

MATEY: MULTISCALE ADAPTIVE FOUNDATION MODELS FOR SPATIOTEMPORAL PHYSICAL SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Accurate representation of the multiscale features in spatiotemporal physical systems using vision transformer (ViT) architectures requires extremely long, computationally prohibitive token sequences. To address this issue, we propose an adaptive tokenization scheme which dynamically adjusts the token sizes based on local features. Moreover, we present a set of spatiotemporal attention schemes, where the temporal or axial spatial dimensions are decoupled, and evaluate their computational and data efficiencies. We assess the performance of the proposed multiscale adaptive model, MATEY, in a sequence of experiments. The results show that adaptive tokenization achieves improved accuracy without significantly increasing token sequence length, but the improvement deteriorates in more complex data configurations. Compared to a full spatiotemporal attention scheme or a scheme that decouples only the temporal dimension, we find that fully decoupled axial attention is less efficient and expressive, requiring more training time and model weights to achieve the same accuracy. Finally, we demonstrate in two fine-tuning tasks featuring different physics that models pretrained on PDEBench data outperform the ones trained from scratch, especially in the low data regime with frozen attention.

1 INTRODUCTION

Developing foundation models for physical systems is vital for energy generation, earth sciences, and power and propulsion systems. These models offer faster solutions than physics-based simulations and can generalize better across multiple systems than single-purpose AI approaches. However, their application to physical systems, often characterized by multiple sub-processes at different scales, is still in the early stages. For instance, fluid flowing around a cylinder creates a von Kármán vortex street, a highly dynamic flow with rapidly evolving vortices. Accurate solutions of such multiscale systems requires a very high resolution representation to capture the most complex features across space and time. However, for scientific machine learning as for modeling and simulation, using very high resolutions to achieve accurate solutions incurs significant computational cost. This is particularly true for developing foundation models using vision transformer (ViT)-based architectures, as using the standard self-attention mechanism for extremely long spatiotemporal sequences can become prohibitively computationally expensive.

Efficient representation of multiscale features in high-resolution inputs has been an active research topic in computer vision. Three broad approaches can be characterized. First, multiscale models like Swin Transformer (Liu et al., 2021) and MViTv2 (Li et al., 2022) introduce multiple stages with decreasing resolution and increasing feature dimension for efficient hierarchical representations. Second, computational techniques have been developed which facilitate training on long sequences (e.g., sequence parallelism across GPUs (Jacobs et al., 2023)) or reduce the effective sequence length in the attention kernel (e.g., decomposing attention along axial directions (Ho et al., 2019)). Third, the actual sequence length can be directly shortened by pruning and merging tokens ((Haurum et al., 2023; Meng et al., 2022; Yin et al., 2022; Bolya & Hoffman, 2023)), though this strategy may lead to critical information loss (Liu et al., 2024).

These techniques have recently been adopted in scientific machine learning (sciML) for physical systems. For example, the atmosphere foundation model Aurora (Bodnar et al., 2024) uses Swin Transformer, while axial attention is applied by MPP (McCabe et al., 2023). Despite the progress,

054 computational constraints remain a bottleneck, as existing approaches do not yet handle high-fidelity
 055 solutions of applications such as computational fluid dynamics, in which input sequences can eas-
 056 ily exceed billions of tokens. More efficient algorithms are needed to enable the development of
 057 foundation models for multiscale multiphysics systems.

058 In this work, we develop a multiscale adaptive foundation model, MATEY (see Figure 1), that pro-
 059 vides two key algorithmic contributions to address the challenges posed by spatiotemporal physical
 060 systems. First, inspired by the adaptive mesh refinement (AMR) technique, we introduce an adaptive
 061 tokenization method that dynamically adjusts patch sizes across the system based on local features,
 062 which provides as much as a $2\times$ reduction in compute for similar or higher accuracy. Second,
 063 we present a set of spatiotemporal attention schemes based on the axial attention (Ho et al., 2019)
 064 that differ in their decomposition of long spatiotemporal sequences and identify the cost in time-to-
 065 accuracy for axial attention. Finally, we assess the fine-tuning performance of models pretrained on
 066 PDEBench (Takamoto et al., 2022) in two highly out-of-distribution settings, colliding thermals and
 067 magnetohydrodynamics (MHD), that include additional physical variables not included in pretrain-
 068 ing and observe the pretrained models outperforming random initialized models.

070 2 RELATED WORK

072 **Scientific foundation models.** Several research directions have been explored for building founda-
 073 tion models for physical systems, including multiple physics pretraining (McCabe et al., 2023)
 074 with PDEBench data, input augmentation with PDE system configurations (Hang et al., 2024), ro-
 075 bust pretraining schemes (Hao et al., 2024), fine-tuning effectiveness investigations (Subramanian
 076 et al., 2024), and data-efficient multiscale ViT architectures (Herde et al., 2024). While these work
 077 made remarkable progress, they do not directly address the issue of token sequence length, which
 078 becomes a computation bottleneck when applying ViTs to high dimension or high resolution data.

080 **Multiscale ViTs.** While most multiscale ViTs achieve hierarchical representations via multi-stage
 081 attention blocks at different resolutions, e.g., MViTv2 (Li et al., 2022) and Swin Transformer (Liu
 082 et al., 2021), there are a few focusing on tokenization schemes, e.g., (Yin et al., 2022; Fan et al.,
 083 2024; Zhang et al., 2024; Havtorn et al., 2023). Among these, the single-stage MSViT with dynamic
 084 mixed-scale tokenization (Havtorn et al., 2023), which leverages a learnable gating neural network
 085 for token refining, is most related to our work. This approach requires a tailored gate loss func-
 086 tion and an adaptive trimming scheme to handle the high overhead cost, which in return hurts gate
 087 training accuracy. In contrast, the tokenization scheme in MATEY adaptively adjusts the patch sizes
 088 directly based on local feature scales, which is simpler and more direct.

089 **Axial attentions.** The quadratic scaling nature of attention makes it computationally prohibitive
 090 for extremely long token sequences from multidimensional systems. To address this challenge, (Ho
 091 et al., 2019) proposed the axial attention, which decomposes the full attention into a sequence of
 092 attention operations along each axis. It reduces the attention cost from $\mathcal{O}(N^{2d})$ to $\mathcal{O}(N^{d+1})$, for a
 093 given d -dimensional system with N^d tokens. ViViT (Arnab et al., 2021) factorized the spatiotem-
 094 poral attention into spatial- and temporal-dimensions for video classification. (McCabe et al., 2023)
 095 applied the axial attention in the Axial ViT (AViT) for spatiotemporal solutions of physical sys-
 096 tems. While these spatiotemporal attention schemes can reduce the sequence length and hence the
 097 attention cost, their impact on accuracy in physical systems is unclear.

099 3 MATEY, EXPLAINED

101 We propose multiscale adaptive foundation models, MATEY, to predict two-dimensional spatiotem-
 102 poral solutions of multiple physical systems. The architecture of MATEY is illustrated in Fig-
 103 ure 1. Given a sequence of T past solutions of some physical system at time t , MATEY predicts
 104 the solution at a future time $t + t_{\text{lead}}$ by learning from sequences of solutions for multiple physical
 105 systems. Specifically, MATEY learns a model \mathbf{f}_w such that $\mathbf{u}_{t+t_{\text{lead}}} \approx \mathbf{f}_w(\mathbf{u}_{t-T+1}, \dots, \mathbf{u}_t; t_{\text{lead}})$
 106 by training parameters w to minimize the loss of the prediction from the solution sequence
 107 $\mathbf{U} = [\mathbf{u}_{t-T+1}, \dots, \mathbf{u}_t]$ against the future solution with a lead time $\mathbf{u}_{t+t_{\text{lead}}}$. In the following para-
 graphs, we give detailed descriptions for each component in MATEY.

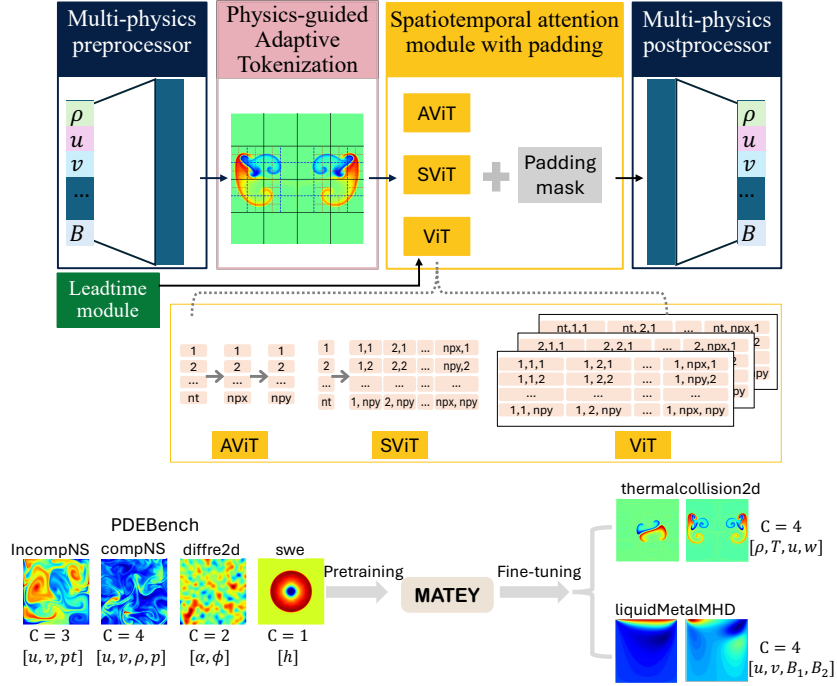


Figure 1: MATEY: multiscale adaptive foundation models for spatiotemporal physical systems.

Multi-physics preprocessor, postprocessor, and training. To accommodate multiple physical systems with different sets of variables at different spatial resolutions, we adopt the multi-physics preprocessor and postprocessor used in the MPP work (McCabe et al., 2023). For system k with C_k variables, the preprocessor first embeds solutions $\mathbf{u}_t(x, y) \in \mathbb{R}^{C_k}$ to a latent space $\mathbb{R}^{C_{\text{uni}}}$, where $C_{\text{uni}} \gg C_k$ is shared among all systems. The resulting embedded solution passes through a convolutional block in the tokenization module and is converted into patch sequences in $\mathbb{R}^{C_{\text{emb}}}$, which are further passed to the attention block and then mapped back to \mathbb{R}^{C_k} by the postprocessor, to predict the solution fields of the C_k variables. To handle solutions with different resolutions, we follow the approach in MPP by performing system-based sampling in the training process and fusing information from samples from different systems via gradient accumulation. We employ a convolutional neural network (CNN) in the tokenization module and 2D transposed convolutional blocks in the multi-physics postprocessor to convert between patch sequences and spatiotemporal solution fields. The CNN block performs the conversion of the embedded solutions in the latent space $\mathbb{R}^{C_{\text{uni}}}$ to the patch sequences in $\mathbb{R}^{C_{\text{emb}}}$, while the transposed convolutional module converts the sequences back to the solution fields. Specifically, the preprocessor embeds the solution of system k , $\mathbf{U}_k \in \mathbb{R}^{T \times H \times W \times C_k}$ into the unified latent representation $\mathbf{U} \in \mathbb{R}^{T \times H \times W \times C_{\text{uni}}}$, which is tokenized into sequences $\mathbf{Z}_p^0 \in \mathbb{R}^{nt \times np_x \times np_y \times C_{\text{emb}}}$, where $nt = T/p_t$, $np_x = H/p_x$, and $np_y = W/p_y$ with prescribed patch size $[p_t, p_x, p_y]$ in the temporal and spatial dimensions. On the other hand, the postprocessor decodes the attention output $\mathbf{Z}_p^L \in \mathbb{R}^{np_x \times np_y \times C_{\text{emb}}}$ to the prediction $\mathbf{u}_{\text{pred}} \in \mathbb{R}^{H \times W \times C_k}$ for system k . In our work, we keep $p_t = 1$ and $C_{\text{uni}} = C_{\text{emb}}/4$.

Attention mechanisms — AViT, SViT, and ViT. The standard ViT attention mechanism takes into account the attention across the entire set of spatiotemporal dimensions, which results in a high attention cost when extremely long spatiotemporal token sequences (e.g., from high-resolution spatiotemporal data) are considered. To address this issue, various factorized attention mechanisms have been proposed, such as AViT (Ho et al., 2019; McCabe et al., 2023) and a spatio-temporal decoupled attention (Arnab et al., 2021), referred to as SViT here. These attention mechanisms mainly consist of the same multihead self attention (MHSA) and feed forward multi-layer perceptron (MLP) but differ in their attention block architecture. When L attention blocks are cascaded, the

standard attention block in ViT is given as

$$\begin{aligned} \widehat{\mathbf{Z}}_p^0 &= [\mathbf{z}_1^0, \mathbf{z}_2^0, \dots, \mathbf{z}_N^0] + \mathbf{E}_{\text{pos}}, \\ \mathbf{Z}_p^1 &= \text{MLP}(\widetilde{\mathbf{Z}}_p^1) + \widetilde{\mathbf{Z}}_p^1, & \widetilde{\mathbf{Z}}_p^1 &= \text{MHSA}(\widehat{\mathbf{Z}}_p^0) + \widehat{\mathbf{Z}}_p^0 + \text{MLP}(t_{\text{lead}}), \\ \mathbf{Z}_p^\ell &= \text{MLP}(\widetilde{\mathbf{Z}}_p^\ell) + \widetilde{\mathbf{Z}}_p^\ell, & \widetilde{\mathbf{Z}}_p^\ell &= \text{MHSA}(\mathbf{Z}_p^{\ell-1}) + \mathbf{Z}_p^{\ell-1}, \quad \ell = 2, \dots, L \end{aligned} \quad (1)$$

where $[\mathbf{z}_1^0, \dots, \mathbf{z}_N^0]$ denotes the full spatiotemporal token sequence of length N with each token $\mathbf{z}_i^0 \in \mathbb{R}^{C_{\text{emb}}}$, \mathbf{E}_{pos} is a positional embedding term, and each MHSA and MLP is followed by an InstanceNorm1d module. In ViT, the token sequence is composed of full spatiotemporal patches, meaning $N = nt \cdot npx \cdot npy$, resulting in an overwhelming costs of $\mathcal{O}((nt \cdot npx \cdot npy)^2)$ operations for attention. In contrast, SViT decouples the attention into $npx \cdot npy$ time-attention blocks and nt space-attention blocks cascaded sequentially, as in “MHSA_{time} \rightarrow MHSA_{space} \rightarrow MLP”, which reduces the MHSA cost to $npx \cdot npy \cdot \mathcal{O}(nt^2) + nt \cdot \mathcal{O}((npx \cdot npy)^2)$. AViT further decomposes the space-attention in SViT into two axial directions following the same approach, which leads to a cost of $npx \cdot npy \cdot \mathcal{O}(nt^2) + nt \cdot npy \cdot \mathcal{O}(npx^2) + nt \cdot npx \cdot \mathcal{O}(npy^2)$. The decomposition approach taken in both AViT and SViT neglects some spatiotemporal correlations and thus gives shorter token sequence length for each attention blocks, at the cost of introducing additional attention blocks. These extra attention blocks moderately increase the model size, as shown in Table II. Note that within the same size category, AViT and ViT are larger than ViT due to the additional MHSA, while AViT and ViT have similar sizes because AViT reuses the same attention for different spatial directions. In MATEY, we implement the three attention mechanisms, AViT, SViT, and ViT, and evaluate their performance on test problems to study how the lost spatiotemporal correlations affect the quality of the solution and to assess the impact of decoupled attentions with additional attention blocks on the learning efficiency for multi-physics foundation models.

Adaptive tokenization. Smaller patch sizes are preferred for better representation accuracy, as ViTs can capture long-range correlations between patches well but lack inductive biases within patches. However, features in physical systems often cross multiple length scales and exhibit strong spatiotemporal inhomogeneities. Consequently, constant patch sizes that are small enough to provide good accuracy in the necessary regions of such systems result in impractically long token sequence lengths over the entire domain. To address this issue, we propose an adaptive ViT that dynamically adjusts the tokenization patch sizes according to local physical features. To maximize expressiveness, we start with coarse patching and identify the most complex patches in each sample based on a simple metric, such as the variance of local features. The identified patches are further refined to the sub-token-scale (STS) to improve representation accuracy in these regions. Adaptive patch size leads to patches at varying length across samples, which are handled with padding mask. Patch position and patch area bias are represented following the embedding method in (Bodnar et al., 2024).

For a given solution $\mathbf{u}_i \in \mathbb{R}^{H \times W \times C}$ and an initial coarse patch size $[p_{x_1}, p_{y_1}]$, the patch sequence is refined adaptively based on local patch variance with two parameters, $[p_{x_{\text{ref}}}, p_{y_{\text{ref}}}]$ and γ_{ref} , as shown in Figure 2. The resulting STS tokens can be incorporated in two ways. In the first approach, referred to as “Adap_Mul” (for adaptive multi-resolution tokenization), we consider the coarse and STS tokens as separate sequences, passing through the attention blocks serially. In the second approach, referred to as “Adap_Mix” (for adaptive mixed-resolution tokenization), we append the sequence of STS tokens directly to the end of the sequence of coarse tokens. While the second approach leads to relatively longer token sequences, it has the potential benefit of better capturing cross-scale correlations than the decoupled first approach.

Pretraining and fine-tuning. We pretrain the models on PDEBench data, which include five basic 2D systems: incompressible flows, compressible flows, turbulent flows, reaction-diffusion systems, and shallow water equations. We consider two fine-tuning cases: 1) colliding thermals between a cold and a warm bubbles from MiniWeather simulations (Norman, 2020) and 2) lid-driven cavity MHD flows (Fambri et al., 2023). As discussed in detail in Appendix A.1, these fine-tuning datasets were selected to be meaningfully out-of-distribution, not only in flow regime but also in including thermal and electromagnetic components that are not represented at all in the pretraining data. Training was performed on the Frontier and Perlmutter supercomputers at the Oak Ridge Leadership Computing Facility (OLCF) and National Energy Research Scientific Computing Center (NERSC), respectively, using distributed data parallelism.

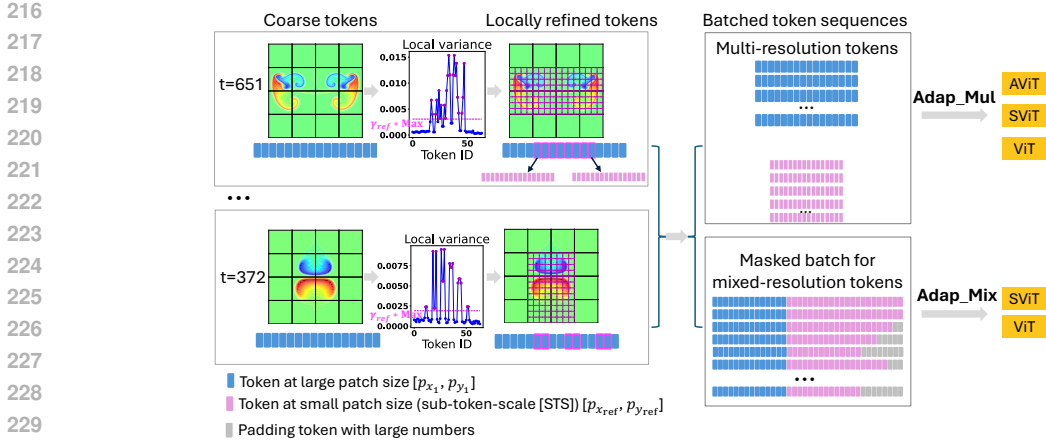


Figure 2: Adaptive tokenization that dynamically adjusts patch sizes based on local features. There are three essential parameters: $[p_{x_1}, p_{y_1}]$, $[p_{x_{ref}}, p_{y_{ref}}]$ and γ_{ref} . The parameter $[p_{x_1}, p_{y_1}]$ denotes the coarse patch size to start with, $[p_{x_{ref}}, p_{y_{ref}}]$ represents the refined patch size, and $\gamma_{ref} \in [0, 1]$ determines which patches to refine. We select patches with local variances greater than γ_{ref} times the maximum variance across all patches.

4 EXPERIMENTS

We design three experiments to evaluate 1) the performance of three spatiotemporal attention schemes, AViT, SViT, and ViT, 2) the impact of adaptive tokenization, and 3) the effectiveness of pretrained models on two fine-tuning tasks that feature physics different from the pretraining data.

4.1 SPATIOTEMPORAL ATTENTION SCHEMES

We evaluate AViT, SViT, and ViT for three model sizes: Tiny (Ti), Small (S), and Base (B) with 3, 6, 12 heads and hidden dimension $C_{emb} = 192, 384, \text{ and } 768$, respectively (Touvron et al., 2022), as shown in Table 1, on the colliding thermals dataset. In the same size category, AViT and SViT are about 30% larger than ViT due to the additional attention block. More details about the experiment are presented in Appendix A.2

Table 1: Number of model parameters in AViT, SViT, and ViT for three model sizes, Tiny, Small, and Base, detailed in Section 4.1. ViT results in about 30% fewer model parameters than AViT and SViT because the latter two require additional attention blocks.

	Tiny	Small	Base
AViT	7.5M	29.9M	119.3M
SViT	7.6M	30.0M	119.3M
ViT	5.8M	22.8M	90.9M

Figure 3 compares the final test error, defined as the normalized root-mean-square error (NRMSE), and the training time, represented as GPU hour per step, for the nine models. For the same size category, SViT (green) achieves the lowest error, followed by ViT (blue), and then AViT (red). In terms of training time, SViT takes longer than AViT, while ViT is the least expensive one. ViT processes longer token sequences and hence is expected to have a higher single-unit attention cost, whereas AViT and SViT have multiple attention units with shorter token sequence length. The results reported in Figure 3 show that the ViT has the lowest cost, which implies that the number of attention blocks plays a more important role than the token sequence length in terms of training cost in this example. This observation is due to the fact that the spatiotemporal token sequence length ($16 \times 8 \times 8$) in this example is relatively short. We expect ViT to become more expensive than AViT and SViT when more refined or higher dimensional solutions are considered, in which longer token sequences are required. In general, we find that SViTs and ViTs are more expressive

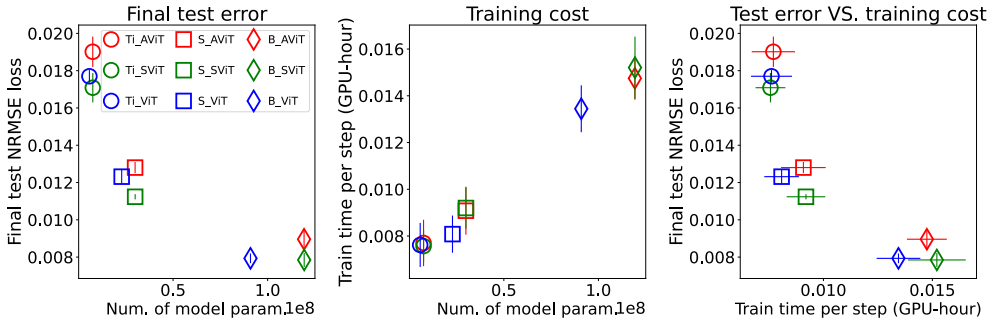


Figure 3: Learning efficiency of AViT, SViT, and ViT at three model sizes regarding final predictive error and training time cost: SViT and ViT are observed to be more expressive and computationally efficient than AViT in the experiment, as they require fewer model parameters and less training time to achieve the same test accuracy.

and computationally efficient than AViTs, in that they achieve the lower predictive errors with fewer model parameters and less training time.

4.2 ADAPTIVE TOKENIZATION

We start the evaluation of our adaptive tokenization methods in a single collision trajectory between two thermal bubbles. Figure 4 compare the temperature contours of the true solution at $t = 590$ with the predicted solutions from Ti-SViT models at constant patch sizes: $ps=16 \times 16$ and $ps=32 \times 32$ and adaptive tokenization (Adap_Mul with $p_{x_1} = p_{y_1} = 32$, $p_{x_{ref}} = p_{y_{ref}} = 16$, and $\gamma_{ref} = 0.2$). The predicted solution from $ps=32 \times 32$ exhibits abrupt changes with clear edges for the local structures inside the patches, while the finer resolution model at $ps=16 \times 16$ captures smoother, finer structures but requiring more patches. In contrast, our adaptive tokenization method (Adap_Mul) capture smooth, fine structures comparable to $ps=16 \times 16$ while requiring much fewer FLOPs, as shown in Figure 5.

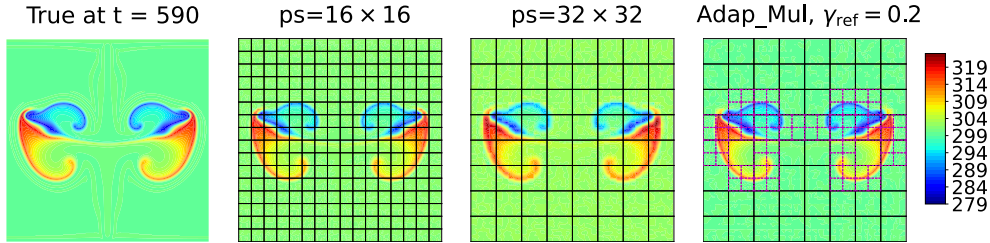


Figure 4: Predicted temperature contours at $t = 590$ from Ti-SViT models with constant patch sizes $ps=16 \times 16$ and $ps=32 \times 32$ and adaptive tokenization (Adap_Mul with $p_{x_1} = p_{y_1} = 32$, $p_{x_{ref}} = p_{y_{ref}} = 16$, and $\gamma_{ref} = 0.2$). Adap_Mul predicts smoother, finer local structures that are overlooked in $ps=32 \times 32$, similar to the more expensive $ps=16 \times 16$.

Figure 5 shows the final NRMSE loss versus the floating-point operations (measured in TeraFLOPs) for 15 models, including the three in Figure 4. We aim to evaluate the adaptive tokenization methods coupled with (Ti-) AViT, SViT, and ViT versus using these models with three fixed patch resolutions: patch sizes $ps=8 \times 8$, 16×16 , and 32×32 . For the same attention scheme with a fixed patch size, as expected, increasing resolution leads to lower errors but also substantially increases the training cost, particularly for ViT (triangles). ViT shows fewer FLOPs than AViT (squares) and SViT (circles) with shorter sequences (blue), consistent with the time measure in Figure 3, but it significantly surpasses the other two at the finest resolution with longer sequences (red). In contrast, our adaptive scheme (magenta markers, Adap_Mul with $p_{x_{ref}} = 16$), which starts with uniform 32×32 patches and locally refines to 16×16 on selected patches, achieves comparable accuracy to uniform 16×16 patches with SViT and ViT. Moreover, Adap_Mul with $p_{x_{ref}} = 16$ obtains this accuracy level at reduced FLOPs. As this reduced cost depends significantly on the spatiotemporal attention, the speedup is modest for

SViT but becomes more significant for ViT, being more than $2\times$ more efficient than constant patch sizes for ViT.

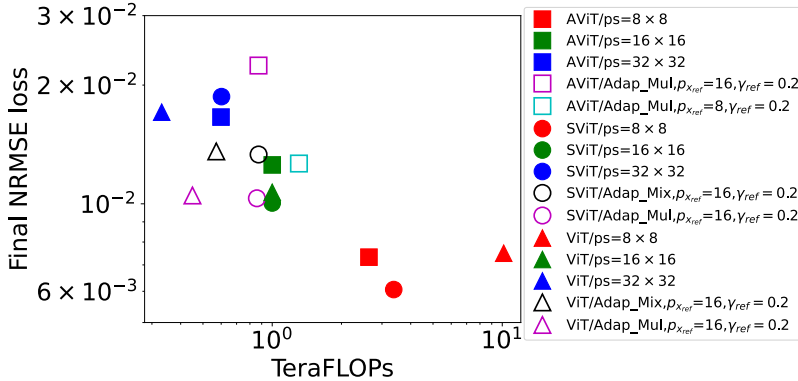


Figure 5: Final NRMSE loss for the three attention schemes (AViT, SViT, and ViT) with adaptive tokenization and constant patch sizes against estimated TeraFLOPs at the Tiny size. The flops values are available in Table A1. Adaptive tokenization methods (empty markers, $p_{x_1} = 32$) achieves comparable accuracy to $ps=16 \times 16$ (green solid markers) but with $2.2\times$ fewer FLOPs in ViT (triangles) and $1.16\times$ fewer FLOPs in SViT (circles), respectively. In contrast, AViT (squares), under the same setting, achieves higher error (magenta square with $p_{x_{ref}} = 16$) or lower error (cyan square with $p_{x_{ref}} = 8$) but with refined patches, making it less suitable for adaptive tokenization.

While the results with adaptive tokenization are positive on a single trajectory from the colliding thermals dataset, the accuracy improvement deteriorates when applied to more complex settings with multiple trajectories that involve varying initial bubble locations and temperature differences. Figure 6 compares the final test errors of Ti-SViT with constant patch sizes: $ps=32 \times 32$, $ps=16 \times 16$, and adaptive tokenization with a few parametric settings. In general, the adaptive errors are between errors of the two reference cases but still noticeably higher than the error from $ps=16 \times 16$. The adaptive error decreases with lower γ_{ref} values and is ideally expected to converge to the model error with the constant fine patch size. However, we encounter training instability issues when further reducing γ_{ref} . Addressing these stability issues is a focus of future work.

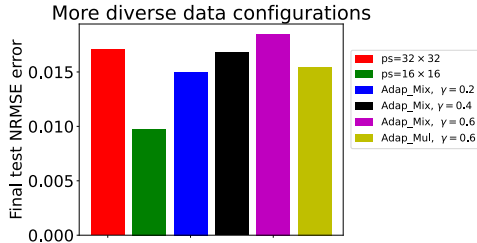


Figure 6: Final NRMSE test loss from constant patch sizes ($ps=32 \times 32$ and $ps=16 \times 16$) and adaptive tokenization for diverse data configurations with initial bubble locations and intensities varying dramatically. Adaptive tokenization (with $p_{x_1} = 32$, $p_{x_{ref}} = 8$) shows varying accuracy improvement with different parameter settings but the improvement is not as optimal as the single trajectory case in Figure 5.

4.3 EFFECTIVENESS OF PRETRAINING IN COLLIDING THERMALS AND MHD FINE-TUNING TASKS

We examine the transferrability of pretrained models to fine-tuning systems with distinct physics and different set of variables, as in Table A2. Specifically, we aim to address three broad questions:

1. Is pretraining effective when the downstream tasks have a distinct set of physical variables?

2. How does limited fine-tuning of non-attention blocks compare to full fine-tuning?
3. How does fine-tuning data size affect convergence?

To address these three questions, we design a sets of experiments, starting from models pretrained on PDEBench or randomly initialized models (*_INIT), and fine-tune them on colliding thermals and MHD datasets with distinct physical variables. For fine-tuning each model, we either allow all model parameters to be tunable ('ALL') or freeze the attention blocks and limit training to the preprocessor, the tokenization module, and the postprocessor ('PREPOST'). Finally, for each initial model and fine-tuning configuration, we train four models with increasing amounts of fine-tuning data.

For the colliding thermals dataset, Figure 7 compares the test loss with full and limited fine-tuning using pretrained and randomly initialized models. The different training data sizes ranging from one set of colliding thermals time-trajectory to 24 sets of trajectories. The fine-tuning task is to predict the solution of the physical system at a lead time of t_{lead} uniformly sampled between 1 and 50 steps. An example of the true and predicted solutions in these four training configurations is illustrated in Figure 8.

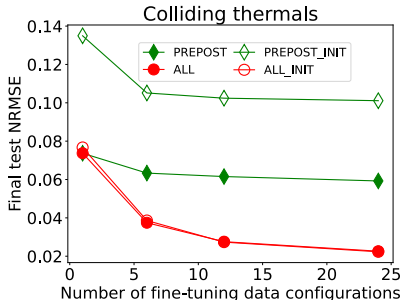


Figure 7: NRMSE loss for test set at different training data sizes in fine-tuning of colliding thermals at a maximum lead time of 50 steps, with full ('ALL') and limited ('PREPOST') fine-tuning using pretrained and randomly initialized models (*_INIT).

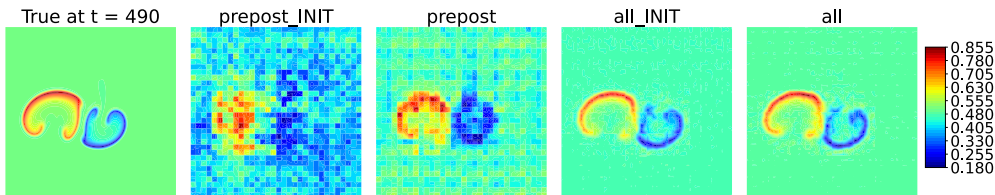
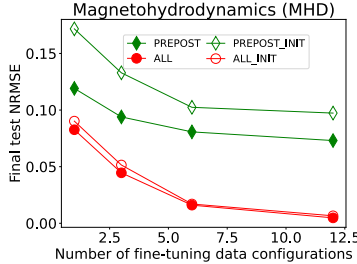


Figure 8: Temperature contours of true solution vs predicted solutions from four fine-tuned models (on 12 trajectories) at $t = 490$ from Ti-SViT models for a lead time of 40 in the collision of two thermal bubbles.

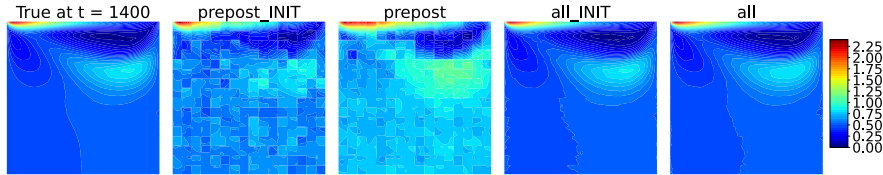
For the limited fine-tuning test with the colliding thermals dataset, the pretrained models achieve significantly lower error than starting from scratch with randomly initialized parameters. Moreover, while this advantage persists as the number of fine-tuning data increases, it is most pronounced in the low data configuration of learning from a single trajectory. Indeed, we find that limited fine-tuning with the pretrained models generalizes well even when learning from one trajectory, seeing only moderate improvements when run on the largest dataset size considered. Overall, the lower converged error from pretrained models suggests the frozen attention blocks clearly learned transferable knowledge during pretraining. For full fine-tuning, the accuracy is much better than limited fine-tuning as a result of the model being more expressive. The difference between the pretrained and randomly initialized models is much lower, being minor in the case of a single data configuration during training and vanishing as the amount of data increases.

For the MHD dataset, Figure 9 shows the final test NRMSE errors in lid-driven cavity flows after fine-tuning against data sizes when starting from pretrained and randomly initialized models for

432 limited and full fine-tuning. The training dataset sizes used for fine-tuning range from 1 to 12 simula-
 433 tion configurations, with each configuration including approximately 1900 samples. The fine-tuning
 434 task is to predict the flow solution at a lead time of t_{lead} uniformly sampled between 1 and 100 steps.
 435 Contour plots from the true solution and the predicted solution from each training configuration are
 436 depicted in Figure 10.



437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447 Figure 9: NRMSE loss for test set at different training data sizes in fine-tuning of lid-driven cavity
 448 MHD flows dataset at a maximum lead time of 100 steps, with full (‘ALL’) and limited (‘PREPOST’)
 449 fine-tuning using pretrained and randomly initialized models (‘*_INIT’).



451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461 Figure 10: Contours of true horizontal magnetic field values B_x vs predicted solutions from four
 462 fine-tuned models (on 12 trajectories) at $t = 1400$ from Ti-SViT models for a lead time of 80 in
 463 lid-driven cavity MHD flows.

462 Overall, the fine-tuning performance is a result of model expressibility, training data size, and the
 463 similarity between training and testing tasks. As with the colliding thermals dataset, pretrained
 464 models outperformed the randomly initialized models for both full and limited fine-tuning runs.
 465 However, the reduced expressibility of the limited fine-tuning configuration consistently shows an
 466 accuracy gap, even with more training data, as they cannot fully represent the data complexity. In
 467 contrast, full fine-tuning leads to more expressive models that can capture all training data informa-
 468 tion when trained on limited data but often show high test errors; as more training data is provided,
 469 they generalize well and lead to a convergent improved test error. In our fine-tuning, the randomly
 470 initialized models perform well in testing even with a single data configuration, likely due to the
 471 similarity between training and testing tasks. Future work will explore more challenging scenarios
 472 with increased heterogeneity within the fine-tuning data.

472 While studies like McCabe et al. (2023) have demonstrated impressive outperformance from fine-
 473 tuning of pretrained models versus randomly initialized models, these fine-tuning tests were per-
 474 formed on data that, while distinct, was fully governed by physical equations and characterized by
 475 physical variables that were represented in the training data. Yet for a model that aims to be founda-
 476 tional for multiphysical systems, we argue that assessing model performance in more realistic
 477 settings, where equations like Navier-Stokes are coupled with those from other domains of physics,
 478 is a more informative test of the effectiveness of pretraining. Accordingly, we assess fine-tuning per-
 479 formance on physical systems that incorporate fluid flows, which are well-represented in PDEBench,
 480 with thermodynamics and electromagnetism, which are not. As reasonably anticipated, we find that
 481 advantages of pretraining are reduced in this more complex setting.

482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539
 540
 541
 542
 543
 544
 545
 546
 547
 548
 549
 550
 551
 552
 553
 554
 555
 556
 557
 558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570
 571
 572
 573
 574
 575
 576
 577
 578
 579
 580
 581
 582
 583
 584
 585
 586
 587
 588
 589
 590
 591
 592
 593
 594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824
 825
 826
 827
 828
 829
 830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865
 866
 867
 868
 869
 870
 871
 872
 873
 874
 875
 876
 877
 878
 879
 880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000

486 decoupled spatiotemporal attention scheme, such as AViT and SViT, provides an intriguing balance
 487 of computational and data efficiency versus the standard ViT approach. Yet using SViT alone does
 488 not sufficiently address the computational challenges associated with attention for high spatial res-
 489 olutions. Second, we instead suggest that our adaptive tokenization scheme provides a promising
 490 approach for working with high resolution data. This sort of adaptivity has the potential to be both
 491 flexible and expressive enough to deal with the dynamic and sparse nature of the multiscale features
 492 in physical data. Third, we suggest an alternative path to evaluate foundation models for multi-
 493 scale physical systems that focuses on fine-tuning problems involving out-of-distribution physics
 494 governed by different equations with distinct sets of physical variables. In two such settings, collid-
 495 ing thermals and magnetohydrodynamics, we find that while pretraining does provide an advantage,
 496 its impact is much more muted compared to fine-tuning on the same set of variables, suggesting
 497 additional effort is required to obtain truly foundational models in this space.

498 REFERENCES

- 500 Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid.
 501 Vivit: A video vision transformer, 2021. URL <https://arxiv.org/abs/2103.15691>,
 502
- 503 Cristian Bodnar, Wessel P Bruinsma, Ana Lucic, Megan Stanley, Johannes Brandstetter, Patrick
 504 Garvan, Maik Riechert, Jonathan Weyn, Haiyu Dong, Anna Vaughan, et al. Aurora: A foundation
 505 model of the atmosphere. *arXiv preprint arXiv:2405.13063*, 2024.
- 506 Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *Proceedings of the*
 507 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp.
 508 4599–4603, June 2023.
- 509 F. Fambri, E. Zampa, S. Busto, L. Río-Martín, F. Hindenlang, E. Sonnendrücker, and M. Dumbser.
 510 A well-balanced and exactly divergence-free staggered semi-implicit hybrid finite volume / finite
 511 element scheme for the incompressible mhd equations. *Journal of Computational Physics*, 493:
 512 112493, 2023. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2023.112493>. URL <https://www.sciencedirect.com/science/article/pii/S0021999123005880>,
 513
- 514 Qihang Fan, Quanzeng You, Xiaotian Han, Yongfei Liu, Yunzhe Tao, Huaibo Huang, Ran He, and
 515 Hongxia Yang. Vitar: Vision transformer with any resolution, 2024. URL <https://arxiv.org/abs/2403.18361>,
 516
- 517 Zhou Hang, Yuezhou Ma, Haixu Wu, Haowen Wang, and Mingsheng Long. Unisolver: Pde-
 518 conditional transformers are universal pde solvers. *arXiv preprint arXiv:2405.17527*, 2024.
 519
- 520 Zhongkai Hao, Chang Su, Songming Liu, Julius Berner, Chengyang Ying, Hang Su, Anima Anand-
 521 kumar, Jian Song, and Jun Zhu. Dpot: Auto-regressive denoising operator transformer for large-
 522 scale pde pre-training. *arXiv preprint arXiv:2403.03542*, 2024.
 523
- 524 Joakim Bruslund Haurum, Sergio Escalera, Graham W. Taylor, and Thomas B. Moeslund. Which to-
 525kens to use? investigating token reduction in vision transformers. In *Proceedings of the IEEE/CVF*
 526 *International Conference on Computer Vision (ICCV) Workshops*, pp. 773–783, October 2023.
 527
- 528 Jakob Drachmann Havtorn, Amélie Royer, Tijmen Blankevoort, and Babak Ehteshami Bejnordi.
 529 Msvit: Dynamic mixed-scale tokenization for vision transformers. In *Proceedings of the*
 530 *IEEE/CVF International Conference on Computer Vision*, pp. 838–848, 2023.
- 531 Maximilian Herde, Bogdan Raonić, Tobias Rohner, Roger Käppeli, Roberto Molinaro, Emmanuel
 532 de Bézenac, and Siddhartha Mishra. Poseidon: Efficient foundation models for pdes. *arXiv*
 533 *preprint arXiv:2405.19101*, 2024.
- 534 Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional
 535 transformers. *arXiv preprint arXiv:1912.12180*, 2019.
 536
- 537 Sam Ade Jacobs, Masahiro Tanaka, Chengming Zhang, Minjia Zhang, Shuaiwen Leon Song,
 538 Samyam Rajbhandari, and Yuxiong He. DeepSpeed Ulysses: System optimizations for enabling
 539 training of extreme long sequence transformer models, 2023. URL <https://arxiv.org/abs/2309.14509>.

- 540 Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and
541 Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and
542 detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-*
543 *tion*, pp. 4804–4814, 2022.
- 544 Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue
545 Huang, Hanchi Sun, Jianfeng Gao, Lifang He, and Lichao Sun. Sora: A review on back-
546 ground, technology, limitations, and opportunities of large vision models, 2024. URL <https://arxiv.org/abs/2402.17177>.
- 547
548
- 549 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
550 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
551 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 552 Michael McCabe, Bruno Régaldo-Saint Blancard, Liam Holden Parker, Ruben Ohana, Miles
553 Cranmer, Alberto Bietti, Michael Eickenberg, Siavash Golkar, Geraud Krawezik, Francois
554 Lanusse, et al. Multiple physics pretraining for physical surrogate models. *arXiv preprint*
555 *arXiv:2310.02994*, 2023.
- 556 Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-
557 Nam Lim. Adavit: Adaptive vision transformers for efficient image recognition. In *Proceedings*
558 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12309–12318,
559 2022.
- 560 Matt Norman. `mnorman/miniWeather`, September 2024. URL <https://github.com/mnorman/miniWeather>, original-date: 2018-07-24T02:29:06Z.
- 561
562
- 563 Matthew R Norman. `miniweather`. Technical report, Oak Ridge National Laboratory (ORNL), Oak
564 Ridge, TN (United States), 2020.
- 565 Shashank Subramanian, Peter Harrington, Kurt Keutzer, Wahid Bhimji, Dmitriy Morozov,
566 Michael W Mahoney, and Amir Gholami. Towards foundation models for scientific machine
567 learning: Characterizing scaling and transfer behavior. *Advances in Neural Information Process-*
568 *ing Systems*, 36, 2024.
- 569 Makoto Takamoto, Timothy Praditia, Raphael Leiteritz, Daniel MacKinlay, Francesco Alesiani,
570 Dirk Pflüger, and Mathias Niepert. Pdebench: An extensive benchmark for scientific machine
571 learning. *Advances in Neural Information Processing Systems*, 35:1596–1611, 2022.
- 572
573
- 574 Hugo Touvron, Matthieu Cord, Alaeldin El-Nouby, Jakob Verbeek, and Hervé Jégou. Three things
575 everyone should know about vision transformers. In *European Conference on Computer Vision*,
576 pp. 497–515. Springer, 2022.
- 577 Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-vit:
578 Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF conference on*
579 *computer vision and pattern recognition*, pp. 10809–10818, 2022.
- 580 Enzhi Zhang, Isaac Lyngaas, Peng Chen, Xiao Wang, Jun Igarashi, Yuankai Huo, Mohamed Wahib,
581 and Masaharu Munetomo. Adaptive patching for high-resolution image segmentation with trans-
582 formers, 2024. URL <https://arxiv.org/abs/2404.09707>.
- 583
584
585
586
587
588
589
590
591
592
593