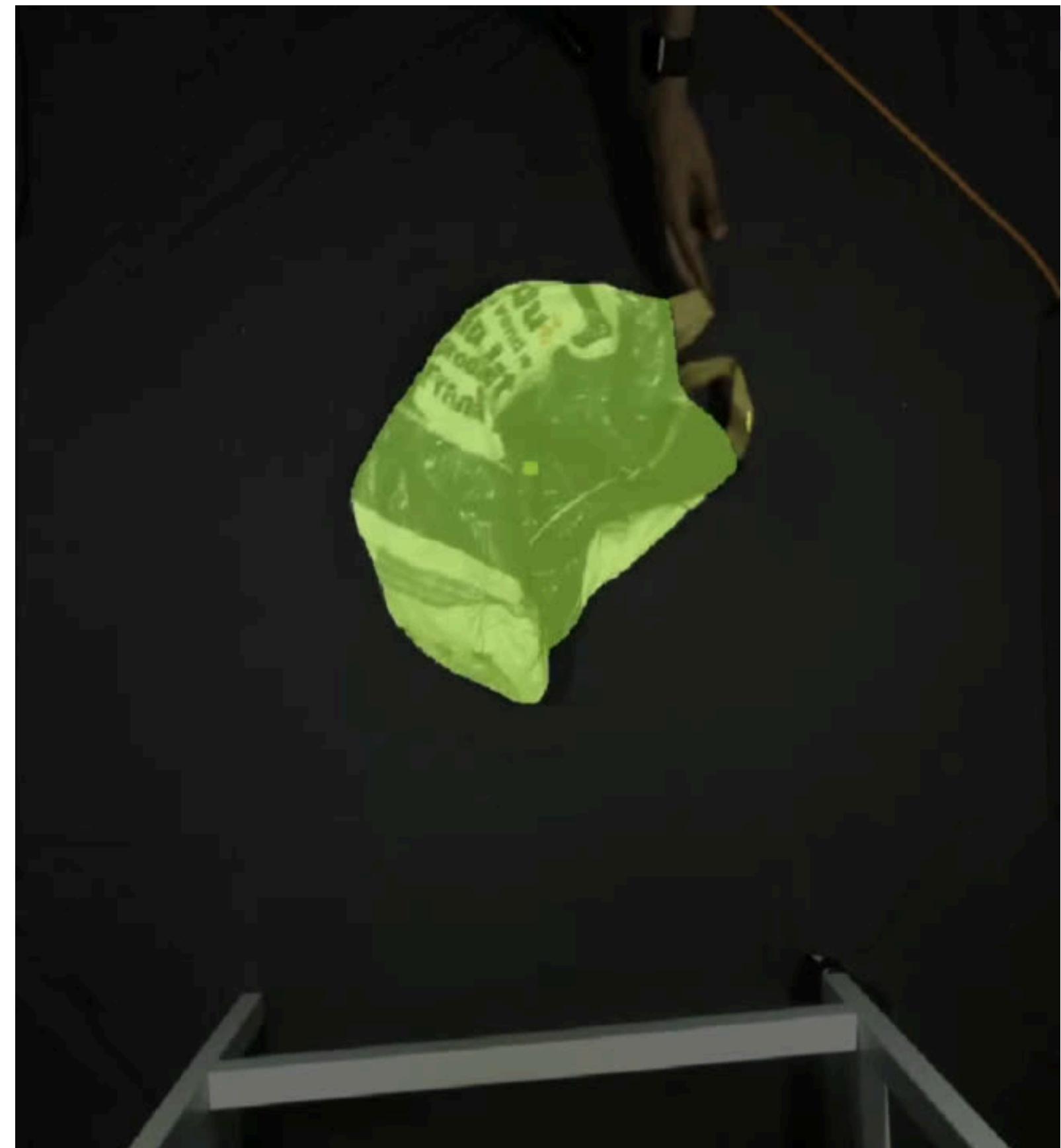
Left: Teacher (MobileSAM), Right: Student (EfficientViT-b0) showing successful knowledge transfer

Mask Visualizer



```
ta-abdi@wasoria:~/Desktop/ML_STUDY/nanosam$ python3 sega.py
di/Desktop/ML_STUDY/nanosam/nanosam/utils/predictor.py:84: UserWarning: The given NumPy array is not writable, and PyTorch does not support non-writable tensors. This means
ult in undefined behavior. You may want to copy the array to protect its data or make it writable before converting it to a tensor. This type of warning will be suppressed f
triggered internally at /opt/conda/conda-bld/pytorch_1682343998658/work/torch/csrc/utils/tensor_numpy.cpp:206.)
sized = torch.from_numpy(image_np_resized).permute(2, 0, 1)
4:48] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
5:00] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
di/Desktop/ML_STUDY/nanosam/nanosam/utils/predictor.py:84: UserWarning: The given NumPy array is not writable, and PyTorch does not support non-writable tensors. This means
ult in undefined behavior. You may want to copy the array to protect its data or make it writable before converting it to a tensor. This type of warning will be suppressed f
triggered internally at /opt/conda/conda-bld/pytorch_1682343998658/work/torch/csrc/utils/tensor_numpy.cpp:206.)
ords = torch.tensor([points]).float().cpu()
5:00] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
ta-abdi@wasoria:~/Desktop/ML_STUDY/nanosam$ python3 sega.py
di/Desktop/ML_STUDY/nanosam/nanosam/utils/predictor.py:84: UserWarning: The given NumPy array is not writable, and PyTorch does not support non-writable tensors. This means
ult in undefined behavior. You may want to copy the array to protect its data or make it writable before converting it to a tensor. This type of warning will be suppressed f
triggered internally at /opt/conda/conda-bld/pytorch_1682343998658/work/torch/csrc/utils/tensor_numpy.cpp:206.)
sized = torch.from_numpy(image_np_resized).permute(2, 0, 1)
7:20] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
7:38] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
di/Desktop/ML_STUDY/nanosam/nanosam/utils/predictor.py:84: UserWarning: The given NumPy array is not writable, and PyTorch does not support non-writable tensors. This means
ult in undefined behavior. You may want to copy the array to protect its data or make it writable before converting it to a tensor. This type of warning will be suppressed f
triggered internally at /opt/conda/conda-bld/pytorch_1682343998658/work/torch/csrc/utils/tensor_numpy.cpp:206.)
ords = torch.tensor([points]).float().cpu()
7:38] [TRT] [W] Using default stream in enqueueV3() may lead to performance issues due to additional calls to cudaStreamSynchronize() by TensorRT to ensure correct synchroni
stream instead.
```

A screenshot of a computer monitor displaying a window titled "WALE_NANO". Inside the window, there is a dark image of several plastic bottles scattered on a surface. A single bottle in the lower-left foreground has a red rectangular bounding box drawn around it. Another bottle in the center has a green semi-transparent mask overlaid on it. In the top-left corner of the window, there is some text: "Image Encoder: mobilevit_s_image_encoder.engine" and "Mask Decoder: mobile_sam_mask_decoder.engine". At the bottom of the window, there is a small text box containing a coordinate list: "[0.6779406070709229, 0.6987832188606262, 0.8143357038497925, 0.7737200260162354]". The background of the monitor shows the terminal command line from which the application was run.





References

1. **ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness.** Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., & Brendel, W. (2019).
2. **Tinyvit: Scaling down vision transformers for efficient transfer learning.** Microsoft Research Blog, 2022.
<https://www.microsoft.com/en-us/research/blog/tinyvit-scaling-down-vision-transformers-for-efficient-transfer-learning/>
3. **Onnx: Open neural network exchange.** ONNX Official Website, 2023. URL <https://onnx.ai/>.
4. **Ketan Doshi. Transformers explained visually** - multi-head attention, deep dive. Ketan Doshi Blog, 2023. URL <https://ketanhdoshi.github.io/transformers-explained-visually>.
5. **Adversarially robust distillation**, Micah Goldblum, Liam Fowl, Soheil Feizi, and Tom Goldstein In Proceedings of the 37th International Conference on Machine Learning (ICML), 2020.
6. **Cross-modal vision transformer for efficient multimodal learning:** Yuhang Gong, Wenhao Jiang, Tiejun Huang, and Yonghong Tian. Xmodalvit: . IEEE Transactions on Multimedia, 24:2373– 2383, 2022.
7. **Knowledge distillation: A survey**, Jianping Gou et al. . International Journal of Computer Vision, 129(6):1789–1819, 2021.
8. **Robust estimation of a location parameter.** Peter J. Huber. The Annals of Mathematical Statistics, 35(1):73–101, 1964.
10. **Efficientvit: Memory efficient vision transformer with cascaded group attention:** Xinyu Liu, Houwen Peng, Ningxin Zheng, Yuqing Yang, Han Hu, and Yixuan Yuan. . Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023. URL <https://arxiv.org/abs/2305.07027>.
11. **Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer** Manoj Kumar et al. . arXiv preprint arXiv:2110.02178, 2022.
12. **NVIDIA. Tensorrt: High-performance deep learning inference.** NVIDIA Developer Blog, 2023. URL <https://developer.nvidia.com/tensorrt>.
13. **Segment anything model (sam)** Meta AI Research. Meta AI Blog, 2023. URL <https://ai.facebook.com/blog/segment-anything-model-sam/>.
14. **Mobilenetv2: Inverted residuals and linear bottlenecks.** Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 4510–4520, 2018a. URL <https://arxiv.org/abs/1801.04381>.
15. **Mobilesam: Lightweight segment anything model**, Kyunghee University Research Team. arXiv preprint arXiv:2304.09867, 2023a.
16. **Exploring mobilevit: Architecture and performance.** Hugo Touvron et al. E Medium, 2022. URL <https://medium.com/@hugotouvron/exploring-mobilevit-architecture-and-performance-9e7f5b0f354a>.
17. **Tinyvit: Fast pretraining distillation for small vision transformers**, Bin Xiao et al. arXiv preprint arXiv:2207.10633, 2022.
18. **Your classifier is secretly an energy-based model and you should treat it like one.** Jianmin Zhang et al.. International Conference on Learning Representations (ICLR), 2020.

Thank You