

Reversible Vision Transformers

Karttikeya Mangalam^{*,2} Haoqi Fan¹ Yanghao Li¹
 Chao-Yuan Wu¹ Bo Xiong¹ Christoph Feichtenhofer^{*,1} Jitendra Malik^{1,2}

¹Meta AI, FAIR

²UC Berkeley

Abstract

We present *Reversible Vision Transformers*, a memory efficient architecture design for visual recognition. By decoupling the GPU memory requirement from the depth of the model, *Reversible Vision Transformers* enable scaling up architectures with efficient memory usage. We adapt two popular models, namely *Vision Transformer* and *Multiscale Vision Transformers*, to reversible variants and benchmark extensively across both model sizes and tasks of image classification, object detection and video classification. *Reversible Vision Transformers* achieve a reduced memory footprint of up to **15.5×** at roughly identical model complexity, parameters and accuracy, demonstrating the promise of reversible vision transformers as an efficient backbone for hardware resource limited training regimes. Finally, we find that the additional computational burden of recomputing activations is more than overcome for deeper models, where throughput can increase up to **2.3×** over their non-reversible counterparts. Full code and trained models are available at <https://github.com/facebookresearch/slowFast>. A simpler, easy to understand and modify version is also available at <https://github.com/karttikeya/minREV>.

1. Introduction

The deep learning revolution in computer vision has rested on the bedrock of high performance hardware accelerators. Fueled by special purpose AI accelerators, the compute requirements for state-of-the-art models are growing exponentially. However, compute is only half the story. The other, and often overlooked half, is memory bandwidth bottleneck, which has been difficult to proportionally scale as compared to peak accelerator FLOPs [54]. In particular, the peak accelerator FLOPs have been increasing at a rate of $\sim 3.1\times$ every 2 years [21, 62]. However, peak bandwidth only scales at a rate of $\sim 1.4\times$ every 2 years. This disparity is exacerbated in transformers, which have been doubling

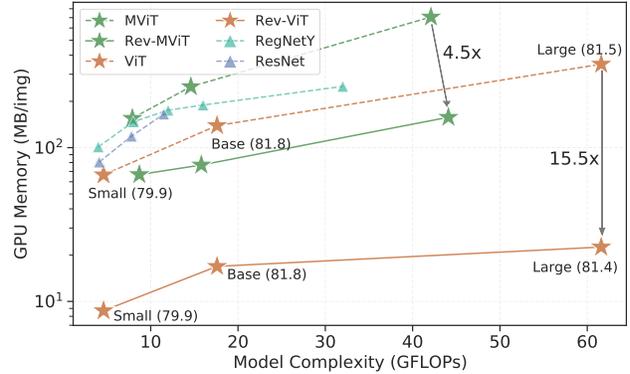


Figure 1. **Reversible Vision Transformers** are more memory-efficient, yet powerful *reversible counterparts* of state-of-the-art Vision Transformer (ViT) [15] and Multiscale Vision Transformer (MVIT) [18] architectures with varying model complexity. Numbers in parentheses denote top-1 ImageNet performance. ResNet [28] and RegNet [58] are only shown for reference. For detailed discussion please refer to §4.1.

in required compute roughly every three months for the past three years, resulting in a so-called memory wall [21] where both the overall model performance as well as the training speed have become tightly memory-bound [34].

As such, for bandwidth bound models, trading compute for memory through re-computation could actually be more efficient than using work-optimal algorithms [70, 71]. In the case of training neural network models, this can be achieved by re-computing activations instead of storing and then loading them from DRAM [31]. Besides training speed, scaling vision transformers in depth naturally hits the GPU memory capacity, especially in memory starved regimes such as video recognition where state-of-the-art models are often limited to batch size 1 due to high memory footprint of intermediate activations.

We propose *Reversible Vision Transformers*, a family of expressive visual recognition architectures with very favorable activation memory footprints (Figure 1) compared to their non-reversible variants. By trading-off GPU activation caching with efficient on-the-fly activation re-computation, reversible vision transformers effectively *decouple* the activation *memory* growth from the *depth* of the model.

*Equal technical contribution.

While the natural language processing community has performed some early exploration of reversible transformers for machine translation [38], these techniques focus on *longer sequence lengths* rather than depth.

Our experiments show that a straightforward adaptation of vision transformers to reversible architectures *fails* to scale for *deeper* models because of training convergence instabilities due to internal sub-block residual connections.

In this work, we reconfigure the residual paths in Vision Transformers (ViT) [15] and Multiscale Vision Transformers (MViT) [18] to overcome this issue. We further find that reversible structures have stronger inherent regularization and therefore, we use a lighter augmentation recipe (repeated augmentation, augmentation magnitude and stochastic depth) and lateral connections between residual blocks.

We benchmark extensively across image recognition tasks such as image classification and object detection as well as video classification, across all of which, reversible vision transformers have competitive performance to their non-reversible counterparts suffering negligible to no performance decay. Moreover, reversible models have extremely favorable per-image memory footprint, saving $15.5\times$ on the ViT-Large model and $4.5\times$ on the MViT-Large model with reversible training.

In summary, our contributions are three-fold.

(i) We propose Reversible Vision Transformer (**Rev-ViT**) and Reversible Multiscale Vision Transformers (**Rev-MViT**), memory efficient reversible adaptations of state-of-the-art visual recognition backbones.

(ii) We observe reversible transformers to have a stronger inherent regularization than vanilla networks. Hence, we develop new training recipes by adapting the original recipes with different repeated augmentations, augmentation magnitudes and drop path rate to match the performance of their non-reversible counterparts.

(iii) We benchmark our models across several tasks: image classification, object detection and action recognition, across accuracy, memory, maximum training batch size and model complexity. In particular, at matched complexity (FLOPs/parameters) and final accuracy, Rev-ViT-B and Rev-ViT-L train with per image memory footprints that are $8.2\times$ and $15.5\times$ lighter than ViT-B and ViT-L respectively. Further, we show how *deep* reversible networks can achieve up to $2-4\times$ throughput than their vanilla counterparts.

2. Related Work

Transformers are a popular network structure that were first proposed for natural language applications [68] and now are widely used in all areas of deep learning such as Reinforcement Learning [7], Speech [41], Music [32], multi-modal learning [35] and recently, in traditional vision tasks [15] as well. Since their introduction, Vision Transformers have experienced enthusiastic adoption and have

been applied to several visual recognition tasks [15, 63, 64] using priors such as multi-scale feature hierarchies [18, 24, 49, 69, 76] and local structure modelling [9, 14, 49]. Further, vision transformers have also been generalized for action recognition and detection in videos [1, 3, 18, 49, 52, 53].

However, a crucial problem with scaling up transformer models is the growth of required GPU memory with depth. This linear growth in memory is prohibitive to the development of very deep models since the batch size needs to be reduced considerably to be able to accommodate storing the intermediate activations on GPU. This problem is exacerbated in video models which process very large input tensors and are often trained with batch size 1 even for shallower depths. A potential systems-level solution to scale up conventional transformer architectures is model parallelism [10] that puts different parts of the model on different GPUs. However in practice, it is quite slow and requires special high bandwidth network infrastructure because of huge across device traffic.

In this work, we use Vision Transformers [15] and Multiscale Vision Transformers [18] as our base models and propose their reversible transformer version that decouple the memory requirement from depth of the model. This facilitates saving GPU memory and allows training with much higher batch size, and consequently, to preserve or even *increase* training throughput of deep non-reversible models.

Reversible Architectures are a family of neural network architectures that are based on the NICE [12, 13] reversible transformation model which are the precursors of the modern day generative flow based image generation architectures [30, 37]. Based on the NICE invertible transformations, Gomez *et al.* [22] propose a Reversible ResNet architecture that employs the reversible transformation [12] for memory-efficient image classification in ResNets [27]. An interesting line of work builds upon the Reversible ResNets ideas proposing better reversible CNN models using ODE characterizations [6, 39, 59], momentum [39, 59], layer-wise inversion [25], fourier transform based inversion [20] and fixed point iteration based inversion [2, 60]. Reversible CNNs have been applied to several traditional image tasks such as compression [46], reconstruction [43], retrieval [42], and denoising [33, 47] as well as to compressed sensing [61], compact resolution [75], image to image translation [67], remote sensing [56], medical image segmentation [55, 74] and MRI reconstruction [57]. Reversible transformation have also been adapted to other networks such as RNNs [51], Unet [4, 16], Masked Convolutional Networks [60] and 1000-layer deep Graph Neural Networks [40]. Some early attempts have also been made to adapt the reversible transformation to the NLP domain, initiated by Kiatev *et al.* [38] and built upon in [78, 79] for machine translation.

However, word-level input partitioning contains much richer semantic content than patch level image partitioning and NLP transformers tend to be shallower in depth but wider in channel dimension. For example, Kiatchev *et al.* [38] focus on expanding on the input sequence dimension rather than model depth and with no benchmarking on maximum batch-size, peak GPU memory and training throughput.

Our experiments show that a naïve adaption of reversible vision transformers performs poorly for deeper (≥ 8 blocks) models. This work is the first to propose Reversible Vision Transformers, adapt it to two state-of-the-art transformer networks, namely, ViT and MViT. Furthermore, this work is the first use of a reversible backbone for object detection and video classification, which tends to be one of the most memory starved domains of visual recognition.

3. Approach

We first present a brief overview of the reversible transformation (§3.1.1) and its benefits in neural network training (§3.1.3). We then present our proposed Reversible Vision Transformer (§3.2) its two residual stream structure (§3.2.1 and associated constraints (§3.1.3. This is followed by our proposed Reversible Multiscale Vision Transformer architecture (§3.3) and its sub-blocks (§3.3.2 and §3.3.1) that allow end-to-end reversible training.

3.1. Reversible Block Structure

The reversible transformer is composed of a stack of reversible blocks that follow the structure of the reversible transformation to allow analytic invertibility of outputs.

3.1.1 Reversible Transformation

Consider a transformation T_1 that transforms an input tensor I partitioned into two d dimensional tensors, $[I_1; I_2]$ into the output tensor O also similarly partitioned into tensors, $[O_1; O_2]$ with an arbitrary differentiable function $F(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as follows:

$$\mathbf{I} = \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} \xrightarrow{T_1} \begin{bmatrix} O_1 \\ O_2 \end{bmatrix} = \begin{bmatrix} I_1 \\ I_2 + F(I_1) \end{bmatrix} = \mathbf{O}$$

Note that the above transformation T_1 allows an inverse transformation T_1' such that $T_1' \circ T_1$ is an identity transform. Also, consider an analogous transposed transformation T_2 using the function $G(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as follows:

$$\mathbf{I} = \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} \xrightarrow{T_2} \begin{bmatrix} O_1 \\ O_2 \end{bmatrix} = \begin{bmatrix} I_1 + G(I_2) \\ I_2 \end{bmatrix} = \mathbf{O}$$

Similar to T_1 , T_2 also allows an inverse transform T_2' . Now consider the composition $T = T_2 \circ T_1$ that transforms both the partitions of the input vector \mathbf{I} and is obtained as,

$$\mathbf{I} = \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} \xrightarrow{T} \begin{bmatrix} O_1 \\ O_2 \end{bmatrix} = \begin{bmatrix} I_1 + G(I_2 + F(I_1)) \\ I_2 + F(I_1) \end{bmatrix} = \mathbf{O} \quad (1)$$

Naturally, T affords the inverse transform $T' = T_1' \circ T_2'$ that follows $T'(T(\mathbf{I})) = \mathbf{I}$. Note that the inverse transform T' queries the functions F and G exactly once and hence has the same computational cost as the forward transform T .

3.1.2 Vanilla networks require caching activations

Consider the back-propagation mechanism. Given a computation graph node, \mathcal{M} , its children nodes $\{\mathcal{N}_j\}$, and the gradients of the children node with respect to final loss $\left\{ \frac{d\mathcal{L}}{d\mathcal{N}_j} \right\}$, the back-propagation algorithm uses the chain rule to calculate the gradient with respect to \mathcal{M} as,

$$\frac{d\mathcal{L}}{d\mathcal{M}} = \sum_{\mathcal{N}_j} \left(\frac{\partial f_j}{\partial \mathcal{M}} \right)^T \frac{d\mathcal{L}}{d\mathcal{N}_j}$$

where f_j denotes the function computing node \mathcal{N}_j from its parents, \mathcal{M} being one of them. The jacobian $\frac{\partial f_j}{\partial \mathcal{M}}$, requires calculating the partial gradient of the f_j output with respect to the current node \mathcal{M} .

Now consider the simplest possible neural network layer $f(X) = W^T X$, where X is an intermediate activation inside the network. Applying the above described back-propagation algorithm to compute the derivative with respect to parent nodes, and using the output Y as the sole child node, \mathcal{N}_j , we get,

$$\frac{d\mathcal{L}}{dW} = \left(\frac{d\mathcal{L}}{dY} \right) X^T \quad \frac{d\mathcal{L}}{dX} = W \frac{d\mathcal{L}}{dY}$$

Thus, because of the function jacobian, the backpropagation algorithm requires intermediate activations during the forward pass to be available in the backward pass to compute the gradients with respect to the weights.

Typically, this is achieved by caching the intermediate activations on GPU memory for use in the backward pass. This allows fast gradient computation at the cost of extra memory. Further, the sequential nature of the network requires the activations for all the layers to be cached in before the loss gradients are calculated and the cached memory is freed. This dependence significantly affects the peak memory usage which thus becomes linearly dependent on the network depth D .

3.1.3 Learning without caching activations

As noted in §3.1.1, an input transformed with the reversible transformation T allows recalculating the input from the output of the transformation. Hence, a network composed of such reversible transformations *does not need to store intermediate activations* since they can be recomputed easily in the backward pass from the output. However the reversible transformation T places an important constraint on the property of the learnt function.

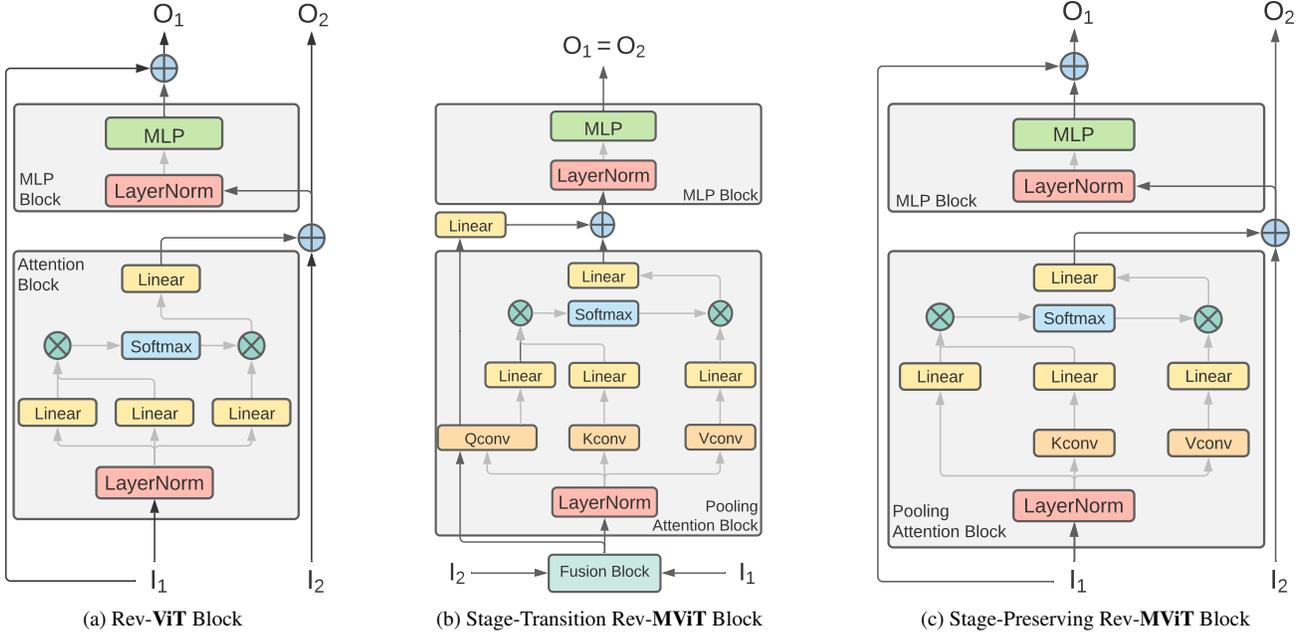


Figure 2. **Reversible ViT** is a two-residual-stream architecture composed of a stack of Reversible ViT blocks (a) that transforms the inputs I_1 and I_2 with the ViT design [15], but in our reversible fashion. **Reversible MViT** is a two-residual-stream architecture as well, made up of a stack of two type of blocks – (b) The stage-transition blocks that act as coupling between the residual streams as well as perform channel upsampling and resolution downsampling and (c) the stage-preserving blocks that form the majority of the computational graph and propagate information preserving input feature dimension.

Equidimensional Constraint. As mentioned in §3.1.1, the functions F and G need to be equidimensional in input and output spaces. Hence, the feature dimensions need to remain constant under T . While this constraint is an obstruction for other vision architectures such as ResNets [27] that require a change of feature dimensions, it is easily satisfied in the Vision Transformer architecture [15] which maintains a constant feature dimension throughout the layers.

3.2. Reversible Vision Transformers

3.2.1 Adapting ViT to Two-Residual-Streams

Fig. 2a shows the reversible transformation T adapted to the Vision Transformer architecture [15]. The input consists of two partitioned tensors I_1 and I_2 that are transformed as per the equation 3.1.1 maintaining reversibility. This leads to a *two-residual-stream* architecture where each of the inputs I_1 and I_2 maintain their own residual streams while mixing information with each other using functions F and G . Following ViT [15], we use the Multi-head attention and the MLP subblocks as functions F and G respectively.

3.2.2 Boundary Conditions

As the ViT architecture only uses a single residual stream, the architecture needs to be modified to support the two-residual-stream design (§3.2.1). We propose the following:

1. Initiation. We keep the stem intact and send the patchification output activations to I_1 and I_2 . Note that this de-

sign choice is different from [23] which proposes to split in halves along the channel dimensions.

2. Termination. The two residual paths need to be fused before the final classifier head to preserve information. We propose to layer-normalize the inputs first, followed by concatenation, to reduce the fusion computational overhead.

3.2.3 Reconfiguring Residual Connections

Residual connections play a key role for signal propagation in deep networks [27]. The reversible transform T itself also depends crucially on the residual connections between the two streams to maintain reversibility. Interestingly, we observe a key relationship between the residual connections and signal propagation in Reversible Vision Transformer.

Note that while it is common practice for neural network blocks to be wrapped around a residual connection for better gradient flow [27], there is no such connection for either the I_1 or I_2 inputs. Specifically, internal residual connections around the MLP and attention sub-blocks for both the I_1 and I_2 streams are absent. Instead, the residual connections for each residual stream flows through the other stream, operating through the inherent skip connection present in the reversible transformation T (§3.1.1). We find these *internal* skip connections *detrimental* to training convergence for deeper models while bringing no additional gain for shallower models and choose to omit them entirely for reversible vision transformer blocks.

3.3. Reversible Multiscale Vision Transformers

The recently proposed MViT architecture develops a feature hierarchy inside the model by *downsampling* the visual resolution and *upsampling* the channel dimension. It obtains state-of-the-art results on both image and video classification benchmarks. To showcase the flexibility of the reversible design, we adapt the MViT model to Reversible Multiscale Vision Transformers. We propose to compose the Reversible MViT architecture in the same structure as the MViT model but using two different layers – the *Stage Transition* and the *Stage-Preserving* blocks.

3.3.1 Stage-Transition Block

Figure 2b depicts the architecture of the proposed stage-transition block. The stage-transition block closely follows the design of the resolution upsampling blocks in MViT [18] with the following crucial modifications:

Lateral Connections. The residual streams I_1 and I_2 are fused with lateral connections at the start of the stage-transition block. This allows efficient computation of the resolution downsampling and feature upsampling without repeat computation in each stream separately.

Feature Upsampling. MViT performs feature upsampling in the last MLP block before the resolution upsampling block. We propose to move the channel upsampling stage inside the pooling attention sub-block of the stage-transition block. Specifically, we propose to upsample the Query, Key and Value vectors in the linear layer following the pooling channel-wise convolutional layers (Figure 2b and 2c). This was the dual benefit of (A) allowing all feature dimension changes to take place in sync inside the same block and allowing other blocks to keep feature dimensions intact, a virtue of reversible architectures (§3.1.3) and (B) saving additional computation from being used in the prior MLP and pooling layers. We follow the same boundary conditions at the stage-transition blocks as in the reversible vision transformer architecture (§3.2.2).

3.3.2 Stage-Preserving Block

Figure 2c shows the reversible transformation T (§3.1.1) adapted to the Multiscale Vision Transformer architecture [18]. The design closely resembles that of the reversible vision transformer block (Figure 2a) with the addition of multi-head pooling attention [18]. Note that even though the attention uses pooling on key and value tensors, thereby changing the sequence length, the output dimensions are still preserved. Hence, the stage-preserving block still follows the equidimensional constraint (§3.1.3)

and hence can be made fully reversible and learnt without caching activations.

Since each stage-transition block changes the spatiotemporal resolution, they occur only a limited number of times in the entire MViT network. In other words, the majority of the computation as well as memory usage is performed within the stage-preserving blocks and is fully reversible. We follow the same residual connection circuit (§3.2.3) as in Reversible Vision Transformer blocks for both the stage-transition and the stage-preserving blocks.

4. Results

Datasets. We benchmark both the Reversible Vision Transformer and the Reversible Multiscale Vision Transformer architectures extensively across image classification (ImageNet [11]), video classification (Kinetics 400 [36] & Kinetics 600 [5]) and object detection (MS-COCO [45]). Across all the benchmarks, we observe significant memory savings by using the reversible architecture with negligible to no accuracy change. All presented results and ablations are trained from random initialization, except for COCO where we initialize from ImageNet weights.

4.1. Image Classification

Settings. We benchmark our proposed models on image classification on the ImageNet-1K dataset [11] with $\sim 1.28\text{M}$ images among 1000 classes. We follow training recipes [17] for both ViT [15] and MViT [18] models with certain crucial adaptations (§6). All models are trained from random initialization without EMA for 300 epochs except for ViT-L and Rev-ViT-L which follow a 200 epoch training recipe. Training details are in Supplementary.

Results. Table 1 shows the results for Reversible Vision and Reversible Multiscale Vision Transformers across different models and FLOP regimes. We benchmark all the models on a single 16 GB V100 GPU under 224×224 image size and otherwise identical conditions. The maximum batch size is obtained as the highest number of images in a batch than can train without running out of GPU memory. The memory per image is measured as the peak GPU memory each image occupies during training.

We note that Reversible Vision Transformers match the FLOP and parameter specifications of their non-reversible counterparts exactly owing to the parsimonious design of the reversible vision transformer block (§3.2). The Reversible Multiscale Vision Transformer has slightly higher FLOPs due to the stage-transition (§3.3.1) stages while still being very GPU memory efficient owing to the stage-preserving (§3.3.2) stages.

model	Acc	Memory (MB/img)	Maximum Batch Size	GFLOPs	Param (M)
ResNet-101 [29]	76.4	118.7	112	7.6	45
ResNet-152 [29]	77.0	165.2	79	11.3	60
RegNetY-4GF [58]	80.0	101.1	136	4.0	21
RegNetY-12GF [58]	80.3	175.2	75	12.1	51.8
RegNetY-32GF [58]	80.9	250.2	46	32.3	32.3
Swin-T [48]	81.3	-	-	4.5	29
ViT-S [63]	79.9	66.5	207	4.6	22
Rev-ViT-S	79.9	8.8 ↓ 7.5 ×	1232 ↑ 5.9 ×	4.6	22
ViT-B [63]	81.8	129.7	95	17.6	87
Rev-ViT-B	81.8	17.0 ↓ 7.6 ×	602 ↑ 6.3 ×	17.6	87
RegNetY-8GF [58]	81.7	147.2	91	8.0	39
CSWin-T [14]	82.7	-	-	4.3	23
Swin-S [48]	83.0	-	-	8.7	50
ViT-L	81.5	349.3	26	61.6	305
Rev-ViT-L	81.4	22.6 ↓ 15.5 ×	341 ↑ 13.1 ×	61.6	305
MViT-B-16 [18]	82.8	153.6	89	7.8	37
Rev-MViT-B-16	82.5	66.8 ↓ 2.3 ×	157 ↑ 1.8 ×	8.7	39

Table 1. **Comparison to prior work on ImageNet-1K classification.** All memory and maximum batch size are on 224×224 input resolution on a 16G V100 GPU. **Rev-ViT** and **Rev-MViT** match performance across different FLOP regimes at a fraction of the per-input GPU memory cost.

Increasing memory savings with depth. In Table 1, we observe that our Rev-ViT matches the performance of vanilla ViT to very close fidelity across all model variants (Small, Base and Large) and FLOP regimes. Since the memory used per image is linearly dependent on the depth of the model for vanilla networks (§3.1.2), the memory gains of the reversible model increases as the network scales in depth. Notably, while the Reversible ViT-S already enjoys an impressive memory saving of about **86.8%** (equivalent to a **7.6**× reduction) with respect to the vanilla ViT-S model, the gain increases further to **15.5**× or, about **93%** memory savings for the Reversible ViT-L model.

Equivalently, the saved memory can be used to increase the training batch size where we observe a similar trend as well. While reversible ViT-S model achieves a **6.1**× increase in batch size on the ViT-S model, the effect is more for ViT-L model where the maximum batch size increases by **14.3**× jumping from a small 24 image per batch to 344 images. This is a very favorable trend, since it is indeed the deeper models that hit the *GPU memory wall* [21].

Further, hierarchical vision transformers such as MViT also enjoy a memory saving of about **52.1%** without suffering any significant drop in performance. The memory savings in Rev-MViT are smaller compared to the ViT variants because of the stage-transition blocks in hierarchical models (§3.3.1) that require storing the input activations due to the non-reversible nature of pooling attention stemming from the feature dimension change [18].

model	top-1	Mem Max (GB)	BS	GFLOPs × views	Param
Two-Stream I3D [5]	71.6	-	-	$216 \times \text{NA}$	25.0
R(2+1)D [66]	72.0	-	-	152×115	63.6
Two-Stream R(2+1)D [66]	73.9	-	-	304×115	127.2
Oct-I3D + NL [8]	75.7	-	-	$28.9 \times 3 \times 10$	33.6
ip-CSN-152 [65]	77.8	-	-	$109 \times 3 \times 10$	32.8
SlowFast 4×16 , R50 [19]	75.6	-	-	36.1×30	34.4
SlowFast 8×8 , R101 [19]	77.9	-	-	106×30	53.7
SlowFast 8×8 +NL [19]	78.7	-	-	$116 \times 3 \times 10$	59.9
ViT-B-VTN-IN-1K [52]	75.6	-	-	$4218 \times 1 \times 1$	114.0
ViT-B-VTN-IN-21K [52]	78.6	-	-	$4218 \times 1 \times 1$	114.0
MViT-B-16, 16×4	78.4	1.27	10	$70.5 \times 1 \times 5$	36.6
Rev-MViT-B-16, 16×4	78.5	0.64	20	$64 \times 1 \times 5$	34.9

Table 2. **Comparison to prior work on Kinetics-400 video classification.** Single view inference cost is reported along with used number of views ($\text{FLOPs} \times \text{view}_{\text{space}} \times \text{view}_{\text{time}}$). Memory (Mem) reported in Gigabytes per input clip. Maximum Batch Size (Max BS) measured as the maximum possible single GPU batch size. All measurements are performed on a single 16G V100 GPU.

model	top-1	Mem Max (GB)	BS	GFLOPs × views	Param
SlowFast 16×8 +NL [19]	81.8	-	-	$234 \times 3 \times 10$	59.9
X3D-XL	81.9	-	-	$48.4 \times 3 \times 10$	11.0
ViT-B-TimeSformer-IN-21K [3]	82.4	-	-	$1703 \times 3 \times 1$	121.4
ViT-L-ViViT-IN-21K [1]	83.0	-	-	$3992 \times 3 \times 4$	310.8
MViT-B-16, 16×4	81.3	-	-	$70.3 \times 1 \times 5$	36.6
MViT-B-16, 32×3	83.4	-	-	$170 \times 1 \times 5$	36.8
MViT-B-24, 32×3	83.8	4.40	2	$236 \times 1 \times 5$	52.9
Rev-MViT-B-24, 32×3	83.7	1.64	7	$223 \times 1 \times 5$	51.8

Table 3. **Comparison to prior work on Kinetics-600 video classification.** Results under same settings as Kinetics-400 in Table 2.

4.2. Video Classification

Settings. We also benchmark Rev-MViT-B models on action classification on Kinetics-400 [36] and Kinetics-600 [5] datasets. All the models are trained from scratch with training recipes adapted from [18].

Results. Table 2 and 3 present the results on action recognition task on the Kinetics-400 [36] and Kinetics-600 [5] respectively. For action recognition, we benchmark our adapted Reversible MViT model and report top-1 and top-5 performance for both datasets. Similar to the image classification benchmark, we observe that the Reversible MViT models closely match the accuracy for their non reversible counterparts at a fraction of the memory cost.

Increasing video model batch sizes. We note that our adapted Reversible MViT forms a very competitive video recognition model. Specifically, for both Kinetics-400 (Table 2) and Kinetics-600 (Table 3) the reversible models

Model	AP ^{box}	AP ^{mask}	Memory(GB)	GFLOPs	Param (M)
Res50 [28]	41.0	37.1	-	260	44
Res101 [28]	42.8	38.5	-	336	63
X101-64 [73]	44.4	39.7	-	493	101
PVT-L [69]	44.5	40.7	-	364	81
MViT-B	48.2	43.9	18.9	668	57
Rev-MViT-B	48.0	43.5	10.9	683	58

Table 4. **Comparison on MS-COCO object detection.** Rev-MViT achieves competitive performance to MViT across all metrics at $1.7\times$ lower memory footprint.

match the overall accuracy very closely at only **51.5%** and **37.2%** of the memory cost respectively.

This allows a batch size increase of $2\times$ on the 16 layer, 70.5 GFLOPs Kinetics-400 MViT-B-16 model and of $3.5\times$ on the 24 layer, 236 GFLOPs Kinetics-600 MViT-B-24 model, a very beneficial result for large video models which are often severely memory limited and trained with very small batch sizes (Table 3). Moreover, due to the more efficient design of stage-transition blocks in Rev-MViT (§3.3.1), *i.e.* bringing the dimension upsampling operation inside the pooling attention instead of being performed in the prior MLP stage [18], the Rev-MViT are also slightly more parameter and FLOP efficient on both Kinetics.

4.3. Object Detection

We benchmark the proposed Rev-ViT-B and Rev-MViT-B models on object detection on MS-COCO [45] as well. All the models are trained on 118K training images and evaluated on the 5K validation images. We take the ViT-B and MViT-B backbones pre-trained on IN and use the standard Mask R-CNN [26] as the detection framework. All models are trained with a standard $3\times$ schedule (36 epochs). For MViT, we integrate the multi-scale backbone with the feature pyramid network [44]. Referring to Table 4 we observe that, the Rev-MViT-B model closely matches the AP performance on MViT-B at only 54.8% of the memory cost.

4.4. Ablations

Stronger Inherent Regularization. Across different models and datasets, we find that at the same FLOP and parameter specifications, the reversible models tend to have stronger inherent regularization than their non-reversible counterparts. Hence, training recipes for reversible vision transformers have lighter repeated augmentations, smaller augmentation magnitudes and consequently, higher weight decay. Table 5 shows the effects of these recipe changes on Rev-ViT-B. We also observe similar effects on other Rev-ViT and Rev-MViT models where a modified training recipe with lighter augmentations and higher weight decay play a crucial role in matching performance.

Lateral Fusion Strategies. The stage-transition blocks employ residual stream fusion blocks for mixing information between I_1 and I_2 (§3.3.1). We explore several fusion

Training Improvement	Train Acc	Top-1 ImageNet Acc
Naïve Rev-ViT-B	15.3	12.1
+ Re-configuring residual streams	82.1	77.2
+ Repeated Augmentation	84.9	80.6
+ Lighter Augmentation magnitude	93.2	81.0
+ Stronger Stochastic Depth	92.0	81.4
+ Higher weight decay	91.0	81.8
Rev-ViT-B	91.0	81.8

Table 5. **Rev-ViT-B Training Recipe.** We observe that reversible transformers tend to have a stronger inherent regularization and require a lighter augmented training recipe for peak performance.

Stage-Transition Fusion	Termination Fusion	Train Acc	Top-1 Acc
Max	Norm \rightarrow Concat	78.1	81.7
Concat	Norm \rightarrow Concat	79.1	82.0
$2\times$ -MLP	Norm $\rightarrow 2\times$ -MLP	80.2	81.8
$2\times$ -MLP + 0.2 dp	Norm $\rightarrow 2\times$ -MLP \rightarrow 0.5dp	77.1	81.2
$2\times$ -MLP	Norm \rightarrow 1-layer	53.6	82.1
$2\times$ -MLP	Norm \rightarrow 1-layer \rightarrow 0.2dp	64.0	82.4
Norm $\rightarrow 2\times$ -MLP	Norm \rightarrow Concat	79.4	82.3
Norm $\rightarrow 2\times$ -MLP	Norm \rightarrow 1-layer \rightarrow 0.2dp \rightarrow Norm	78.3	82.3
$4\times$ -MLP	Norm \rightarrow Concat	80.4	82.3
$2\times$ -MLP	Concat \rightarrow Norm	80.5	82.2
$2\times$ -MLP	Norm \rightarrow Concat	80.1	82.5

Table 6. **Lateral Fusion Strategies.** Residual streams I_1 and I_2 are fused in state-transition blocks (§3.3.1) as well as on termination (§3.2.2) before the network head. We find fusion strategy to play a key role for ReV-MViT performance. Rev-MViT-B uses a 2-layer MLP with $2\times$ hidden dimensions in stage-transition blocks (gray). Please see Section 4.4 for details.

strategies in Table 6 using a combination of: (A) $n\times$ -MLP: Two layer perceptrons with n times the hidden dimension and GELU activations. (B) $0.n$ dp: $n \times 10$ percent dropout on output activations. (C) Simple operators such as channel-wise maximum of I_1 and I_2 activations, and channel-wise concatenation of tensor.

Lateral connections in stage-transition stages allows effective information mixing between the residual streams and hence increases network capacity. For example, compared to concatenation, $2\times$ -MLP increases to training accuracy by 1% and also the top-1 performance by 0.5%. However an even heavier strategy, such as $4\times$ -MLP widens the generalization gap and promotes over-fitting. Note that the training accuracy is often lower than the top-1 performance because of training data augmentations.

Re-configuring residual connections. As discussed in §3.2.3, the reversible vision transformer design removes the skip connections that are commonly used inside the Attention and MLP blocks (Figure 2). Specifically, for all of the reversible blocks in Rev-MViT and Rev-ViT, the inputs I_1 and I_2 do not have residual connections that allow residual signal propagation by directly skipping their respective functions (MLP for I_2 and Attention for I_1). Instead their

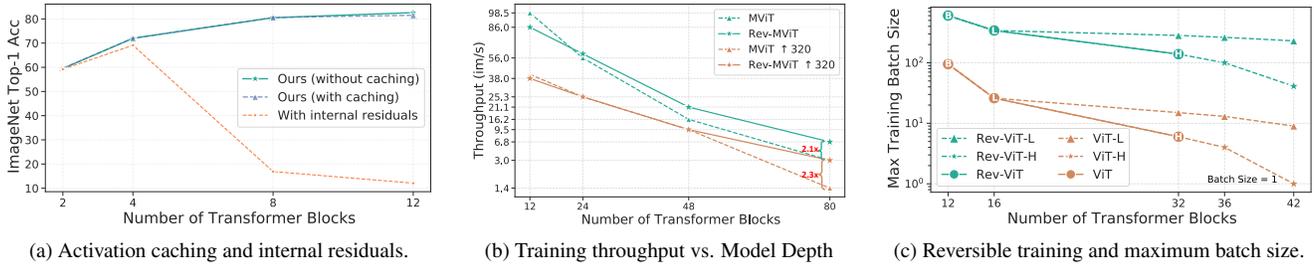


Figure 3. **Ablation Experiments.** (a): We observe that (1) Learning without activation caching does not hurt reversible accuracy for Rev-ViT-B of varying depths and (2) Internal residual connections train well for shallow models but the training *diverges* for deeper models. (b): Rev-MViT has higher *throughput* for higher input resolution and deeper MViT models increasing up to **2.3**× at 224 resolution for 80 layers. (c): We benchmark the maximum batch size for Rev-ViT Base (B), Large (L) and Huge (H) and their non-reversible counterparts.

residuals are performed via the other residual stream operating through the reversible transform T (§3.1.1).

We ablate this design choice for the ViT architecture in Figure 3a. We vary the model depth without changing any other model dimensions and observe the performance of the two reversible variants *i.e.*, with and without internal skip connections. We note that while the naïve version with internal skip connections trains stably for shallow models, for deeper models the accuracy drops significantly. On the other hand, Rev-MViT scales well with depth, just as accurate as the other variant at shallower depths but significantly better with deeper models.

Effect of learning without caching activations. Figure 3a also compares the image classification performance of the Rev-ViT-B architecture trained with and without caching activations. This allows us to disentangle the effect of the proposed residual configurations necessary for reversible vision transformer from any artefacts that might result from learning without caching activations. However, for all depths we notice the Rev-ViT-B performance trained without caching activations to closely track the performance of the same architecture trained with caching. The slight difference at depth 12 results might stem from the training recipe being adapted for the actual Rev-ViT-B architecture trained without activations.

Model size and Input Resolution vs. Throughput. Figure 3b shows the training throughput comparisons for different models sizes at 224 and 320 input resolutions. We note that while for smaller models such as, MViT-B with a depth 12 layers, Rev-MViT-B has a slightly smaller training throughput (98.5 vs. 86.0), the additional re-computation burden of intermediate activations is easily overcome at both higher resolution training as well as for deeper models. In particular, at 224 resolution, the 24-layer and 48-layer Rev-MViT models have similar throughput as the MViT models increasing up to **2.1**× higher throughput at 384

resolution and **2.3**× higher throughput at 384 resolution for the 80 layer depth models. Further, the rate of memory increase for deeper models is much lower for reversible variants than vanilla networks, allowing scaling to much deeper models without any additional training infrastructure burden or memory requirement like with gradient checkpointing or model parallelism.

Maximum batch-size. We benchmark the maximum possible batch size for Rev-ViT Base (B), Large (L) and Huge (H) and their non-reversible counterparts in Fig.3c. We extrapolate the trend (denoted by ----) to larger models by scaling ViT-L and ViT-H in depth (keeping other model dimensions constant) and benchmark the maximum batch size for their reversible counterparts.

Model size vs. GPU memory footprint. Figure 1 plots the GPU Memory footprint for both Rev-ViT and Rev-MViT family of models as well as for several other prior networks such as MViT [18], ViT [15], ResNets and RegNetY [58]. We note that at fixed GFLOPs, reversible variants are extremely memory efficient going up to **4.5**× for MViT and **15.5**× for ViT surpassing prior convolutional variants by orders of magnitude.

5. Conclusion

We present Reversible Vision Transformers, memory-efficient reversible architectural adaptations of ViT and MViT models. We benchmark across several tasks, such as image classification, object detection and video classification and across several metrics, such as model complexity, throughput, accuracy and memory usage. Given any specification, our Rev-ViT and Rev-MViT match the accuracy of non-reversible variants at a tiny fraction of the memory cost while maintaining similar training throughput for smaller models and up to **2.3**× higher throughput for larger models. Specifically, we observe that the Rev-ViT and the Rev-

MViT models achieve upto $15.5\times$ and $4.5\times$ lighter memory footprint than ViT and MViT models respectively.

Follow-Up Work

Chen *et al.* [77] apply the proposed Reversible Vision Transformer architecture design for Reversible Swin Transformers and apply the model for memory-efficient temporal action localization. Temporal action localization involves detecting precise temporal frame positions for the start and end boundaries of an action and hence needs to be performed on a densely sampled video. Dense frame sampling causes GPU memory overheads during training that prohibit finetuning the backbone end-to-end on the temporal action localization task (TAL). Reversible backbone alleviate the memory overhead and allow efficient end-to-end TAL training, thereby providing significant localization performance boost.

Concurrently, [80] proposes a training procedure for reversible transformers that allows speeding up training while ensuring exact replication of the original computation. In particular, [80] proposes to stagger the activation re-computation one transformer block ahead of the gradient computation using the recomputed activations of the previous block. This allows the activations to be available for gradient calculation of the next block, as soon as the previous block finishes. Hence, the gradient calculation step does not need to wait for activation recomputation, *effectively* hiding the latency of the burden of re-computation behind latency of gradient calculations, thus effectively speeding up training. This requires maintaining separate cuda work-streams that process the above two steps asynchronously. Depending on the hardware and computation size, there can be significant speedups from such operator parallelization.

Acknowledgements

The authors would like to thank Harshayu Girase for help with benchmarking models, Amir Gholami, Ajay Jain and Nikita Kiatev for helpful research discussions and reference suggestions, Ajay Jain, Matthew Tancik and Hang Gao for writing discussions and Shubh Gupta, Suzie Petryk, Hang Gao, Abhinav Agarwal, Medhini Narasimhan and Amur Ghosh for proofreading the manuscript.

Appendix

A. Architecture Details

Reversible Vision Transformers Table 7 shows the architectures for all the Reversible Vision Transformer Models. All models closely follow the original ViT architectures [15] in matched performance, parameters, FLOPs and much lower memory footprint (Table 1). Output sizes denote the tensor shapes of the two residual streams at the

end of each reversible Vision Transformer block. Note that even though the intermediate activations are twice the non-reversible variant, the actual memory needed is much lower because of memory reuse in reversible training. Further, the FLOPs are matched since each layer is performed only one of the two streams.

Reversible Multiscale Vision Transformers Table 8 shows the architecture for the Rev-MViT-B model for image classification. The backbone is made-up of two stages – Stage-transition blocks that increase the channel capacity and down-sample the resolution and the reversible Stage-preserving blocks that perform the majority of computation without changing feature dimensions. Similar to Rev-ViT, the output sizes of both the streams are denoted. Fusion blocks operate on Y_1 and Y_2 together, hence operate with computationally light operations (Table 6).

B. Training Settings

ImageNet. Table 9 shows the training recipes for ViT-L and Rev-ViT-L models presented in Table 1. Note that ViT-L is quite heavy with 61.6 GFLOPs and hence we adopt a shorter 200 epochs recipe for faster experiment cycle for developing Rev-ViT-L. Smaller ViT models – ViT-S and ViT-B – are trained according to the Data efficient transformers [63] and are all trained for 300 epochs. Hence, the accuracy difference between ViT-L which achieves 81.5% while ViT-B achieves 81.8% overall. MViT-B model follows the 300 epochs recipe as well proposed in [18].

Kinetics-400 & Kinetics-600. We follow the recipes proposed in [18] to train the Rev-MViT-B architecture (Table 8) following crucial modifications shown in Table 5.

MS-COCO. For object detection experiments, we adopt the Mask R-CNN [26] object detection framework in Detectron2 [72]. We follow the same training settings from [48], AdamW optimizer [50] ($\beta_1, \beta_2 = 0.9, 0.999$, base learning rate $1.6e-4$ for base size of 64, and weight decay of 0.1), and 3x schedule (36 epochs). The drop path rate is set as 0.4. We use PyTorch’s automatic mixed precision during training.

References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proc. ICCV*, 2021. 2, 6
- [2] Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, pages 573–582. PMLR, 2019. 2

stage	operators	output sizes
data		224×224
patch	$1 \times 16 \times 16$, 384 stride $1 \times 16 \times 16$	$384 \times 14 \times 14$
rev	$\begin{bmatrix} \mathbf{F} : \text{MHA}(384) \\ \mathbf{G} : \text{MLP}(1536) \end{bmatrix} \times 12$	$\begin{bmatrix} Y_1 : 384 \times 14 \times 14 \\ Y_2 : 384 \times 14 \times 14 \end{bmatrix}$

(a) **Rev-ViT-S** with **4.6G** FLOPs, **22M** param, **8.8MB/img** memory, and **79.9%** top-1 accuracy.

stage	operators	output sizes
data		224×224
patch	$1 \times 16 \times 16$, 768 stride $1 \times 16 \times 16$	$768 \times 14 \times 14$
rev	$\begin{bmatrix} \mathbf{F} : \text{MHA}(768) \\ \mathbf{G} : \text{MLP}(3072) \end{bmatrix} \times 12$	$\begin{bmatrix} Y_1 : 768 \times 14 \times 14 \\ Y_2 : 768 \times 14 \times 14 \end{bmatrix}$

(b) **Rev-ViT-B** with **17.6G** FLOPs, **87M** param, **17MB/img** memory, and **81.8%** top-1 accuracy.

stage	operators	output sizes
data		224×224
patch	$1 \times 16 \times 16$, 1024 stride $1 \times 16 \times 16$	$1024 \times 14 \times 14$
rev	$\begin{bmatrix} \mathbf{F} : \text{MHA}(1024) \\ \mathbf{G} : \text{MLP}(4096) \end{bmatrix} \times 24$	$\begin{bmatrix} Y_1 : 1024 \times 14 \times 14 \\ Y_2 : 1024 \times 14 \times 14 \end{bmatrix}$

(c) **Rev-ViT-L** with **61.6G** FLOPs, **305M** param, **22.6MB/img** memory, and **81.4%** top-1 accuracy.

Table 7. **Reversible Vision Transformer Architectures:** Rev-ViT are reversible adaption of ViT with exactly matched FLOPs, parameters and accuracy under identical conditions but with much lower GPU memory footprints.

stage	operators	output sizes
data		224×224
cubification	7×7 , 96 stride 4×4	$96 \times 56 \times 56$
Stage-Preserving	$\begin{bmatrix} \mathbf{F} : \text{MHPA}(96) \\ \mathbf{G} : \text{MLP}(384) \end{bmatrix} \times 1$	$\begin{bmatrix} Y_1 : 96 \times 56 \times 56 \\ Y_2 : 96 \times 56 \times 56 \end{bmatrix}$
Stage-Transition	$\begin{bmatrix} \text{FUSION}(192) \\ \text{MHPA}(192) \\ \text{MLP}(768) \end{bmatrix} \times 1$	$192 \times 28 \times 28$
Stage-Preserving	$\begin{bmatrix} \mathbf{F} : \text{MHPA}(192) \\ \mathbf{G} : \text{MLP}(768) \end{bmatrix} \times 1$	$\begin{bmatrix} Y_1 : 192 \times 28 \times 28 \\ Y_2 : 192 \times 28 \times 28 \end{bmatrix}$
Stage-Transition	$\begin{bmatrix} \text{FUSION}(384) \\ \text{MHPA}(384) \\ \text{MLP}(1536) \end{bmatrix} \times 1$	$384 \times 14 \times 14$
Stage-Preserving	$\begin{bmatrix} \mathbf{F} : \text{MHPA}(384) \\ \mathbf{G} : \text{MLP}(1536) \end{bmatrix} \times 10$	$\begin{bmatrix} Y_1 : 384 \times 14 \times 14 \\ Y_2 : 384 \times 14 \times 14 \end{bmatrix}$
Stage-Transition	$\begin{bmatrix} \text{FUSION}(768) \\ \text{MHPA}(768) \\ \text{MLP}(3072) \end{bmatrix} \times 1$	$768 \times 7 \times 7$
Stage-Preserving	$\begin{bmatrix} \mathbf{F} : \text{MHPA}(768) \\ \mathbf{G} : \text{MLP}(3072) \end{bmatrix} \times 1$	$\begin{bmatrix} Y_1 : 768 \times 7 \times 7 \\ Y_2 : 768 \times 7 \times 7 \end{bmatrix}$

Table 8. **Rev-MViT-B** with **8.7G** FLOPs, **39M** param, **66.8MB/img** memory, and **82.5%** top-1 accuracy is reversible adaption of MViT-B architecture [15].

[3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding?

Training Hyperparameter	ViT-B	Rev-ViT-B
Learning Rate	1e-4	7e-5
Random augment Repeats (N)	1	2
Random augment Magnitude (M)	9	7
Optimizer Momentum	(0.9, 0.95)	(0.9, 0.999)
Weight Decay	0.3	0.3
Batch Size	4096	4096
Epochs	200	200
Label Smoothing	0.1	0.1
Drop Path Rate	0.2	0.2
Mixup	0.8	0.8
Cutmix	1.0	1.0

Table 9. **Training Recipe for ViT-L and Rev-ViT-L**

In *Proc. ICCV*, 2021. 2, 6

- [4] Robin Brügger, Christian F Baumgartner, and Ender Konukoglu. A partially reversible u-net for memory-efficient volumetric image segmentation. In *International conference on medical image computing and computer-assisted intervention*, pages 429–437. Springer, 2019. 2
- [5] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proc. CVPR*, 2017. 5, 6
- [6] Bo Chang, Lili Meng, Eldad Haber, Lars Ruthotto, David Begert, and Elliot Holtham. Reversible architectures for arbitrarily deep residual neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018. 2
- [7] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *arXiv preprint arXiv:2106.01345*, 2021. 2
- [8] Yunpeng Chen, Haoqi Fan, Bing Xu, Zhicheng Yan, Yanis Kalantidis, Marcus Rohrbach, Shuicheng Yan, and Jiashi Feng. Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution. In *Proc. CVPR*, 2019. 6
- [9] Zhengsu Chen, Lingxi Xie, Jianwei Niu, Xuefeng Liu, Longhui Wei, and Qi Tian. Visformer: The vision-friendly transformer. In *Proc. ICCV*, 2021. 2
- [10] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’auelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks. *Advances in neural information processing systems*, 25:1223–1231, 2012. 2
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, pages 248–255. Ieee, 2009. 5
- [12] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014. 2
- [13] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016. 2
- [14] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Bain-

- ing Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. *arXiv preprint arXiv:2107.00652*, 2021. 2, 6
- [15] 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. In *Proc. ICLR*, 2021. 1, 2, 4, 5, 8, 9, 10
- [16] Christian Etmann, Rihuan Ke, and Carola-Bibiane Schönlieb. iunets: Fully invertible u-nets with learnable up-and downsampling. *arXiv preprint arXiv:2005.05220*, 2020. 2
- [17] Haoqi Fan, Yanghao Li, Bo Xiong, Wan-Yen Lo, and Christoph Feichtenhofer. Pyslowfast. <https://github.com/facebookresearch/slowfast>, 2020. 5
- [18] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proc. ICCV*, 2021. 1, 2, 5, 6, 7, 8, 9
- [19] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. SlowFast networks for video recognition. In *Proc. ICCV*, 2019. 6
- [20] Marc Finzi, Pavel Izmailov, Wesley Maddox, Polina Kirichenko, and Andrew Gordon Wilson. Invertible convolutional networks. In *Workshop on Invertible Neural Nets and Normalizing Flows, International Conference on Machine Learning*, 2019. 2
- [21] Amir Gholami, Zhewei Yao, Kim Sehoon, Michael W. Mahoney, and Kurt Keutzer. Ai and memory wall. *RiseLab Medium Post*, 2021. 1, 6
- [22] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 2211–2221, 2017. 2
- [23] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 2211–2221, 2017. 4
- [24] Benjamin Graham, Alaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Herve Jegou, and Matthijs Douze. LeViT: A vision transformer in ConvNet’s clothing for faster inference. In *Proc. ICCV*, 2021. 2
- [25] Tristan Hascoet, Quentin Fevrier, Weihao Zhuang, Yasuo Ariki, and Tetsuya Takiguchi. Layer-wise invertibility for extreme memory cost reduction of cnn training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Proc. ICCV*, 2017. 7, 9
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proc. CVPR*, 2015. 2, 4
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016. 1, 7
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Proc. ECCV*, 2016. 6
- [30] Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *International Conference on Machine Learning*, pages 2722–2730. PMLR, 2019. 2
- [31] Mark Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, pages 10–14. IEEE, 2014. 1
- [32] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinulescu, and Douglas Eck. Music transformer. In *ICLR*, 2019. 2
- [33] Jun-Jie Huang and Pier Luigi Dragotti. Winnet: Wavelet-inspired invertible network for image denoising. *arXiv preprint arXiv:2109.06381*, 2021. 2
- [34] Andrei Ivanov, Nikoli Dryden, Tal Ben-Nun, Shigang Li, and Torsten Hoefler. Data movement is all you need: A case study on optimizing transformers. *arXiv preprint arXiv:2007.00072*, 2020. 1
- [35] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver: General perception with iterative attention. *arXiv preprint arXiv:2103.03206*, 2021. 2
- [36] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv:1705.06950*, 2017. 5, 6
- [37] Diederik P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *arXiv preprint arXiv:1807.03039*, 2018. 2
- [38] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*. 2, 3
- [39] Duo Li and Shang-Hua Gao. m-revnet: Deep reversible neural networks with momentum. *arXiv preprint arXiv:2108.05862*, 2021. 2
- [40] Guohao Li, Matthias Müller, Bernard Ghanem, and Vladlen Koltun. Training graph neural networks with 1000 layers. *arXiv preprint arXiv:2106.07476*, 2021. 2
- [41] Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis with transformer network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6706–6713, 2019. 2
- [42] Shanshan Li, Qiang Cai, Zhuangzi Li, Haisheng Li, Naiguang Zhang, and Jian Cao. Attention-aware invertible hashing network. In *International Conference on Image and Graphics*, pages 409–420. Springer, 2019. 2
- [43] Shaohui Li, Ziyang Zheng, Wenrui Dai, Junni Zou, and Hongkai Xiong. Rev-ae: A learned frame set for image reconstruction. In *ICASSP 2020-2020 IEEE Interna-*

- tional Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1823–1827. IEEE, 2020. 2
- [44] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 7
- [45] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Proc. ECCV*, 2014. 5, 7
- [46] Kang Liu, Dong Liu, Li Li, Ning Yan, and Houqiang Li. Semantics-to-signal scalable image compression with learned revertible representations. *International Journal of Computer Vision*, pages 1–17, 2021. 2
- [47] Yang Liu, Zhenyue Qin, Saeed Anwar, Pan Ji, Dongwoo Kim, Sabrina Caldwell, and Tom Gedeon. Invertible denoising network: A light solution for real noise removal. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13365–13374, 2021. 2
- [48] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 6, 9
- [49] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proc. CVPR*, 2022. 2
- [50] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018. 9
- [51] Matthew MacKay, Paul Vicol, Jimmy Ba, and Roger Grosse. Reversible recurrent neural networks. *arXiv preprint arXiv:1810.10999*, 2018. 2
- [52] Daniel Neimark, Omri Bar, Maya Zohar, and Dotan Aselsmann. Video transformer network. *arXiv preprint arXiv:2102.00719*, 2021. 2, 6
- [53] Mandela Patrick, Dylan Campbell, Yuki M Asano, Ishan Misra Florian Metzke, Christoph Feichtenhofer, Andrea Vedaldi, Jo Henriques, et al. Keeping your eye on the ball: Trajectory attention in video transformers. *NIPS*, 2021. 2
- [54] David A Patterson. Latency lags bandwidth. *Communications of the ACM*, 47(10):71–75, 2004. 1
- [55] Mihir Pendse, Vithursan Thangarasa, Vitaliy Chiley, Ryan Holmdahl, Joel Hestness, and Dennis DeCoste. Memory efficient 3d u-net with reversible mobile inverted bottlenecks for brain tumor segmentation. In *International MICCAI Brainlesion Workshop*, pages 388–397. Springer, 2020. 2
- [56] Bas Peters, Eldad Haber, and Keegan Lensink. Fully reversible neural networks for large-scale surface and sub-surface characterization via remote sensing. *arXiv preprint arXiv:2003.07474*, 2020. 2
- [57] Patrick Putzky and Max Welling. Invert to learn to invert. *Advances in Neural Information Processing Systems*, 32:446–456, 2019. 2
- [58] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proc. CVPR*, June 2020. 1, 6, 8
- [59] Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Momentum residual neural networks. *arXiv preprint arXiv:2102.07870*, 2021. 2
- [60] Yang Song, Chenlin Meng, and Stefano Ermon. Mintnet: Building invertible neural networks with masked convolutions. *arXiv preprint arXiv:1907.07945*, 2019. 2
- [61] Bingfeng Sun and Jian Zhang. Invertible image compressive sensing. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, pages 548–560. Springer, 2021. 2
- [62] Yifan Sun, Nicolas Bohm Agostini, Shi Dong, and David Kaeli. Summarizing cpu and gpu design trends with product data. *arXiv preprint arXiv:1911.11313*, 2019. 1
- [63] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In *icml*, 2021. 2, 6, 9
- [64] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proc. ICCV*, 2021. 2
- [65] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *Proc. ICCV*, 2019. 6
- [66] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proc. CVPR*, 2018. 6
- [67] Tycho FA van der Ouderaa and Daniel E Worrall. Reversible gans for memory-efficient image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4720–4728, 2019. 2
- [68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 2
- [69] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proc. ICCV*, 2021. 2, 7
- [70] Samuel Williams, Andrew Waterman, and David Patterson. Roofline: an insightful visual performance model for multi-core architectures. *Communications of the ACM*, 52(4):65–76, 2009. 1
- [71] Samuel Webb Williams. *Auto-tuning performance on multicore computers*. University of California, Berkeley, 2008. 1
- [72] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 9
- [73] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proc. CVPR*, 2017. 7
- [74] Kashu Yamazaki, Vidhiwar Singh Rathour, and T Le. Invertible residual network with regularization for effective medical image segmentation. *arXiv preprint arXiv:2103.09042*, 2021. 2
- [75] Jieming Yang, Hongwei Ge, Jinlong Yang, and Yubing Tong. Image compact-resolution and reconstruction using reversible network. *IET Image Processing*, 14(16):4376–4384, 2020. 2

- [76] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [2](#)
- [77] Chen Zhao, Shuming Liu, Karttikeya Mangalam, and Bernard Ghanem. Re²tal: Rewiring pretrained video backbones for reversible temporal action localization. *arXiv preprint arXiv:2211.14053*, 2022. [9](#)
- [78] Yuekai Zhao, Shuchang Zhou, and Zhihua Zhang. Multi-split reversible transformers can enhance neural machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 244–254, 2021. [2](#)
- [79] Zaixiang Zheng, Hao Zhou, Shujian Huang, Jiajun Chen, Jingjing Xu, and Lei Li. Duplex sequence-to-sequence learning for reversible machine translation. *arXiv preprint arXiv:2105.03458*, 2021. [2](#)
- [80] Tyler Zhu. Speeding up reversible vision transformers. <http://bit.ly/3J6Q0Cb>, 2022. [9](#)