

308 A Implementation Details

309 **Model Architecture.** We employ the ViT-B/16 version of the Segment Anything Model (SAM) as our
310 base architecture [20], comprising 12 transformer layers. To integrate CLIP capabilities, we append a
311 lightweight CLIP head consisting of 3 transformer layers to the SAM backbone. The patch token
312 outputs from this CLIP head undergo a pooling layer to produce an image-level embedding, akin to the
313 role of the CLS token output in ViT models. We adopt max-pooling since we observe that it can lead
314 to better zero-shot classification and semantic segmentation performance of SAM-CLIP than average
315 pooling. It is noteworthy that max-pooling has been found to be able to encourage the learning of
316 spatial visual features [38]. With the pooling layer, the CLIP head can output an embedding for the
317 whole image, which can be aligned with a text embedding just like the original CLIP model [37].

318 **Dataset Preparation.** For the CLIP distillation, we merge images from several datasets: CC3M [44],
319 CC12M [4], YFCC-15M [37] (a curated subset of YFCC-100M [47] by OpenAI) and ImageNet-
320 21k [41]. This forms our $\mathcal{D}_{\text{CLIP}}$ containing 40.6M unlabeled images. For the SAM self-distillation,
321 we sample 5.7% subset from the SA-1B dataset to form \mathcal{D}_{SAM} , which originally comprises 11M
322 images and 1.1B masks. We randomly select 1% of $\mathcal{D}_{\text{CLIP}}$ and \mathcal{D}_{SAM} as validation sets. Overall, we
323 have 40.8M images for training, which we term as Merged-41M in this work.

324 **Training.** As we discussed in Sec. 2 the training is conducted in two phases to optimize convergence,
325 in a “*probing then full finetuning*” style. The first stage of CLIP-head probing takes 20 epochs on
326 $\mathcal{D}_{\text{CLIP}}$, while the backbone is kept frozen. Here, the teacher model is the OpenCLIP [18] ViT-L/14
327 trained on the DataComp-1B dataset [12]. In the second stage (16 epochs), we unfreeze the backbone
328 $\text{Enc}_{\text{SAM-CLIP}}$ and proceed with joint fine-tuning together with $\text{Head}_{\text{CLIP}}$ and Head_{SAM} , incorporating
329 both CLIP and SAM distillation losses at the ratio of 1:10. The original SAM ViT-B model serves
330 as the teacher in SAM loss. Further, the learning rates applied to $\text{Enc}_{\text{SAM-CLIP}}$ and Head_{SAM} are 10
331 times smaller than that of $\text{Head}_{\text{CLIP}}$ in order to reduce the forgetting of the original SAM abilities.
332 Besides, we adopt a mixed input resolution strategy for training. A notable difference between SAM
333 and CLIP is their pretraining resolution. SAM is trained and works best on 1024px resolution while
334 often lower resolutions (e.g., 224/336/448px) are adopted for CLIP training and inference [37, 7, 45].
335 Hence, we employ variable resolutions of 224/448px for the CLIP distillation via the variable batch
336 sampler approach of [31], while SAM distillation utilizes a 1024px resolution in accordance with
337 SAM’s original training guidelines [20]. In every optimization step, we form a batch of 2048 images
338 from $\mathcal{D}_{\text{CLIP}}$ and 32 images (each with 32 mask annotations) from \mathcal{D}_{SAM} and perform training in a
339 multi-task fashion.

340 **Resolution Adaption.** After the two training stages, SAM-CLIP can accomplish CLIP tasks (e.g.,
341 zero-shot classification) using the CLIP-head under 224/336/448px, and run inference with the
342 SAM-head under 1024px. However, if one wants to apply the two heads together on a single input
343 image for certain tasks (we present a demo of this in Sec. A.3), it would be inefficient to pass the
344 image twice to the image encoder with two resolutions for the two heads respectively. To remedy this
345 issue, we adapt the CLIP head for 1024px input using a very short and efficient stage of fine-tuning:
346 freezing the image encoder and only finetuning the CLIP-head with $\mathcal{L}_{\text{CLIP}}$ for 3 epochs (it is the
347 same as the first stage of training, which is also CLIP-head probing) under variable resolutions of
348 224/448/1024px. *Note:* resolution upscaling strategies are prevalent in CLIP training: [37, 45, 22]
349 show it is more efficient than training with high resolution from the beginning.

350 A.1 Zero-Shot Evaluations

351 **CLIP Task: Zero-Shot Image Classification.** To examine the CLIP-related capabilities of
352 SAM-CLIP, we evaluate it with zero-shot image classification on ImageNet [8], ImageNet-v2 [39] and
353 Places365 [54], under image resolution of 224x. We use the text templates as CLIP [37] utilizing the
354 textual embeddings from the text encoder of SAM-CLIP (which is kept frozen from our CLIP teacher)
355 to perform zero-shot classification without any finetuning. The evaluation results are presented in
356 Table I. Employing a ViT-B architecture, our model achieves zero-shot accuracy comparable to the
357 state-of-the-art CLIP ViT-B models pretrained on LAION-2B [43] and DataComp-1B [12] (both
358 released by [18]), over the three datasets. These results validate the efficacy of our merging approach
359 in inheriting CLIP’s capabilities.

360 **SAM Task: Zero-Shot Instance Segmentation.** For the SAM component of SAM-CLIP, we evaluate
361 its performance in instance segmentation, a task at which the original SAM model excels [20], with

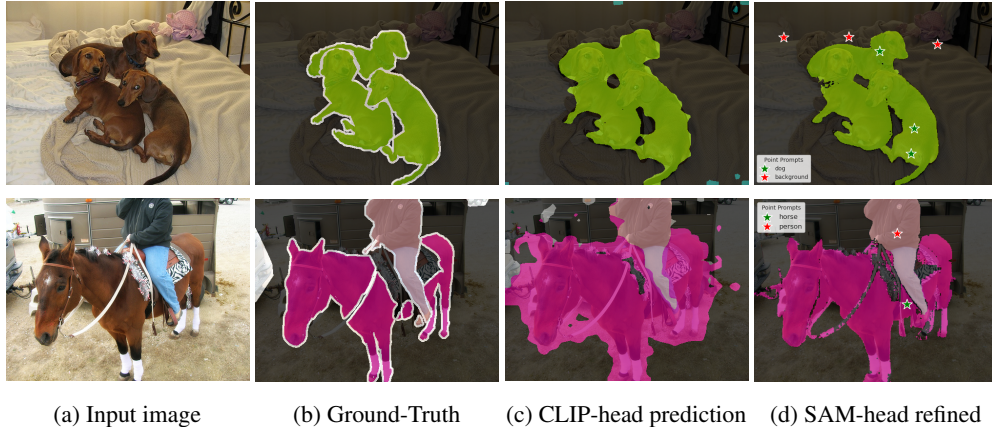


Figure 5: Demo on zero-shot semantic segmentation. Passing an input image through the image encoder, $\text{Head}_{\text{CLIP}}$ can predict a semantic segmentation mask, and Head_{SAM} can refine it to a more fine-grained mask with auto-generated geometric prompts.

Table 1: Zero-shot evaluations on classification and instance segmentation tasks, comparing SAM-CLIP with state-of-the-art models that use the ViT-B architecture. SAM-CLIP demonstrates minimal forgetting compared to the baseline FMs on their original tasks.

Model	Training Data	0-Shot Classification (%)			0-Shot Instance Seg. (mAP)	
		ImageNet	ImageNet-v2	Places-365	COCO	LVIS
SAM [20]	SA-1B	-	-	-	41.2	36.8
CLIP [37]	OpenAI-400M	68.3	62.6	42.2	-	-
CLIP [7]	LAION-2B	71.1	61.7	43.4	-	-
CLIP [12]	DataComp-1B	73.5	65.6	43.0	-	-
SAM-CLIP (Ours)	Merged-41M	72.4	63.2	43.6	40.9	35.0

362 COCO [26] and LVIS [14] datasets. Following the original practices of [20], we first generate object
 363 detection bounding boxes using a ViT-Det model (ViT-B version) [23]. These bounding boxes act as
 364 geometric prompts for SAM’s prompt encoder, which then predicts masks for each object instance.
 365 The evaluation results of SAM-CLIP and the original SAM ViT-B are provided in Table 1 (both
 366 under 1024px resolution), showing that SAM-CLIP is very close to SAM on the two benchmarks, not
 367 suffering from catastrophic forgetting during training.

368 **Zero-Shot Transfer to Semantic Segmentation.** We extend our evaluation to (text-prompted) zero-
 369 shot semantic segmentation over 5 datasets, Pascal VOC [10], Pascal Context [33], ADE20k [55],
 370 COCO-Stuff [2] and COCO-Panoptic [19, 26]. We adopt a common evaluation protocol for this
 371 task: i) each input image is resized to 448×448 px and pass to the image encoder and CLIP-head
 372 of SAM-CLIP to obtain 28×28 patch features; ii) OpenAI’s 80 pre-defined CLIP text templates
 373 are employed to generate textual embeddings for each semantic class, and these embeddings act as
 374 mask prediction classifiers and operate on the patch features from the CLIP head; iii) we linearly
 375 upscale the mask prediction logits to match the dimensions of the input image. Evaluation results of
 376 SAM-CLIP and previous zero-shot models over the five datasets are demonstrated in Fig. 2. Notably,
 377 SAM-CLIP establishes new state-of-the-art performance on all 5 datasets, with a significant margin
 378 over past works.

379 A.2 Head-Probing Evaluations on Learned Representations

380 By merging the SAM and CLIP models, we anticipate that the resultant model will inherit advantages
 381 at the representation level from both parent models. Specifically, SAM excels at capturing low-
 382 level spatial visual details pertinent to segmentation tasks, while CLIP specializes in high-level
 383 semantic visual information encompassing the entire image. We hypothesize that the merged model
 384 combines these strengths, thereby enhancing its utility in broad range of downstream vision tasks. To

Table 2: Zero-shot semantic segmentation performance comparison with recent works. ([†]SegCLIP is trained on COCO data, so it is not zero-shot transferred to COCO-Stuff.)

Model	Arch	Training Data	0-Shot Semantic Segmentation (mIoU %)				
			Pascal VOC	Pascal-Context	ADE20k	COCO-Stuff	COCO-Panoptic
GroupViT [49]	ViT-S	Merged-26M	52.3	22.4	-	24.3	-
ViewCo [40]	ViT-S	Merged-26M	52.4	23.0	-	23.5	-
ViL-Seg [27]	ViT-B	CC12M	37.3	18.9	-	18.0	-
OVS [50]	ViT-B	CC4M	53.8	20.4	-	25.1	-
CLIPpy [38]	ViT-B	HQITP-134M	52.2	-	13.5	-	25.5
TCL [3]	ViT-B	CC3M+CC12M	51.2	24.3	14.9	19.6	-
SegCLIP [28]	ViT-B	CC3M+COCO	52.6	24.7	8.7	26.5 [†]	-
SAM-CLIP	ViT-B	Merged-41M	60.6	29.2	17.1	31.5	28.8

Table 3: Head probing evaluations on semantic segmentation datasets, comparing our model with SAM and CLIP that use the ViT-B architecture. Avg is the average evaluation results of three heads.

Model	Training Data	Pascal VOC				ADE20k			
		Linear	DeepLabv3	PSPNet	Avg	Linear	DeepLabv3	PSPNet	Avg
SAM	SA-1B	46.6	69.9	71.2	62.6	26.6	32.8	36.2	31.9
CLIP	DataComp-1B	70.7	78.9	79.7	76.4	36.4	39.4	40.7	38.8
SAM-CLIP	Merged-41M	75.0	80.3	81.3	78.8	38.4	41.1	41.7	40.4

Table 4: Composing both CLIP and SAM heads of SAM-CLIP for zero-shot semantic segmentation on Pascal VOC.

Method	Resolution	mIoU
CLIP head only	448px	60.6
CLIP+SAM heads	1024px	66.0

Table 5: Linear probing evaluations on image classification datasets with ViT-B models.

Model	Linear Probing	
	ImageNet	Places365
SAM	41.2	41.5
CLIP (DataComp1B)	81.3	55.1
CLIP (LAION-2B)	79.6	55.2
SAM-CLIP	80.5	55.3

385 investigate this hypothesis, we conduct head-probing (i.e., learn a task specific head with a frozen
 386 image backbone) evaluations on SAM, CLIP, and SAM-CLIP, utilizing different segmentation head
 387 structures (linear head, DeepLab-v3 [5] and PSPNet [53]) across two semantic segmentation datasets,
 388 Pascal VOC and ADE20k. The results are presented in Table 3. We observe that SAM representations
 389 do not perform as well as those of CLIP for tasks that require semantic understanding, even for
 390 semantic segmentation task. However, SAM-CLIP outperforms both SAM and CLIP across different
 391 head structures and datasets, thereby confirming its superior visual feature representation capabilities.

392 Besides, we apply linear probing to these models for image classification tasks on two datasets,
 393 ImageNet and Places365. Results in Table 5 show that SAM-CLIP attains comparable performance
 394 with CLIP, implying that the image-level representation of SAM-CLIP is also well-learned. All head
 395 probing evaluation results are visualized in Figure 3 to deliver messages more intuitively.

396 A.3 Composing Both CLIP and SAM Heads for Better Segmentation

397 Given that SAM-CLIP is a multi-task model with SAM and CLIP heads, one would naturally ask if
 398 the two heads can work together towards better performance on some tasks. Here, we showcase that a
 399 simple composition of the CLIP and SAM heads can lead to better zero-shot semantic segmentation.
 400 Specifically, we resize the input image to 1024px and pass it through $Enc_{SAM-CLIP}$, and use the CLIP
 401 head to generate low-resolution mask prediction (32×32) using text prompts. Then, we generate
 402 some point prompts from the mask prediction (importance sampling based on the mask prediction
 403 confidence), and pass the mask prediction and point prompts together to the prompt encoder module
 404 as geometric prompts. Finally, $Head_{SAM}$ takes embeddings from both the prompt encoder and the
 405 image encoder to generate high-resolution mask predictions (256×256) as shown in Figure 2 (right).
 406 Examples of this pipeline are shown in Figure 5. One can clearly observe that the refined segmentation
 407 by the SAM-head is more fine-grained.

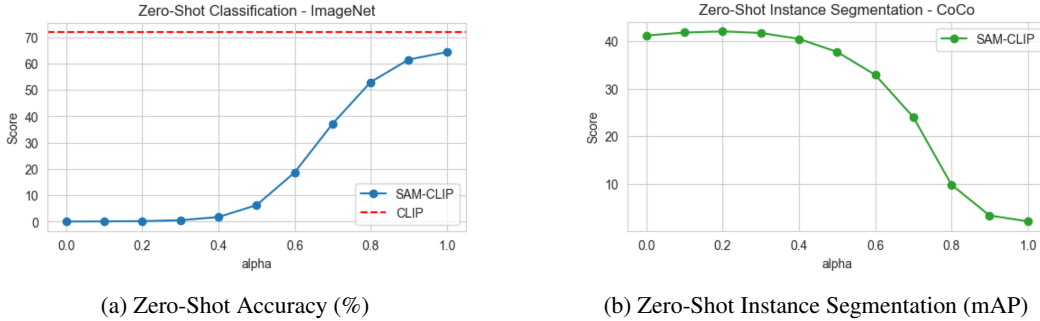


Figure 6: Wise-FT [48] to a CLIP-distilled SAM ViT-B model. The red dashed line marks the performance of the CLIP teacher model.

408 Note that this pipeline requires *only one forward pass* on $\text{Enc}_{\text{SAM-CLIP}}$ with 1024px resolution. For
 409 fair comparison, in Table 1 and Figure 1 we report SAM-CLIP zero-shot segmentation performance
 410 with 448px resolution using $\text{Head}_{\text{CLIP}}$ only. Using our high-resolution pipeline we obtain further
 411 gain in zero-shot semantic segmentation as shown in Table 4

412 B Weight Averaging

413 Weight averaging is a straightforward post-processing method proven to mitigate forgetting across a
 414 variety of fine-tuning tasks. Specifically, Wise-FT [48] proposes linearly interpolating the pretrained
 415 and fine-tuned parameters using a coefficient α . In this study, we explore the application of Wise-FT
 416 in our setup. We focus exclusively on CLIP distillation applied to SAM ViT-B (serving as the
 417 student model), with a CLIP ViT-B/16 model acting as the teacher model. The model is trained on
 418 ImageNet-21k for 20 epochs. It is evident that the fine-tuned student model ($\alpha = 1$) gains zero-shot
 419 classification capabilities at the expense of forgetting its original zero-shot instance segmentation
 420 abilities. Upon applying Wise-FT to the fine-tuned model, we observe an inherent tradeoff between
 421 learning and forgetting. Notably, no optimal point exists where both high classification accuracy
 422 ($> 60\%$ on ImageNet) and a high mAP (> 35 mAP on COCO) are achieved simultaneously.