# GTP-ViT: Efficient Vision Transformers via Graph-based Token Propagation

Xuwei Xu[1,3], Sen Wang[1], Yudong Chen[1,3], Yanping Zheng[2,3], Zhewei Wei[2], Jiajun Liu[3,1]

[1] The University of Queensland, Australia
[2] Renmin University of China, China
[3] CSIRO Data61, Australia

{xuwei.xu, sen.wang, yudong.chen}@uq.edu.au, {zhengyanping, zhewei}@ruc.edu.cn
jiajun.liu@csiro.au

## Abstract

*Vision Transformers (ViTs) have revolutionized the field of computer vision, yet their deployments on resource-constrained devices remain challenging due to high computational demands. To expedite pre-trained ViTs, token pruning and token merging approaches have been developed, which aim at reducing the number of tokens involved in the computation. However, these methods still have some limitations, such as image information loss from pruned tokens and inefficiency in the token-matching process. In this paper, we introduce a novel **G**raph-based **T**oken **P**ropagation (**GTP**) method to resolve the challenge of balancing model efficiency and information preservation for efficient ViTs. Inspired by graph summarization algorithms, GTP meticulously propagates less significant tokens' information to spatially and semantically connected tokens that are of greater importance. Consequently, the remaining few tokens serve as a summarization of the entire token graph, allowing the method to reduce computational complexity while preserving essential information of eliminated tokens. Combined with an innovative token selection strategy, GTP can efficiently identify image tokens to be propagated. Extensive experiments have validated GTP's effectiveness, demonstrating both efficiency and performance improvements. Specifically, GTP decreases the computational complexity of both DeiT-S and DeiT-B by up to 26% with only a minimal 0.3% accuracy drop on ImageNet-1K without finetuning, and remarkably surpasses the state-of-the-art token merging method on various backbones at an even faster inference speed. The source code is available at* [https://github.com/Ackesnal/GTP-ViT](https://github.com/Ackesnal/GTP-ViT).

## 1. Introduction

In recent years, Vision Transformer (ViT) has rapidly emerged as the leading backbone for various computer vision tasks, demonstrating remarkable performance in image
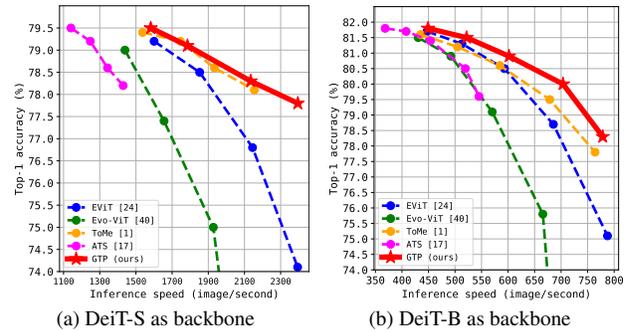


Figure 1. **Comparisons among token reduction methods, taking DeiT models [38] as backbones.** Our GTP presents the best trade-off between model efficiency and performance.

classification [13, 15, 29, 45], object detection [6, 28, 29] and segmentation [6, 43]. Despite its impressive accomplishments in the computer vision domain, the high computational cost hinders the applicability of ViT on devices with constrained computing resources. As a result, improving ViT's computational efficiency has become a growing area of interest in ViT research.

Various approaches have been explored to alleviate the heavy computational burden on ViT, such as integrating self-attention with convolution [7, 16, 37, 41, 44] and designing regional self-attention [3, 14, 29]. In contrast to the methods that put forward novel efficient architectures for ViT, token pruning techniques [10, 18, 21, 25, 31, 34, 42] are proposed to expedite pre-established ViT models. In particular, token pruning methods first measure the importance of each token and then discard the insignificant ones, aiming to gradually reduce the number of tokens involved in the computation. While improving the model efficiency, pruning image tokens inevitably leads to an irreversible information loss of the removed tokens and subsequently compromises performance, especially when a large number of tokens are eliminated. Besides, token pruning methods necessitate further finetuning to prevent a significant performance drop, which also increases their computational cost. Regarding the defects
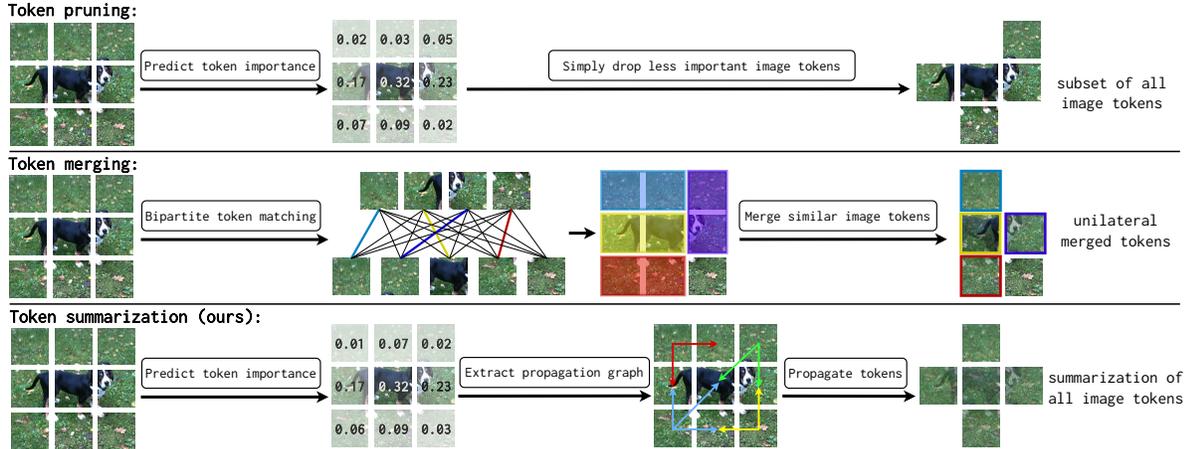
**Token pruning:**
Predict token importance → 0.02 0.03 0.05 / 0.17 0.32 0.23 / 0.07 0.09 0.02 → Simply drop less important image tokens → subset of all image tokens

**Token merging:**
Bipartite token matching → Merge similar image tokens → unilateral merged tokens

**Token summarization (ours):**
Predict token importance → 0.01 0.07 0.02 / 0.17 0.32 0.23 / 0.06 0.09 0.03 → Extract propagation graph → Propagate tokens → summarization of all image tokens

Figure 2. **Comparisons among existing token pruning [21, 25, 34] (top), token merging [1] (middle) and our token summarization (bottom) methods.** Both token pruning and token summarization can efficiently measure the importance of each token and determine which tokens should be discarded, providing a computational advantage over token merging. However, only token merging and token summarization successfully preserve the information of eliminated tokens.

of token pruning, a recent study [1] suggests token merging as a solution to preserve information and avoid finetuning. However, token merging incurs a token-matching process per layer with computational complexity proportional to the feature dimensions and the square of the number of tokens, making it less efficient than token pruning.

Limitations of the existing methods highlight a research challenge: **how to effectively balance the model efficiency and information preservation for ViT models**. Meanwhile, **how to enhance the efficiency of pre-trained ViT models while achieving minimal performance drop without finetuning** also leaves an open research question, especially in computing resource-constrained environments.

In this work, we present a novel Graph-based Token Propagation (GTP) approach to address these challenges. Motivated by graph summarization techniques [24, 27, 35] that target generating condensed representations of a graph, we redefine the problem of token removal and information preservation as a token summarization task where remaining tokens encapsulate the information of the removed ones.

First, we propose an innovative and efficient token selection strategy for GTP to measure the importance of each token. The token selection strategy is based on the regeneration difficulty and broadcasting ability of each token, which can be easily drawn from the attention map that has already been calculated in the self-attention module. Consequently, our token selection strategy is more efficient than pairwise token matching in [1], and can perform on ViTs without the [CLS] token while [10, 18, 21, 25, 42] all depend on it. Second, inspired by the message-passing mechanism in Graph Neural Networks [2, 20, 39], GTP constructs an image token graph and distributes the information of eliminated tokens to their neighbours in the graph. As a result, GTP preserves token information through multilateral relationships among

tokens while the existing token merging method [1] concentrates exclusively on one-to-one matching and merging. The kept tokens eventually make up a smaller representation of the image without abandoning much information. Figure 2 illustrates the comparison between our method and previous approaches. Moreover, we observe that directly discarding tokens in a pre-trained ViT tends to yield a smoother attention map after softmax activation. To resolve this issue, our approach integrates attention map sparsification as an anti-oversmoothing mechanism.

Our contributions are as follows: 1) we introduce a novel and efficient token selection strategy (Section 3.2.1); 2) we design a graph-based token propagation method to summarize the whole image, preserving information of removed tokens (Sections 3.2.2 and 3.2.3); 3) we sparsify the attention map to enforce tokens to focus on significant information (Section 3.2.4). Extensive experiments have demonstrated the effectiveness of GTP. Remarkably, taking pre-trained DeiT-B [38] as the backbone, GTP achieves 28% real inference speed up at the cost of merely a 0.3% accuracy drop without finetuning, and outperforms state-of-the-art token reduction methods in terms of the trade-off between performance and efficiency, as illustrated in Figure 1.

## 2. Related works

**Efficient Vision Transformers.** Ever since the success of Vision Transformer (ViT) [15], numerous studies have been investigating efficient ViTs. Some approaches devise fast self-attention computations that scale linearly or close to linearly with respect to either input length or feature dimensions [9, 23, 30, 40, 46]. Besides, some combine self-attention layers with efficient convolutional layers [5, 7, 11, 37, 41]. Additionally, regional self-attention methods [4, 8, 14, 28, 29, 43] that calculate self-attention within a constrained area have

been proposed to reduce the computation complexity of global token interactions. Distinct from these methods, our method concentrates on expediting pre-trained ViTs with plug-and-play components instead of proposing new backbone architectures.

**Token pruning and merging.** Leveraging the inherent redundancy among image tokens, many studies have attempted to reduce the number of tokens in ViTs [18, 25, 26, 32, 34, 36, 42, 48]. [34] introduces a predictor module to identify tokens that can be removed without significant performance degradation. [25] removes tokens based on their attention to the [CLS] token. [18] scores and samples important tokens to retain. However, these approaches result in information loss of the discarded tokens and necessitate finetuning from pre-trained models. Apart from the pruning-based methods, [1] proposes a token merging method that maintains the information by merging similar tokens, which does not need further finetuning. Our method builds upon these ideas by introducing a graph-based token propagation technique, addressing the limitations of existing token reduction strategies.

# 3. Methods

## 3.1. Preliminaries

**Vision Transformer.** The vanilla ViT [15] divides an input image into several image patches, which are then projected into image token embeddings. We denote the embedded feature map of an image as $X \in \mathbb{R}^{N \times C}$, where $N$ and $C$ are the number of tokens and the dimension of features, respectively. Each ViT block comprises a multi-head self-attention (MHSA) layer and a feed-forward network (FFN) layer. In the MHSA layer, a layer-normalized feature map is first linearly transformed into Query ($Q$), Key ($K$) and Value ($V$) matrices. Then, ViT calculates the similarity between each pair of tokens by the dot product between Query and Key with a softmax activation as

$$A = \text{softmax}(\frac{QK^\top}{\sqrt{d_K}}), \qquad (1)$$

where $A \in \mathbb{R}^{N \times N}$ is the attention map and $d_K = C$ is the feature dimension of $K$. Additionally, the MHSA layer calculates multiple attention maps to increase diversity.

**Graph Neural Network.** Graph Neural Networks (GNNs) are typically constructed by stacking message-passing layers, during which all nodes in a graph update their representations by aggregating information from neighbours [2]. This mechanism can also be regarded as each node propagating its information to the neighbouring nodes. Graph Convolutional Network (GCN) [20] employs the convolution operation as a node aggregation method on the graph-structured data. Given a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ that consists of a set $\mathcal{V}$ of nodes

and a set $\mathcal{E}$ of edges with its adjacency matrix $\mathcal{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$, GCN updates the node feature map $Z \in \mathbb{R}^{|\mathcal{V}| \times C}$ by

$$\text{GCN}(Z) = \sigma(\mathcal{D}^{-\frac{1}{2}} \mathcal{A} \mathcal{D}^{-\frac{1}{2}} Z\Theta), \qquad (2)$$

where $\Theta$ is the projection weight, $\sigma$ represents a nonlinearity (i.e., ReLU), $\mathcal{D}^{-\frac{1}{2}} \mathcal{A} \mathcal{D}^{-\frac{1}{2}}$ is the symmetrically normalized adjacency matrix, and $\mathcal{D}$ is the degree matrix $\mathcal{D}_{i,i} = \sum_j \mathcal{A}_{i,j}$.

## 3.2. Efficient token propagation

### 3.2.1 Token selection

A swift and effective token selection strategy is essential for identifying which tokens can be propagated and discarded without significantly sacrificing. In our method, we evaluate the importance of a token from two aspects.

**Regeneration difficulty.** We assume that a token is less important than others if it is primarily aggregated by other tokens during the self-attention process. These less important tokens can be dropped since they are more easily to be regenerated by other tokens and their information is less significant in token summarization results. Specifically, the regeneration difficulty score $\gamma_i$ of an image token $x_i$ is calculated by the negative sum of attentions from all other tokens to $x_i$ as

$$\gamma_i = \oplus \left( - \sum_{j \in \{0, \dots, N-1\} \setminus \{i\}} A_{i,j} \right) = \oplus (A_{i,i} - 1), \quad (3)$$

where $\oplus(\cdot)$ is a permutation-invariant aggregator to fuse the values from multiple heads. A greater $\gamma_i$ indicates a more important token $x_i$. Since we only need to know the order of $\gamma$s corresponding to different image tokens rather than the values themselves, the constant term in the formula can be omitted, thereby $\gamma_i = \oplus(A_{i,i})$. Consequently, the regeneration difficulty scores $\Gamma_X$ for all the image tokens in a feature map $X$ can be directly obtained from the main diagonal of its attention map $A$ as $\Gamma_X = [\gamma_0, \dots, \gamma_{N-1}] = [\oplus(A_{0,0}), \dots, \oplus(A_{N-1,N-1})] = diag(\oplus(A))$.

**Broadcasting ability.** Despite the regeneration difficulty, an image token is also indispensable if it considerably contributes to other tokens in the self-attention computation. We quantify the broadcasting ability of a token $x_i$ by adding up the attention scores from this token to all other tokens and denote the score as $\psi_i$:

$$\psi_i = \oplus \left( \sum_{j \in \{0, \dots, N-1\} \setminus \{i\}} A_{j,i} \right). \qquad (4)$$

The broadcasting ability score $\psi$ reflects the significance of a token's role in broadcasting information to other tokens in ViT. Specifically, we use $\Psi_X$ to denote the broadcasting abilities for all the image tokens in a feature map $X$, where $\Psi_X = [\psi_0, \dots, \psi_{N-1}]$.
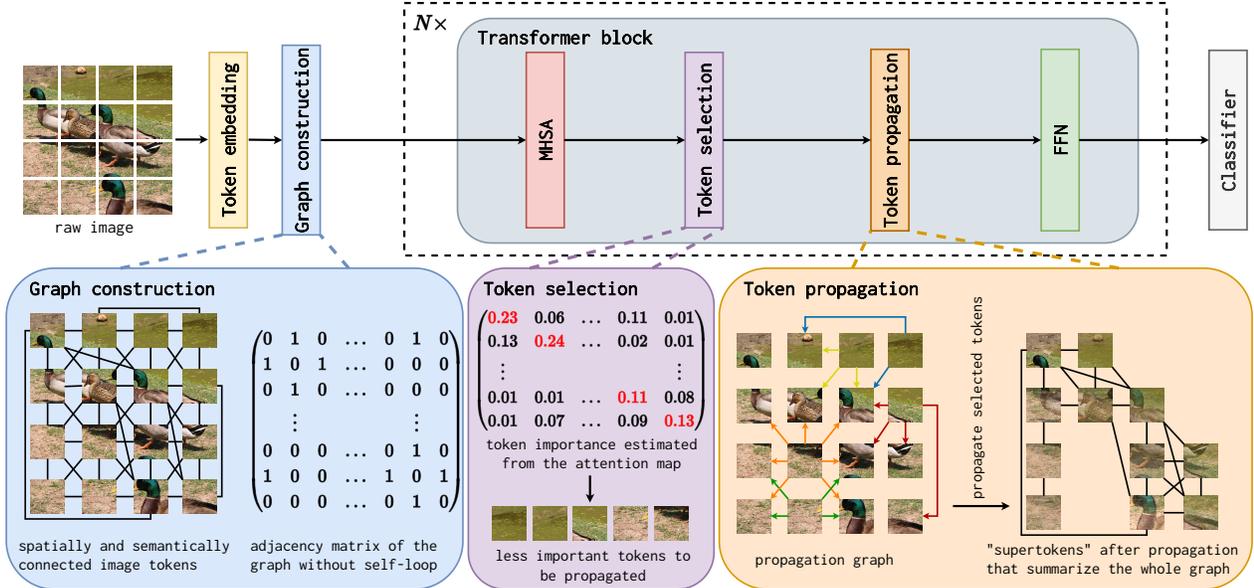
Figure 3. **Graph-based Token Propagation (GTP) visualization.** GTP constructs a graph of image tokens after the token embedding layer **only once**. Within each transformer block, GTP utilizes the attention map computed in the MHSA layer to estimate the importance score for each image token. Next, it propagates less significant tokens to important tokens w.r.t. a subgraph that only contains edges from propagated tokens to kept tokens. As a result, the remaining tokens form a condensed graph representation of the entire image.

**Token selection.** Taking account of both regeneration difficulty score $\Gamma$ and broadcasting ability score $\Psi$, we keep $N - P$ tokens with the largest $\Gamma \times \Psi$ values and propagate the rest $P$ tokens. The propagated tokens are denoted by $\boldsymbol{X}^p \in \mathbb{R}^{P \times C}$ while the kept tokens are denoted by $\boldsymbol{X}^k \in \mathbb{R}^{(N-P) \times C}$, where $N$ and $P$ represent the total number of tokens and the number of propagated tokens, respectively. In addition, we exclude the [CLS] token from the token selection procedure and retain it by default. In particular, we choose $max(\cdot)$ as the $\oplus(\cdot)$ by default.

**Analysis.** Our token selection strategy offers three key advantages. First, in contrast to [21, 34], our strategy does not introduce additional parameters. Second, unlike [25], our method can operate without the [CLS] token, which can extend our method to ViTs that do not use the [CLS] token. Moreover, our strategy is computationally efficient, as it does not necessitate the calculation of pairwise similarities among tokens, making it faster than [1] in practice. We have compared different token selection strategies in Figure 5, including mixed strategy of both regeneration difficulty and broadcasting ability (*MixedAttn*), solely regeneration difficulty (*DiagAttn*), solely broadcasting ability (*BroadAttn*), [CLS] token attention (*CLSAttn*) [25], cosine similarity between tokens (*CosSim*) [1] and random selection (*Random*).

### 3.2.2 Sparse graph construction

GTP regards the image tokens as nodes in a graph and constructs a sparse graph based on spatial and semantic relationships between tokens. Notably, the token graph is con-

structed subsequent to the token embedding layer **only once** and remains static throughout the network, eliminating the need for repeated construction in each layer.

**Spatial graph.** Since each image token corresponds to a region on the raw image, we can simply generate a spatial graph with respect to the tokens' original locations on the raw image. The adjacency matrix $\mathcal{A}^{\text{spatial}}$ of the spatial graph is defined as

$$\mathcal{A}^{\text{spatial}}_{i,j} = \begin{cases} 1 & \text{if } \boldsymbol{x}_i \text{ and } \boldsymbol{x}_j \text{ are adjacent and } i \neq j \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

which enables GTP to capture the spatial information of image tokens in a graph representation. The adjacency matrix $\mathcal{A}^{\text{spatial}}$ is fixed for all images.

**Semantic graph.** While the spatial graph reflects the spatial connections among tokens, capturing their semantic connections is also essential. We leverage the cosine similarity to measure the semantic affinity between tokens $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ in the initial feature map $\boldsymbol{X}_0$ as

$$\text{CosSim}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \frac{\boldsymbol{x}_i \boldsymbol{x}_j^\top}{\|\boldsymbol{x}_i\| \cdot \|\boldsymbol{x}_j\|}. \quad (6)$$

Then, the adjacency matrix $\mathcal{A}^{\text{semantic}}$ of the semantic graph is defined as

$$\mathcal{A}^{\text{semantic}}_{i,j} = \begin{cases} 1 & \text{if } \text{CosSim}(\boldsymbol{x}_i, \boldsymbol{x}_j) \geq T_i \text{ and } i \neq j \\ 0 & \text{otherwise,} \end{cases} \quad (7)$$

where $T_i$ represents the $M^{\text{th}}$ ($M \ll N$) largest cosine similarity value of node $\boldsymbol{x}_i$ to other nodes. $T_i$ serves as a threshold to ensure that each token $\boldsymbol{x}_i$ has a maximum of $M$ edges.

Distinct from the spatial graph, the semantic graph provides image-specific relationships for token propagation.

**Mixed graph.**    Next, we generate a mixed graph that effectively represents both the spatial and semantic relationships among tokens, by integrating the spatial graph and semantic graph. The adjacency matrix $\mathcal{A}$ of mixed graph is simply the union of $\mathcal{A}^{\text{spatial}}$ and $\mathcal{A}^{\text{semantic}}$:

$$\mathcal{A}_{i,j} = \begin{cases} 1 & \text{if } \mathcal{A}_{i,j}^{\text{spatial}} = 1 \text{ or } \mathcal{A}_{i,j}^{\text{semantic}} = 1 \\ 0 & \text{otherwise.} \end{cases} \quad (8)$$

Note that none of the three graphs contains self-loops. In GTP, the graph structure is only used to propagate information from eliminated tokens to the remaining tokens. Therefore, a token chosen for elimination never needs to gather information from itself. Follow Equation 2, we symmetrically normalize the graph as:

$$\hat{\mathcal{A}} = \mathcal{D}^{-\frac{1}{2}} \mathcal{A} \mathcal{D}^{-\frac{1}{2}}, \quad (9)$$

where $\mathcal{D}$ is the diagonal degree matrix defined as $\mathcal{D}_{i,i} = \sum_j \mathcal{A}_{i,j}$. Notably, the token graph is constructed **only once** before the Transformer blocks.

**Implementation optimizations.**    We offer a detailed introduction to implementation optimizations for sparse graph propagation and a fast algorithm for determining the $M^{\text{th}}$ largest value in the supplementary material. In the following experiments, we set $M = 8$ unless otherwise noted. We also provide a study on the choice of $M$ in the supplementary material.

### 3.2.3    Token summarization

Motivated by the message-passing mechanism in GNNs where a node distributes its information to neighbouring nodes, we put forward the token summarization process in which an image token propagates its feature to spatially and semantically connected tokens. In each layer, GTP broadcasts the propagated tokens $X^p$ to the kept tokens $X^k$ by

$$X^s = X^k + \alpha \hat{\mathcal{A}}^p X^p, \quad (10)$$

where $\alpha$ is a hyperparameter controlling the magnitude of propagated token features. The term $\hat{\mathcal{A}}^p \in \mathbb{R}^{(N-P) \times P}$ is extracted from the normalized adjacency matrix $\hat{\mathcal{A}}$, such that the row and column indices correspond to the kept and propagated tokens, respectively. $X^s$ is the summarization of image tokens in the current layer and participates in subsequent computations. GTP implements the token propagation process immediately after the MHSA module on a layer-by-layer basis. After the token propagation procedure, we only maintain the normalized adjacency matrix $\hat{\mathcal{A}}^s \in \mathbb{R}^{(N-P) \times (N-P)}$ for the remaining tokens $X^s$.

### 3.2.4    Attention sparsification

**Proportional attention.**    After reducing the number of tokens, the vanilla softmax outputs become smoother, which

could negatively impact performance [1,47]. To address this issue, we introduce the *proportional attention* from [1] into GTP. The proportional attention is computed as

$$A = \text{softmax}(\frac{QK^\top}{\sqrt{d_K}} + \log s), \quad (11)$$

where $s \in \mathbb{R}^{N \times 1}$ represents the size of each token. Furthermore, we use $s^k$ and $s^p$ to denote the sizes for kept tokens and propagated tokens, respectively. The size of a kept token is dynamically updated according to the number of tokens that it has summarized:

$$s^s = s^k + \alpha \hat{\mathcal{A}}^p s^p. \quad (12)$$

**Attention map sparsification.**    In addition to proportional attention, we refine the attention map by filtering out trivial attention values. In particular, we maintain the largest $\theta N^2$ values in the attention map and assign a zero value to the rest $(1 - \theta) N^2$ elements, where $N$ is the number of tokens and $\theta \in [0, 1]$ represents the attention sparsity. Attention map sparsification helps to concentrate token attention on the most significant signals, thereby relieving the smoothness of the attention map and enhancing model performance.

## 4. Experiments

### 4.1. Implementation settings

All the experiments in this section focus on the image classification task using the ImageNet-1K dataset [12], which contains approximately 1.28 million training images and 50 thousand validation images. We report the top-1 accuracy on the validation set as the main performance metric. For finetuned models, we utilize the same data augmentation and training recipe as implemented in DeiT [38] and finetune for only 30 epochs. The base and minimum learning rates for finetuning are set to $10^{-5}$ and $10^{-6}$, respectively. We measure the inference speed for GTP and all other compared models on the same NVIDIA A6000 GPU with fixed batch size 128 unless noted otherwise. We ensure the PyTorch and CUDA versions are the same for all the models.

### 4.2. Main result

We first apply our GTP on pre-trained DeiT-S [38], DeiT-B [38], LV-ViT-S [19] and LV-ViT-B [19] without additional finetuning and present the performance for various numbers of propagated tokens in Table 1. These four models are popular ViT backbones for token reduction methods. Table 1 demonstrates the capability of GTP to expedite ViT without necessitating finetuning. In particular, when propagating 8 tokens per layer (i.e., $P = 8$), GTP achieves 25% real throughput speed-up (1581.3 image/s vs 1268.3 image/s) and 26% fewer computational complexity (3.4 GMACs vs 4.6 GMACs) with an insignificant accuracy decrease of 0.3% (79.5% vs 79.8%) compared to full-size DeiT-S model. Even for the more complex model, DeiT-B, GTP still accomplishes a 26% reduction in computational complexity (13.1GMACs
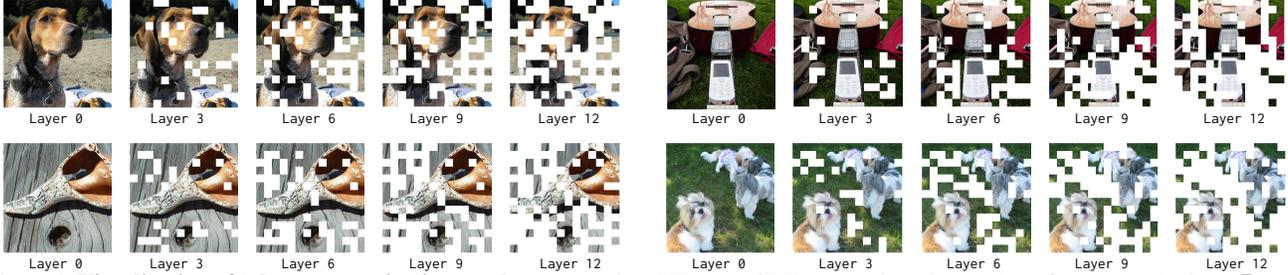
Figure 4. **Visualization of token summarization results.** We employ GTP on DeiT-B [38] and set the number of propagated tokens $P$ to 8. Unlike existing token pruning models that focus primarily on eliminating less significant background tokens, GTP ensures the retention of certain background tokens, thereby providing a summarized representation of the original image.

| Backbone | # Prop. | $P=0$ | $P=1$ | $P=2$ | $P=3$ | $P=4$ | $P=5$ | $P=6$ | $P=7$ | $P=8$ | $P=9$ | $P=10$ | $P=11$ | $P=12$ | $P=13$ | $P=14$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DeiT-S [38] | Acc. (%) | 79.8 | 79.9 | 79.8 | 79.8 | 79.8 | 79.8 | 79.7 | 79.7 | 79.5 | 79.4 | 79.2 | 79.1 | 78.8 | 78.6 | 78.3 |
| | GMACs | 4.6 | 4.5 | 4.3 | 4.2 | 4.0 | 3.9 | 3.7 | 3.6 | 3.4 | 3.3 | 3.2 | 3.0 | 2.9 | 2.7 | 2.6 |
| | img/s | 1265.3 | 1236.7 | 1259.3 | 1325.5 | 1369.7 | 1424.2 | 1471.0 | 1511.8 | 1581.9 | 1645.0 | 1724.9 | 1784.8 | 1874.6 | 1965.5 | 2134.3 |
| LV-ViT-S [19] | Acc. (%) | 81.8 | 81.9 | 81.9 | 81.8 | 81.8 | 81.7 | 81.6 | 81.6 | 81.5 | 81.4 | 81.1 | 80.9 | 80.7 | 80.4 | 80.0 |
| | GMACs | 17.6 | 17.0 | 16.5 | 15.9 | 15.4 | 14.8 | 14.2 | 13.7 | 13.1 | 12.6 | 12.1 | 11.6 | 11.0 | 10.4 | 9.8 |
| | img/s | 408.8 | 405.1 | 418.2 | 432.3 | 449.1 | 462.3 | 480.1 | 500.2 | 521.8 | 544.4 | 577.5 | 601.8 | 630.0 | 659.6 | 703.6 |
| LV-ViT-S [19] | Acc. (%) | 83.3 | 83.0 | 82.8 | 82.6 | 82.4 | 82.4 | 82.2 | 82.1 | 81.9 | 81.7 | 81.5 | 80.8 | 80.0 | - | - |
| | GMACs | 6.6 | 6.4 | 6.1 | 5.9 | 5.7 | 5.4 | 5.2 | 5.0 | 4.8 | 4.6 | 4.4 | 4.1 | 3.9 | - | - |
| | img/s | 940.9 | 870.6 | 902.3 | 940.7 | 978.0 | 1038.2 | 1081.5 | 1159.7 | 1208.8 | 1258.7 | 1338.6 | 1419.2 | 1513.3 | - | - |
| LV-ViT-M [19] | Acc. (%) | 84.0 | 83.8 | 83.7 | 83.6 | 83.4 | 83.3 | 83.2 | 83.0 | 82.8 | 82.5 | - | - | - | - | - |
| | GMACs | 12.7 | 12.1 | 11.5 | 10.9 | 10.3 | 9.7 | 9.1 | 8.5 | 8.0 | 7.4 | - | - | - | - | - |
| | img/s | 524.8 | 496.9 | 521.3 | 549.9 | 582.5 | 624.3 | 667.4 | 715.7 | 770.6 | 833.6 | - | - | - | - | - |

Table 1. **GTP main results on ImageNet-1K *without finetuning*.** In this table, we report the best top-1 accuracy for various numbers of propagated tokens $P$ among different hyperparameter settings. $P=0$ represents the full-size backbone model. Note that LV-ViT-S and LV-ViT-M [19] can reduce at most 12 and 9 tokens per layer, respectively.

vs 17.6GMACs) and approximately 28% improvement in inference speed (521.8 image/s vs 408.8 image/s) with a mere 0.3% drop in accuracy (81.5% vs 81.8%). We also visualize some token summarization examples in Figure 4.

## 4.3. Comparisons with state-of-the-art methods

In Tables 2 and 3, we present comparisons of GTP against token pruning and token merging methods, including DynamicViT [34], EViT [25], ATS [18], Evo-ViT [42], Tri-Level [22] and ToMe [1]. We present the top-1 accuracy, computational complexity (measured in GMACs), and inference speed (measured in images per second) for comparison. We compare these benchmarks since they have released their official source codes so that we can reproduce the results for these models both with and without finetuning for various computational complexities. More implementation details are provided in the table captions.

Table 2 displays both finetuned and finetune-free results on DeiT-S. At a similar inference speed, GTP can match the performance of token pruning methods with finetuning and outperform them when a larger number of tokens are eliminated. For instance, at the same computational complexity of 2.6GMACs, GTP exceeds EViT by 0.2% top-1 accuracy (79.1% vs 78.9%), which reflects GTP's ability to preserve information. It is worth noting that when a substantial number of tokens are dropped, token pruning methods would suffer a dramatic accuracy drop without finetuning. Table 3 presents the finetune-free results on DeiT-B, where GTP
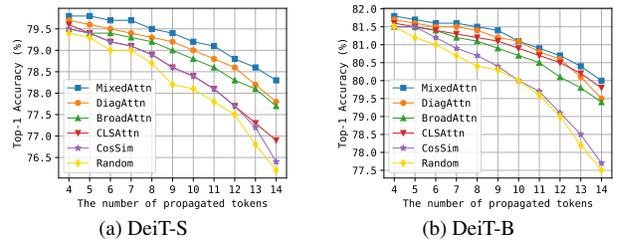


(a) DeiT-S          (b) DeiT-B

Figure 5. **Comparisons of different token selection strategies.** We apply different token selection strategies with GTP and report the top-1 accuracy for various numbers of propagated tokens ($P$).

surpasses all the compared models at a similar inference speed. The naive elimination of tokens leads to a considerable performance decline in token pruning methods. For instance, when reducing the computational complexity of DeiT-B to 8.8GMACs, EViT can only obtain 75.1% top-1 accuracy, which is worse than its performance on DeiT-S (76.8%) with merely 2.6GMACs complexity. On the contrary, GTP reaches 78.3% accuracy, surpassing EViT by a significant 3.2% top-1 accuracy. The only exception, ATS, maintains the performance at the cost of extremely slower inference speed. This indicates that token pruning methods struggle to preserve the information of pruned tokens while remaining efficient and are ineffective without finetuning. In contrast, GTP achieves the best trade-off between model performance and efficiency without finetuning, which is also shown in Figure 1.

| Method | #Param | Approx. 3.5 GMACs | | | Approx. 3.0 GMACs | | | Approx. 2.6 GMACs | | | Approx. 2.3 GMACs | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | w/ F Acc (%) | w/o F Acc (%) | Speed (img/s) | w/ F Acc (%) | w/o F Acc (%) | Speed (img/s) | w/ F Acc (%) | w/o F Acc (%) | Speed (img/s) | w/ F Acc (%) | w/o F Acc (%) | Speed (img/s) |
| DyViT [34] | 22.8M | 79.6 | 74.0 | 1478.4 (×1.30) | 79.3 | 67.4 | 1700.1 (×1.36) | 78.5 | 58.3 | 1980.9 (×1.47) | 77.5 | 51.7 | 2217.7 (×1.48) |
| EViT [25] | 22.1M | **79.8** | 79.2 | **1600.5 (×1.40)** | 79.5 | 78.5 | **1852.3 (×1.48)** | 78.9 | 76.8 | **2144.5 (×1.60)** | 78.5 | 74.1 | **2393.8 (×1.67)** |
| Evo-ViT [42] | 22.1M | 78.4 | 79.0 | 1439.3 (×1.26) | 78.2 | 77.4 | 1655.4 (×1.33) | 78.0 | 75.0 | 1927.9 (×1.44) | 77.7 | 72.1 | 2016.7 (×1.41) |
| Tri-Level [22] | 22.1M | 79.5 | 67.6 | 1345.9 (×1.08) | 79.1 | 67.6 | 1551.2 (×1.24) | 78.8 | 67.6 | 1793.8 (×1.34) | 78.1 | 67.6 | 2013.2 (×1.41) |
| ToMe [1] | 22.1M | 79.7 | 79.4 | 1536.7 (×1.35) | 79.4 | **79.2** | 1746.0 (×1.40) | 78.9 | **78.6** | 1934.0 (×1.44) | 78.4 | 78.1 | 2154.7 (×1.51) |
| ATS [18] | 22.1M | 79.7 | **79.5** | 1140.6 (×1.00) | **79.7** | **79.2** | 1248.1 (×1.00) | 79.0 | **78.6** | 1343.0 (×1.00) | **78.6** | **78.2** | 1429.4 (×1.00) |
| GTP (ours) | 22.1M | 79.7 | **79.5** | 1581.9 (×1.39) | 79.5 | 79.1 | 1784.8 (×1.43) | **79.1** | 78.3 | 2134.3 (×1.59) | **78.6** | 77.8 | 2304.7 (×1.61) |

Table 2. **Comparisons with state-of-the-art methods, taking DeiT-S [38] as the backbone.** "w/ F" and "w/o F" represent the performance with and without 30-epoch finetuning, respectively. We categorize the performance in terms of computational complexity. For example, *Approx. 3.5GMACs* stands for the computational complexity at about 3.5 GMACs, which is equivalent to the keep ratio at 0.8 for DynamicViT [34], EViT [25] and Tri-Level [22], selection ratio at 0.7 for Evo-ViT [42], ATS block from layers 7 to 11 for ATS [18] and the number of reduced tokens at 8 per layer for ToMe [1] and our GTP. More details are provided in the supplementary material. We leverage the slowest inference speed in each category as the baseline (i.e. ×1.00) for speed comparisons. To ensure fairness, we reproduce the finetuned results for these models using their officially released codes. Figure 1a shows visualized comparisons. Bold font means better.

| Method | #Param | Approx. 15.3 GMACs | | Approx. 13.1 GMACs | | Approx. 11.6 GMACs | | Approx. 9.8 GMACs | | Approx. 8.8 GMACs | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) |
| DyViT [34] | 89.5M | 79.9 | 420.0 (×1.14) | 77.7 | 489.0 (×1.20) | 75.5 | 579.7 (×1.28) | 69.5 | 676.2 (×1.30) | 50.2 | 784.2 (×1.44) |
| EViT [25] | 86.6M | 81.7 | 447.2 (×1.21) | 81.3 | 513.0 (×1.26) | 80.5 | 593.2 (×1.31) | 78.7 | 685.4 (×1.32) | 75.1 | **787.4 (×1.44)** |
| Evo-ViT [42] | 86.6M | 81.5 | 430.6 (×1.17) | 80.9 | 492.2 (×1.21) | 79.1 | 570.2 (×1.26) | 75.8 | 665.6 (×1.28) | 60.6 | 736.6 (×1.35) |
| Tri-Level [22] | 86.6M | 64.6 | 398.9 (×1.08) | 64.6 | 466.7 (×1.14) | 64.6 | 539.7 (×1.19) | 64.6 | 622.9 (×1.20) | 64.6 | 707.7 (×1.30) |
| ToMe [1] | 86.6M | 81.6 | 435.1 (×1.18) | 81.2 | 504.9 (×1.24) | 80.6 | 584.5 (×1.29) | 79.5 | 678.6 (×1.31) | 77.8 | 764.0 (×1.40) |
| ATS [18] | 86.6M | **81.8** | 368.3 (×1.00) | **81.7** | 407.9 (×1.00) | **81.4** | 453.7 (×1.00) | **80.5** | 519.6 (×1.00) | **79.6** | 545.2 (×1.00) |
| GTP (ours) | 86.6M | **81.8** | **449.1 (×1.22)** | 81.5 | **521.8 (×1.28)** | 80.9 | **601.8 (×1.33)** | 80.0 | **703.6 (×1.35)** | 78.3 | 778.5 (×1.43) |

Table 3. **Comparisons with state-of-the-art methods *without finetuning*, taking DeiT-B [38] as the backbone.** Due to the computing resource limitation, we only present the finetune-free results for DeiT-B. All the formats presented in this table align with the format in Table 2. Figure 1b shows visualized comparisons. Bold font means better.

## 4.4. Ablation studies

### 4.4.1 Token selection strategy

In Figure 5, we compare different token selection strategies, including mixed strategy of both regeneration difficulty and broadcasting ability (*MixedAttn*), solely regeneration difficulty (*DiagAttn*), solely broadcasting ability (*BroadAttn*), [CLS] token attention (*CLSAttn*) [25], cosine similarity between tokens (*CosSim*) [1] and random selection (*Random*). Figure 5 demonstrates that our token selection strategy consistently outperforms the other methods w.r.t. different numbers of eliminated tokens. Besides, simply adopting the regeneration difficulty score (i.e., *DiagAttn*) can achieve close or even higher performance than that of the traditional [CLS] attention. Considering that many new ViT architectures do not contain the [CLS] token, our method illustrates a potential solution for token pruning methods on them.

### 4.4.2 Graph type

We study the token graph types in the token summarization process and report their best top-1 accuracy in Table 4. Graph type *None* represents for only selecting and removing tokens without propagation, which is analogous to layer-by-layer token pruning. For semantic graphs, we set the number of semantically connected tokens ($M$) to 8 by default, thereby creating a graph of equivalent size to the spatial graph.

Primarily, we note that relying solely on a semantic graph can yield substantial performance as using the mixed graph in most scenarios. It indicates that the semantic relationship among tokens is more important than the spatial relationship.

| Backbone | Graph type | Acc (%) | | | |
|---|---|---|---|---|---|
| | | $P=4$ | $P=8$ | $P=11$ | $P=14$ |
| DeiT-S [38] | Mixed | **79.8** | **79.5** | **79.1** | **78.3** |
| | Spatial | **79.8** | **79.5** | 78.9 | 78.1 |
| | Semantic | **79.8** | 79.4 | 79.0 | 78.2 |
| | None | 79.6 | 79.2 | 78.7 | 77.6 |
| DeiT-B [38] | Mixed | 81.7 | **81.5** | 80.9 | **80.0** |
| | Spatial | **81.8** | 81.4 | 80.9 | 79.8 |
| | Semantic | **81.8** | **81.5** | 80.9 | **80.0** |
| | None | 81.7 | 81.4 | 80.7 | 79.6 |

Table 4. **Ablation study on graph types.**

| Backbone | Sparsity | Prop. attn. | Acc (%) | | | |
|---|---|---|---|---|---|---|
| | | | $P=4$ | $P=8$ | $P=11$ | $P=14$ |
| DeiT-S [38] | 1.0 | n/a | 79.7 | 79.2 | 78.7 | 77.6 |
| | 0.9 | √ | 79.7 | 79.3 | 79.0 | 78.2 |
| | 0.8 | | 79.7 | 79.4 | 79.0 | 78.1 |
| | 0.7 | | **79.8** | 79.4 | **79.1** | 78.2 |
| | 0.6 | | **79.8** | **79.5** | 79.0 | **78.3** |
| | 0.5 | | **79.8** | 79.4 | 79.0 | 78.0 |
| | best | × | 79.7 | 79.3 | 78.8 | 77.8 |
| DeiT-B [38] | 1.0 | n/a | 81.7 | 81.4 | 80.7 | 79.8 |
| | 0.9 | √ | 81.7 | 81.4 | **80.9** | **80.0** |
| | 0.8 | | 81.7 | **81.5** | **80.9** | 79.9 |
| | 0.7 | | 81.7 | **81.5** | 80.9 | 79.9 |
| | 0.6 | | **81.8** | 81.4 | 80.9 | 79.9 |
| | 0.5 | | **81.8** | **81.5** | 80.9 | **80.0** |
| | best | × | **81.8** | **81.5** | 80.9 | 79.9 |

Table 5. **Ablation studies on attention sparsity.** Attention sparsity at 1.0 represents the result without the attention sparsification. *Prop. attn.* stands for using proportional attention. *Best* stands for the best attention sparsity for each number of propagated tokens $P$.

Secondly, we observe that the effectiveness of graph propagation signifies when the number of eliminated tokens $P$ increases. For instance, on DeiT-B, graph propagation can only enhance the top-1 accuracy by 0.1% when $P=4$ or $P=8$, but this accuracy difference escalates to 0.4% when

| Method | P = 0 | | P = 10 | | P = 20 | | P = 30 | | P = 40 | | P = 50 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) | Acc (%) | Speed (img/s) |
| ToMe [1] | 85.8 | 48.3 | 85.8 | 53.5 | 85.7 | 59.7 | 85.5 | 67.2 | 85.3 | 75.9 | 84.9 | 86.8 |
| GTP (ours) | | | **85.9** | **57.6 (+7.6%)** | **85.8** | **63.7 (+6.7%)** | **85.7** | **74.6 (+11.0%)** | **85.5** | **83.2 (+9.6%)** | **85.3** | **95.5 (+10.0%)** |

Table 6. **GTP performance on ViT with more tokens.** We validate GTP's effectiveness and efficiency on ViT-B-Patch8 [15], which has 785 tokens. Due to our GPU memory constraints, the inference speeds are tested with batch size 32.

| Backbone | # Prop. | P = 0 | P = 1 | P = 2 | P = 3 | P = 4 | P = 5 | P = 6 | P = 7 | P = 8 | P = 9 | P = 10 | P = 11 | P = 12 | P = 13 | P = 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. (%) | 84.0 | 84.1 | 84.0 | 83.8 | 83.8 | 83.7 | 83.7 | 83.7 | 83.6 | 83.5 | 83.4 | 83.3 | 83.1 | 82.7 | 82.3 |
| ViT-Medium-GAP [15] | GMACs | 10.6 | 10.4 | 10.1 | 9.9 | 9.6 | 9.3 | 9.0 | 8.8 | 8.5 | 8.3 | 8.0 | 7.8 | 7.5 | 7.2 | 7.0 |
| | img/s | 535.0 | 523.7 | 536.8 | 552.7 | 566.7 | 581.5 | 598.0 | 615.3 | 635.0 | 656.5 | 670.8 | 748.7 | 778.9 | 805.6 | 833.6 |

Table 7. **GTP results on ViT model without the [CLS] token *without finetuning*.**

$P = 14$. Meanwhile, graph propagation benefits small models more than large models. For example, graph propagation can increase the top-1 accuracy by 0.7% on DeiT-S when $P = 14$, in contrast to only 0.4% for DeiT-B, which is a significant improvement in the field of expediting ViTs. We think this is because larger ViT models are more robust when confronted with image information loss.
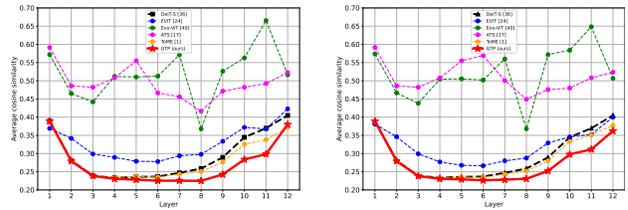
### 4.4.3 Attention sparsification
We conduct an ablation of the proportional attention utilized in the GTP method, as detailed in Table 5. For DeiT-S, we observe that proportional attention consistently enhances performance. For example, it increases the best top-1 accuracy of DeiT-S by 0.5% when $P = 14$. However, it becomes less effective on DeiT-B with only trivial improvements. We then explore the attention sparsity presenting the top-1 accuracy for different attention sparsities. Similar to the results of proportional attention, we find that attention sparsification is less impactful for larger models. For example, when removing 14 tokens per layer, a proper attention sparsity can increase the top-1 accuracy for DeiT-S by 0.7%, much higher than the accuracy increase on DeiT-B, which is only 0.2%. We think this is also because larger models are more robust than their smaller counterparts and already concentrate on the most significant parts of an image. In other words, the attention map of DeiT-B is already fairly sparse.

### 4.5. Anti-oversmoothing
Figure 6 illustrates the trend of average cosine similarity between image tokens in each layer for various token reduction models. A higher average cosine similarity indicates a more severe oversmoothing problem where all the remaining tokens tend to be similar. Oversmoothing would lead to performance degradation both in GCN and ViT. Our GTP can mitigate the oversmoothing problem and yield lower similarities between image tokens, which is one of the key factors to our outstanding performance.

### 4.6. Performance on ViT with more tokens
The token selection and fusion strategies of our GTP are more efficient than those of ToMe [1]. We provide a theoretical comparison of computational complexities in the supplementary material. In this section, we present empirical comparisons between GTP and ToMe, taking ViT-B-Patch8 [15] as the backbone. ViT-B-Patch8 contains 765 tokens,


(a) Approx. 2.6 GMACs    (b) Approx. 3.0 GMACs

Figure 6. **Average cosine similarity.** We calculate the average cosine similarity among image tokens in each layer for various token reduction methods finetuned on DeiT-S.

significantly more than ViT-B or DeiT-B which holds only 197 tokens. Experimental results in Table 6 demonstrate that with more tokens in the backbone ViT, our GTP achieves better performance and around 10% faster inference speeds.

### 4.7. Performance on ViTs without the [CLS] token
As introduced in section 3.2.1, GTP selects tokens without the need for [CLS] token. In Table 7, we present GTP's results on ViT-Medium-GAP [15] where global average pooling is used instead of the [CLS] token. It is worth noting that [18, 22, 25, 42] cannot work with this backbone.

### 4.8. Additional experiments
In the appendix, we provide additional ablation studies on the hyperparameter $\alpha$ in token propagation and the number of semantic neighbours $M$. We also provide GTP's performance on large ViT models and a theoretical analysis of the computational complexity.

## 5. Conclusion
In this work, we treat the challenge of expediting ViTs for computing resource-constrained environments by reducing tokens as a token summarization task. We present a graph-based token propagation (GTP) method to address this issue without the need for finetuning. GTP constructs a sparse graph representation for image tokens and strategically selects less important tokens for propagation based on their regeneration difficulty and broadcasting ability. Then, GTP propagates the information of insignificant tokens into the other remaining tokens, which constitutes a condensed token representation. Extensive experiments have demonstrated the efficacy and efficiency of GTP. We hope our work can inspire future research in token reduction for ViTs.

# References

[1] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. In *ICLR*, 2023. 2, 3, 4, 5, 6, 7, 8, 10, 11, 12, 13

[2] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021. 2, 3

[3] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *ICCV*, 2021. 1

[4] Chun-Fu Chen, Rameswar Panda, and Quanfu Fan. Regionvit: Regional-to-local attention for vision transformers. In *ICLR*, 2022. 2

[5] Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan, and Zicheng Liu. Mobile-former: Bridging mobilenet and transformer. In *CVPR*, 2022. 2

[6] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *ICLR*, 2023. 1

[7] Zhengsu Chen, Lingxi Xie, Jianwei Niu, Xuefeng Liu, Longhui Wei, and Qi Tian. Visformer: The vision-friendly transformer. In *ICCV*, 2021. 1, 2

[8] Zhiyang Chen, Yousong Zhu, Chaoyang Zhao, Guosheng Hu, Wei Zeng, Jinqiao Wang, and Ming Tang. Dpt: Deformable patch-based transformer for visual recognition. In *ACMMM*, 2021. 2

[9] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *NeurIPS*, 2021. 2

[10] Zheng Chuanyang, Zheyang Li, Kai Zhang, Zhi Yang, Wenming Tan, Jun Xiao, Ye Ren, and Shiliang Pu. Savit: Structure-aware vision transformer pruning via collaborative optimization. In *NeurIPS*, 2022. 1, 2

[11] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. In *NeurIPS*, 2021. 2

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 5, 11

[13] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *ECCV*, 2022. 1

[14] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *CVPR*, 2022. 1, 2

[15] 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 *ICLR*, 2021. 1, 2, 3, 8, 10, 12

[16] Stéphane d'Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *ICML*, 2021. 1

[17] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *CVPR*, 2023. 10, 12

[18] 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 *ECCV*, 2022. 1, 2, 3, 6, 7, 8, 11

[19] Zi-Hang Jiang, Qibin Hou, Li Yuan, Daquan Zhou, Yujun Shi, Xiaojie Jin, Anran Wang, and Jiashi Feng. All tokens matter: Token labeling for training better vision transformers. In *NeurIPS*, 2021. 5, 6

[20] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017. 2, 3

[21] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *ECCV*, 2022. 1, 2, 4

[22] Zhenglun Kong, Haoyu Ma, Geng Yuan, Mengshu Sun, Yanyue Xie, Peiyan Dong, Xin Meng, Xuan Shen, Hao Tang, Minghai Qin, et al. Peeling the onion: Hierarchical reduction of data redundancy for efficient vision transformer training. In *AAAI*, 2023. 6, 7, 8, 11

[23] Hai Lan, Xihao Wang, Hao Shen, Peidong Liang, and Xian Wei. Couplformer: Rethinking vision transformer with coupling attention. In *WACV*, 2023. 2

[24] Kristen LeFevre and Evimaria Terzi. Grass: Graph structure summarization. In *ICDM*, 2010. 2

[25] Youwei Liang, GE Chongjian, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Evit: Expediting vision transformers via token reorganizations. In *ICLR*, 2021. 1, 2, 3, 4, 6, 7, 8, 11

[26] Yue Liu, Christos Matsoukas, Fredrik Strand, Hossein Azizpour, and Kevin Smith. Patchdropout: Economizing vision transformers using patch dropout. In *WACV*, 2023. 3

[27] Yike Liu, Tara Safavi, Abhilash Dighe, and Danai Koutra. Graph summarization methods and applications: A survey. *CSUR*, 2018. 2

[28] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 1, 2

[29] 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 *ICCV*, 2021. 1, 2

[30] 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 *NeurIPS*, 2021. 2

[31] Dmitrii Marin, Jen-Hao Rick Chang, Anurag Ranjan, Anish Prabhu, Mohammad Rastegari, and Oncel Tuzel. Token pooling in vision transformers. In *WACV*, 2021. 1

[32] 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 *CVPR*, 2022. 3

[33] PyTorch. Torch.sparse, pytorch 2.0 documentation. `https://pytorch.org/docs/stable/sparse.html`, 2023. 10

[34] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. In *NeurIPS*, 2021. 1, 2, 3, 4, 6, 7, 11

[35] Matteo Riondato, David García-Soriano, and Francesco Bonchi. Graph summarization with quality guarantees. *DMKD*, 2017. 2

[36] Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: Adaptive space-time tokenization for videos. In *NeurIPS*, 2021. 3

[37] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *CVPR*, 2021. 1, 2

[38] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICLR*, 2021. 1, 2, 5, 6, 7, 10, 11, 12

[39] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018. 2

[40] 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 *ICCV*, 2021. 2

[41] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, 2021. 1, 2

[42] Yifan Xu, Zhijie Zhang, Mengdan Zhang, Kekai Sheng, Ke Li, Weiming Dong, Liqing Zhang, Changsheng Xu, and Xing Sun. Evo-vit: Slow-fast token evolution for dynamic vision transformer. In *AAAI*, 2022. 1, 2, 3, 6, 7, 8, 11

[43] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal attention for long-range interactions in vision transformers. In *NeurIPS*, 2021. 1, 2

[44] Jinsu Yoo, Taehoon Kim, Sihaeng Lee, Seung Hwan Kim, Honglak Lee, and Tae Hyun Kim. Enriched cnn-transformer feature aggregation networks for super-resolution. In *WACV*, 2023. 1

[45] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2021. 1

[46] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In *ICCV*, 2021. 2

[47] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021. 5

[48] Zhuofan Zong, Kunchang Li, Guanglu Song, Yali Wang, Yu Qiao, Biao Leng, and Yu Liu. Self-slimmed vision transformer. In *ECCV*, 2022. 3

## A. Implementation optimizations

First, we note that both the spatial graph and semantic graph are sparse graphs, with graph sparsity $\frac{8}{N}$ and $\frac{M}{N}$, respectively. For ViT [15] and DeiT [38] models, the total number of image tokens $N$ is usually 196, which indicates less than 5% non-trivial values in the two adjacency matrices. As a result, the mixed graph is also a sparse graph whose sparsity is no more than $\frac{8+M}{N}$ (less than 7% on average). Therefore, we can store the graph in sparse tensors [33] and perform sparse matrix multiplication to accelerate the graph propagation in Sections 3.2.2 and 3.2.3. Besides, for batched inputs that the sparse matrix multiplication does not support, we can use the scatter reduction operation to avoid the dense matrix multiplication. Second, the threshold $T_i$ for sparsifying the semantic graph can be determined by finding the $M^{\text{th}}$ largest value in the unsorted array $\mathcal{A}_i^{\text{semantic}}$ with a complexity $O(N)$ rather than sorting the whole array with a complexity $O(N \log N)$.

## B. Experiment settings

We provide the hyperparameter settings for the compared methods in Tables 8 and 9. These hyperparameters are used to control the reduced computational complexity for the backbone ViT, ensuring fair comparisons.

## C. Additional experiments

### C.1. Larger ViT backbones

In the main submission, we have validated GTP's effectiveness on small to medium-sized ViT backbones. Moreover, we employ GTP on two large-size ViT backbones: ViT-L [15] and EVA-L [17], which represent ViT backbones with and without the [CLS] token, respectively. Since ViT-L and EVA-L both have 24 layers, the maximum number of propagated tokens $P$ per layer is limited to 8. We also reproduce the state-of-the-art ToMe [1] on these models and present the comparisons in Table 10. It is worth noting that ToMe consumes considerable GPU memory for computing the cosine similarity in each layer. We omit the evaluation on larger models (-H, -G) due to hardware constraints.

From Table 10, we point out that our GTP outperforms ToMe in both model performance and efficiency. Such performance difference is even more significant on ViT backbones without the [CLS] token. For instance, GTP achieves 85.4% top-1 accuracy at 212.0 image/second when taking EVA-L as the backbone, **surpassing ToMe's 80.1% top-1 accuracy by a significant 5.3%** at a similar inference speed.

| Backbone | Method | Hyperparameter | Approximate computational complexity | | | |
|---|---|---|---|---|---|---|
| | | | 3.5 GMACs | 3.0 GMACs | 2.6 GMACs | 2.3 GMACs |
| DeiT-S [38] | DynamicViT [34] | base keep ratio $\rho$ | $\rho = 0.8$ | $\rho = 0.7$ | $\rho = 0.6$ | $\rho = 0.5$ |
| | EViT [25] | token keeping rate $k$ | $k = 0.8$ | $k = 0.7$ | $k = 0.6$ | $k = 0.5$ |
| | Evo-ViT [42] | selection ratio $p$ | $p = 0.7$ | $p = 0.5$ | $p = 0.4$ | $p = 0.3$ |
| | Tri-Level [22] | token keep ratio $R_T$ | $R_T = 0.8$ | $R_T = 0.7$ | $R_T = 0.6$ | $R_T = 0.5$ |
| | ToMe [1] | token reduce $r$ per layer | $r = 8$ | $r = 11$ | $r = 14$ | $r = 16$ |
| | ATS [18] | ATS block layers | layer 7 to 11 | layer 6 to 11 | layer 5 to 11 | layer 3 to 11 |
| | GTP (ours) | propagated tokens $P$ per layer | $P = 8$ | $P = 11$ | $P = 14$ | $P = 16$ |

Table 8. **Hyperparameter settings for the baseline methods, taking DeiT-S as the backbone.**

| Backbone | Method | Hyperparameter | Approximate computational complexity | | | | |
|---|---|---|---|---|---|---|---|
| | | | 15.3 GMACs | 13.1 GMACs | 11.6 GMACs | 9.8 GMACs | 8.8 GMACs |
| DeiT-B [38] | DynamicViT [34] | base keep ratio $\rho$ | $\rho = 0.9$ | $\rho = 0.8$ | $\rho = 0.7$ | $\rho = 0.6$ | $\rho = 0.5$ |
| | EViT [25] | token keeping rate $k$ | $k = 0.9$ | $k = 0.8$ | $k = 0.7$ | $k = 0.6$ | $k = 0.5$ |
| | Evo-ViT [42] | selection ratio $p$ | $p = 0.8$ | $p = 0.7$ | $p = 0.5$ | $p = 0.4$ | $p = 0.3$ |
| | Tri-Level [22] | token keep ratio $R_T$ | $R_T = 0.9$ | $R_T = 0.8$ | $R_T = 0.7$ | $R_T = 0.6$ | $R_T = 0.5$ |
| | ToMe [1] | token reduce $r$ per layer | $r = 4$ | $r = 8$ | $r = 11$ | $r = 14$ | $r = 16$ |
| | ATS [18] | ATS block layers | layer 9 to 11 | layer 8 to 11 | layer 6 to 11 | layer 3 to 11 | layer 1 to 11 |
| | GTP (ours) | propagated tokens $P$ per layer | $P = 4$ | $P = 8$ | $P = 11$ | $P = 14$ | $P = 16$ |

Table 9. **Hyperparameter settings for the baseline methods, taking DeiT-B as the backbone.**
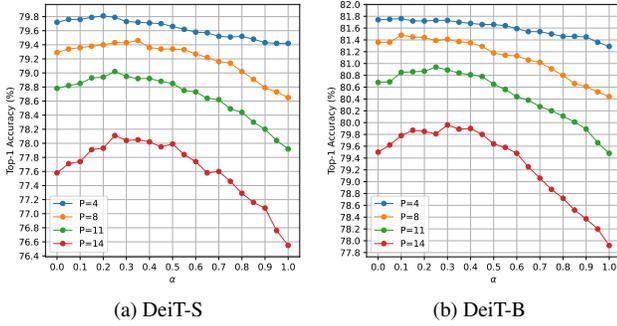


(a) DeiT-S  (b) DeiT-B

Figure 7. **Top-1 accuracy of GTP on ImageNet-1K [12] for various $\alpha$s.** We evaluate the performance of GTP on both DeiT-S and DeiT-B [38] *without finetuning* w.r.t. different $\alpha$. For fair comparisons, we employ the same graph type for propagation and set the attention sparsity static at 0.6 for DeiT-S and 0.5 for DeiT-B, respectively. The findings in this experiment are consistent across different settings.

## C.2. Graph propagation hyperparameter $\alpha$

In the graph summarization process $X^s = X^k + \alpha \hat{\mathcal{A}}^p X^p$, we use $\alpha$ to control the magnitude of information broadcast from propagated tokens $X^p$ to kept tokens $X^k$. In this section, we investigate the performance of GTP with respect to different $\alpha$ and visualize the results in Figure 7. In general, as $\alpha$ increases within the range of $[0, 1]$, the corresponding accuracy first rises and then declines, reaching its peak between 0.2~0.4 for DeiT-S and 0.1~0.3 for DeiT-B. We explain this phenomenon from two aspects. First, when $\alpha$ is approaching 0, the propagated information becomes trivial, leading to a situation where the propagated tokens' information is not preserved. When $\alpha = 0$, this process merely prunes tokens.

Secondly, as $\alpha$ increases, the propagated information gradually dominates the original information of the remaining tokens, which results in an over-smoothing problem and subsequently hinders performance.

## C.3. The number of graph neighbours

We study the influence of the number of semantic neighbours $M$ on model performance and plot the accuracy in Figure 8. For fair comparisons, we apply GTP on DeiT-S and DeiT-B with static attention sparsity and alpha at 0.5 and 0.2, respectively. Figures 8(a) and 8(c) illustrate the results obtained by only employing the semantic graph for token propagation, with respect to different $M$. It can be observed that when the number of propagated tokens $P$ is small (e.g., $P = 4$ or $P = 8$), increasing the semantic neighbours would first slightly improves the accuracy and then converges. However, when $P$ becomes large (e.g., $P = 14$), increasing the semantic neighbours may lead to a performance drop. This can be attributed to the aggregation of redundant information, where one kept token encapsulates an excessive number of propagated tokens that may not be semantically close to it. Figures 8(b) and 8(d) show that integrating the semantic graph with the spatial graph stabilizes the trend of accuracy, indicating the significance of the spatial relationship in token summarization.

## C.4. Computational complexity comparison

As stated in the Introduction section, ToMe [1] encounters a computational bottleneck in the pair-wise token matching process, whose computational complexity is proportional to the feature dimensions and the square of the number of tokens. Compared with ToMe, GTP demonstrates faster infer-

| Backbone | Method | $P=0$ | | $P=1$ | | $P=2$ | | $P=3$ | | $P=4$ | | $P=5$ | | $P=6$ | | $P=7$ | | $P=8$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) | Acc. (%) | Speed (img/s) |
| ViT-L [38] | ToMe [1] | 85.8 | 123.4 | **85.8** | 122.0 | 85.7 | 132.1 | 85.5 | 142.2 | 85.3 | 153.8 | 85.0 | 169.0 | 84.7 | 184.8 | 84.2 | 204.0 | **83.7** | 228.5 |
| | GTP (ours) | 85.8 | 125.4 | **85.8** | **134.4** | **85.8** | **144.9** | **85.8** | **157.4** | **85.5** | **172.0** | **85.3** | **188.3** | **85.0** | **208.8** | **84.3** | **234.0** | **83.7** | 234.0 |
| EVA-L [17] | ToMe [1] | 87.7 | 123.6 | 87.7 | 123.6 | 87.6 | 132.3 | 87.3 | 142.5 | 87.0 | 155.2 | 86.5 | 169.6 | 85.6 | 185.9 | 84.0 | 207.0 | 80.1 | 230.9 |
| | GTP (ours) | 87.9 | 123.1 | **87.9** | **125.6** | **87.8** | **134.6** | **87.8** | **145.3** | **87.7** | **158.0** | **87.5** | **172.5** | **87.2** | **188.9** | **86.5** | **212.0** | **85.4** | **234.9** |

Table 10. **Performance on larger ViT models.** We validate GTP's performance on two large-size ViT models, ViT-L [15] and EVA-L [17], where ViT-L employs the [CLS] token while EVA-L does not have the [CLS] token. Since both models have 24 layers, we can eliminate at most 8 tokens per layer. Bond font means better. GTP constantly outperform ToMe on large ViT models with and without the [CLS] token.
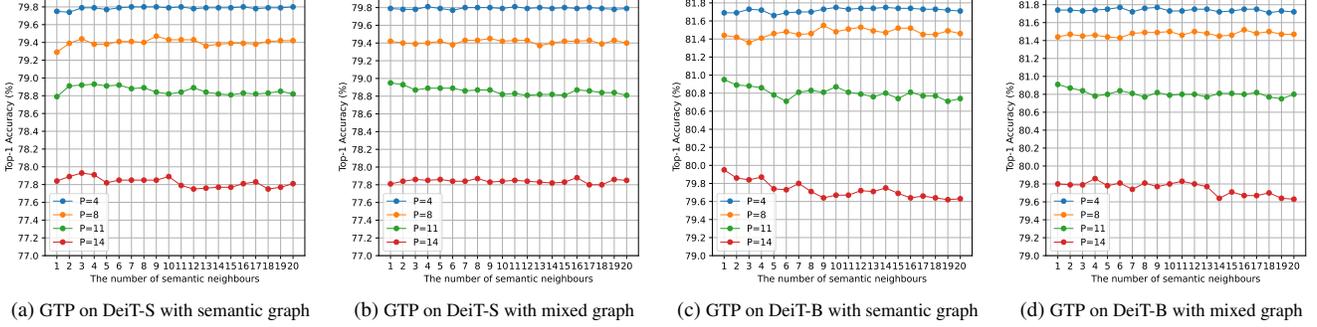


(a) GTP on DeiT-S with semantic graph  (b) GTP on DeiT-S with mixed graph  (c) GTP on DeiT-B with semantic graph  (d) GTP on DeiT-B with mixed graph

Figure 8. **Top-1 accuracy of GTP with various numbers of semantic neighbours $M$.** We evaluate the performance of GTP on both DeiT-S and DeiT-B [38] *without finetuning* w.r.t. different numbers of semantically connected neighbours.

ence speed and accomplishes better information preservation results. In this section, we provide an in-depth analysis of the enhancements in computational efficiency achieved by GTP.

As a plug-and-play component, GTP inserts the token summarization module between the MHSA layer and the FFN layer in each ViT block, which behaves analogously to ToMe. Therefore, when the number of eliminated tokens is the same, the computational complexity of the backbone model for GTP and ToMe should be the same. Consequently, we only consider the additional computational costs (e.g., token matching, token selection and token propagation) introduced by the two models in this analysis. We list the denotations before the theoretical analysis as follows:

$N$ : The total number of tokens in the backbone network.

$N_l$ : The number of remaining tokens in layer $l$, where
$$N_l = N - (l-1)M.$$

$M$ : The number of eliminated tokens in each layer.

$C$ : The dimension of features.

$L$ : The total number of layers.

$H$ : The number of heads.

$$(13)$$

For ToMe, the token matching processing first splits tokens into two sets and then calculates the cosine similarity between each pair of tokens from the two sets. The computational complexity for this process in layer $l$ is $\frac{1}{4}N_l^2C$.

Besides, ToMe merges $M$ tokens in each layer, whose total computational complexity is $MC$ in each layer. As a result, the total additional computational complexity $G_{\text{ToMe}}$ introduced by ToMe is calculated as

$$
\begin{aligned}
G_{\text{ToMe}} &= \sum_{l=1}^{L}\left(\frac{1}{4}N_l^2C + MC\right) \\
&= \frac{1}{4}C\sum_{l=1}^{L}N_l^2 + LMC \\
&= \frac{1}{4}C\sum_{l=1}^{L}(N-(l-1)M)^2 + LMC \\
&= \frac{1}{4}C\sum_{l=1}^{L}\left(N^2 - 2(l-1)NM + (l-1)^2M^2\right) \\
&\quad + LMC \\
&= \frac{1}{4}LN^2C + \frac{1}{4}(L-L^2)NMC \\
&\quad + \left(\frac{1}{12}L^3 + \frac{1}{12}L^2 + \frac{1}{8}L\right)M^2C + LMC
\end{aligned}
$$
$$(14)$$

We then calculate the computational complexity for GTP. GTP first constructs the semantic graph for an input image after the token embedding layer with a computational complexity $N^2C$. Next, it selects tokens with a computational complexity at most $HN_l$ in layer $l$. And finally, the tokens are propagated with computational complexity at most
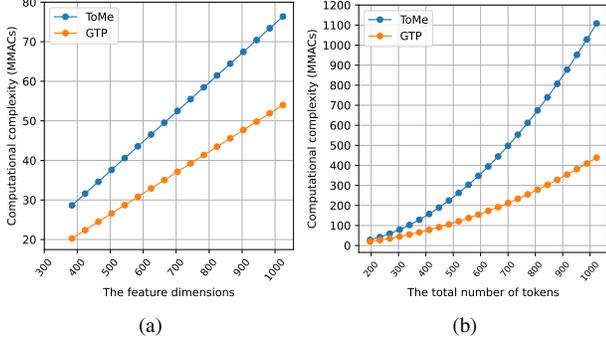
Figure 9. **Comparisons on the additional computational complexities introduced by ToMe [1] and our GTP.** We plot the computational complexity (measured in MMACs) with respect to the dimension of token features in (a) and the total number of tokens in (b).

$(N_l - M)MC$ in layer $l$. Consequently, the total additional computational complexity $G_{\text{GTP}}$ of GTP is

$$
\begin{aligned}
G_{\text{GTP}} &= N^2C + \sum_{l=1}^{L}(HN_l + (N_l - M)MC) \\
&= N^2C + \sum_{l=1}^{L}(H(N - (l-1)M) \\
&\quad + (N - (l-1)M - M)MC) \\
&= N^2C + LHN + LMNC - \frac{1}{2}(L^2 - L)HM \\
&\quad - \frac{1}{2}(L + L^2)M^2C.
\end{aligned}
$$
(15)

Given $N = 197, L = 12, H = 6, C = 384$ and $M = 8$ for DeiT-S, we can get $G_{\text{GTP}} \approx 20.1$MMACs, which is smaller than $G_{\text{ToMe}} \approx 28.3$MMACs. On DeiT-B where $N = 197, L = 12, M = 8, H = 12$ and $C = 768$, we observe that $G_{\text{GTP}} \approx 40.5$MMACs is much smaller than $G_{\text{ToMe}} \approx 57.3$MMACs. Figure 9 illustrates the additional computational complexity change with respect to the total number of tokens $N$ and the feature dimensions $C$. It is obvious that our GTP introduces far less additional computational complexity than ToMe, and the difference signifies when $N$ and $C$ increase. It indicates the efficiency of GTP on large-size ViTs whose feature dimensions may exceed 1024, as well as on ViTs for dense prediction tasks where the number of tokens may be more than 1024.