
Q-SAM2: Accurate Quantization for Segment Anything Model 2

Nicola Farronato

IBM Research - Zurich
ETH Zurich

nicola.farronato@ibm.com

Florian Scheidegger

IBM Research - Zurich
eid@zurich.ibm.com

Mattia Rigotti

IBM Research - Zurich
mrg@zurich.ibm.com

Cristiano Malossi

IBM Research - Zurich
acm@zurich.ibm.com

Michele Magno

ETH Zurich
michele.magno@pbl.ee.ethz.ch

Haotong Qin

ETH Zurich
haotong.qin@pbl.ee.ethz.ch

Abstract

The Segment Anything Model 2 (SAM2) has gained significant attention as a foundational approach for promptable image and video segmentation. However, its expensive computational and memory consumption poses a severe challenge for its application in resource-constrained scenarios. In this paper, we propose an accurate low-bit quantization method for efficient SAM2, termed Q-SAM2. To address the performance degradation caused by the singularities in weight and activation distributions during quantization, Q-SAM2 introduces two novel technical contributions. We first introduce a linear layer calibration method for low-bit initialization of SAM2, which minimizes the Frobenius norm over a small image batch to reposition weight distributions for improved quantization. We then propose a Quantization-Aware Training (QAT) pipeline that applies clipping to suppress outliers and allows the network to adapt to quantization thresholds during training. Our comprehensive experiments demonstrate that Q-SAM2 allows for highly accurate inference while substantially improving efficiency. Both quantitative and visual results show that our Q-SAM2 surpasses existing state-of-the-art general quantization schemes, especially for ultra-low 2-bit quantization. While designed for quantization-aware training, our proposed calibration technique also proves effective in post-training quantization, achieving up to a 66% mIoU accuracy improvement over non-calibrated models.

1 Introduction

The Segment Anything Model 2 (SAM2) has demonstrated outstanding success for promptable visual segmentation [23]. Its capabilities in terms of video segmentation accuracy make this model attractive for applications in robotics and tiny machine learning. However, taking advantage of real-time inference with this model requires high-end computation hardware, such as an NVIDIA A100 GPU. Additionally, edge deployment scenarios typically impose tighter memory and storage constraints (even at the MB level in extreme cases), which make the use of the existing version of SAM2 prohibitive. An effective strategy to reduce the size and computation of the models is to quantize them using fewer bits. Quantization is a widely adopted technique for optimizing deep

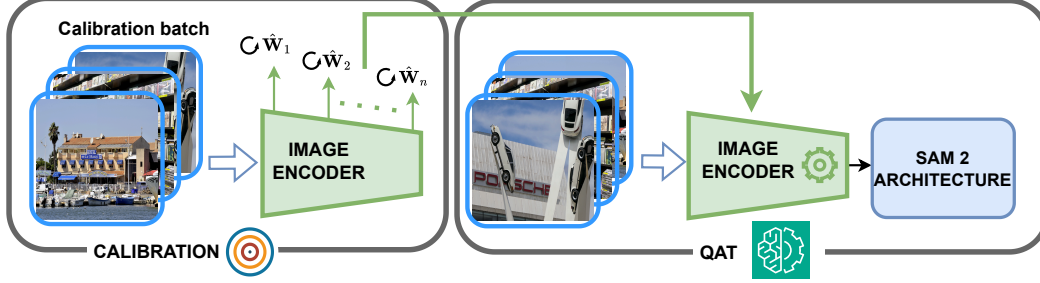


Figure 1: The Q-SAM2 approach. The weight distributions of the linear layers in the image encoder are calibrated to reduce variance. We substitute the original encoder, then quantize and train the entire architecture using a custom QAT pipeline. This calibration process ensures that the initial quantization error is lower than that produced by the original image encoder under the same conditions.

learning models, particularly when deploying them on resource-constrained hardware. By reducing the numerical precision of the model parameters and activations, quantization leads to faster inference, lower power consumption, and smaller model sizes.

For these reasons the demand for post-training quantization (PTQ) methods has grown substantially in recent years, particularly due to the increasing size and complexity of large language models (LLMs) [5, 26, 1, 17]. These models often consist of billions of parameters, making them difficult to deploy efficiently without compression techniques. PTQ offers a practical solution by enabling model compression without requiring retraining, thus reducing computational costs and accelerating deployment [10]. Quantization Aware Training (QAT), on the other hand, offers a more accurate solution by simulating quantization effects during the training phase, allowing the model to adapt to reduced numerical precision [11, 24]. QAT enables the model to recover accuracy while achieving aggressive compression, but is inapplicable when resources are limited or the data set is massive.

While no previous work focuses on the SAM2 architecture, some authors have proposed PTQ algorithms for the original Segment Anything Model [12]. The complexity of quantizing the SAM model was studied in [18, 16], which can efficiently obtain per-channel quantized models with few data samples and limited computation resources. But although these PTQ methods offer solutions without retraining, they inherently lack the adaptability to fine-tune model parameters. Thus, they potentially lead to suboptimal accuracy performance when lower-bit quantization is required, especially considering the quantization for a more compact SAM2 and challenging video segmentation tasks.

We find that quantizing SAM2 is nontrivial and presents a challenge due to outliers in weight and activation distributions. Specifically, the linear layer weight matrices in the image encoder exhibit values that deviate significantly from the center of the distribution, introducing outliers that hinder effective quantization. In addition, the activation outputs also exhibit heavy-tailed distributions, further degrading the performance when uniform quantization schemes are applied.

To address these challenges and meet the objective of ultra-low bit quantization, specifically 2-bit precision, we propose Q-SAM2, a QAT pipeline to achieve low-bit resolution. Figure 1 introduces the complete approach that consists of two steps. In the initial calibration step, we reduce the variance of the weight distributions in the linear layers of the image encoder by applying a closed-form update derived from a small batch of calibration data. This operation preserves the encoder’s generalization capability during inference while effectively narrowing the weight distributions, a property that is critical for the subsequent quantization-aware training phase. When combined with distribution clipping, this leads to a more efficient use of the available quantization levels, ensuring that each bit contributes meaningfully.

Our main contributions are the following.

- We introduce Q-SAM2, to the best of our knowledge the first quantized model based on the SAM2 architecture, demonstrating strong performance even under ultra-low bit schemes.
- A novel light calibration method that reduces the variance of weight distributions in the SAM2 image encoder using a small batch of images. This procedure provides better initialization for quantization, significantly lowering the initial quantization error.

- A tailored quantization-aware training (QAT) pipeline that leverages the calibrated encoder and applies clipping to weights and activations based on their distributions. This enables effective ultra-low-bit quantization using only a small portion of the full training dataset.

The remainder of this paper is organized as follows. Section 2 describes the background and related work, Section 3 shows the Q-SAM2 pipeline, Section 4 presents our experimental setup and results, and Section 5 concludes our findings.

2 Background and Related Work

2.1 Segment Anything 2 (SAM2)

The Segment Anything Model 2 (SAM2)[23] is an open-source foundation model by Meta AI for promptable segmentation in both images and videos[19]. Unlike the original SAM [12], which targets static images, SAM2 extends segmentation to video, addressing challenges like temporal consistency, motion, and occlusion. SAM2 employs a hierarchical transformer encoder with a streaming memory mechanism to retain spatio-temporal context across frames, enabling consistent mask predictions without re-prompting. It supports prompts such as points, boxes, or text and achieves real-time, state-of-the-art performance on high-end hardware (e.g., NVIDIA A100).

For training, SAM2 leverages the CC BY 4.0 licensed SA-V dataset¹, which contains 51,000 video clips and 600,000 segmentation masks. This dataset was built using a model-in-the-loop approach, significantly reducing the amount of manual annotation required while ensuring high-quality labels.

In this work, our main focus is the SAM2 image encoder, which accounts for more than 90% of the total parameters (in the *base_plus* model) and poses a significant challenge for ultra-low-bit quantization (see Section 3).

2.2 Model Quantization

Quantization is the process of approximating a continuous or high-resolution discrete domain by a finite set of representative levels, to optimize computational efficiency and memory usage with controlled performance degradation. The quantization process can be formulated as:

$$\mathbf{X}_q = \text{clip} \left(\text{round} \left(\frac{\mathbf{X}}{s} \right) + z, 0, 2^b - 1 \right). \quad (1)$$

Where \mathbf{X} is a high-precision input matrix or tensor, typically represented in floating-point (e.g., FP32 or FP16), and z is the zero-point ensuring that the real value zero is mapped exactly to the quantized value z . Then the number of bits b defines the dynamic range of the quantized values. For unsigned quantization, the range is $[0, 2^b - 1]$; for signed quantization, it is $[-2^{b-1}, 2^{b-1} - 1]$. Proper clipping is required to avoid overflow. The scaling factor s and the zero-point z can be a scalar ($s \in \mathbb{R}$) for tensor quantization, a vector ($\mathbf{s} \in \mathbb{R}^d$) for per-axis quantization, or a tensor ($\mathbf{s} \in \mathbb{R}^{d_1 \times d_2 \times \dots}$) for group quantization. Finer granularity (vector or tensor scales) typically improves accuracy by adapting to local variations but increases storage and computational overhead compared to global quantization. We focus on per-axis uniform quantization, which offers a good balance between efficiency and complexity, avoiding the overhead of group-wise methods.

2.3 Quantization Aware Training (QAT)

In quantization-aware training (QAT), the quantization process is simulated during model training by inserting fake quantization operations based on equation 1. The scale s and zero point z are treated as learnable parameters or are dynamically adapted during training, allowing the model to compensate for quantization-induced distortions. Formally, QAT introduces non-differentiable quantization operations into the optimization loop, typically handled using straight-through estimators (STE) [3] to enable gradient-based updates.

In contrast, post-training quantization (PTQ) applies quantization after training without modifying the original weights. In PTQ, s and z are statically computed from calibration data, such as by minimizing the reconstruction error or matching dynamic ranges, and are not adjusted through gradient descent.

¹<https://ai.meta.com/datasets/segment-anything-video/>

Consequently, QAT generally achieves higher accuracy, especially under aggressive quantization regimes (e.g., low-bit settings), at the cost of increased training complexity and time, while PTQ offers a faster and simpler pipeline but may suffer from larger performance degradation.

Several QAT methods have been proposed to improve accuracy in low-precision regimes. One of the earliest and most influential methods is DoReFa-Net [28], which introduces quantization of weights, activations, and gradients using deterministic functions. However, its use in modern transformer architectures is limited due to the inability to quantize models with large activation ranges [14].

PACT [4] introduced a learnable clipping parameter for activations, enabling better control of dynamic range. Nevertheless, this introduces additional parameters and gradients, complicating the deployment. LSQ [8] proposes learning the quantization step size via backpropagation, achieving strong accuracy but requiring gradient computation for scale factors, which are not always hardware-friendly.

3 Method

3.1 Challenges in quantizing the SAM2 Image Encoder

The quantization of the Hiera-based image encoder presents additional challenges compared to conventional transformer architectures such as LLMs [5]. Hiera employs a hierarchical architecture that progressively merges tokens, enabling multi-scale feature extraction and improved efficiency for dense visual tasks, in contrast to the standard Vision Transformer (ViT) [7] that processes flat sequences of image patches with uniform global attention and no inherent spatial hierarchy. Our analysis outlines the differences and limitations of current algorithms before proposing our method.

As with many transformer-based architectures [5, 25], some activations in this model have very large magnitude values. The occurrence of such outliers is limited to a small number of channels, where they tend to appear repeatedly for each input image. Many recent works in LLMs focused on reducing the difficulty of quantizing activation by taking advantage of the mostly uniform weight distribution which allows one to trade off outliers in activations with outliers in weight distributions [26, 17, 1].

However, in the case of the SAM2 model the approach of focusing only on activation outliers turns out to be ineffective, as the weight distribution of the image encoder is characterized by pronounced outliers and heavy-tailed behavior, due to a small number of weights with substantially larger magnitudes compared to the majority. This irregularity complicates quantization, as standard methods struggle to represent extreme values under limited precision, leading to severe information loss or degraded performance. The largest outliers are concentrated in early layers, causing errors to propagate through the network. Empirically, we observed that reducing the bit-width of weights to 4 or 2 bits, sharply increased quantization error, causing significant degradation in model performance until the inability to segment even the easiest instances.

To address the problem of weights as well as activations outliers, we analyzed the impact of weight value clipping as a potential mitigation strategy using a simple quantization strategy. By truncating extreme values at fixed thresholds, value clipping reduces the dynamic range, enabling quantization to allocate more precision where most values are concentrated.

Although clipping initially resulted in a notable loss in model accuracy, further analysis revealed that the core predictive behavior was preserved, suggesting that outliers were not essential for model effective functioning. We observed the same behavior for activations that remained more difficult to quantize on low bits due to the higher values range.

Building on this analysis, eliminating the tail values allowed us to concentrate the quantization effort on the central part of the distribution. To take advantage of this effect, we designed a pipeline composed of a calibration phase (Section 3.2), followed by clipping and quantization. The role of the calibration process is to optimize the weight distributions of the original network to reduce error in the subsequent quantization-aware training (QAT) phase (Section 3.3), where clipping is applied on the weight and activation distributions.

Figure 2 compares a standard 3-bit quantization approach to our proposed method, applied to a linear layer weight distribution. The Gaussian curve represents the distribution of weight values, and the Gaussian shape represents the distribution of the weight values, while the bins represents the quantization levels. In Figure 2a, five out of eight quantization levels encode fewer than 1,000 weights each, with one level representing only a single value. Therefore, even when the model is

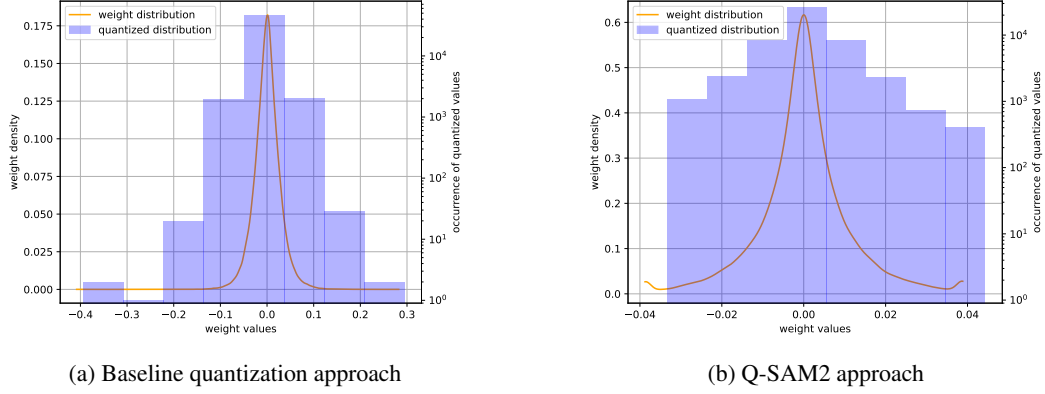


Figure 2: Weight distributions and quantized distributions of linear layer "blocks.0.mlp.layers.1" using a baseline quantization approach and our method, Q-SAM2, in a 3-bit configuration. In (a) outliers extend the dynamic range, causing most weights to be mapped to a small subset of levels. In (b) we compressed the distribution, yielding a more balanced mapping across quantization levels, thus reducing quantization error.

optimized through QAT, the initial performance drop, due to the high quantization error, cannot be easily recovered by finetuning. On the other hand, Figure 2b represents the weight distribution using our method, after calibration and clipping. The quantization error in this case may be very large for the few outliers but clearly low for most of the values around the mean, resulting in more effective utilization of quantization levels and lower overall information loss.

3.2 Linear Layers Calibration for Initializing Quantized SAM2

Calibration is a central argument to PTQ weight optimization. AdaRound [20] learns adaptive rounding offsets to minimize layer-wise reconstruction error without retraining. BrecQ [13] solves a block-wise reconstruction problem using second-order information to align quantized and full-precision outputs. GPTQ [9] minimizes the output reconstruction error using second-order loss approximations for block-wise calibration. MagR [27] smooths outliers by minimizing the L_∞ norm, reducing weight value range.

In contrast, the use of a preprocessing step for weight calibration in QAT has not been thoroughly investigated [14]. In this work, we show that introducing such a step can significantly improve performance, particularly at very low bit-widths (below 4 bits).

Since we apply clipping later, our calibration focuses on reducing weight variance in each linear layer, concentrating values around the mean to improve quantization. This enhances resolution where it matters most, lowering the quantization error while preserving key information. We use a low-impact, closed-form method based on Frobenius norm minimization over a small calibration set, as the model will be trained afterward. This approach is justified by the fact that, when the mean of the weights is small, the Frobenius norm of a matrix closely approximates its variance up to a constant factor, making it a suitable proxy for measuring overall weight magnitude.

In principle, any modification of a fully trained linear layer weight matrix risks degrading the accuracy of the model. A straightforward but crucial observation is that the accuracy degradation in the final model is directly related to the discrepancy between the *output* of the linear layers before and after quantization, rather than the quantization error of the matrix elements themselves. Any changes of \mathbf{W} will therefore need to balance the compression gains due to quantization with the resulting errors at the output of the linear layer.

Based on this insight, we propose a quantization procedure aimed at optimizing this trade-off through an initial calibration step that projects the trained weights onto the subspace spanned by a calibration batch of data samples, which ensures that the layer's outputs are preserved on these known inputs. In particular, we propose to perform this projection by learning the *Moore-Penrose Pseudoinverse* between the input activations and outputs of the linear layer while running a forward pass of the original trained model on the calibration samples. The use of the Moore-Penrose Pseudoinverse is

motivated by its known properties in terms of 1) minimizing the residual with the calibration output activations, and 2) minimizing the Frobenius norm of the obtained weights. In other words, this simultaneously achieves that the dynamic range of the weight distribution is minimized making them easier to quantize, while at the same time minimizing the change of output of the linear layer (at least on the calibration samples).

While the residual between pre- and post-calibration output activations is guaranteed to be minimized only on the calibration samples, the Moore-Penrose Pseudoinverse has additional properties that make it suitable as a preprocessing step before quantization. The fact that it allows us to control the singular values of the calibrated weight matrix (see later the discussion on Tikhonov regularization) implies that the calibration process “cleans” its spectrum by eliminating components that are irrelevant, unstable, or poorly supported by the calibration data, leading to a simpler and more robust solution. As a result, the modified layer maintains essential functionality on the calibration inputs while promoting smoother and potentially more general behavior elsewhere, which is particularly important for robustness under quantization or other perturbations.

Calibration data selection is critical: larger batches improve faithfulness to the original model, while smaller ones reduce memory and runtime costs. Moreover, choosing representative frames is equally important because stacking ill-conditioned inputs in the calibration set can collapse the key directions learned during training, harming generalization. As a result, the layer may perform well on calibration samples, but fail on unseen or diverse data. To prevent this, calibration inputs must be diverse, well-conditioned, and representative of the operational domain of the model.

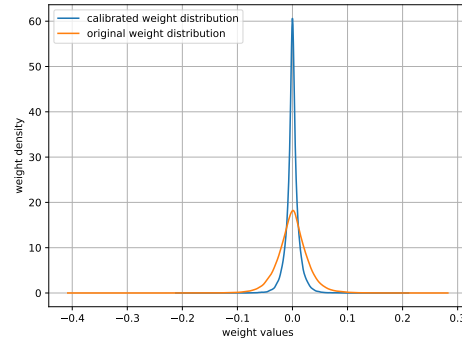


Figure 3: Weight distributions of the linear layer “blocks.0.mlp.layers.1” before and after calibration defined by equation 3.

The condition number of activations reflects calibration input quality: low values indicate stable, rich subspaces. However, stacking inputs requires care, as a few high-condition-number inputs can destabilize the overall matrix and impair projection quality. Ideally, all inputs should be individually stable, but precomputing condition numbers is infeasible in large datasets. We address this by using regularization to remove ill-conditioning, allowing random input selection.

To define our method more formally, let us recall the definition of a linear layer $f(\mathbf{X}, \mathbf{W})$. It maps an input tensor $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}}$, where n is the batch size and d_{in} is the input feature dimension, to an output tensor $\mathbf{Y} \in \mathbb{R}^{n \times d_{\text{out}}}$, where d_{out} is the output feature dimension, through an affine transformation defined as $f(\mathbf{X}, \mathbf{W}) = \mathbf{X}\mathbf{W}^\top + \mathbf{b}$, with $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and $\mathbf{b} \in \mathbb{R}^{d_{\text{out}}}$. We can consider $\mathbf{b} = \mathbf{0}$ for the following analysis, without loss of generality.

Given a batch of data \mathbf{X} , in order to obtain a minimum Frobenius norm weight matrix $\hat{\mathbf{W}}$, we can use the Moore-Penrose pseudoinverse [21]. For $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{Y} \in \mathbb{R}^{n \times p}$, among all matrices $\mathbf{W} \in \mathbb{R}^{p \times d}$ satisfying $\mathbf{Y} = \mathbf{X}\mathbf{W}^\top$, the unique matrix minimizing the Frobenius norm $\|\mathbf{W}\|_F$ is

$$\mathbf{W}^\top = \mathbf{X}^\dagger \mathbf{Y}, \quad (2)$$

where \mathbf{X}^\dagger denotes the Moore–Penrose pseudoinverse of \mathbf{X} .

To mitigate the instability caused by poorly conditioned calibration samples, we apply Tikhonov regularization during the projection. From equation (2) the new formulation is derived:

$$\mathbf{W}^\top = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (3)$$

where \mathbf{I} is the $d \times d$ identity matrix, and $\lambda \in \mathbb{R}^+$ is a cutoff hyper-parameter. To improve generalization, we set λ larger than the smallest singular value, preventing the projection from overfitting unstable directions in the calibration data. In Appendix A.1 we present analysis of our calibration method and the role of λ in equation 3.

Figure 3 shows the weight distribution of a real layer of the image encoder. The image compare the original floating point distribution before (orange) and after (blue) calibration, defined in equation

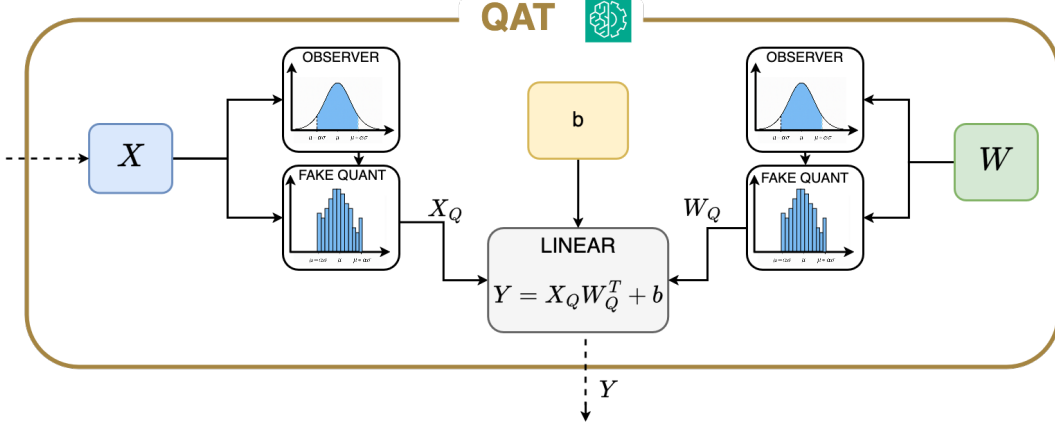


Figure 4: Representation of a linear layer during the forward pass of the quantization aware training procedure. Weight and activation distributions are analyzed by the observer, which will provide the clipping values, the scale, and the zero point to the fake quantization.

3. By reducing the variance of the weight distribution, our method compresses the dynamic range, which in turn reduces the quantization error.

3.3 Quantization-Aware Training for Finetuning Quantized SAM2

The calibrated image encoder is then ready for the quantization step, depicted in Figure 4. During quantization-aware training (QAT), we apply clipping to weights and activations when simulating quantization through the fake quantization operation. The clipping is driven by an observer that, based on the statistics of the corresponding weight or activation, clips the values based on the variance of that distribution. In particular, from equation 1 we can deduce the customized quantization:

$$\mathbf{X}_q = \text{clip} \left(\text{round} \left(\frac{\text{clip}(\mathbf{X}, \mu_{\mathbf{X}} - \alpha \cdot \sigma_{\mathbf{X}}, \mu_{\mathbf{X}} + \alpha \cdot \sigma_{\mathbf{X}})}{s} \right) + z, -2^{b-1}, 2^{b-1} - 1 \right), \quad (4)$$

where $\mu_{\mathbf{X}} \in \mathbb{R}^n$ can be the mean value of \mathbf{X} if using per-tensor quantization ($n = 1$), or the average value of the input channels if using per-channel quantization ($n = d_{\text{in}}$). In the same way $\sigma_{\mathbf{X}} \in \mathbb{R}^n$ is the standard deviation; finally $\alpha \in \mathbb{R}^{>0}$ is the cutoff multiplication factor.

4 Experiments

4.1 Experimental Setup

4.1.1 Task, Dataset, and Metrics

Our experimental setup resembled the one in the original SAM2 work [23]. We calibrated and trained the models on a subset of SA-1B [12] and SA-V [23] datasets. SA-1B dataset is composed by 11M images and more than 1B masks, while SA-V contains 51k videos and 643K spatio-temporal segmentation masks (i.e., masklets). To reduce training time, we used a smaller subset of the dataset and still obtained high accuracy, demonstrating that effective performance can be achieved without full retraining. From the original data, we subsample at random $\sim 8\%$ and $\sim 45\%$, respectively.

We evaluated our models on the semi-supervised Visual Object Segmentation (VOS) task, which involves segmenting a target object across a video sequence given its ground-truth mask in the first frame. Evaluation was conducted on SA-V val, SA-V test (both subsets of SA-V), and the MOSE val from the MOSE dataset [6], using the accuracy metric $J\&F$ [22]. For instance segmentation, we used the COCO2017 validation set [15], reporting mean Intersection over Union (mIoU) over all 5,000 images.

Model	Bit-Width (W/A)	Method	$J\&F$		mIoU	
			SA-V val	SA-V test	MOSE val	COCO2017
B+	FP	-	78.1	78.2	73.7	59.1
	2/8	Baseline Q-SAM2 (ours)	65.6 68.8	66.9 69.0	67.1 67.5	49.0 49.6
	2/4	Baseline Q-SAM2 (ours)	52.3 62.2	56.8 65.0	55.3 62.3	39.0 44.7
S	FP	-	77.7	76.6	73.5	59.7
	2/8	Baseline Q-SAM2 (ours)	65.4 66.3	66.9 68.7	67.3 68.9	47.0 47.7
	2/4	Baseline Q-SAM2 (ours)	56.8 58.7	59.9 61.5	59.1 60.9	41.2 42.1
T	FP	-	75.0	75.0	70.9	57.9
	2/8	Baseline Q-SAM2 (ours)	64.0 63.8	66.7 67.7	67.8 66.5	45.5 46.4
	2/4	Baseline Q-SAM2 (ours)	53.5 55.3	57.1 59.4	58.8 60.7	40.3 40.7

Table 1: Evaluation of the semi-supervised VOS and instance segmentation tasks. For each SAM2 model, we presented the reference full precision values and two bit-width configurations. Our solution is then compared with the MinMax [11] QAT algorithm as the baseline.

4.1.2 Implementation

Our starting models are the checkpoints available from the SAM2 repository [19], named *sam2.1_hiera_tiny(T)*, *sam2.1_hiera_small(S)*, *sam2.1_hiera_base_plus(B+)*. The calibration was then performed with 50 input images, chosen uniformly at random from the SA-1B dataset. The images are forwarded through the original image encoder, and the activations are recorded. The activations are then stacked and the calibrated weight matrices are obtained by applying equation 3. To assess the performance of our method, we compared our pipeline with respect to the classical MinMax method, described in [11]. For a fair comparison, the training procedure is the same for each experiment. The implementation of the QAT pipelines is based on PyTorch 2.6.0 and Cuda 12.4, and all experiments are performed on the NVIDIA A100 80GB GPU. The details of our training pipeline, and additional training ablations for our setup are presented in the Appendix A.2, A.3.

4.2 Ultra-Low-Bit Quantization Performance

Table 1 summarizes the results of our experiments. The table shows that our method outperforms the classical MinMax method in both semi-supervised VOS and instance segmentation for all experiments considered for *B+* and *S* models. Furthermore, under the W2A4 configuration, the improvement is even more pronounced, with the *B+* model achieving up to a 18.9% gain over the baseline. Figure 5 shows qualitative results using the W2A4 configuration for the *B+* model. Our method produces more precise masks than the baseline, closely approaching the output quality of the full-precision model. For the tiny model (*T*), the performance differences are less pronounced, especially in the W2A8 configuration. This is likely due to its smaller number of parameters, which leads to fewer outliers compared to larger models. In addition, the activations are better conditioned, reducing the

Method	Bit-Width (W)	$J\&F$			mIoU
		SA-V val	SA-V test	MOSE val	COCO2017
FP	-	78.1	78.2	73.7	59.1
Baseline	3	34.4	33.6	29.9	11.8
	4	69.2	70.2	71.4	53.4
Baseline + Calibration	3	36.6	35.8	33.1	13.4
	4	70.9	70.0	71.8	54.2
HQQ	3	38.1	39.7	37.3	16.6
	4	71.0	72.5	72.3	55.0
HQQ + Calibration	3	41.2	42.2	42.5	27.6
	4	71.0	72.6	73.1	54.8

Table 2: Evaluation of Q-SAM2 calibration under PTQ with weight-only quantization using B+ model.

overall influence of calibration. This highlights a key limitation: when the bit-width is sufficient to represent weights and activations, or in smaller networks with fewer outliers, our method performs similarly or slightly worse than the baseline due to calibration-induced errors and information loss from clipping. In Appendix A.4 other qualitative results (both static images and video) are presented.

4.3 Calibration in Post Training Quantization

We evaluated our calibration procedure on two PTQ methods with weight-only quantization at 3 and 4 bits using B+ model: a standard MinMax quantizer, and HQQ [2], a Hessian-aware method popular on Hugging Face for efficient deployment. This evaluation investigates the effectiveness of our calibration method in the PTQ setting. As reported in Table 2, calibration consistently improves performance across both quantization strategies. In the 3-bit configuration, the improvement is particularly significant, with HQQ achieving over 66% mIoU. In the 4-bit setting, our method still outperforms the non-calibrated baseline in most cases, although the performance gap is less pronounced.

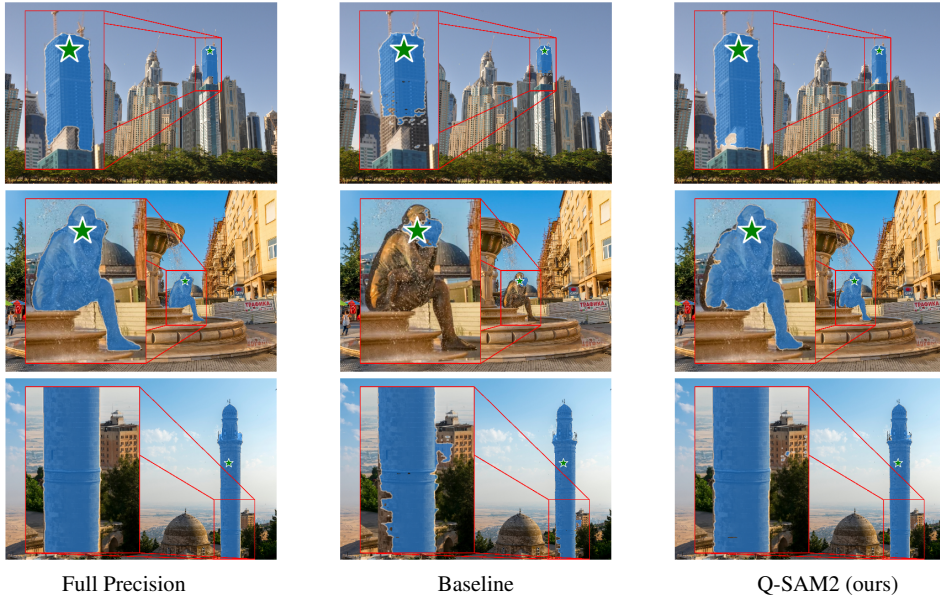


Figure 5: Qualitative results on W2A4 configuration for B+ model for instance segmentation.

5 Conclusions

We presented the first quantized Segment Anything Model 2, Q-SAM2. Using a novel calibration procedure, we were able to reduce the variance of the weights distribution in the linear layer of the image encoder. A quantization aware training procedure was then developed achieving ultra-low 2 bit quantization, using clipping and two proposed losses. Our results significantly surpassed the classical MinMax method for the semi-supervised VOS and instance segmentation tasks, especially when the compression is very high. While QAT generally succeeds in recovering accuracy after quantization, our results highlight the importance of a proper weight pre-processing step, especially under aggressive compression settings. Moreover, our calibration method demonstrated strong generalization even with minimal data, and its effectiveness across different architectures presents a promising direction for future research.

Acknowledgement

We acknowledge funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie actions HORIZON-MSCA-2022-DN-01 call (Grant agreement ID: 101119554).

References

- [1] Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. *arXiv preprint arXiv:2404.00456*, 2024.
- [2] Hicham Badri and Appu Shaji. Half-quadratic quantization of large machine learning models, November 2023.
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013.
- [4] Y. Choi, M. El-Khamy, and J. Lee. Pact: Parameterized clipping activation for quantized neural networks. In *arXiv preprint arXiv:1805.06085*, 2018.
- [5] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [6] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, Philip HS Torr, and Song Bai. MOSE: A new dataset for video object segmentation in complex scenes. In *ICCV*, 2023.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [9] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: Accurate post-training compression for generative pretrained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [10] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-power computer vision*, pages 291–326. Chapman and Hall/CRC, 2022.
- [11] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2018.

- [12] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
- [13] Jianfei Li, Shupeng Zhang, Kaiqi Liu, and Zhen Han. Brecq: Pushing the limit of post-training quantization by block reconstruction. *arXiv preprint arXiv:2102.05426*, 2021.
- [14] Yuhang Li, Mingzhu Shen, Jian Ma, Yan Ren, Mingxin Zhao, Qi Zhang, Ruihao Gong, Fengwei Yu, and Junjie Yan. Mqbench: Towards reproducible and deployable model quantization benchmark, 2022.
- [15] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.
- [16] Xiaoyu Liu, Xin Ding, Lei Yu, Yuanyuan Xi, Wei Li, Zhijun Tu, Jie Hu, Hanting Chen, Baoqun Yin, and Zhiwei Xiong. Pq-sam: Post-training quantization for segment anything model. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol, editors, *Computer Vision – ECCV 2024*, pages 420–437, Cham, 2025. Springer Nature Switzerland.
- [17] Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinqant: Llm quantization with learned rotations, 2025.
- [18] Chengtao Lv, Hong Chen, Jinyang Guo, Jinyang Guo, Jinyang Guo, Yifu Ding, Xianglong Liu, Xianglong Liu, and Xianglong Liu. Ptq4sam: Post-training quantization for segment anything. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15941–15951, 2024.
- [19] Meta AI Research. SAM2 github repository. <https://github.com/facebookresearch/sam2>, 2024. Accessed June 12, 2025.
- [20] Markus Nagel, Mart van Baalen, Titus Blankevoort, and Max Welling. Up or down? adaptive rounding for post-training quantization. In *International Conference on Machine Learning (ICML)*, pages 7197–7206, 2020.
- [21] R. Penrose. A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society*, 51(3):406–413, 1955.
- [22] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 DAVIS challenge on video object segmentation. *CoRR*, abs/1704.00675, 2017.
- [23] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [24] Hao Wu, Patrick Judd, Xiaojie Zhang, Mikhail Isaev, and Paulius Micikevicius. Integer quantization for deep learning inference: Principles and empirical evaluation. *arXiv preprint arXiv:2004.09602*, 2020.
- [25] Junyi Wu, Haoxuan Wang, Yuzhang Shang, Mubarak Shah, and Yan Yan. Ptq4dit: Post-training quantization for diffusion transformers. In *NeurIPS*, 2024.
- [26] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. SmoothQuant: Accurate and efficient post-training quantization for large language models. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- [27] Aozhong Zhang, Naigang Wang, Yanxia Deng, Xin Li, Zi Yang, and Penghang Yin. Magr: Weight magnitude reduction for enhancing post-training quantization. *arXiv preprint arXiv:2406.00800*, 2024.

- [28] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, and He Wen. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. In *arXiv preprint arXiv:1606.06160*, 2016.

A Technical Appendices and Supplementary Material

This appendix provides more details to support the findings presented in the main paper. In Section A.1, we analyze the bias-variance trade-off introduced by the regularization hyperparameter λ , Section A.3 outlines our training configuration, Section A.2.1 presents training loss curves, Section A.2.2 discusses the differences between our training procedure and the original SAM2 model [23], and Section A.3 presents additional training ablations for our setup. Finally, Section A.4 includes qualitative segmentation results to illustrate the effectiveness of our approach.

A.1 Calibration

To better analyze the behavior of our calibration method introduced in Section 3.2, we implemented Equation 3 within the linear layer *cut.blocks.21.mlp.layers.1* of the image encoder in the B+ model. Calibration was performed by applying the regularized pseudoinverse formulation to batches of 50 images while systematically varying the regularization hyperparameter λ . For each value of λ , we computed a calibrated weight matrix $\hat{\mathbf{W}}_\lambda$.

To evaluate the effect of calibration, we compared the outputs of the original and calibrated layers using a separate evaluation batch of 10 images, denoted as \mathbf{X}_{eval} . The output discrepancy was quantified using the L_2 norm between the original projection $\mathbf{X}_{\text{eval}}\mathbf{W}^\top$ and the calibrated output $\mathbf{X}_{\text{eval}}\hat{\mathbf{W}}_\lambda^\top$. This metric provides a direct measure of the reconstruction error induced by the calibration process.

Figure 6a illustrates how the L_2 error varies with λ , revealing a clear bias-variance trade-off. For small regularization values ($\lambda < 10^{-3}$), the error is high and exhibits large variance across different calibration batches, indicating sensitivity to batch selection due to ill-conditioned activations. This suggests that when regularization is too weak, the inversion step amplifies noise, degrading output fidelity. On the opposite end, overly strong regularization ($\lambda > 1$) significantly alters $\hat{\mathbf{W}}_\lambda$, reducing its ability to generalize and match the behavior of the original layer. In the intermediate range ($10^{-3} \leq \lambda \leq 1$), the reconstruction error is minimized, indicating an optimal trade-off that stabilizes calibration without overfitting to the batch.

To further assess the impact of calibration on quantization, we analyzed the quantization error of $\hat{\mathbf{W}}_\lambda$, as shown in Figure 6b. The error was computed as $\|\hat{\mathbf{W}}_\lambda - \hat{\mathbf{W}}_{\lambda_Q}\|$, where $\hat{\mathbf{W}}_{\lambda_Q}$ is the quantized version of the calibrated weight matrix using uniform quantization. The orange line in the figure represents the quantization error of the original weight matrix \mathbf{W} . We observe that counterintuitively for low λ values, the quantization error is extremely high due to the poor conditioning of $\hat{\mathbf{W}}_\lambda$, which leads to numerical instability and large weight magnitudes. As λ increases, particularly beyond 5×10^{-2} , the quantization error of the calibrated weights becomes consistently lower than that of the original weights. This suggests that regularization not only improves numerical stability but also produces weight matrices that are better structured for quantization, exhibiting smaller dynamic ranges.

The selection of λ is not straightforward. One option is to choose the value that minimizes the reconstruction error in Figure 6a, which provides a good match to the original layer outputs. However, this choice does not necessarily lead to the lowest quantization error, and its impact on initialization quality for quantized training may be limited. Conversely, setting $\lambda \approx 1$ improves the quantization behavior but can result in higher L_2 output error, compromising the layer’s ability to generalize. This highlights the need to balance fidelity to the original outputs with improved robustness under quantization when selecting the regularization strength.

To select an appropriate value of λ in practice, we adopt a rule based on the singular value spectrum of the input matrix. If the activation matrix is ill-conditioned (i.e., with a condition number greater than 10^2) but not effectively rank deficient, we set λ to five times the smallest non-zero singular value. In the case of low-precision formats such as `bf16`, matrices with condition numbers exceeding 10^3 can already exhibit numerical instability. When the matrix is more severely ill-conditioned, with a condition number significantly greater than 10^3 , it is likely to be effectively rank deficient. In this case, we identify the smallest singular value σ_* such that $\sigma_{\max}/\sigma_* \leq 10^3$, and set λ to five times σ_* .

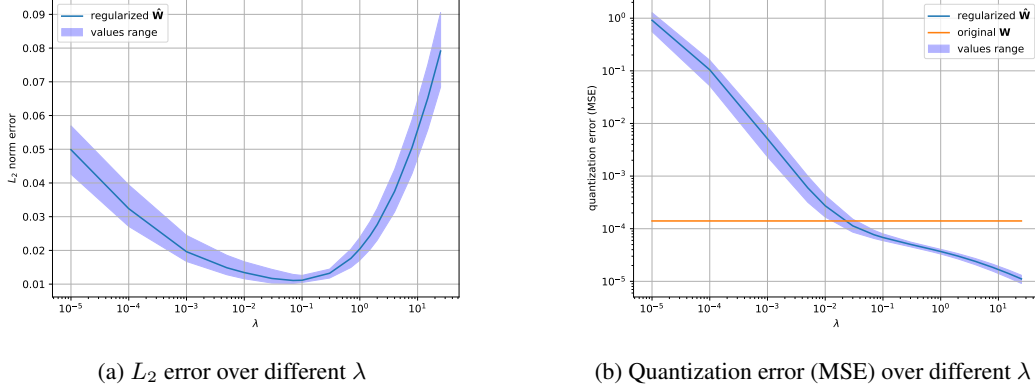


Figure 6: Calibration error for the linear layer *cut.blocks.21.mlp.layers.1* from the image encoder of the B+ model. In subfigure (a), we compare the output of the original linear layer with the output produced using the calibrated weights $\hat{\mathbf{W}}_\lambda$, measured via the L_2 norm, across different values of the hyperparameter λ . In subfigure (b), we report the quantization error of $\hat{\mathbf{W}}_\lambda$ for varying λ , relative to the quantization error of the original weights \mathbf{W} .

As a general rule, depending on the conditioning characteristics of the specific network and input data, λ is chosen between two and five times the selected singular value. This strategy stabilizes the pseudoinverse computation while preserving the most meaningful directions in the activation space.

A.2 Training

To ensure fairness, we used a consistent training and evaluation setup for both our Q-SAM2 and the baseline approach [11]. The training procedure follows the one proposed in the SAM2 [23] paper where the authors referred as the "full training" step. The details are reported in Table 3 where we highlight the configurations and hyperparameters that are different from the original work. All trained models use the same image and video data, with a fixed random seed to ensure reproducibility.

Compared to the original SAM2 training setup, we introduced several modifications to accommodate hardware constraints and reduce training time. First, we reduced the maximum number of masks per frame from 64 to 60 to fit within GPU memory limits on an NVIDIA A100 80GB. We also removed the warm-up phase, reduced the amount of input data, and limited training to a single epoch. Ablation studies on dataset size and number of epochs are presented in Section A.3.

All training experiments were conducted using 8 NVIDIA A100 80GB GPUs, with a total training time of less than 35 hours for the B+ architecture.

Furthermore, we adjusted the learning rate from 4×10^{-5} to 1×10^{-5} and decreased the effect of the gradient scaler by halving its initial scale value.

A.2.1 Losses Plots

We adopt the standard loss of SAM2 [23] as $\mathcal{L}_{\text{SAM2}}$, which is composed of a linear combination of focal and dice losses for mask prediction, mean absolute error (MAE) loss for IoU prediction, and cross-entropy loss for object prediction. The weights are set in a fixed ratio of 20 : 1 : 1 : 1, respectively.

Figures 7 and 8 present the loss curves for the W2A4 and W2A8 quantization settings. Specifically, Figure 7 shows the results for the B+ model and one configuration of the S model (W2A4), while Figure 8 includes the second S configuration (W2A8) together with both configurations for the T model.

In nearly all settings and across all loss components, Q-SAM2 demonstrates lower training losses compared to the MinMax baseline. This consistently improved performance indicates that our calibration-based initialization plays a critical role in improving the subsequent QAT phase. In

Configuration	Value
data	SA-1B(8%) SA-V(45%)
resolution	1024
precision	bfloat16
epochs	1
optimizer	AdamW
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.999$
gradient	clipping type: L_2 , max: 0.1
weight decay	0.1
learning rate (lr)	img. enc.: $1e-5$, other: $3.0e-4$
lr schedule	cosine
warmup	no warmup
layer-wise decay	0.8 (T, S), 0.9 (B+)
image augmentation	hflip, resize to 1024 (square)
video augmentation	hflip, affine (deg: 25, shear: 20), colorjitter (b: 0.1, c: 0.03, s: 0.03, h: null), grayscale (0.05), per frame colorjitter (b: 0.1, c: 0.05, s: 0.05, h: null), mosaic-2x2 (0.1)
batch size	256
drop path	0.1 (T, S), 0.2 (B+)
mask losses (weight)	focal (20), dice (1)
IoU loss (weight)	L_1 (1)
occlusion loss (weight)	cross-entropy (1)
max. masks per frame	image: 60, video: 3
# correction points	7
global attn. blocks	5-7-9 (T), 7-10-13 (S), 12-16-20 (B+)

Table 3: Hyperparameters and details of our Q-SAM2 and MinMax baseline [11] QAT training for the three image encoder sizes B+,S,T.

addition, our model converges faster as the difference in loss values, particularly for the dice and IoU components, increases progressively with each training step.

A.2.2 Comparison with SAM2 training pipeline

The Segment Anything Model 2 (SAM2) is trained in two main stages followed by a fine-tuning procedure. The first stage, pre-training, involves training on static images using the SA-1B dataset. The second stage, called full training, uses a combination of images from SA-1B, videos from SA-V, and additional internal data that is not publicly available. Finally, to improve segmentation performance on long video sequences, the authors fine-tuned the model by increasing the number of input video sequences from 8 to 16. To accommodate this within GPU memory constraints, the image encoder is frozen during this step.

In our work, we reproduced the second stage, called "full training", focusing on a full architecture training with static image and video data, while reducing both the data volume and the number of training steps with respect to the original paper. As demonstrated in Section A.2.1, our method learns faster than the baseline and achieves accurate results with fewer data, thus reducing training time. However, a direct comparison with the original SAM2 model remains difficult due to the differences in training procedures. We plan to explore a fully aligned comparison with the complete SAM2 pipeline in future work.

A.3 Training data ablations

To justify the use of a single training epoch in our setup, we conducted a set of experiments varying both the number of epochs and the dataset size. Table 4 summarizes the results. We evaluated three dataset sizes: 10%, 50%, and 100%, where the percentages refer to the subset described in Section 4. In particular, training for ten epochs on 10% of the data produces performance comparable to training for one epoch on 50% of the data, while requiring only half the training time in the latter case.

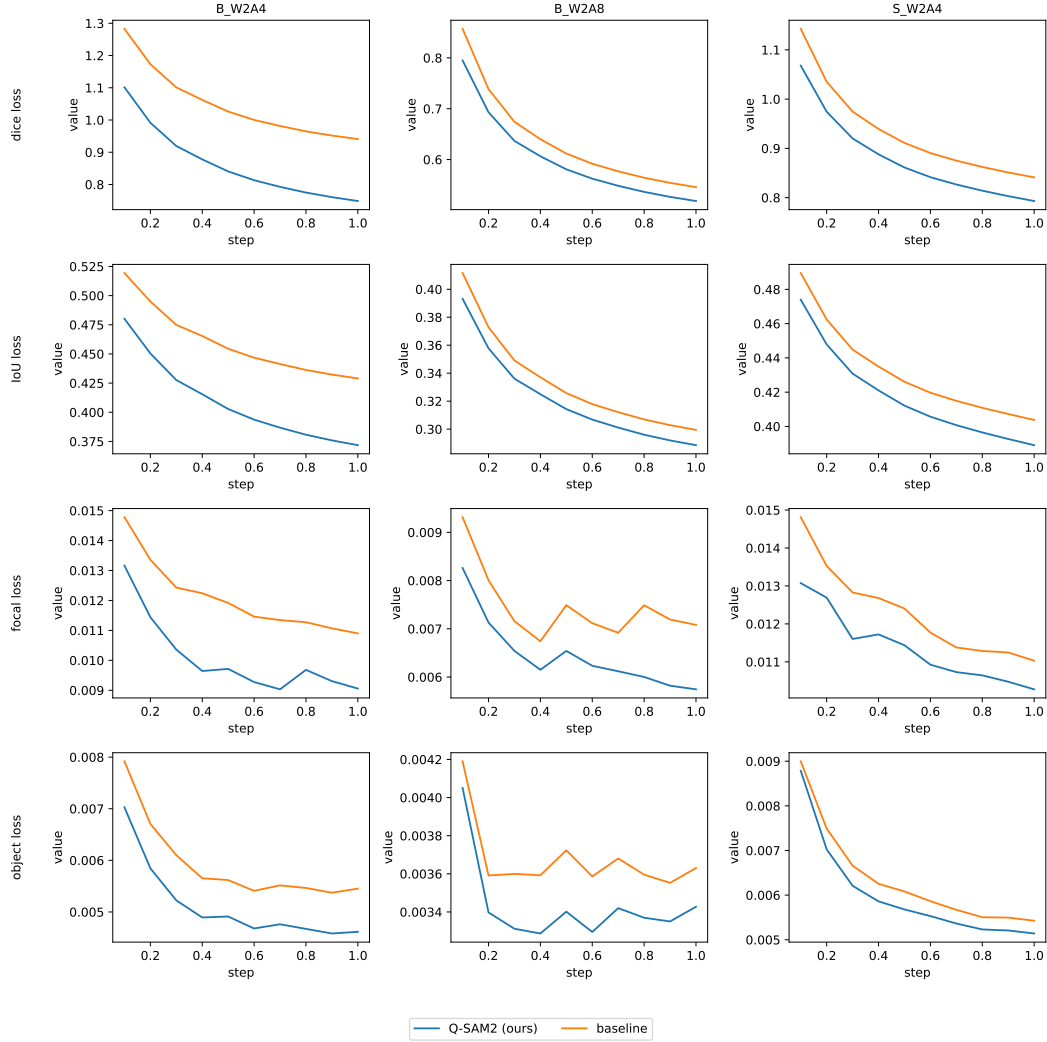


Figure 7: Training loss curves for image encoder size B+,S under two configurations 2W4A and 2W8A. Each sub-plot contains a comparison between Q-SAM2 and the baseline for the four losses: dice, IoU, focal, and occlusion. Our solution outperforms the baseline in all losses for all configurations.

Data ratio [%]	Normalized training time	Epochs	SA-V test $J\&F$
100	1.00	1	69.0
50	0.50	1	66.3
10	0.10	1	65.7
10	1.00	10	66.2

Table 4: Ablation study on the number of training epochs and dataset size. Results show that training for one epoch on 50% of the data achieves comparable performance to training for ten epochs on 10%, while significantly reducing total training time.

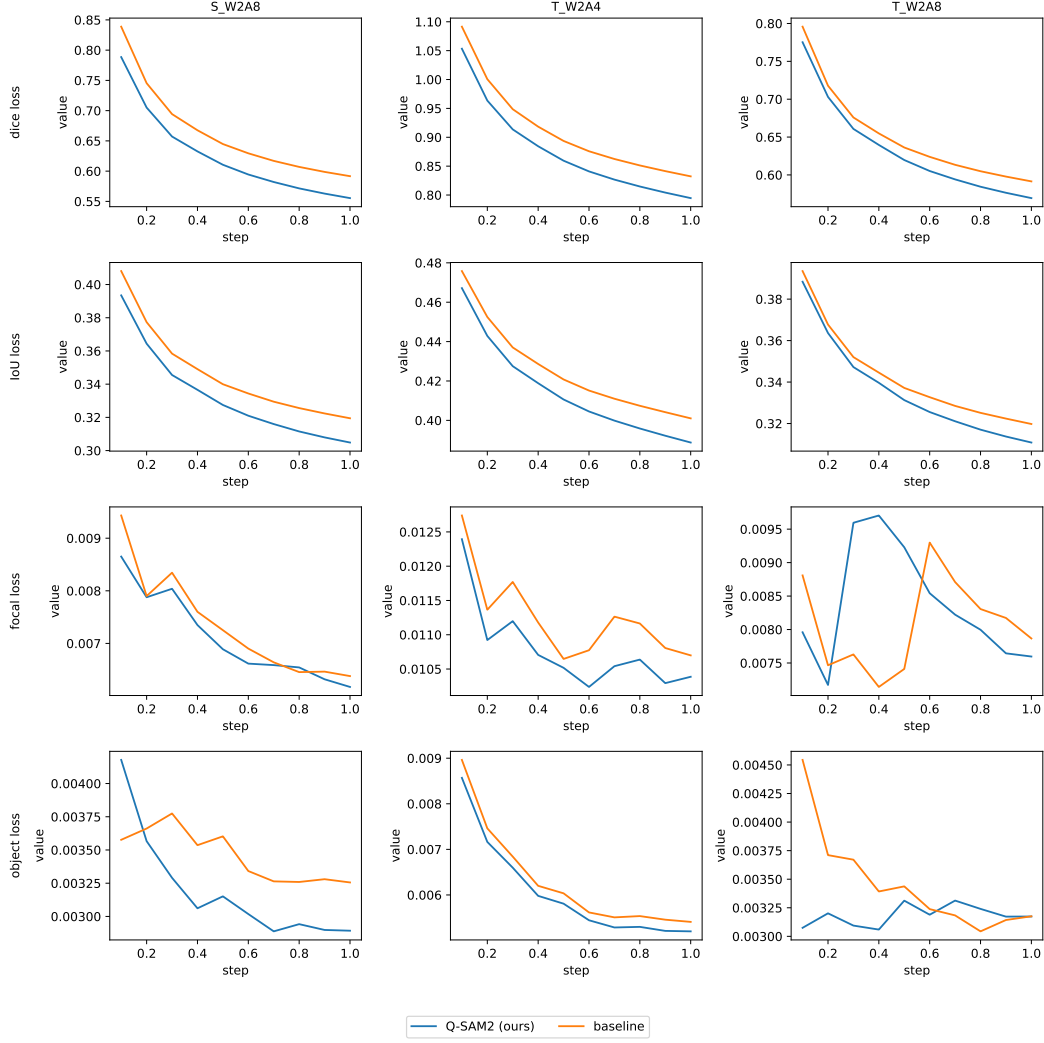


Figure 8: Continuation of the training loss curves from Figure 7 for image encoder size S,T under two configurations 2W4A and 2W8A. Each sub-plot contains a comparison between Q-SAM2 and the baseline for the four losses: dice, IoU, focal, and occlusion. Our solution outperforms the baseline in all losses for almost all configurations.

A.4 Qualitative Results

We have extensively benchmarked the original variants of SAM2, the MinMax baseline, and our algorithm Q-SAM2. The variation includes different architectures, prompts, and quantization configurations for a fair comparison in real-world scenarios. Figure 9, Figure 10, and Figure 11 show the impact of adding more input points as prompt. We observe that our algorithm provides cleaner masks compared to the MinMax algorithm with the same quantization scheme. We also note, that by providing one prompt point only, typically all variants of the algorithms tend to provide a clean mask. However, and that is a well-studied behavior of the original SAM model, this mask might be only a part of the user’s intended selection, such as the front of the racing car as in Figure 10, or a single garlic bulb as in Figure 11. Henceforth, interesting examples include those that require multiple prompts to select the entire object. Those 3-point or 5-point prompts trigger differences among the quantized models. Compared against MinMax our 2-bit models require fewer prompt points to recover a larger fraction of the entire object.

Figure 12, Figure 13, and Figure 14 show the entire family of models including SAM2 base plus, small, and tiny architectures at full precision, as well as all variants of MinMax and our Q-SAM2



Figure 9: Qualitative results of SAM2, MinMax, and QSAM2(ours) on the base plus architecture. From left to right, we add 1, 3, and 5 prompt points.

models with 2-bit quantization for weights and 4 and 8 bits for activations. More bits help to improve the results. The quantization has a strong impact on the results, for example, the MinMax baseline with the W2A4 quantization fails to segment the entire sheep by using only one prompt. Q-SAM2 demonstrates a similar behavior for the S and T architecture, however, for the B+ architecture our low W2A4 quantization operates comparable to full precision results.

Finally, Figures 15 and 16 present a qualitative comparison on a video from the SA-V test set for the semi-supervised Visual Object Segmentation (VOS) task. The models included in this comparison are SAM2.1, the MinMax baseline, and our proposed Q-SAM2. The first column of Figure 15 shows the initial ground-truth mask used to prompt the models, as required in the semi-supervised VOS setup. Although all models generally succeed in segmenting the target object, Q-SAM2 produces cleaner masks with fewer border artifacts compared to the baseline.



Figure 10: Qualitative results of SAM2, MinMax, and QSAM2(ours) on the base plus architecture. From left to right, we add 1, 3, and 5 prompt points.

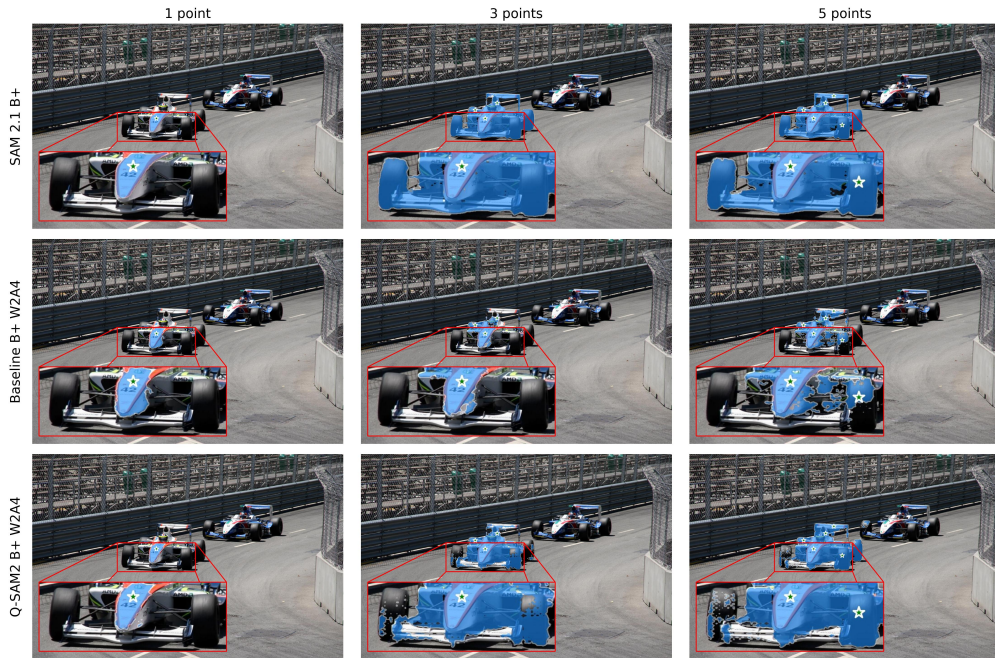


Figure 11: Qualitative results of SAM2, MinMax, and QSAM2(ours) on the base plus architecture. From left to right, we add 1, 3, and 5 prompt points.

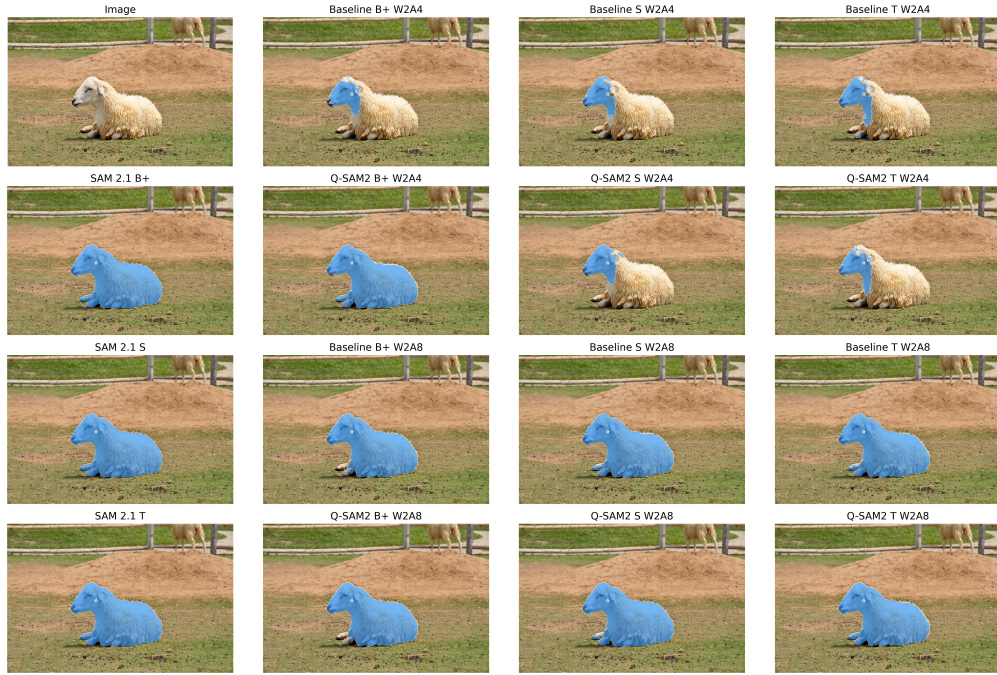


Figure 12: Qualitative results of variants of models.



Figure 13: Qualitative results of variants of models.

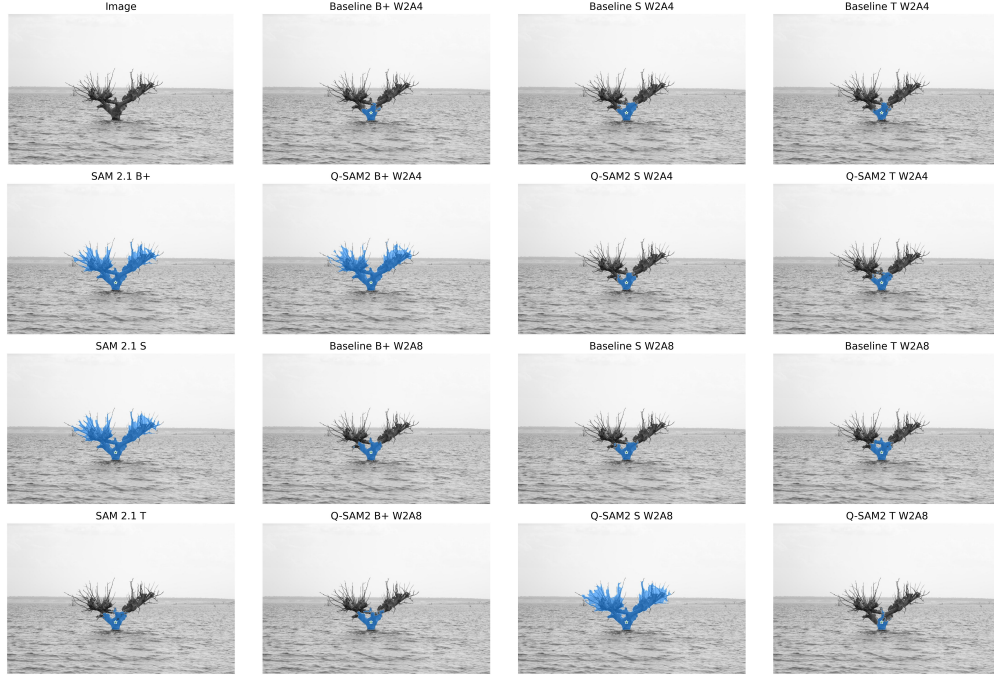


Figure 14: Qualitative results of variants of models.

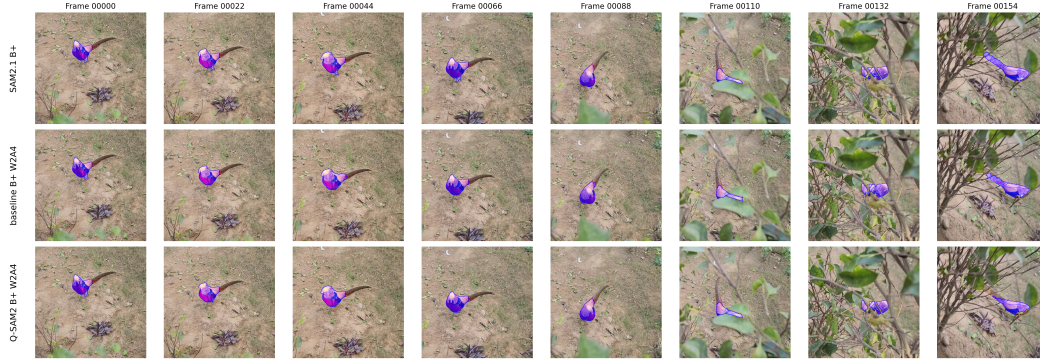


Figure 15: Qualitative results of variants of models SAM2, MinMax, and QSAM2(ours) on the base plus architecture for a video sequence. The first column shows the input mask provided in the first frame, while the remaining columns display the predicted segmentations for each model across subsequent frames.



Figure 16: Continuation of qualitative results from Figure 15, showing additional frames from the same video sequence.