

# Eventful Transformers: Leveraging Temporal Redundancy in Vision Transformers

Matthew Dutson, Yin Li, and Mohit Gupta  
 University of Wisconsin–Madison  
 {dutson, yin.li, mgupta37}@wisc.edu

## Abstract

Vision Transformers achieve impressive accuracy across a range of visual recognition tasks. Unfortunately, their accuracy frequently comes with high computational costs. This is a particular issue in video recognition, where models are often applied repeatedly across frames or temporal chunks. In this work, we exploit temporal redundancy between subsequent inputs to reduce the cost of Transformers for video processing. We describe a method for identifying and re-processing only those tokens that have changed significantly over time. Our proposed family of models, Eventful Transformers, can be converted from existing Transformers (often without any re-training) and give adaptive control over the compute cost at runtime. We evaluate our method on large-scale datasets for video object detection (ImageNet VID) and action recognition (EPIC-Kitchens 100). Our approach leads to significant computational savings (on the order of 2-4x) with only minor reductions in accuracy.

## 1. Introduction

Transformers, initially designed for language modeling [57], have been recently explored as an architecture for vision tasks. Vision Transformers [16] have achieved impressive accuracy across a range of visual recognition problems, attaining state-of-the-art performance in tasks including image classification [16], video classification [1, 2, 18], and object detection [8, 37, 40, 61].

One of the primary drawbacks of vision Transformers is their high computational cost. Whereas typical convolutional networks (CNNs) consume tens of GFlops per image [7], vision Transformers often require an order of magnitude more computation, up to hundreds of GFlops per image. In video processing, the large volume of data further amplifies these costs. High compute costs preclude vision Transformers from deployment on resource-constrained or latency-critical devices, limiting the scope of this otherwise

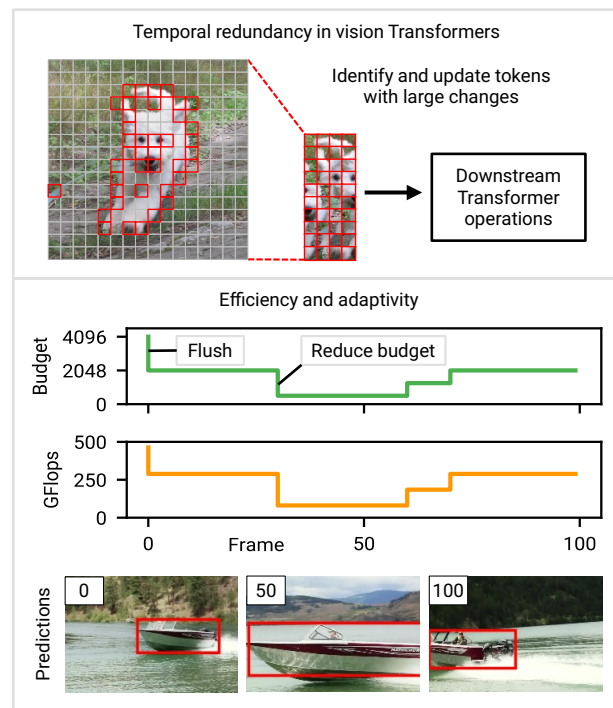


Figure 1. **Eventful Transformers.** Our method exploits temporal redundancy between subsequent model inputs. **(Top)** Within each Transformer block, we identify and update only those tokens with significant changes over time. Image: [5]. **(Bottom)** In addition to improving efficiency, our method gives fine-grained control over the compute cost at runtime. “Budget” refers to parameter  $r$  as described in Section 4.3. “Flush” refers to the initialization of all tokens on the first time step. This example shows the ViTDet [37] object detection model on a video from the VID [54] dataset.

exciting technology. In this paper, we present one of the first methods to use *temporal redundancy between subsequent inputs* to reduce the cost of vision Transformers when applied to video data.

**Temporal redundancy.** Consider a vision Transformer that is applied frame-by-frame or clip-by-clip to a video sequence. This Transformer might be a simple frame-wise

model (e.g., an object detector) or an intermediate step in some spatiotemporal model (e.g., the first stage of the factorized model from [1]). Unlike in language processing, where one Transformer input represents a complete sequence, we consider Transformers applied to several distinct inputs (frames or clips) over time.

Natural videos contain significant temporal redundancy, with only slight differences between subsequent frames. Despite this fact, deep networks (including Transformers) are commonly computed “from scratch” on each frame. This approach is wasteful, discarding all potentially relevant information from previous inferences. Our key intuition is that we can reuse intermediate computations from earlier time steps to improve efficiency on redundant sequences.

**Adaptive inference.** For vision Transformers (and deep networks in general), the inference cost is typically fixed by the architecture. However, in real-world applications, the available resources may vary over time (e.g., due to competing processes or variations in power supply). As such, there is a need for models whose computational cost can be modified at runtime [45]. In this work, adaptivity is one of our primary design objectives; we design our method to allow real-time control over the compute cost. See Figure 1 (bottom portion) for an example where we vary the compute budget throughout a video.

**Challenges and opportunities.** There are past works exploring temporal redundancy [17, 23, 48] and adaptivity [44, 58, 66] for CNNs. However, these methods are generally incompatible with vision Transformers, owing to substantial architectural differences between Transformers and CNNs. Specifically, Transformers introduce a new primitive, self-attention, that does not conform to the assumptions of many CNN-based methods.

Despite this challenge, vision Transformers also represent a unique opportunity. In CNNs, it is difficult to translate sparsity improvements (i.e., the sparsity gained by considering temporal redundancy) into concrete speedups. Doing so requires imposing significant restrictions on the sparsity structure [23] or using custom compute kernels [48]. In contrast, the structure of Transformer operations (centered on manipulating token vectors) makes it easier to translate sparsity into reduced runtime using standard operators.

**Eventful Transformers.** We propose Eventful Transformers, a new class of Transformer that leverages temporal redundancy between inputs to enable efficient, adaptive inference. The term “Eventful” is inspired by event cameras [4, 39], sensors that produce sparse outputs based on scene changes. Eventful Transformers track token-level changes over time, selectively updating the token representations and self-attention maps on each time step. Blocks in an Eventful Transformer include gating modules that allow controlling the number of updated tokens at runtime.

Our method can be applied to off-the-shelf models (generally without re-training) and is compatible with a wide range of video processing tasks. Our experiments demonstrate that Eventful Transformers, converted from existing state-of-the-art models, significantly reduce computational costs while largely preserving the original model’s accuracy. We publicly release our code, which includes PyTorch modules for building Eventful Transformers. See our project page: [visionlab.com/project/eventful-transformers](https://visionlab.com/project/eventful-transformers).

**Limitations.** We demonstrate wall-time speedups on both the CPU and GPU. However, our implementation (based on vanilla PyTorch operators) is likely sub-optimal from an engineering standpoint. With additional effort to reduce overhead (e.g., implementing a fused CUDA kernel for our gating logic), we are confident that the speedup ratios could be further improved. Our method also involves some unavoidable memory overheads. Perhaps unsurprisingly, reusing computation from previous time steps requires maintaining some tensors in memory. These memory overheads are relatively modest; see Section 6 for further discussion.

## 2. Related Work

**Efficient Transformers.** Several past works improve the efficiency of Transformers. Many of these methods focus on reducing the quadratic complexity of self-attention, often using low-rank or sparse approximations [11, 13, 22, 31, 32, 33, 41, 52, 53, 60]. In this work, we consider standard self-attention (with windowing in some cases). Our approach is orthogonal to the above methods.

**Selecting and summarizing tokens.** Some recent works improve the efficiency of vision Transformers by exploiting spatial redundancy within each input. Many of these methods prune or fuse tokens based on a salience measure [19, 21, 38, 46, 49]. A notable example is the Adaptive Token Sampling (ATS) algorithm [19], which has an adaptive computation cost and does not require re-training. Other spatial redundancy methods include adaptive token pooling [3, 42], hierarchical pooling [47], learned tokenization [55], and progressive token sampling [69].

Unlike these works, which consider *spatial* redundancy, our method targets *temporal* redundancy. This makes our work complementary to these approaches. A single model can leverage both spatial and temporal redundancy by only updating tokens that are both salient *and* not temporally repetitive. We illustrate this compatibility in our experiments by building a simple proof-of-concept model that combines temporal and spatial redundancy.

Another related work is Spatiotemporal Token Selection (STTS) [59], which exploits spatiotemporal redundancy for video inference. STTS is intended for models with explicit temporal reasoning that take an entire video as input. In contrast, our method is designed for models that are repet-

itively applied to frames or clips. Compared to STTS, our method covers a wider range of architectures and tasks.

**Temporal redundancy between inputs.** There has been recent work on exploiting inter-frame temporal redundancy in CNNs [9, 17, 23, 48]. While we draw some inspiration from these methods, directly applying them to vision Transformers is not feasible due to significant architectural differences between CNNs and Transformers.

There is limited existing research on exploiting temporal redundancy between subsequent vision Transformers inputs. To our knowledge, the only past work in this area is the Spatiotemporal Gated Transformers (STGT) method [36]. There are two noteworthy differences between STGT and our work. Most notably, STGT only considers temporal redundancy within token-level operations (*e.g.*, token-wise linear transforms), and not within the self-attention operator. Our method accelerates all major Transformer components, including self-attention. Further, STGT uses lossy gating logic that leads to accuracy degradation on long sequences with gradual changes. Our method avoids this issue by employing an improved, reference-based gating mechanism.

**Adaptive neural networks.** Many existing methods add adaptivity to deep CNNs [12, 20, 28, 44, 56, 58, 62, 64, 66, 67]. However, due to architectural differences (*e.g.*, the use of relative position embeddings in Transformers), these methods (*e.g.*, those based on input resizing) often do not translate to vision Transformers.

There has been some recent work on adaptive vision Transformers [43, 63, 68]. These works leverage redundancy within a single input, whereas we consider redundancy between inputs. Unlike our method, these approaches generally require re-training or fine-tuning the model.

**Efficient neural networks.** There is a substantial body of work on improving the efficiency of deep networks. Some works propose efficient CNN architectures [27, 30, 70]. Others use reduced-precision arithmetic [14, 29, 50] or pruning [24, 25, 34, 35]. Our method is loosely connected to pruning; it can be viewed as adaptively pruning redundant tokens on each time step.

### 3. Background: Vision Transformers

In this section, we describe the basic elements of a vision Transformer (see [16] for more details) and define the notation we use throughout the rest of the paper.

A vision Transformer consists of a sequence of Transformer blocks. The input to each block is a list of  $N$ ,  $D$ -dimensional token vectors; we denote this as  $\mathbf{x} \in \mathbb{R}^{N \times D}$ . Before the first Transformer block, a vision Transformer maps each image patch to a token vector using a linear transform. Positional embedding [57] can be injected before the first block [16] or at every block [37].

**A Transformer block.** A Transformer block maps input  $\mathbf{x} \in \mathbb{R}^{N \times D}$  to output  $\mathbf{z} \in \mathbb{R}^{N \times D}$ , according to

$$\mathbf{y} = \text{MSA}(\text{LN}(\mathbf{x})) + \mathbf{x}, \quad (1)$$

$$\mathbf{z} = \text{MLP}(\text{LN}(\mathbf{y})) + \mathbf{y}, \quad (2)$$

where “MSA” denotes multi-headed self-attention. “MLP” is a token-wise multilayer perceptron with two layers and one GELU nonlinearity. “LN” denotes layer normalization.

**Multi-headed self-attention (MSA).** The self-attention operator first applies three linear transforms  $W_q, W_k, W_v \in \mathbb{R}^{D \times D}$  to its input  $\mathbf{x}' = \text{LN}(\mathbf{x})$ .

$$\mathbf{q} = \mathbf{x}'W_q \quad \mathbf{k} = \mathbf{x}'W_k \quad \mathbf{v} = \mathbf{x}'W_v. \quad (3)$$

$\mathbf{q}$ ,  $\mathbf{k}$ , and  $\mathbf{v}$  are the “query,” “key,” and “value” tensors, respectively. In practice,  $W_q, W_k, W_v$  are often fused into a single transform  $W_{qkv} = [W_q, W_k, W_v]$ . These transforms may include a bias; we omit the bias here for brevity.

The self-attention operator then computes a normalized similarity matrix (attention matrix)  $A \in \mathbb{R}^{N \times N}$  between the tokens of  $\mathbf{q}$  and  $\mathbf{k}$ .

$$A = \text{Softmax} \left( \frac{\mathbf{q}\mathbf{k}^T}{\sqrt{D}} \right). \quad (4)$$

Softmax normalization is applied along rows of the matrix.

The MSA output  $\mathbf{y}'$  is an attention-weighted sum of the value tokens  $\mathbf{v}$ , followed by a linear projection  $W_p$ .

$$\mathbf{y}' = (A\mathbf{v})W_p. \quad (5)$$

*Multi-headed self-attention* (as opposed to single-headed self-attention) splits  $\mathbf{q}$ ,  $\mathbf{k}$ , and  $\mathbf{v}$  into  $H$  tensors of shape  $\mathbb{R}^{N \times (D/H)}$  and applies self-attention in parallel across these  $H$  heads. Before applying  $W_p$ , the results of all heads are concatenated into a tensor with shape  $\mathbb{R}^{N \times D}$ .

**Windowed attention.** Standard MSA has a complexity of  $\mathcal{O}(N^2)$  (quadratic in the number of tokens). To reduce this cost, many vision Transformers adopt *windowed* attention. Windowed attention constrains the attention computation to local windows. Information can be exchanged between windows by shifting the windows between blocks [40] or by interleaving global attention [37].

### 4. Eventful Transformers

Our goal is to accelerate vision Transformers for video recognition, in the situation where a Transformer is applied repetitively across frames or chunks of frames (*e.g.*, for video object detection or video action recognition, respectively). Our key idea is to exploit temporal redundancy by re-using computation from previous time steps. In this section, we describe how to modify Transformer blocks to add temporal redundancy awareness.

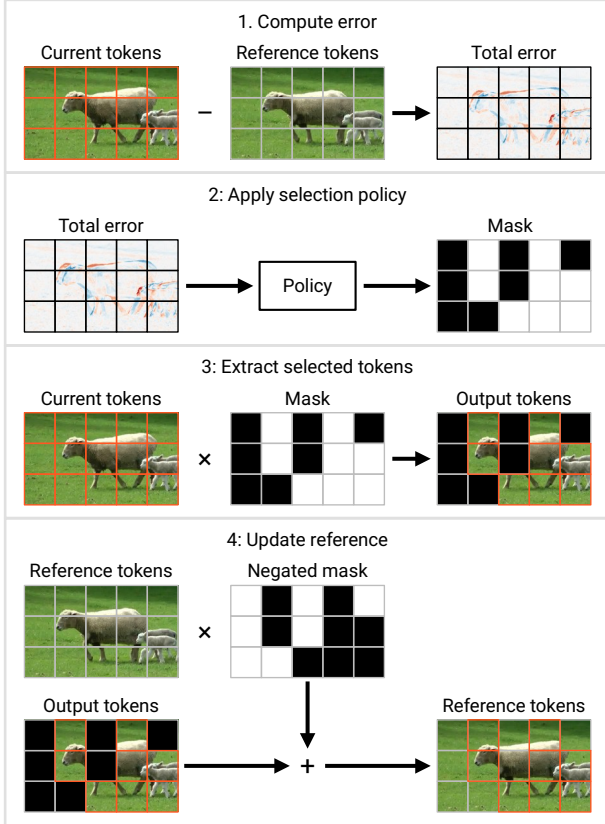


Figure 2. **Token Gating.** A gating module compares incoming tokens against a stored reference. If the difference between a token and its reference is large, then the token is selected to be updated. See Section 4.3 for details on selection policies. Images are from the VID [54] dataset.

In Section 4.1, we present a token-gating module that monitors temporal changes and determines which tokens to update. In Section 4.2, we integrate our token gating logic into a Transformer block, creating a redundancy-aware Eventful Transformer block. In Section 4.3, we explore policies for selecting which tokens to update.

#### 4.1. Token Gating: Detecting Redundancy

In this subsection, we propose two modules: token gates and token buffers. These modules allow us to identify and update only those tokens that have changed significantly since their last update.

**Gate module.** A gate selects  $M \leq N$  of its input tokens to send to downstream layers for re-computation. The gate maintains a set of *reference tokens* in memory, which we denote as  $\mathbf{u} \in \mathbb{R}^{N \times D}$ . The reference tensor contains the value of each token on the time step it was most recently updated. On each time step, tokens are compared against their references; those that deviate significantly from their reference are selected for an update.

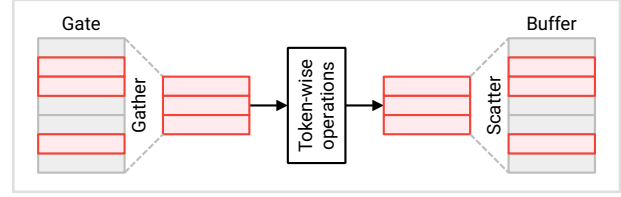


Figure 3. **Accelerating token-wise operations.** The gate reduces the number of active tokens from  $N$  to  $M$ . Subsequent token-wise operations operate on a smaller tensor and therefore have a lower computational cost (proportional to  $M$ ).

Let  $\mathbf{c} \in \mathbb{R}^{N \times D}$  denote the current input to the gate. On each time step, we update the gate’s state and determine its output according to the following procedure (see Figure 2):

1. Compute the total error  $\mathbf{e} = \mathbf{u} - \mathbf{c}$ .
2. Apply a *selection policy* to the error  $\mathbf{e}$ . A selection policy returns a binary mask  $\mathbf{m}$  (equivalently, a list of token indices) indicating which  $M$  tokens should be updated.
3. Extract the tokens selected by the policy. In Figure 2, we depict this as the product  $\mathbf{c} \times \mathbf{m}$ ; in practice, we achieve this with a “gather” operation along the first axis of  $\mathbf{c}$ . We denote the gathered tokens as  $\tilde{\mathbf{c}} \in \mathbb{R}^{M \times D}$ . The gate returns  $\tilde{\mathbf{c}}$  as its output.
4. Update the references for selected tokens. In Figure 2, we depict this as  $\mathbf{u} \leftarrow \mathbf{e} \times (\sim \mathbf{m}) + \mathbf{c} \times \mathbf{m}$ ; in practice, we apply a “scatter” operation from  $\tilde{\mathbf{c}}$  into  $\mathbf{u}$ .

On the first time step, the gate updates all tokens (initializing  $\mathbf{u} \leftarrow \mathbf{c}$  and returning  $\tilde{\mathbf{c}} = \mathbf{c}$ ).

**Buffer module.** A buffer module maintains a state tensor  $\mathbf{b} \in \mathbb{R}^{N \times D}$  that tracks the most recent known value for each of its input tokens. When receiving a new input  $f(\tilde{\mathbf{c}}) \in \mathbb{R}^{M \times D}$ , the buffer scatters the tokens from  $f(\tilde{\mathbf{c}})$  into their corresponding locations in  $\mathbf{b}$ . It then returns the updated  $\mathbf{b}$  as its output. See Figure 3.

We pair each gate with a subsequent buffer. One simple usage pattern is as follows. The gate output  $\tilde{\mathbf{c}} \in \mathbb{R}^{M \times D}$  is passed to a series of token-wise operations  $f(\tilde{\mathbf{c}})$ . The resulting tensor  $f(\tilde{\mathbf{c}}) \in \mathbb{R}^{M \times D}$  is then passed to a buffer, which restores the full shape  $\mathbb{R}^{N \times D}$ .

#### 4.2. Building Redundancy-Aware Transformers

In this subsection, we propose a modified Transformer block that exploits temporal redundancy. Figure 4 shows our design for an Eventful Transformer block. Our method accelerates token-wise operations (e.g., the MLP), as well as the query-key and attention-value multiplications (Equations 4 and 5, respectively).

**Token-wise operations.** Many of the operations in a Transformer block are token-wise, meaning they do not involve information exchange between tokens. These include the

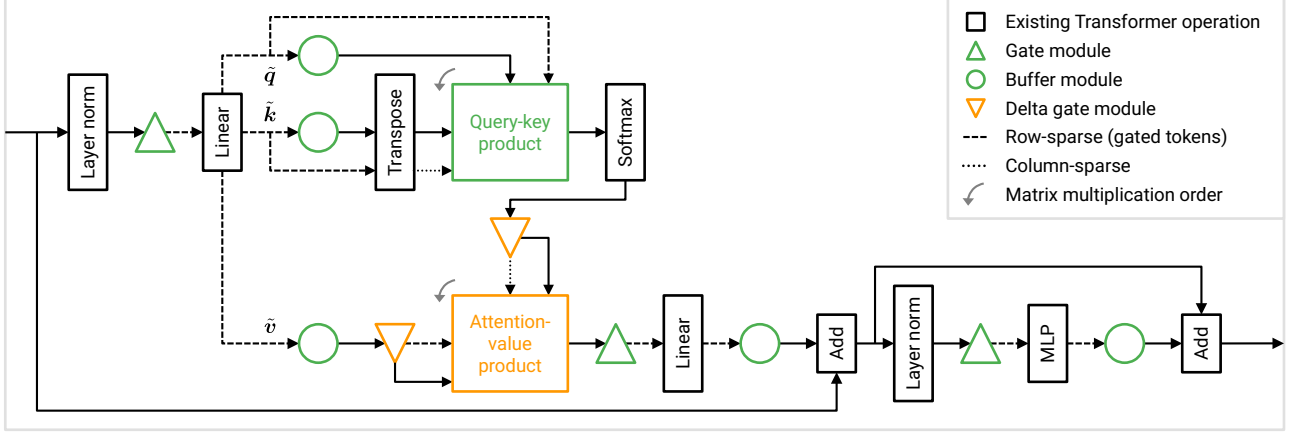


Figure 4. **An Eventful Transformer block.** To exploit temporal redundancy, we strategically apply token gating throughout the block and compute a modified, sparse self-attention update. Rectangles are standard Transformer components (see Section 3). For clarity, we have omitted some minor operations (e.g., scaling after the first matrix multiplication) from this figure.

MLP and the linear transforms in the MSA. We can save computation in token-wise operations by skipping those tokens not selected by the gate. Due to token-wise independence, this does not change the result of the operation for the selected tokens. See Figure 3.

Specifically, we place a gate-buffer pair around each contiguous sequence of token-wise operations, including the  $W_{qkv}$  transform (Equation 3), the  $W_p$  transform (Equation 5), and the MLP. Note that we add buffers before the skip connections (Equations 1 and 2) to ensure that the tokens of the two addition operands are correctly aligned.

The cost of a token-wise operation is proportional to the number of tokens. A gate reduces the number of tokens from  $N$  to  $M$ . This, in turn, reduces the computational cost of downstream token-wise operations by a factor of  $N/M$ .

**The query-key product.** We now consider the query-key product  $B = \mathbf{q}\mathbf{k}^T$  (part of Equation 4). Writing this matrix multiplication explicitly, we have

$$B_{ij} = \sum_p \mathbf{q}_{ip} (\mathbf{k}^T)_{pj}. \quad (6)$$

Element  $B_{ij}$  needs to be updated if either (a) there is a change in the  $i^{\text{th}}$  row of  $\mathbf{q}$ , or (b) there is a change in the  $j^{\text{th}}$  column of  $\mathbf{k}^T$ . Due to the gate that we inserted before the  $W_{qkv}$  transform (shown in Figure 4), only some rows of  $\mathbf{q}$  and some columns of  $\mathbf{k}^T$  have changed. Therefore, we only need to recompute a subset of the elements of  $B$ .

Let  $\tilde{\mathbf{x}}' \in \mathbb{R}^{M \times D}$  denote the output of the gate before the  $W_{qkv}$  transform. We define  $\tilde{\mathbf{q}} = \tilde{\mathbf{x}}'W_q$  and  $\tilde{\mathbf{k}} = \tilde{\mathbf{x}}'W_k$  (following Equation 3). Let  $\mathbf{q}$  and  $\mathbf{k}$  denote the outputs of the  $\tilde{\mathbf{q}}$  and  $\tilde{\mathbf{k}}$  buffers (shown in Figure 4).  $\tilde{\mathbf{q}}$  contain  $\tilde{\mathbf{k}}$  the subset of tokens from  $\mathbf{q}$  and  $\mathbf{k}$  that are being updated.

Figure 5 depicts our method for sparsely updating  $B$ . The product  $\tilde{\mathbf{q}}\tilde{\mathbf{k}}^T$  contains the elements of  $B$  that need to

be updated due to a change in  $\tilde{\mathbf{q}}$ . We compute  $\tilde{\mathbf{q}}\tilde{\mathbf{k}}^T$ , then scatter the result row-wise into the old  $B$  (the value of  $B$  from the last time step). We use an analogous approach for the  $\tilde{\mathbf{k}}$ -induced updates; we compute  $\mathbf{q}\tilde{\mathbf{k}}^T$  and scatter the result column-wise into  $B$ .

The overall cost of these updates is  $2NMD$ , compared to a cost of  $N^2D$  to compute  $B$  from scratch. Note that the cost of our method is proportional to  $M$ , the number of tokens selected by the gate. We save computation when  $M < N/2$  (when we update fewer than half of the tokens).

The above method for updating  $B$  involves some redundant computation. Some elements of the first scattered matrix  $\tilde{\mathbf{q}}\tilde{\mathbf{k}}^T$  are also present in the second matrix  $\mathbf{q}\tilde{\mathbf{k}}^T$ . These overlapping elements are computed twice. Eliminating this redundancy would reduce the cost of our method to  $NMD \leq N^2D$ . This could be achieved by removing the tokens in  $\tilde{\mathbf{q}}$  from  $\mathbf{q}$  before computing  $\mathbf{q}\tilde{\mathbf{k}}^T$ . We would then scatter the result by indexing along both axes of  $B$ . We leave this as an optimization for future implementations.

**The attention-value product.** We now describe a method for updating the attention-value product  $A\mathbf{v}$  (part of Equation 5). Writing this multiplication explicitly, we have

$$(A\mathbf{v})_{ij} = \sum_p A_{ip} \mathbf{v}_{pj}. \quad (7)$$

Because of the gate before the  $W_{qkv}$  transform, only some rows (tokens) of  $\mathbf{v}$  change on each time step. However, there are some updated values in every *column* of  $\mathbf{v}$ . Therefore, every element of  $A\mathbf{v}$  will change on each time step. This means we cannot use the same strategy that we used for  $B$ , where we only updated some of the output elements.

Instead, we propose a delta-based update strategy. Let  $A_o$  and  $\mathbf{v}_o$  denote the last known values for  $A$  and  $\mathbf{v}$ . Let  $A_\Delta$  and  $\mathbf{v}_\Delta$  denote changes in  $A$  and  $\mathbf{v}$ . Define  $A_n =$

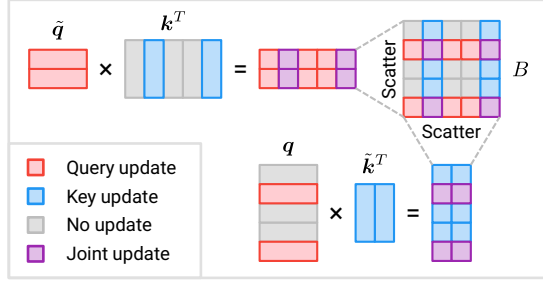


Figure 5. **The query-key product.** We reduce the cost of computing  $B = \mathbf{q}\mathbf{k}^T$  by only updating a subset of its elements. We first compute changes induced by updated rows in  $\mathbf{q}$  (top-left), then compute changes induced by updated columns in  $\mathbf{k}^T$  (bottom).

$A_o + A_\Delta$  and  $\mathbf{v}_n = \mathbf{v}_o + \mathbf{v}_\Delta$ . We can compute the updated attention-value product  $A_n \mathbf{v}_n$  as

$$\begin{aligned}
 A_n \mathbf{v}_n &= (A_o + A_\Delta)(\mathbf{v}_o + \mathbf{v}_\Delta) \\
 &= A_o \mathbf{v}_o + A_o \mathbf{v}_\Delta + A_\Delta \mathbf{v}_o + A_\Delta \mathbf{v}_\Delta \\
 &= A_o \mathbf{v}_o + (A_o + A_\Delta) \mathbf{v}_\Delta + A_\Delta (\mathbf{v}_o + \mathbf{v}_\Delta) - A_\Delta \mathbf{v}_\Delta \\
 &= A_o \mathbf{v}_o + A_n \mathbf{v}_\Delta + A_\Delta \mathbf{v}_n - A_\Delta \mathbf{v}_\Delta. \quad (8)
 \end{aligned}$$

Therefore, on each time step, we can update  $A\mathbf{v}$  by adding  $A_n \mathbf{v}_\Delta + A_\Delta \mathbf{v}_n - A_\Delta \mathbf{v}_\Delta$  to the previous result  $A_o \mathbf{v}_o$ .

We obtain  $A_\Delta$ ,  $\mathbf{v}_\Delta$ ,  $A_o$ , and  $\mathbf{v}_o$  using *delta gate modules*. Delta gates are similar to the gates defined in Section 4.1, with one difference: instead of returning  $\tilde{\mathbf{c}}$ , a delta gate returns  $\mathbf{u}$  and  $\tilde{\mathbf{e}}$  (where  $\tilde{\mathbf{e}}$  is the result of gathering the selected indices from  $\mathbf{e}$ ).  $\mathbf{u}$  represents the effective current value of the gate’s output, corresponding to  $A_n$  or  $\mathbf{v}_n$  in Equation 8.  $\tilde{\mathbf{e}}$  represents the amount of change on the current time step, corresponding to  $A_\Delta$  or  $\mathbf{v}_\Delta$  in Equation 8.

Figure 6 illustrates our approach for efficiently computing the three delta terms in Equation 8. We remove the columns of  $A_n$  that correspond to zero rows in  $\mathbf{v}_\Delta$  (these columns will always be multiplied by zero). Let  $\tilde{A}_n$  denote  $A_n$  with these columns removed. We remove rows of  $\mathbf{v}_n$  analogously to produce  $\tilde{\mathbf{v}}_n$ . We then compute

$$\tilde{A}_n \tilde{\mathbf{v}}_\Delta + \tilde{A}_\Delta \tilde{\mathbf{v}}_n - \tilde{A}_\Delta \tilde{\mathbf{v}}_\Delta, \quad (9)$$

adding the result to the previous value of  $A\mathbf{v}$ .

The product  $\tilde{A}_\Delta \tilde{\mathbf{v}}_\Delta$  assumes the columns of  $\tilde{A}_\Delta$  are correctly aligned with the rows of  $\tilde{\mathbf{v}}_\Delta$ . We achieve this alignment by forcing the  $A$  gate to select the same indices as the  $\mathbf{v}$  gate. Using a separate policy in the  $A$  gate would be possible, but would require a re-alignment operation before computing  $\tilde{A}_\Delta \tilde{\mathbf{v}}_\Delta$ . Further, forcing alignment allows us to eliminate a multiplication by rearranging Equation 9 as

$$\tilde{A}_n \tilde{\mathbf{v}}_\Delta + \tilde{A}_\Delta (\tilde{\mathbf{v}}_n - \tilde{\mathbf{v}}_\Delta). \quad (10)$$

Equation 10 has a cost of  $2MND$  (assuming the addition has a negligible cost), compared to  $N^2D$  for a standard multiplication. We see savings when  $M < N/2$ .

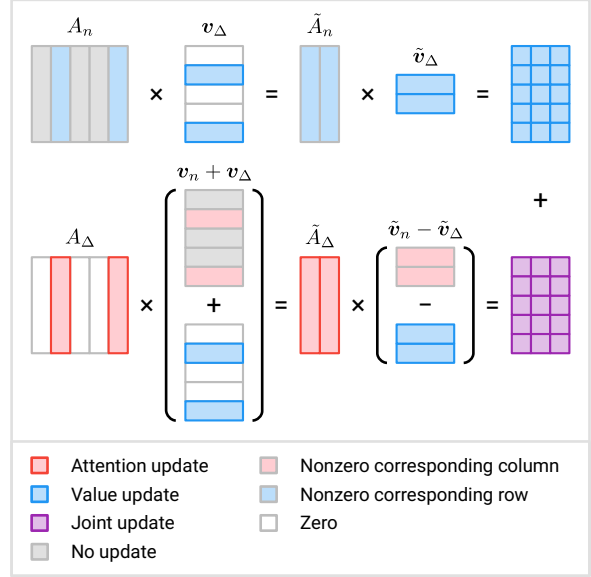


Figure 6. **The attention-value product.** We propose a delta-based strategy for sparsely updating the product  $A\mathbf{v}$ . We reduce the cost of each sub-product by cutting rows and columns that do not contribute to the result (due to a zero multiplication).

### 4.3. Token Selection Policies

An important design choice for an Eventful Transformer is the token selection policy. Given a gate error tensor  $\mathbf{e}$ , a policy generates a mask  $\mathbf{m}$  indicating which tokens should be updated. We now discuss the design of selection policies.

**Top- $r$  policy.** This policy selects the  $r$  tokens whose error  $\mathbf{e}$  has the largest norm (we use the L2 norm). The top- $r$  policy is lightweight and has a single parameter that can be easily tuned by hand. Varying  $r$  gives direct control over the model’s computation cost. These properties make the top- $r$  policy a good fit for applications with tight (potentially time-varying) computational constraints. We use a top- $r$  policy in our main experiments.

**Threshold policy.** This policy selects all tokens where the norm of the error  $\mathbf{e}$  exceeds a threshold  $h$ . A threshold policy is input-adaptive; the number of tokens selected depends on the amount of change in the scene. This input adaptivity can potentially lead to a better accuracy-cost tradeoff. However, the best value for the threshold  $h$  depends on the distribution of token vectors (which varies across layers) and is difficult to decide. In addition, a threshold policy does not give a fixed compute cost. This policy is likely better suited to applications with more flexible resources, where achieving the best possible accuracy-cost tradeoff is critical.

**Other policies.** More sophisticated token selection policies could lead to an improved accuracy-cost tradeoff. For example, we could use a learned policy (e.g., a lightweight

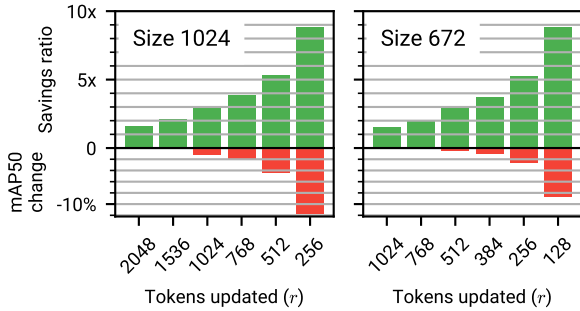


Figure 7. **Video object detection results.** Computation savings ratio (positive axis) and relative reductions in mAP50 score (negative axis) for our method. Results are for the ViTDet model [37] on the VID [54] dataset. See the supplement for tables.

policy network). However, training the policy’s decision mechanism might be challenging, due to the general non-differentiability of the binary mask  $m$ . Another idea is to use an importance score (e.g., as proposed in [19]) to inform the selection. We leave these ideas as potential topics for future work.

## 5. Experiments

In this section, we present our experiments and results. We evaluate our method for video object detection (Section 5.1) and video action recognition (Section 5.2). We show additional analysis in Section 5.3.

### 5.1. Video Object Detection

**Task and dataset.** We test our method on video object detection using the ILSVRC 2015 ImageNet VID dataset [54]. We report results on the validation set, which contains 555 videos with lengths of up to 2895 frames. Following prior works [10, 51], we evaluate the mean average precision (mAP) metric with an IoU threshold of 0.5.

**Implementation details.** We consider the ViTDet model from [37], which we apply to individual frames of an input video. ViTDet combines a plain Transformer backbone (based on ViT-B [16]) with a standard detection head [6, 26]. The backbone consists of 12 blocks with interleaved global and windowed self-attention (blocks 3, 6, 9, and 12 use global attention). Windowed self-attention uses a window size of  $14 \times 14$  tokens ( $224 \times 224$  pixels). Token vectors are 768-dimensional. Self-attention operators have 12 heads and employ learned relative position embeddings.

Before the backbone, the model maps each  $16 \times 16$  image patch to a token vector using a linear transform. The model expects fixed-size inputs (due to resolution-specific position embeddings). Therefore, following from [37], we rescale and pad all video frames to a uniform size (e.g.,  $1024 \times 1024$ ) before applying the model.

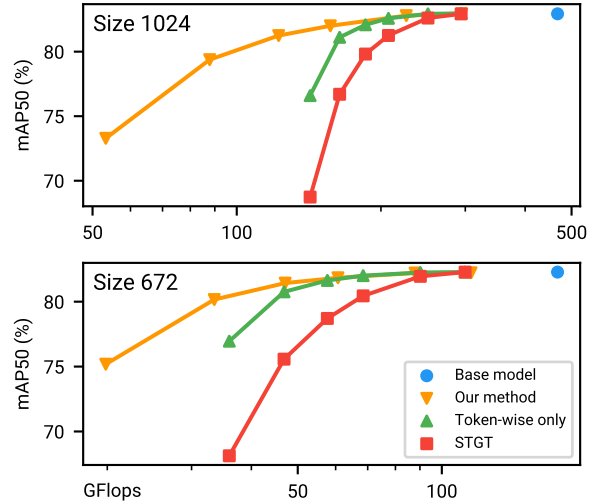


Figure 8. **Video object detection comparison and ablation.** The accuracy-cost tradeoff for our method, compared with STGT [36] and an ablation that only accelerates token-wise operations. See the supplement for tables.

We convert the model to an Eventful Transformer following the method in Section 4. In blocks that use windowed attention, we exploit temporal redundancy only within token-wise operations (not within the query-key or attention-value products). Our complete approach is compatible with windowed attention; however, windowing leads to some implementation challenges (ragged tensor shapes across windows, making batched computation more difficult). Note that for ViTDet, global self-attention represents the bulk of the self-attention compute cost.

**Experiment protocol and baselines.** We fine-tune the original ViTDet weights (trained on COCO) for VID object detection. See the supplement for training parameters. Note that we fine-tune *before* we add temporal redundancy awareness to the model. We train and evaluate at resolution  $1024 \times 1024$ . To understand the effect of token count (which strongly influences compute cost), we also evaluate at resolution  $672 \times 672$ . Rather than training a separate lower-resolution model, we adapt the  $1024 \times 1024$  model by interpolating the learned position embeddings. The resulting adapted model retains most of its accuracy.

We compare against a version of the STGT method [36]. Due to unavailable source code, we were unable to evaluate all components of this method (notably, the use of a learned policy network). Instead, we consider a simplified version that uses the same top- $r$  policy as our method. This setup enables a direct comparison of the core gating and update mechanisms. In addition, we evaluate an ablated version of our approach that only accelerates token-wise operations. We vary the policy  $r$  to explore the accuracy-compute trade-

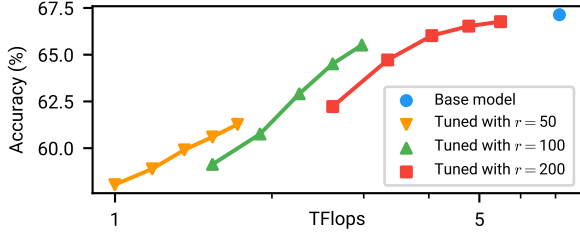


Figure 9. **Video action recognition results.** Our results for action recognition on EPIC-Kitchens 100 using the ViViT model [1]. We report the total TFlops per video (spatial + temporal sub-models). See the supplement for a table containing this data.

off. At resolution 1024, we test  $r = 256, 512, 768, 1024, 1536,$  and  $2048$  (from a maximum of 4096 tokens). At resolution 672, we test  $r = 128, 256, 384, 512, 768,$  and  $1024$  (from a maximum of 1764).

**Results.** Figure 7 shows our results. Our method gives significant savings with only minor reductions in accuracy. For example, at size 1024 with  $r = 768$ , our approach reduces the cost from 467.4 GFlops to 122.3 GFlops (3.8x lower) while reducing the mAP50 score from 82.93 to 81.25 (-1.68% in absolute mAP50). At size 672 with  $r = 384$ , we reduce the cost by 3.7x with a -0.85% change in mAP50.

In these experiments, some tokens correspond to padded space and are therefore “easy” from a temporal redundancy standpoint. However, even in a padding-free deployment (e.g., with a single, known training and inference resolution) our method would still give strong computation savings. For example, consider resolution 1024 with  $r = 768$ . We are skipping 66% of all non-padded tokens here (based on a measured mean padding ratio of 44.6% on VID – corresponding to a  $\sim 16:9$  aspect ratio). This corresponds to a savings of  $>2x$ , with an accuracy drop of only 1.68%. Note that our ViViT experiments (Section 5.2 and supplementary) do not involve padding.

Figure 8 shows the accuracy-compute tradeoff for our method, along with baselines. Our approach gives a considerable improvement in the accuracy-compute tradeoff compared to STGT [36]. Further, adding redundancy awareness to the query-key and attention-value products reduces the cost significantly, especially at low  $r$  values.

## 5.2. Video Action Recognition

**Task and dataset.** We evaluate our method on action recognition using the EPIC-Kitchens 100 dataset [15]. EPIC-Kitchens 100 contains highly dynamic egocentric videos annotated with 97 verb and 300 noun classes. We consider the verb classification task. The training and validation set contains 67217 and 9668 action instances, respectively.

**Implementation details.** We use the ViViT model [1] with factorized spatial and temporal sub-models based on ViT-B.

Table 1. **Adding spatial redundancy to ViTDet.** “Spatial” is a model with pooling in  $k$  and  $v$ . “Spatiotemporal” is a model with both pooling and temporal redundancy awareness.

Variant	$r$	mAP50 (%)	GFlops
Base model	–	82.93	467.4
Spatial	–	80.15	388.1
Spatiotemporal	2048	80.14	217.0
Spatiotemporal	1536	80.07	169.3
Spatiotemporal	1024	79.50	121.0
Spatiotemporal	768	78.69	96.3
Spatiotemporal	512	76.96	70.9
Spatiotemporal	256	71.35	44.5

The spatial sub-model (the bulk of the compute cost) is applied sequentially to 16 2-frame input clips. The outputs of the spatial model are concatenated and passed to the temporal model, which returns a class prediction. The prediction is the average over 12 video views (4 temporal views, each divided into 3 spatial crops). Each view has a shape of  $320 \times 320 \times 32$ . Unlike ViTDet, ViViT adds a class embedding token (see [16]), does not use windowed self-attention, and does not use relative position embeddings.

We convert the spatial model to an Eventful Transformer. Naively replacing the spatial model with an Eventful version leads to a considerable drop in accuracy (about -10% with  $r = 100$ ). We conjecture that the cause is a distribution shift in the inputs to the temporal model (see the supplement for further discussion). We recover most of the lost accuracy by fine-tuning the *non-Eventful* temporal model on the outputs of a frozen Eventful spatial model.

**Experiment protocol.** We start with ViViT pre-trained on EPIC-Kitchens 100 and fine-tune the temporal model as described above (on the EPIC-Kitchens training set). We fine-tune different model variants with policy  $r$  values of 50, 100, and 200 (out of a maximum of 401 tokens). See the supplement for training parameters. We report results using the top-1 accuracy metric, following standard protocol [15].

**Results.** Figure 9 shows our results for the Eventful ViViT model. We evaluate a range of  $r$  values for each of the fine-tuned variants. We test the original fine-tuned  $r$ -value, along with  $\pm 20\%$  and  $\pm 40\%$  of this value. We observe considerable computation savings with only moderate reductions in accuracy. For example, with  $r = 140$ , we reduce the cost by 2.4x while reducing the accuracy by only 1.62%. In addition, the model retains adaptivity despite being fine-tuned with a single- $r$  value, exhibiting a favorable accuracy-compute tradeoff over a range of  $r$ -values.

## 5.3. Spatial Redundancy and Runtime

**Considering spatial redundancy.** Eventful Transformers exploit *temporal* redundancy and thus complement prior works that leverage *spatial* redundancy. Here we present



Table 2. **Runtimes (ms)**. ViTDet runtimes are for the Transformer backbone only. ViViT runtimes include the temporal sub-model.

Model	Size	Variant	$r$	GPU	CPU
ViTDet	1024	Base model	–	86.6	5150
ViTDet	1024	Spatial	–	58.9	3116
ViTDet	1024	Temporal	512	69.9	3570
ViTDet	1024	Spatiotemporal	512	38.1	1682
ViTDet	672	Base model	–	28.3	1492
ViTDet	672	Spatial	–	23.3	1055
ViTDet	672	Temporal	256	21.6	838
ViTDet	672	Spatiotemporal	256	20.8	478
ViViT	320	Base model	–	950	5.45e4
ViViT	320	Temporal	50	545	2.15e4

a simple proof-of-concept experiment that considers spatial redundancy in Eventful Transformers.

Specifically, we adopt a variant of [65], which applies spatial pooling to the self-attention key and value tokens. We apply this method with  $2 \times 2$  pooling to the global self-attention operators in the ViTDet model. We evaluate this method both with and without temporal redundancy awareness. In the temporal-redundancy model, we pool  $k$  and  $v$  after their respective buffers. We pool  $\tilde{k}$  by first pooling the active indices (equivalent to max-pooling the mask  $m$ ), then gathering the buffered  $k$  using the pooled indices.

Table 1 shows our results for resolution 1024 (see the supplement for resolution 672). We see that the spatial and temporal methods are complementary; both meaningfully contribute to reducing the computational cost. See Section 6 for further discussion on spatial redundancy methods.

**Runtime.** We show preliminary runtime results on a CPU (Xeon Silver 4214, 2.2 GHz) and a GPU (NVIDIA RTX 3090). See the supplementary material for experiment details. Table 2 shows our results. Adding temporal redundancy awareness leads to speedups of up to 1.74x on the GPU and 2.48x on the CPU. These results should be seen just as a proof of concept – we are confident that these speedups could be improved with further engineering effort (e.g., by replacing vanilla PyTorch operators with custom kernels or using a high-performance inference framework).

**Visualization of updates.** Figure 10 shows an example video sequence. We visualize the model predictions (top), the token-wise L2 norm of the error  $e$  (middle), and the update mask  $m$  (bottom). We see that larger error values correspond to dynamic regions in the image.

## 6. Discussion

**Memory overhead.** Our method reduces floating point operations at the cost of higher memory usage. Each gate or buffer maintains a reference tensor ( $u$  or  $b$ , respectively).

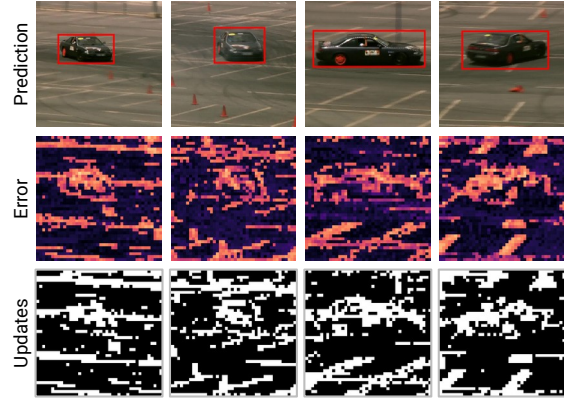


Figure 10. **Visualization of updates.** The error  $e$  and update mask  $m$ , for the pre-QKV gate in block 3 of ViTDet. Video source: [54]

The memory overhead for token gates and buffers is generally modest. For, consider size-1024 ViTDet. The model has 4096 768-dimensional tokens, meaning a token gate or buffer takes 12.6/6.3 MB of memory at full/half precision.

However, gating or buffering the attention matrix  $A$  can require a larger amount of memory. For example, in the global attention layers of the size-1024 ViTDet model, the  $A$  matrix has shape  $4096 \times 4096 \times 12$ . Buffering this  $A$  requires 805/403 MB at full/half precision. Fortunately, the situation dramatically improves if we reduce the number of tokens or use windowed attention (due to the quadratic size of  $A$ ). For example, the  $A$  matrices for the size-672 ViTDet model (1764 tokens) and the ViViT model (301 tokens) occupy 149/75 MB and 4.3/2.2 MB, respectively. In applications where memory is tight and the  $A$  matrix is large, it is possible to save memory by removing temporal redundancy awareness in the query-key and/or attention-value products (each eliminates one  $A$ -shaped state tensor).

**Integration with spatial redundancy methods.** A promising avenue for future work is the further integration of our approach with spatial redundancy methods. Conceptually, these methods summarize a large, redundant set of tokens using a more compact set of tokens. The gating module in an Eventful Transformer assumes that most of its inputs vary smoothly over time. When combining our approach with spatial redundancy methods, we need to ensure that the compact spatial summary is relatively stable. For example, with adaptive token clustering [3, 42], we would need to sort the clusters in a mostly-consistent order.

There are many potentially interesting questions regarding joint modeling of spatial and temporal redundancy. For example, how is the temporal axis different from the spatial one? Should the temporal axis be modeled separately? We leave such questions for future work.

**Acknowledgments.** This research was supported in part by NSF CAREER award #1943149.

## References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. ViViT: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6836–6846, October 2021. [1](#), [2](#), [8](#)
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? arXiv, 2021. [1](#)
- [3] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your ViT but faster. arXiv, 2022. [2](#), [9](#)
- [4] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A 240x180 130 db 3  $\mu$ s latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014. [2](#)
- [5] BuzzFarmers. [flickr.com/photos/buzzfarmers/7318008726](https://www.flickr.com/photos/buzzfarmers/7318008726), 2011. Online, accessed March 2023, CC BY 2.0 license. [1](#)
- [6] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: High quality object detection and instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(5):1483–1498, 2021. [7](#)
- [7] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An analysis of deep neural network models for practical applications. arXiv, 2016. [1](#)
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 213–229, Cham, 2020. Springer International Publishing. [1](#)
- [9] Lukas Cavigelli, Philippe Degen, and Luca Benini. CBinfer: Change-based inference for convolutional neural networks on video data. In *Proceedings of the 11th International Conference on Distributed Smart Cameras*. Association for Computing Machinery, September 2017. [3](#)
- [10] Yihong Chen, Yue Cao, Han Hu, and Liwei Wang. Memory enhanced global-local aggregation for video object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [7](#)
- [11] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. arXiv, 2019. [2](#)
- [12] Ting-Wu Chin, Ruizhou Ding, and Diana Marculescu. AdaScale: Towards real-time video object detection using adaptive scaling. In A. Talwalkar, V. Smith, and M. Zaharia, editors, *Proceedings of Machine Learning and Systems*, volume 1, pages 431–441, 2019. [3](#)
- [13] Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarnlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with Performers. In *International Conference on Learning Representations*, 2021. [2](#)
- [14] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. BinaryConnect: Training deep neural networks with binary weights during propagations. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. [3](#)
- [15] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The EPIC-KITCHENS dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. [8](#)
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. [1](#), [3](#), [7](#), [8](#)
- [17] Matthew Dutton, Yin Li, and Mohit Gupta. Event neural networks. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 276–293, Cham, 2022. Springer Nature Switzerland. [2](#), [3](#)
- [18] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6824–6835, October 2021. [1](#)
- [19] Mohsen Fayyaz, Soroush Abbasi Koohpayegani, Farnoush Rezaei Jafari, Sunando Sengupta, Hamid Reza Vaezi Joze, Eric Sommerlade, Hamed Pirsiavash, and Jürgen Gall. Adaptive token sampling for efficient vision transformers. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 396–414, Cham, 2022. Springer Nature Switzerland. [2](#), [7](#)
- [20] Michael Figurnov, Maxwell D. Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry Vetrov, and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. [3](#)
- [21] Saurabh Goyal, Anamitra Roy Choudhury, Saurabh Raje, Venkatesan Chakaravarthy, Yogish Sabharwal, and Ashish Verma. PoWER-BERT: Accelerating BERT inference via progressive word-vector elimination. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3690–3699. PMLR, July 2020. [2](#)
- [22] Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. Star-transformer. arXiv, 2019. [2](#)
- [23] Amirhossein Habibi, Davide Abati, Taco S. Cohen, and Babak Ehteshami Bejnordi. Skip-convolutions for efficient

- video processing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2695–2704, June 2021. [2](#), [3](#)
- [24] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. [3](#)
- [25] Babak Hassibi and David Stork. Second order derivatives for network pruning: Optimal brain surgeon. In S. Hanson, J. Cowan, and C. Giles, editors, *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann, 1992. [3](#)
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, October 2017. [7](#)
- [27] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. arXiv, 2017. [3](#)
- [28] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Weinberger. Multi-scale dense networks for resource efficient image classification. In *International Conference on Learning Representations*, 2018. [3](#)
- [29] Kyu Yeon Hwang and Wonyong Sung. Fixed-point feedforward deep neural network design using weights +1, 0, and -1. In *2014 IEEE Workshop on Signal Processing Systems (SiPS)*, pages 1–6, 2014. [3](#)
- [30] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and 0.5MB model size. arXiv, 2016. [3](#)
- [31] Sebastian Jaszczur, Aakanksha Chowdhery, Afroz Mohiuddin, LUKASZ KAISER, Wojciech Gajewski, Henryk Michalewski, and Jonni Kanerva. Sparse is enough in scaling transformers. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 9895–9907. Curran Associates, Inc., 2021. [2](#)
- [32] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast autoregressive transformers with linear attention. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5156–5165. PMLR, July 2020. [2](#)
- [33] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020. [2](#)
- [34] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann, 1989. [3](#)
- [35] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient ConvNets. In *International Conference on Learning Representations*, 2017. [3](#)
- [36] Yawei Li, Babak Ehteshami Bejnordi, Bert Moons, Tijmen Blankevoort, Amirhossein Habibian, Radu Timofte, and Luc Van Gool. Spatio-temporal gated transformers for efficient video processing. In *Advances in Neural Information Processing Systems Workshops*, 2021. [3](#), [7](#), [8](#)
- [37] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 280–296, Cham, 2022. Springer Nature Switzerland. [1](#), [3](#), [7](#)
- [38] Youwei Liang, Chongjian GE, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. EViT: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2022. [2](#)
- [39] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128x128 120 db 15  $\mu$ s latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008. [2](#)
- [40] 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 (ICCV)*, pages 10012–10022, October 2021. [1](#), [3](#)
- [41] Jiachen Lu, Jinghan Yao, Junge Zhang, Xiatian Zhu, Hang Xu, Weiguo Gao, Chunjing XU, Tao Xiang, and Li Zhang. SOFT: Softmax-free transformer with linear complexity. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 21297–21309. Curran Associates, Inc., 2021. [2](#)
- [42] Dmitrii Marin, Jen-Hao Rick Chang, Anurag Ranjan, Anish Prabhu, Mohammad Rastegari, and Oncel Tuzel. Token pooling in vision transformers. arXiv, 2021. [2](#), [9](#)
- [43] Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-Nam Lim. AdaViT: Adaptive vision transformers for efficient image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12309–12318, June 2022. [3](#)
- [44] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. AR-Net: Adaptive frame resolution for efficient action recognition. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 86–104, Cham, 2020. Springer International Publishing. [2](#), [3](#)
- [45] Sparsh Mittal. A survey of techniques for approximate computing. *ACM Computing Surveys*, 48(4), March 2016. [2](#)
- [46] Bowen Pan, Rameswar Panda, Yifan Jiang, Zhangyang Wang, Rogerio Feris, and Aude Oliva. IA-RED<sup>2</sup>: Interpretability-aware redundancy reduction for vision transformers. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S.

- Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 24898–24911. Curran Associates, Inc., 2021. [2](#)
- [47] Zizheng Pan, Bohan Zhuang, Jing Liu, Haoyu He, and Jianfei Cai. Scalable vision transformers with hierarchical pooling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 377–386, October 2021. [2](#)
- [48] Mathias Parger, Chengcheng Tang, Christopher D. Twigg, Cem Keskin, Robert Wang, and Markus Steinberger. DeltaCNN: End-to-end cnn inference of sparse frame differences in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12497–12506, June 2022. [2](#), [3](#)
- [49] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. DynamicViT: Efficient vision transformers with dynamic token sparsification. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 13937–13949. Curran Associates, Inc., 2021. [2](#)
- [50] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. XNOR-Net: ImageNet classification using binary convolutional neural networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 525–542, Cham, 2016. Springer International Publishing. [3](#)
- [51] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [7](#)
- [52] Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and Bo Dai. Combiner: Full attention transformer with sparse computation cost. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 22470–22482. Curran Associates, Inc., 2021. [2](#)
- [53] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. arXiv, 2020. [2](#)
- [54] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, April 2015. [1](#), [4](#), [7](#), [9](#)
- [55] Michael S. Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. TokenLearner: What can 8 learned tokens do for images and videos? arXiv, 2021. [2](#)
- [56] Surat Teerapittayanon, Bradley McDanel, and H.T. Kung. BranchyNet: Fast inference via early exiting from deep neural networks. In *23rd International Conference on Pattern Recognition (ICPR)*, pages 2464–2469, 2016. [3](#)
- [57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. [1](#), [3](#)
- [58] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. [2](#), [3](#)
- [59] Junke Wang, Xitong Yang, Hengduo Li, Li Liu, Zuxuan Wu, and Yu-Gang Jiang. Efficient video transformers with spatial-temporal token selection. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 69–86, Cham, 2022. Springer Nature Switzerland. [2](#)
- [60] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. arXiv, 2020. [2](#)
- [61] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 568–578, October 2021. [1](#)
- [62] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E. Gonzalez. SkipNet: Learning dynamic routing in convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. [3](#)
- [63] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 11960–11973. Curran Associates, Inc., 2021. [3](#)
- [64] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. LiteEval: A coarse-to-fine framework for resource efficient video recognition. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. [3](#)
- [65] Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. Nyströmformer: A Nyström-based algorithm for approximating self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14138–14148, May 2021. [9](#)
- [66] Ran Xu, Fangzhou Mu, Jayoung Lee, Preeti Mukherjee, Somali Chaterji, Saurabh Bagchi, and Yin Li. SmartAdapt: Multi-branch object detection framework for videos on mobiles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2528–2538, June 2022. [2](#), [3](#)
- [67] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive networks for efficient inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [3](#)

- [68] Hongxu Yin, Arash Vahdat, Jose M. Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-ViT: Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10809–10818, June 2022. [3](#)
- [69] Xiaoyu Yue, Shuyang Sun, Zhanghui Kuang, Meng Wei, Philip H.S. Torr, Wayne Zhang, and Dahua Lin. Vision transformer with progressive sampling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 387–396, October 2021. [2](#)
- [70] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [3](#)