

Visual Dependency Transformers: Dependency Tree Emerges from Reversed Attention

Mingyu Ding^{13*} Yikang Shen² Lijie Fan³ Zhenfang Chen²
Zitian Chen⁴ Ping Luo¹ Josh Tenenbaum³ Chuang Gan²⁴

¹The University of Hong Kong ²MIT-IBM Watson AI Lab ³MIT ⁴UMass Amherst

Abstract

Humans possess a versatile mechanism for extracting structured representations of our visual world. When looking at an image, we can decompose the scene into entities and their parts as well as obtain the dependencies between them. To mimic such capability, we propose Visual Dependency Transformers (DependencyViT) ¹ that can induce visual dependencies without any labels. We achieve that with a novel neural operator called reversed attention that can naturally capture long-range visual dependencies between image patches. Specifically, we formulate it as a dependency graph where a child token in reversed attention is trained to attend to its parent tokens and send information following a normalized probability distribution rather than gathering information in conventional self-attention. With such a design, hierarchies naturally emerge from reversed attention layers, and a dependency tree is progressively induced from leaf nodes to the root node unsupervisedly.

DependencyViT offers several appealing benefits. (i) Entities and their parts in an image are represented by different subtrees, enabling part partitioning from dependencies; (ii) Dynamic visual pooling is made possible. The leaf nodes which rarely send messages can be pruned without hindering the model performance, based on which we propose the lightweight DependencyViT-Lite to reduce the computational and memory footprints; (iii) DependencyViT works well on both self- and weakly-supervised pretraining paradigms on ImageNet, and demonstrates its effectiveness on 8 datasets and 5 tasks, such as unsupervised part and saliency segmentation, recognition, and detection.

1. Introduction

Humans have a rich mental representation of our surrounding environments. When looking at an image (see Figure 1(a)), we can recognize the scene and also can

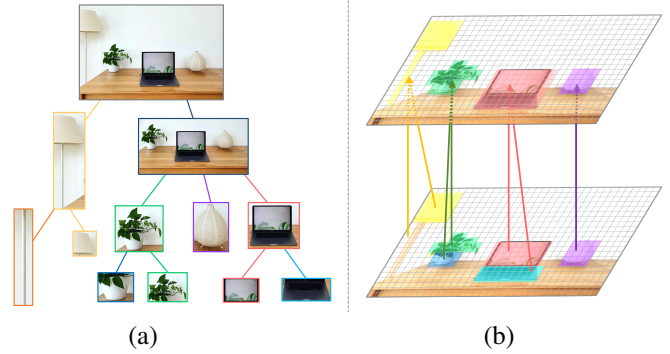


Figure 1. (a) is an example of hierarchical dependency structure. (b) illustrates the dynamic pooling and information aggregation process of DependencyViT.

quickly decompose it into hierarchical elements with dependencies, e.g., a laptop consisting of a screen and a keyboard is placed on the table. This ability to construct dependencies between objects (and/or their parts) serves as the cornerstone of human intelligence, enabling us to perceive, interact, and reason about the world.

From the pre-deeplearning era, many classical image dependency parsing algorithms [29, 31, 83, 87, 100, 119] have been proposed. For example, Bayesian framework [87], And-Or graph [31], and hierarchical probabilistic models [29, 83] for parsing images into their constituent visual patterns. Apart from that, Capsule Network [49, 74] shows the potential to learn geometrically organized parts from images. After that, visual grounding methods [14, 22, 25, 27, 105] try to align the semantic meaning between visual objects and words to distill effective structures for the vision branch from language. Similarly, human-object interaction approaches [46] learn the relationships between two objects, e.g., a boy “holds” an ice cream, from manually annotated labels. Such methods struggle to learn hierarchical visual structures, such as different parts of an object, unless exhaustive and time-consuming manual annotations are provided. Recently, vision-language (VL) grammar induction [89] proposes to extract shared hierarchical object

*This work was done when Mingyu was visiting MIT.

¹<https://github.com/dingmyu/DependencyViT>

dependencies for both vision and language unsupervisedly from image-caption pairs. However, the above works suffer two key issues: 1) the parsing relies heavily on supervision from natural language or human annotations rather than the image itself, and 2) their parsed structures are object-level based on a pre-trained object detection model, like Faster/Mask-RCNN [35, 71], hindering their generalizability in part-level and non-detector scenarios.

This paper answers a question naturally raised from the above issues: can we efficiently induce visual dependencies and build hierarchies from images without human annotations? Currently, visual parsing works mainly lie in semantic and instance segmentation. Unlike detector-based works that rely on pre-trained detectors, they parse the image at the pixel level, which is resource-intensive and costly. Inspired by vision transformers [26] that take image patches as input and leverage self-attention to perform interactions between patches, we propose to build a dependency tree at the patch level. Taking patches as basic elements and building a tree structure based on them has two benefits: 1) it unifies part-level and object-level dependencies, all of which are formulated into subtrees; 2) in the dependency structure, information can be aggregated from leaves to the parent (as shown in Figure 1(b)) to produce a hierarchy of representations for different parts and object along the path.

In practice, it is non-trivial to build the dependency tree with the standard transformer. Although the self-attention mechanism is designed to collect information dynamically from other patches, the number of attention heads constraints the number of tokens that a patch can attend to. 1) However, each parent could have an arbitrary number of children in a dependency tree, while each child only has one parent. Thus it's more straightforward for a node to select its parent instead of selecting the child. 2) Furthermore, the transformer treats each patch equally, it does not distinguish between root and leaf nodes. Contributions for different subtrees should be distinct.

Motivated by the above observations, in this work, we propose a dependency-inspired vision transformer, named Visual Dependency Transformers (DependencyViT). We propose three innovations to the standard self-attention, as shown in Figure 2. Firstly, to form a root-centric dependency parser, we introduce a reversed self-attention mechanism by transposing the adjacency matrix. In this way, leaf nodes can send information to their parents and form hierarchical subtrees. Secondly, we propose a message controller to determine how a node or subtree sends messages. Thirdly, a soft head selector is introduced to generate a unique dependency graph for each layer. As a result, self-attentions in DependencyViT naturally form a dependency tree parser. We did extensive studies in both supervised and self-supervised pretraining to show DependencyViT is capable of capturing either object- or part-level dependencies.

Intuitively, dependency parsing should ease scene understanding, as humans can understand complex scenes at a glance based on visual dependencies. Based on this, we further introduce a lightweight model DependencyViT-Lite by proposing a dynamic pooling scheme, reducing the computational cost largely. Within each subtree, we prune those leaf nodes with the least information received because they have sent information to their parent node. We show the pruned nodes can be retrieved by soft aggregations from their parents, preserving the model capability and dense representation capability.

We make three main contributions. (i) DependencyViT performs visual dependency parsing by reversed attention in self- or weakly-supervised manners. We demonstrate its effectiveness in both part-level and object-level parsing. (ii) We propose a visual dynamic pooling scheme for DependencyViT hence DependencyViT-Lite. The dependency tree can also be progressively built during the pruning process. (iii) Extensive experiments on both self- and weakly-supervised pretraining on ImageNet, as well as five downstream tasks, show the effectiveness of DependencyViT.

2. Related Work

Dependency Parsing in Vision. Unsupervised dependency parsing is a long-standing task in computer vision with many classical image dependency parsing algorithms that have been proposed in the pre-deeplearning era [29, 31, 83, 87, 100, 119]. Dating back to [87] proposed a Bayesian framework for parsing images into their constituent visual patterns. [31, 100, 119] surveyed on stochastic and context sensitive grammar of images with Bayesian framework, And-Or graph and probabilistic models. [29, 83] proposed to use hierarchical probabilistic models for detection and recognition of objects in cluttered environments.

In the deep learning era, a representative accomplishment is Capsule Network [74], where the activity vector of a capsule represents the instantiation parameters of an object part. After that, Stacked Capsule Autoencoders [49] leverages dynamic routing among capsules to automatically discover sub-patterns and recover the compositional relations on the MNIST dataset [21]. There are also works [28, 39, 43, 101, 102, 108] that further extend the composition relations in Capsule Networks and apply them to more tasks, *e.g.*, generative adversarial scenarios. However, it remains challenging to make them work on complex natural images. Most recently, supervised hierarchical semantic segmentation [53, 57, 58] became more popular. There are works to perform human parsing [93, 94] based on human part relations. Recently there are also attempts to perform part segmentation [7, 17, 42, 62] in an unsupervised manner. [79] explored spectral clustering on self-supervised features and pseudo labels on unsupervised saliency detection.

This work provides a *new perspective*, discovering visual dependencies automatically from neural attention in vision transformers. We believe it is of great significance to both the traditional grammar induction field, and the recent vision transformer and multimodal learning research. We provide an initial study that enables a flexible model that can simultaneously work on hierarchical parsing, scene graph, and downstream tasks like detection and segmentation. Furthermore, our model can adaptively induce different kinds of structures conditions on the given task.

Vision Transformers. ViT [26] first applies self-attention directly to a sequence of image patches. Works [15, 30, 38, 72, 73, 82, 85, 92, 98, 113] follows the discipline to stack multiple self-attention layers to model the information across patch tokens. After that hierarchical designs are widely adopted to vision transformers [1, 16, 23, 24, 32, 41, 51, 52, 55, 56, 63, 67, 84, 88, 89, 92, 97, 104, 109, 109–111, 114, 116, 117] for better efficiency and lower memory cost. For example, Swin [63], ViL [114], and HaloNet [88] apply local windows attention to the patch tokens, which reduce the quadratic complexity to linear, but lose the ability of long-range dependency modeling. PVT [92] and CvT [97] perform attention on the squeezed tokens to reduce the computational cost. However, previous transformer models fail to discover object parts in images and resolve their dependencies.

In this work, we focus on efficient transformers for dependency parsing, based on the standard ViT [26]. We propose DependencyViT, a dependency-inspired vision transformer built on reversed self-attention, which captures hierarchies and dependencies between patches automatically. DependencyViT is orthogonal and seamlessly compatible with the SoTA transformer training methods, makes it more attractive than traditional grammar models from a practical perspective. Moreover, we show that the standard ViT layout can be highly efficient with our DependencyViT-Lite and dynamic pooling technique.

3. Method

This work proposes Visual Dependency Transformers (DependencyViT), a dependency-inspired backbone model based on reversed self-attention, capturing dependencies between patches automatically from self- or weakly-supervised signals for vision tasks.

Preliminaries. Let us assume a $\mathbb{R}^{N \times C}$ dimensional visual feature \mathbf{X} , where N is the number of total image patches and C is the number of token dimensions. The number of heads is H . The standard (forward) multi-head self-

attention is defined as:

$$\begin{aligned} \mathcal{A}_f(\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \text{Concat}(\text{head}_1, \dots, \text{head}_H) \mathbf{W}_o \\ \text{where } \text{head}_h &= \text{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h) \\ &= \text{softmax} \left[\frac{\mathbf{Q}_h (\mathbf{K}_h)^T}{\sqrt{C_h}} \right] \mathbf{V}_h \end{aligned} \quad (1)$$

where $\mathbf{Q}_h = \mathbf{X} \mathbf{W}_h^Q$, $\mathbf{K}_h = \mathbf{X} \mathbf{W}_h^K$, and $\mathbf{V}_h = \mathbf{X} \mathbf{W}_h^V$ are $\mathbb{R}^{N \times C_h}$ dimensional visual features of H heads, $\mathbf{X} \in \mathbb{R}^{N \times C}$ denotes the input feature and $\mathbf{W}_h \in \mathbb{R}^{C \times C_h}$ denotes the projection weights of the h_{th} head for $\mathbf{Q}, \mathbf{K}, \mathbf{V}$, $C = C_h * H$, and \mathbf{W}_o is the weight of the output projection. $\mathbf{A}_F = \text{softmax}(\mathbf{Q} \mathbf{K}^T) \in \mathbb{R}^{H \times N \times N}$ is called the attention matrix of the layer. In subsequent sections, we will omit the head dimension and focus on analyzing the attention within a single head.

3.1. Reversed Attention

The standard self-attention mechanism learns the $N \times N$ attention adjacency matrix to exchange information between different image patches. It treats all patches equally and does not follow a tree or graph structure, *i.e.*, it does not distinguish root and leaf nodes. To generate an adjacency graph, let us assume that each node can find its parent node via the $\text{argmax}(\cdot)$ function since the second dimension of the matrix follows a normalized probability distribution. In this case, the forward self-attention works by gathering information from parent nodes following the soft probabilities. All the nodes receive information from others, and eventually, they are dominated by the root node and the structural information of the image is lost. This learning scheme may perform well on visual recognition tasks due to its powerful attentive fusion and interaction capabilities, but it is not based on explicit hierarchical structures and dependencies.

Ideally, to build a dependency tree, we need to identify which patches are child nodes or parent nodes, so that information can be progressively aggregated to the root node. In turn, the root node distributes messages to leaf nodes. We achieve this by proposing reversed self-attention, which simply transposes the adjacency probability matrix so that the child node sends messages to the parent node. Considering each element a_{ij} in the attention matrix \mathbf{A} , we have:

$$a_{ij} = \text{softmax} \left(\left\{ \frac{q_i k_j}{\sqrt{C_h}} \right\}_{j \in [0, N)} \right)_j, \quad (2)$$

where q_i is the i_{th} element of \mathbf{Q} , and k_j is the j_{th} element of \mathbf{K} . Then, after transposing the matrix \mathbf{A} , the information

flow also changes as follows:

$$o_i = \left(\sum_j a_{ij} v_j \right) \mathbf{W}_o \implies o_i = \left(\sum_j a_{ji} v_j \right) \mathbf{W}_o, \quad (3)$$

where o_i denotes the i_{th} output and \mathbf{W}_o is the weight of the output projection. We can see the child node ‘receive’ messages in forward attention but ‘send’ messages in reversed attention following the softmax probability distribution. Each child node has only one parent, but each parent node can have multiple children. In this way, information can be collected progressively from leaf nodes to the root node through multiple reversed attention layers. At the same time, the dependency tree is also built bottom-up, and different subtrees may represent part-level or object-level semantics.

3.2. Dependency Block

Simply applying transposed attention matrices does not guarantee a good dependency graph induced. This is because: (i) The amount of token that a patch can attend to is controlled by multiple attention heads, thus the dependency graph is not unique. (ii) Contributions for different subtrees are not well distinguished. In some downstream tasks like image classification, foreground and background trees should be distinct. To solve the above questions: we further introduce two modules: a head selector and a message controller. An overview of our dependency block is shown in Figure 2.

Head Selector. The head selector \mathbf{P} is used to choose proper reversed attention heads for dependency induction. We obtain it by applying the `softmax()` function on the linear projections of the input tokens: $\mathbf{P} = \text{softmax}(\mathbf{X}\mathbf{W}_p)$, where $\mathbf{W}_p \in \mathbb{R}^{C \times H}$ is the projection weight. By the head selector, we can build dependencies over all attention heads following the learnable soft probabilities and generate a unique dependency graph for each layer.

Message Controller. Similarly, the message controller \mathbf{M} is learnable weights imposed on tokens during reversed self-attention. The goal of \mathbf{M} is to determine the extent to which a node or a subtree sends messages. Specifically, we use two linear projection layers (who have the dimensions $\mathbb{R}^{C \times \frac{C}{2}}$ and $\mathbb{R}^{\frac{C}{2} \times 1}$) with a GELU activation [37] between them to learn it. After that, a `sigmoid()` function is used to get the probability in $[0, 1]$ of sending messages.

Note that the weights learned by the message controller are cumulative across all layers. The message controller \mathbf{M} in the i_{th} layer is computed by $\mathbf{M}_1 \cdot \mathbf{M}_2 \dots \cdot \mathbf{M}_i$, where the subscript represents the index of the layer. We also use \mathbf{M} to weight the pooling to get the final representation over all

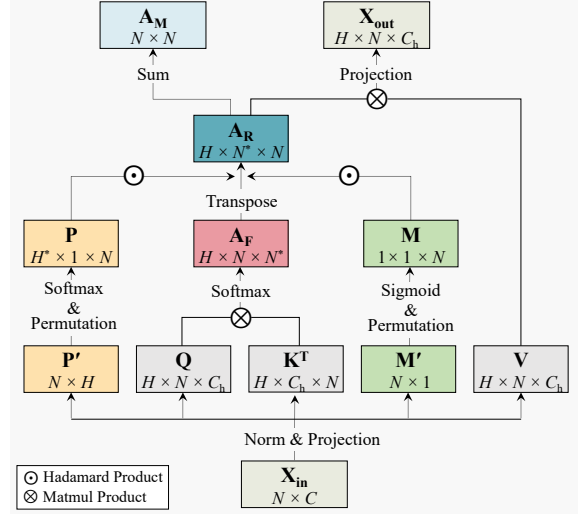


Figure 2. An architecture overview of our proposed reversed attention block in DependencyViT. FeedForward Networks (FFNs) and residual paths are omitted here. The input and output tokens are \mathbf{X}_{in} and \mathbf{X}_{out} with the number of tokens N and token dimensions C , respectively. The number of attention heads is H , and the per-head token dimension is C_h . We obtain the reversed attention matrices \mathbf{A}_R by transposing the forward attention weights \mathbf{A}_F with a head selector \mathbf{P} and a message controller \mathbf{M} imposing on it. After that, the soft dependency mask \mathbf{A}_M is induced by applying summation on \mathbf{A}_R over the head dimension. ‘*’ indicates the dimension that is normalized by softmax probability distribution.

patches. It has two benefits: (i) If a node does not send information in a layer, it keeps the status in subsequent layers, making the induced structure clearer. (ii) Different subtrees are weighted differently, which filters meaningless patches, benefiting downstream tasks like recognition and detection.

In summary, we have the reversed attention $\mathbf{A}_R = \mathbf{A}_F \cdot \mathbf{P} \cdot \mathbf{M}$ with dimensional permutations, where \mathbf{A}_F is the forward attention, as shown in Figure 2. We then compute the soft dependency mask by applying the `sum()` operator on \mathbf{A}_R over the head dimension. The dependency graph and tree structure are then obtained by `argmax()` and the chu-liu-edmonds algorithm [19], respectively.

3.3. Dynamic Pooling based on Dependencies

Our dependency block is able to learn dynamic and comprehensive information flow between patches for dependency induction. Intuitively, with such visual dependencies, scene understanding can be simplified with less computational effort as most of the information can be represented by a few nodes. With this inspiration, we introduce a dynamic visual pooling scheme that reduces the computational cost largely (*i.e.*, FLOPs and GPU memory), and propose a lightweight model DependencyViT-Lite.

Specifically, we rank and prune those leaf nodes which have the least information received, because 1) they are not

the parent of any node and 2) they rarely transmit messages or they have sent enough information to their parent nodes. We progressively prune the leaf nodes with the least messages as the depth of the network increases. In this way, the memory and resource costs are largely reduced. Meanwhile, the tree architecture is still maintained by recording relationships between the pruned nodes and their parents to form a complete tree. Most importantly, DependencyViT-Lite is able to perform dense prediction tasks though some of its tokens are removed. According to the dependency graph, we retrieve those pruned nodes by a soft aggregation from their parents.

3.4. Model Analysis and Protocols

Model Instantiation. In this work, we follow the design strategy of the standard ViT (DeiT) [26, 85]. To show the efficiency and effectiveness of our model, we choose two different model sizes and build DependencyViT-T and DependencyViT-S based on tiny and small ViTs as backbones, respectively. We set the number of attention heads $H = 12$, the number of dependency blocks $L = 12$ with residual paths and FFNs of ratio 4 as in standard ViT. We set the token dimensions $C = \{192, 384\}$ for tiny and small models, respectively. Take an image with an arbitrary resolution, a C -dimensional $16 \times$ down-sampling feature is obtained after the patch embedding layer. There are no overlaps between any of the two patches. Conditional positional encoding is used as in [18]. Based on our observation that the ‘cls’ token passes information between two visual patches and leads to confusion in dependencies, we remove it from our model. For DependencyViT-Lite models, we prune 16% number of nodes (e.g., 32 of 196) at the $\{2, 5, 8, 11\}_{th}$ layers, respectively.

Complexity Analysis. Simply applying the standard global self-attention leads to a complexity of $O(2N^2C + 12NC^2)$, which contains $O(2N^2C)$ for self-attentions, $O(4NC^2)$ for linear projections, and $O(8NC^2)$ for feedforward networks (FFNs). Our head selector and message controller lead to additional costs of $O(NCH)$ and $O(NC)$, respectively, which are much smaller than the costs of other components. In contrast, our DependencyViT-Lite reduces the number of tokens N to $0.3N$ and even smaller through dynamic pooling, which lowers the complexity exponentially (to 10% and even smaller). DependencyViT-Lite can run with batch sizes more than three times that of ViT on a same GPU.

Pretraining Protocols. We apply two different pretraining methods on DependencyViT: weakly-supervised and self-supervised. The first one is supervised pretraining on ImageNet by leveraging the information in class-level labels. The supervision encourages the model to learn high-level object-aware semantics, based on which DependencyViT learns to model object-aware dependencies by gathering information from subtrees to the root node (centered object).

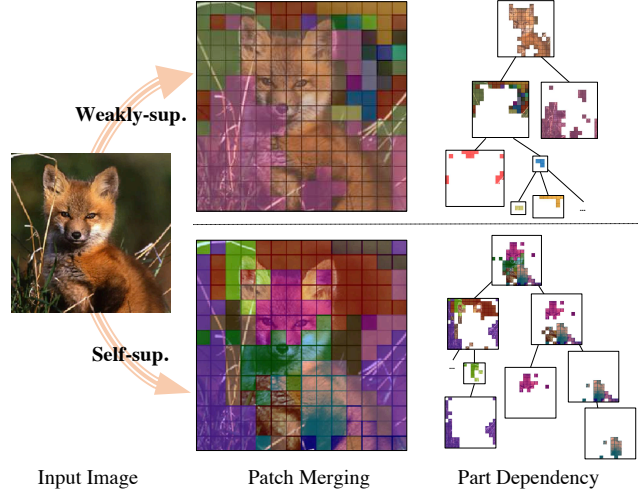


Figure 3. Visualizations of dependency trees parsed by self- and weakly-supervised pretrained DependencyViT, respectively. Patches are aggregated gradually until the root node is formed. To facilitate observation, the background area is not filled to the root node. It can be seen that weakly-supervised DependencyViT focuses more on the whole object, while the self-supervised DependencyViT captures more fine-grained part-aware dependencies.

For self-supervised pretraining, we take inspiration from recent contrastive learning and masked image modeling methods [3, 8, 11, 13, 34, 118] as they can learn both object-level global representations and part-level local features. Specifically, we follow the same pretraining protocol as iBOT [118] (e.g., employ self-distillation and masked image modeling on DependencyViT) and enjoy the benefit of its powerful pretraining capabilities. After pretraining, DependencyViT can establish a dependency tree for an unseen image, containing part-to-part, part-to-object, and object-to-object dependencies.

Figure 3 shows the dependency trees of an image from ImageNet parsed by weakly-supervised and self-supervised pretrained DependencyViT, respectively. It can be seen that weakly-supervised pretrained DependencyViT focuses more on the entire object, while the self-supervised pretrained DependencyViT can capture more fine-grained part-aware dependencies. The parsed dependency tree is expected to help many downstream tasks, such as saliency detection and part segmentation. For more analysis and detailed settings, please refer to Appendix.

4. Experiments

In this section, we conduct extensive experiments to show the effectiveness of DependencyViT and DependencyViT-Lite on visual parsing and recognition. They are: unsupervised part segmentation on the Pascal-Part [12] and Car-Parts [68] datasets; unsupervised saliency detection on the ECSSD [78], DUTS [90] and

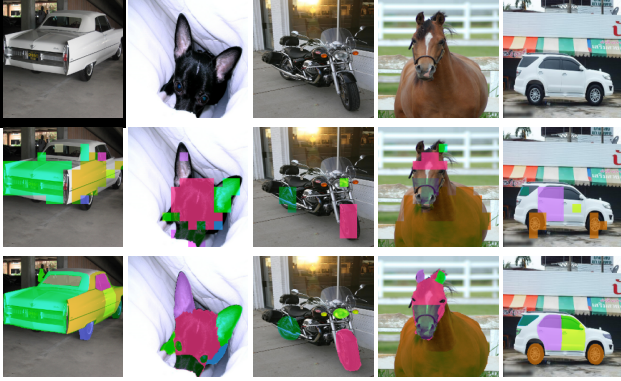


Figure 4. Visualization of part partitioning on the Pascal-Part [12] and Car-Parts [68] datasets. From top to bottom: 1) the original image; 2) our generated part mask of which each color represents a subtree in the hierarchy; 3) the ground truth part segments.

DUT-OMRON [107] datasets; dependency parsing on the COCO dataset [61]; and image classification on ImageNet-1K [20].

4.1. Unsupervised Part Segmentation

To show the effectiveness of DependencyViT on visual dependency parsing, we apply it to the unsupervised part segmentation task without part labels, which is challenging and under-explored as it requires a comprehensive dependency understanding between parts. Considering available part parsing datasets, *e.g.*, Pascal-Part [12] and Car-Parts [68], are of small resolution and data scale, tiny ViT model is enough to work on this situation and further scaling model size up brings no gains. We take DependencyViT-T as our base model.

Both weakly- and self-supervised pretrained models are evaluated. For DependencyViT, we average the learned dependency masks of all layers, and then leverage the Chu-Liu-Edmonds maximum spanning algorithm [19] to generate the dependency tree. After that, we perform matching between all subtrees and part segments by the Hungarian maximum matching algorithm [50] and compute the mean intersection over union (mIoU) and mean accuracy (mAcc) metrics for evaluation. Note that we remove small part regions from evaluation for more reliable results. We take DeiT-Tiny [85] as the baseline. Since there are no explicit dependencies in it, the Naïve solution is to partition the patches in their latent representation space by k-means clustering [64] (k is set to 20 in this paper). To get a stronger baseline, we also build tree structures on DeiT by applying the maximum spanning algorithm on its mean pooled attention map. For self-supervised models, DependencyViT-T is evaluated with the maximum spanning algorithm for dependency tree generation. We take the self-supervised iBOT (tiny DeiT) as a strong baseline for fair comparison.

From the results shown in Table 1, we observe that DependencyViT consistently outperforms the baseline meth-



Figure 5. Visualization for unsupervised saliency detection on the ECSSD [78], DUTS [90] and DUT-OMRON [107] datasets. From top to bottom: 1) the original image; 2) our generated saliency mask; 3) our results post-processed by the bilateral solver [4]; 4) the ground truth part partitions.

ods by a large margin on both weakly- and self-supervised settings and two datasets, demonstrating the effectiveness of our dependency parsing. Self-supervised DependencyViT shows better performance than the weakly-supervised one as it can learn more fine-grained dependencies. We visualize our patch-level part partitioning results in Figure 4.

4.2. Unsupervised Saliency Detection

Besides part-level partitioning, DependencyViT can also work on object-level comprehensions, thanks to its built-in hierarchical dependencies. We evaluate the unsupervised saliency detection results of DependencyViT on three datasets. Except baseline methods [54, 66, 106, 120] that are specifically designed for the task based on pseudo-labels, we evaluate weakly-supervised DeiT for fair comparison. Following [95], we leverage normalized cut [77] on token representations to get the salient area of an image for DeiT. For DependencyViT, the soft dependency mask is added to the representation for better results. Bilateral solver [4] is used as post-processing for segment smoothing.

Figure 5 visualizes the saliency detection map of our method DependencyViT (weakly-sup.). From the figure and Table 2, we see that: 1) DependencyViT is superior to its counterparts, including the pseudo label-based saliency detection methods and DeiT. 2) Weakly-supervised DependencyViT outperforms the self-supervised one, indicating weakly-supervised model is better at modeling object-level semantics. 3) The failure case in the last column of Figure 5 demonstrates how DependencyViT works. The two birds belong to the same semantic category but different objects, hence two subtrees.

To verify the effectiveness of dependency for object-level understanding, we make ablative comparisons by whether adding the soft dependency (+dependency) to the feature representation, see Figure 7). We see that 1) the dependency mask improves the performance significantly, showing the effectiveness of DependencyViT in learning

Table 1. Part segmentation results on the Pascal-Part [12] and Car-Parts [68] datasets. ‘clustering’ indicates applying k-means [64] on feature representations; ‘maximum spanning’ denotes the dependency tree is generated by Chu-Liu-Edmonds maximum spanning algorithm [19].

Method	Pretraining Type	Part Discovery by	Pascal-Part [12]		Car-Parts [68]	
			mIoU (%)	mAcc (%)	mIoU (%)	mAcc (%)
DeiT [85]	weakly-sup.	clustering	7.2	22.6	8.9	29.5
DeiT [85]	weakly-sup.	maximum spanning	18.9	35.5	17.8	37.7
DependencyViT	weakly-sup.	clustering	11.6	31.7	10.9	29.7
DependencyViT	weakly-sup.	maximum spanning	23.2	41.7	22.6	40.0
iBOT [118]	self-sup.	maximum spanning	25.1	44.8	25.7	46.1
DependencyViT	self-sup.	maximum spanning	28.7	47.9	27.0	47.2

Table 2. Unsupervised saliency detection on ECSSD [78], DUTS [90] and DUT-OMRON [107]. Tiny models, token normalized cut [77,95] and bilateral solver [4] post-processing are used for DeiT and DependencyViT.

Method	ECSSD [78]			DUTS [90]			DUT-OMRON [107]		
	$maxF_\beta$ (%)	IoU (%)	Acc. (%)	$maxF_\beta$ (%)	IoU (%)	Acc. (%)	$maxF_\beta$ (%)	IoU (%)	Acc. (%)
DeepUSPS [66]	58.4	44.0	79.5	42.5	30.5	77.3	41.4	30.5	77.9
HS [106]	67.3	50.8	84.7	50.4	36.9	82.6	56.1	43.3	84.3
wCtr [120]	68.4	51.7	86.2	52.2	39.2	83.5	54.1	41.6	83.8
WSC [54]	68.3	49.8	85.2	52.8	38.4	86.2	52.3	38.7	86.5
DeiT [85]	49.3	40.5	72.7	34.2	26.8	72.7	33.2	27.2	71.1
DependencyViT (self-sup.)	62.1	55.0	78.4	43.0	35.9	73.2	32.5	28.0	67.2
DependencyViT (weakly-sup.)	62.0	48.4	83.6	53.8	37.0	87.5	52.0	39.7	88.4

Table 4. Comparison of image classification on ImageNet-1K. All models are trained and evaluated with 224×224 resolution. * denotes the method can not be used for dense predictions.

Model	Hierarchical	Cost	#Params (M)	FLOPs (G)	Top-1 (%)
ResNet-18 [36]	✓	low	11.7	1.8	69.9
ConvMixer-512/16 [86]	×	high	5.4	–	73.7
DeiT-Tiny/16 [85]	×	high	5.7	1.3	72.2
CrossViT-Tiny [10]	×	high	6.9	1.6	73.4
PVT-Tiny [92]	✓	low	13.2	1.9	75.1
DependencyViT-Lite-T	×	low	6.2	0.8	73.7
DependencyViT-T	×	high	6.2	1.3	75.4
ResNet-50 [36]	✓	low	25.0	4.1	76.2
ConvMixer-768/32 [86]	×	high	21.1	–	80.2
DeiT-Small/16 [85]	×	high	22.1	4.5	79.8
CrossViT-Small [10]	×	high	26.7	5.6	81.0
PVT-Small [92]	✓	low	24.5	3.8	79.8
Swin-Tiny [63]	✓	low	28.3	4.5	81.2
CvT-13 [97]	✓	high	20.0	4.5	81.6
DynamicViT-LV-S/0.5 [70]*	×	–	26.9	3.7	82.0
PVTv2-B2 [91]	✓	low	25.4	4.0	82.0
DependencyViT-Lite-S	×	low	24.0	3.0	80.6
DependencyViT-S	×	high	24.0	5.0	82.1

object-level dependencies; and 2) the bilateral solver brings considerable improvement over all models.

4.3. Visualization

As shown in Figure 6, we visualize our visual dependency parsing on images from the COCO dataset, which

does not overlap with the pretraining ImageNet dataset. We can see that the foreground and the background areas are represented by different subtrees, which further construct the overall scene dependency tree. The root subtree is generally an important part of the foreground object.

More downstream experiments, *e.g.*, semantic segmentation on ADE20K [115], object detection on the COCO dataset [61], and video recognition on Kinetics-400 [45], can be found in Appendix.

4.4. Visual Recognition

We show DependencyViT can work as a visual backbone for recognition and its downstream tasks. Two different model configurations, *i.e.*, DependencyViT and DependencyViT-Lite, are evaluated and compared with many counterparts. We make the following summaries from Table 4. 1) DependencyViT outperforms all counterparts, *e.g.*, 3.2% and 2.3% improvements over DeiT-Tiny and DeiT-small, respectively, indicating dependency parsing is likely to contribute to visual recognition tasks. 2) DependencyViT-Lite is the most efficient one (0.8 GFLOPs only) of all models and shows good performance, demonstrating the effectiveness of our progressively dynamic pooling. Typically, hierarchical transformers are more efficient and save computations for downstream tasks. Our DependencyViT-Lite reduces costs through induced dependencies even using a standard ViT layout. 3) DynamicViT [70] is a pruning-based transformer for the classifi-

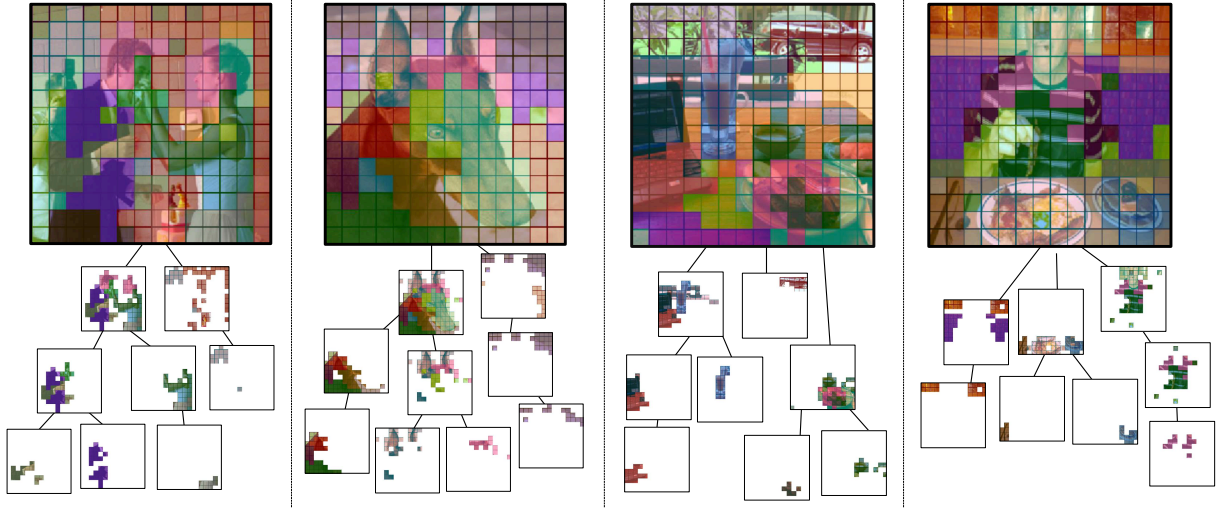


Figure 6. Visualizations of dependency trees parsed by self-supervised DependencyViT (small). Different colors represent different subtrees. Here we ignore the nodes in the small region (less important) and display the main subtrees.

Table 3. Comparisons of image classification on ImageNet-1K. All models (tiny) are trained and evaluated with 224×224 resolution.

Model	Direction	Head Selector	Message Controller	#Params (M)	FLOPs (G)	Top-1 (%)
Baseline (DeiT) [85]	forward	×	×	5.7	1.3	73.3
Forward + P	forward	✓	×	5.7	1.3	73.4
Forward + M	forward	×	✓	6.1	1.3	74.8
Forward + P + M	forward	✓	✓	6.2	1.3	74.8
Reverse + P	reverse	✓	×	5.7	1.3	73.6
Reverse + M	reverse	×	✓	6.1	1.3	74.9
Reverse + P + M (DependencyViT)	reverse	✓	✓	6.2	1.3	75.4
DependencyViT-Lite (forward)	forward	✓	✓	6.2	0.8	71.1
DependencyViT-Lite (reverse)	reverse	✓	✓	6.2	0.8	73.7

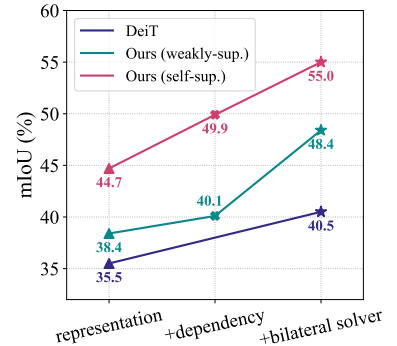
cation task. However, it can not perform dense predictions because the information of its pruned patches is lost. On the contrary, the pruned nodes in our DependencyViT-Lite can be retrieved from their parents for dense predictions, showing the importance of dependency induction.

4.5. Ablation Study

We perform ablation studies on tiny models in Table 3. P denotes the head selector, and M denotes the message controller. We use ‘forward’ and ‘backward’ to indicate the attention direction. We can see that the head selector brings smaller gains than the message controller. And the gains in reverse attention are larger than gains in forward attention. Both the head selector and the message controller are important to dependency induction and the dynamic pooling scheme, *i.e.*, DependencyViT-Lite.

More ablation studies and downstream experiments can be found in Appendix.

Figure 7. Ablative comparisons (tiny) of saliency detection on ECSSD dataset.



5. Conclusion

This paper studies patch-level visual dependency parsing using our proposed DependencyViT. We show that the reversed self-attention mechanism in transformers can naturally capture long-range visual dependencies between image patches. With reversed self-attention, a child node is trained to attend to its parent and send the information to the parent node, and a hierarchical dependency tree can be established automatically. Furthermore, dynamically image pooling is made possible by learned dependencies, *i.e.*, merging child nodes into their corresponding parent nodes, based on which we propose a lightweight model DependencyViT-Lite. Extensive experiments on both self- and weakly-supervised pretraining on ImageNet, as well as five downstream tasks, show the model’s effectiveness.

Limitations. Although our work achieves good performance on many tasks by visual dependency induction, it is an initial study with a fixed patch size and efficient settings

The current patch size limits its performance on small objects. We will explore more and further scale up our model. The proposed approach has no ethical or societal issues on its own, except those inherited from computer vision.

Acknowledgements. Ping Luo is partially supported by the National Key R&D Program of China No.2022ZD0161000 and the General Research Fund of HK No.17200622. Chuang Gan was supported by the MIT-IBM Watson AI Lab, DARPA MCS, DSO grant DSOCO21072, and gift funding from MERL, Cisco, Sony, and Amazon.

References

- [1] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. *NeurIPS*, 34, 2021. [3](#)
- [2] Waleed Ammar, Chris Dyer, and Noah A Smith. Conditional random field autoencoders for unsupervised structured prediction. *Advances in Neural Information Processing Systems*, 27, 2014. [15](#)
- [3] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT pre-training of image transformers. *arXiv:2106.08254*, 2021. [5](#)
- [4] Jonathan T Barron and Ben Poole. The fast bilateral solver. In *European conference on computer vision*, pages 617–632. Springer, 2016. [6, 7](#)
- [5] Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain: a unified image embedding for classes and instances. *arXiv preprint arXiv:1902.05509*, 2019. [17](#)
- [6] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding. *arXiv preprint arXiv:2102.05095*, 2(3):4, 2021. [15](#)
- [7] Sandro Braun, Patrick Esser, and Björn Ommer. Unsupervised part discovery by unsupervised disentanglement. In *DAGM German Conference on Pattern Recognition*, pages 345–359. Springer, 2020. [2](#)
- [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. [5](#)
- [9] Boyu Chen, Peixia Li, Baopu Li, Chuming Li, Lei Bai, Chen Lin, Ming Sun, Junjie Yan, and Wanli Ouyang. Psvit: Better vision transformer via token pooling and attention sharing. *arXiv preprint arXiv:2108.03428*, 2021. [14](#)
- [10] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *ICCV*, 2021. [7](#)
- [11] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised representation learning. *arXiv preprint arXiv:2202.03026*, 2022. [5](#)
- [12] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1971–1978, 2014. [5, 6, 7](#)
- [13] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, pages 9640–9649, 2021. [5](#)
- [14] Zhenfang Chen, Jiayuan Mao, Jiajun Wu, Kwan-Yee Kenneth Wong, Joshua B Tenenbaum, and Chuang Gan. Grounding physical concepts of objects and events through dynamic visual reasoning. In *International Conference on Learning Representations*. [1](#)
- [15] Zitian Chen, Yikang Shen, Mingyu Ding, Zhenfang Chen, Hengshuang Zhao, Erik Learned-Miller, and Chuang Gan. Mod-squad: Designing mixture of experts as modular multi-task learners. *arXiv preprint arXiv:2212.08066*, 2022. [3](#)
- [16] Zhengsu Chen, Lingxi Xie, Jianwei Niu, Xuefeng Liu, Longhui Wei, and Qi Tian. Visformer: The vision-friendly transformer. In *ICCV*, pages 589–598, 2021. [3](#)
- [17] Subhabrata Choudhury, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Unsupervised part discovery from contrastive reconstruction. *Advances in Neural Information Processing Systems*, 34:28104–28118, 2021. [2](#)
- [18] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Conditional positional encodings for vision transformers. *Arxiv preprint 2102.10882*, 2021. [5](#)
- [19] Yoeng-Jin Chu. On the shortest arborescence of a directed graph. *Scientia Sinica*, 14:1396–1400, 1965. [4, 6, 7](#)
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. [6, 16](#)
- [21] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. [2](#)
- [22] Mingyu Ding, Zhenfang Chen, Tao Du, Ping Luo, Josh Tenenbaum, and Chuang Gan. Dynamic visual reasoning by learning differentiable physics models from video and language. *Advances In Neural Information Processing Systems*, 34:887–899, 2021. [1](#)
- [23] Mingyu Ding, Xiaochen Lian, Linjie Yang, Peng Wang, Xiaojie Jin, Zhiwu Lu, and Ping Luo. Hr-nas: Searching efficient high-resolution neural architectures with lightweight transformers. In *CVPR*, 2021. [3](#)
- [24] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. *arXiv preprint arXiv:2204.03645*, 2022. [3, 16](#)
- [25] Mingyu Ding, Yan Xu, Zhenfang Chen, David Daniel Cox, Ping Luo, Joshua B Tenenbaum, and Chuang Gan. Embodied concept learner: Self-supervised learning of concepts and mapping through instruction following. In *Conference on Robot Learning*, pages 1743–1754. PMLR, 2023. [1](#)
- [26] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold,

- Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 3, 5, 16
- [27] Ye Du, Zehua Fu, Qingjie Liu, and Yunhong Wang. Visual grounding with transformers. *arXiv preprint arXiv:2105.04281*, 2021. 1
- [28] Marzieh Edraki, Nazanin Rahnavard, and Mubarak Shah. Subspace capsule network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10745–10753, 2020. 2
- [29] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010. 1, 2
- [30] Benjamin Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. In *ICCV*, pages 12259–12269, 2021. 3
- [31] Feng Han and Song-Chun Zhu. Bottom-up/top-down image parsing with attribute grammar. *IEEE transactions on pattern analysis and machine intelligence*, 31(1):59–73, 2008. 1, 2
- [32] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. *arXiv preprint arXiv:2204.07143*, 2022. 3
- [33] Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. Unsupervised learning of syntactic structure with invertible neural projections. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1292–1302, 2018. 15
- [34] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. 5
- [35] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017. 2, 14, 17
- [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 7, 16
- [37] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 4
- [38] Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. In *ICCV*, pages 11936–11945, 2021. 3
- [39] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with em routing. In *International conference on learning representations*, 2018. 2
- [40] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefer, and Daniel Soudry. Augment your batch: Improving generalization through instance repetition. In *CVPR*, pages 8129–8138, 2020. 17
- [41] Zilong Huang, Youcheng Ben, Guozhong Luo, Pei Cheng, Gang Yu, and Bin Fu. Shuffle transformer: Rethinking spatial shuffle for vision transformer. *arXiv preprint arXiv:2106.03650*, 2021. 3
- [42] Wei-Chih Hung, Varun Jampani, Sifei Liu, Pavlo Molchanov, Ming-Hsuan Yang, and Jan Kautz. Scops: Self-supervised co-part segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 869–878, 2019. 2
- [43] Ayush Jaiswal, Wael AbdAlmageed, Yue Wu, and Premkumar Natarajan. Capsulegan: Generative adversarial capsule network. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 2
- [44] 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. *NeurIPS*, 34, 2021. 17
- [45] Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 7
- [46] Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, and Hyunwoo J Kim. Hotr: End-to-end human-object interaction detection with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 74–83, 2021. 1
- [47] Dan Klein and Christopher D Manning. Corpus-based induction of syntactic structure: Models of dependency and constituency. In *Proceedings of the 42nd annual meeting of the association for computational linguistics (ACL-04)*, pages 478–485, 2004. 15
- [48] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Bin Ren, Minghai Qin, Hao Tang, and Yanzhi Wang. Spvit: Enabling faster vision transformers via soft token pruning. *arXiv preprint arXiv:2112.13890*, 2021. 14
- [49] Adam Kosiorek, Sara Sabour, Yee Whye Teh, and Geoffrey E Hinton. Stacked capsule autoencoders. *Advances in neural information processing systems*, 32, 2019. 1, 2
- [50] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 6
- [51] Changlin Li, Tao Tang, Guangrun Wang, Jiefeng Peng, Bing Wang, Xiaodan Liang, and Xiaojun Chang. Bossnas: Exploring hybrid cnn-transformers with block-wisely self-supervised neural architecture search. In *ICCV*, pages 12281–12291, 2021. 3
- [52] Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *arXiv preprint arXiv:2201.09450*, 2022. 3
- [53] Liulei Li, Tianfei Zhou, Wenguan Wang, Jianwu Li, and Yi Yang. Deep hierarchical semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1246–1257, 2022. 2
- [54] Nianyi Li, Bilin Sun, and Jingyi Yu. A weighted sparse coding framework for saliency detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5216–5223, 2015. 6, 7

- [55] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Improved multiscale vision transformers for classification and detection. *arXiv: Computer Vision and Pattern Recognition*, 2021. 3
- [56] Yawei Li, Kai Zhang, Jie Zhang Cao, Radu Timofte, and Luc Van Gool. Localvit: Bringing locality to vision transformers. *arXiv preprint arXiv:2104.05707*, 2021. 3
- [57] Zhiheng Li, Wenxuan Bao, Jiayang Zheng, and Chenliang Xu. Deep grouping model for unified perceptual parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4053–4063, 2020. 2
- [58] Xiaodan Liang, Hongfei Zhou, and Eric Xing. Dynamic-structured semantic propagation network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 752–761, 2018. 2
- [59] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 2117–2125, 2017. 15
- [60] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017. 16, 17
- [61] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. 6, 7, 14, 17
- [62] Shilong Liu, Lei Zhang, Xiao Yang, Hang Su, and Jun Zhu. Unsupervised part segmentation through disentangling appearance and shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8355–8364, 2021. 2
- [63] 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. 3, 7, 14, 15, 16, 17
- [64] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982. 6, 7
- [65] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 17
- [66] Tam Nguyen, Maximilian Dax, Chaithanya Kumar Mumtadi, Nhung Ngo, Thi Hoai Phuong Nguyen, Zhongyu Lou, and Thomas Brox. Deepusps: Deep robust unsupervised saliency prediction via self-supervision. *Advances in Neural Information Processing Systems*, 32, 2019. 6, 7
- [67] Zizheng Pan, Bohan Zhuang, Jing Liu, Haoyu He, and Jianfei Cai. Scalable visual transformers with hierarchical pooling. *arXiv preprint arXiv:2103.10619*, 2021. 3
- [68] Kitsuchart Pasupa, Phongsathorn Kittiworapanya, Napasin Hongngern, and Kuntpong Woraratpanya. Evaluation of deep learning algorithms for semantic segmentation of car parts. *Complex & Intelligent Systems*, pages 1–13, May 2021. 5, 6, 7
- [69] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992. 17
- [70] 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. 7, 14
- [71] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 2
- [72] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *NeurIPS*, 34, 2021. 3
- [73] Michael S. Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: What can 8 learned tokens do for images and videos? *arXiv: Computer Vision and Pattern Recognition*, 2021. 3
- [74] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *Advances in neural information processing systems*, 30, 2017. 1, 2
- [75] Yikang Shen, Shawn Tan, Sordani Alessandro, Li Peng, Jie Zhou, and Aaron Courville. Unsupervised dependency graph network. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022. 16
- [76] Yikang Shen, Yi Tay, Che Zheng, Dara Bahri, Donald Metzler, and Aaron Courville. Structformer: Joint unsupervised induction of dependency and constituency structure from masked language modeling. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7196–7209, 2021. 16
- [77] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000. 6, 7
- [78] Jianping Shi, Qiong Yan, Li Xu, and Jiaya Jia. Hierarchical image saliency detection on extended cssd. *IEEE transactions on pattern analysis and machine intelligence*, 38(4):717–729, 2015. 5, 6, 7
- [79] Gyungin Shin, Samuel Albanie, and Weidi Xie. Unsupervised salient object detection with spectral cluster voting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3971–3980, 2022. 2
- [80] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, page 1100612. International Society for Optics and Photonics, 2019. 17
- [81] Valentin I Spitkovsky, Hiyan Alshawhi, Angel Chang, and Dan Jurafsky. Unsupervised dependency parsing without gold part-of-speech tags. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 1281–1290, 2011. 15
- [82] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *CVPR*, pages 16519–16529, 2021. 3

- [83] Erik B Sudderth, Antonio Torralba, William T Freeman, and Alan S Willsky. Learning hierarchical models of scenes, objects, and parts. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 2, pages 1331–1338. IEEE, 2005. 1, 2
- [84] Shitao Tang, Jiahui Zhang, Siyu Zhu, and Ping Tan. Quadtree attention for vision transformers. *arXiv preprint arXiv:2201.02767*, 2022. 3
- [85] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357. PMLR, 2021. 3, 5, 6, 7, 8, 14, 15, 16, 17
- [86] Asher Trockman and J Zico Kolter. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022. 7, 17
- [87] Zhuowen Tu, Xiangrong Chen, Alan L Yuille, and Song-Chun Zhu. Image parsing: Unifying segmentation, detection, and recognition. *International Journal of computer vision*, 63(2):113–140, 2005. 1, 2
- [88] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *CVPR*, pages 12894–12904, 2021. 3
- [89] Bo Wan, Wenjuan Han, Zilong Zheng, and Tinne Tuytelaars. Unsupervised vision-language grammar induction with shared structure modeling. In *International Conference on Learning Representations*, 2021. 1
- [90] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 136–145, 2017. 5, 6, 7
- [91] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt2: Improved baselines with pyramid vision transformer. *arXiv preprint arXiv:2106.13797*, 2021. 7
- [92] 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. 3, 7, 14, 15, 16
- [93] Wenguan Wang, Zhijie Zhang, Siyuan Qi, Jianbing Shen, Yanwei Pang, and Ling Shao. Learning compositional neural information fusion for human parsing. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5703–5713, 2019. 2
- [94] Wenguan Wang, Hailong Zhu, Jifeng Dai, Yanwei Pang, Jianbing Shen, and Ling Shao. Hierarchical human parsing with typed part-relation reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8929–8939, 2020. 2
- [95] Yangtao Wang, Xi Shen, Shell Hu, Yuan Yuan, James Crowley, and Dominique Vaufreydaz. Self-supervised transformers for unsupervised object discovery using normalized cut. *arXiv preprint arXiv:2202.11539*, 2022. 6, 7
- [96] Ross Wightman. Pytorch image models. [\(cited on p.\)](https://github.com/rwightman/pytorch-image-models), 2019. 16
- [97] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021. 3, 7, 16
- [98] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *ICCV*, pages 10033–10041, 2021. 3
- [99] Sitong Wu, Tianyi Wu, Haoru Tan, and Guodong Guo. Pale transformer: A general vision transformer backbone with pale-shaped attention. *arXiv preprint arXiv:2112.14000*, 2021. 14
- [100] Tianfu Wu and Song-Chun Zhu. A numerical study of the bottom-up and top-down inference processes in and-or graphs. *International journal of computer vision*, 93(2):226–252, 2011. 1, 2
- [101] Edgar Xi, Selina Bing, and Yang Jin. Capsule network performance on complex data. *arXiv preprint arXiv:1712.03480*, 2017. 2
- [102] Canqun Xiang, Lu Zhang, Yi Tang, Wenbin Zou, and Chen Xu. Ms-capsnet: A novel multi-scale capsule network. *IEEE Signal Processing Letters*, 25(12):1850–1854, 2018. 2
- [103] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, pages 418–434, 2018. 15, 17
- [104] Yufei Xu, Qiming Zhang, Jing Zhang, and Dacheng Tao. Vitae: Vision transformer advanced by exploring intrinsic inductive bias. *NeurIPS*, 34, 2021. 3
- [105] Hongwei Xue, Yupan Huang, Bei Liu, Houwen Peng, Jianlong Fu, Houqiang Li, and Jiebo Luo. Probing inter-modality: Visual parsing with self-attention for vision-and-language pre-training. *Advances in Neural Information Processing Systems*, 34, 2021. 1
- [106] Qiong Yan, Li Xu, Jianping Shi, and Jiaya Jia. Hierarchical saliency detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1155–1162, 2013. 6, 7
- [107] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via graph-based manifold ranking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3166–3173, 2013. 6, 7
- [108] Jinyu Yang, Peilin Zhao, Yu Rong, Chaochao Yan, Chunyuan Li, Hehuan Ma, and Junzhou Huang. Hierarchical graph capsule network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10603–10611, 2021. 2
- [109] Qihang Yu, Yingda Xia, Yutong Bai, Yongyi Lu, Alan L Yuille, and Wei Shen. Glance-and-gaze vision transformer. *NeurIPS*, 34, 2021. 3
- [110] Tan Yu, Gangming Zhao, Ping Li, and Yizhou Yu. Boat: Bilateral local attention vision transformer. *arXiv preprint arXiv:2201.13027*, 2022. 3, 14

- [111] Li Yuan, Qibin Hou, Zihang Jiang, Jiashi Feng, and Shuicheng Yan. Volo: Vision outlooker for visual recognition. *arXiv preprint arXiv:2106.13112*, 2021. 3
- [112] Wang Zeng, Sheng Jin, Wentao Liu, Chen Qian, Ping Luo, Ouyang Wanli, and Xiaogang Wang. Not all tokens are equal: Human-centric visual analysis via token clustering transformer. *arXiv preprint arXiv:2204.08680*, 2022. 14
- [113] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. *arXiv: Computer Vision and Pattern Recognition*, 2021. 3
- [114] 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. 3, 16
- [115] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, pages 633–641, 2017. 7, 14, 15, 17
- [116] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiao Chen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021. 3
- [117] Jingkai Zhou, Pichao Wang, Fan Wang, Qiong Liu, Hao Li, and Rong Jin. Elsa: Enhanced local self-attention for vision transformer. *arXiv: Computer Vision and Pattern Recognition*, 2021. 3
- [118] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. Ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021. 5, 7
- [119] Song-Chun Zhu, David Mumford, et al. A stochastic grammar of images. *Foundations and Trends® in Computer Graphics and Vision*, 2(4):259–362, 2007. 1, 2
- [120] Wangjiang Zhu, Shuang Liang, Yichen Wei, and Jian Sun. Saliency optimization from robust background detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2814–2821, 2014. 6, 7

Appendix

Overview

In this appendix, we supplement the main paper by providing more thorough evaluations and empirical analyses to back up our claims. We also include more detailed descriptions of our experiments to help readers better understand our paper.

This appendix is organized as follows.

- In Section A, we give the notations used in this work.
- In Section B, we benchmark our models on two dense prediction downstream tasks.
- In Section C, we introduce detailed analysis to our model, including the relationship to pruning-based transformers, the comparison between reversed attention and forward attention, possible applications on video recognition, and some ablation studies.
- In Section D, we detail the training configurations and implementation details for each downstream task.

A. Notations

We provide the notations shown in Table 5 for this work.

B. Downstream Tasks

We benchmark our models on two dense prediction downstream tasks. All the model training follows common practices and protocols, as in [85, 92].

Semantic segmentation. In Table 6, we show the performance of our models on ADE20K [115] against several powerful counterparts. Considering DeiT [85] is the baseline that can be apple-to-apple comparable to us, we pre-train DeiT and our models on ImageNet-1K and produce the results of them under the same setting. We can see that: our DependencyViT consistently outperforms its counterparts including Swin [63]; and even DependencyViT-Lite surpasses the baseline PVT [92] by a large margin. Notably, the backbone model for DependencyViT-Lite only costs 1/3 computations (see the numbers in parentheses of the table) of our DependencyViT, showing its efficiency.

Object detection and instance segmentation. We benchmark our models on object detection with COCO 2017 [61] based on Mask R-CNN [35]. Table 7 show the detection and instance segmentation results. The results of DeiT and our models are implemented by us under the same setting. We observe substantial gains across all settings and metrics compared with several CNN and transformer baselines. Surprisingly, the backbone FLOPs consumption of DependencyViT-Lite-T is 3.5 GFLOPs, costing only 1.5% of the entire network.

C. Analysis

In this section, we introduce detailed analysis to our model.

C.1. Relation to Pruning-based Methods

Our work is related to dynamic-merged [99, 110] or pruning-based [9, 48, 70, 112] vision transformers. For example, DynamicViT [70] is a pruning-based transformer by optimizing a learnable weight for each token through Gumbel-Softmax.

However, the above methods mainly focus on the image classification tasks. They can not perform dense predictions because the information of their pruned patches is lost. On the contrary, pruning in a tree structure preserves the information lost by explicitly learned structures. As shown in the main paper, the pruned nodes in our DependencyViT-Lite can be retrieved from their parents for dense predictions, showing the importance of dependency induction.

C.2. Reversed attention vs. Forward one

Though forward attention well models the information interaction between patches, it mainly focuses on the task-specific region rather than the entire image, *e.g.*, the foreground region for the image classification task. This is because forward self-attention works through “gathering information”, thus the information in the background region that does not contribute to the recognition task is to a large extent suppressed and not gathered. The observation is evidenced by many previous works.

However, for our reversed self-attention, all the patches are get attended, *e.g.*, a subtree will be generated for the background area. The background information is kept because we do not prune any parent nodes. We then use the message controller to filter the useless information out for the final image recognition. Therefore, reversed attention has better generalization when extended to dense prediction tasks such as semantic segmentation, which is empirically validated by our experiments.

C.3. Pruning ratio

We also show DependencyViT-Lite with different pruning ratios by keeping the remaining token number as 32, 64, and 128. The results are shown in Table 8. We can see that when we keep 128 tokens, the performance drop is minor relative to the full DependencyViT. The performance gap could be larger when more tokens are pruned.

C.4. Dynamic Pruning on Video Recognition

We evaluate the models on the validation sets of Kinetics-400 (K400). Kinetics-400 consists of 240K training videos and 20K validation videos that span 400 human action categories. The results can be found in Table 9.

Table 5. Notations and their corresponding representations for DependencyViT.

Notation	Representation	Notation	Representation
\mathbf{X}	Input	\mathbf{A}_F	Forward Attention Map
\mathbf{P}	Head Selector	\mathbf{A}_R	ReverseAttention Map
\mathbf{M}	Message Controller	N	number of patches
\mathbf{Q}	Query	H	number of heads
\mathbf{K}	Key	C	token dims
\mathbf{V}	Value	C_h	token dims per head
\mathbf{W}	Projections		

Table 6. Comparison with SoTA methods for semantic segmentation on ADE20K [115] val set. Single-scale evaluation is used. FLOPs are measured by 512×2048 . Considering the segmentation head UperNet [103] is heavy, while the network backbone occupies only a small part of the computation, we mark the GFLOPs of the backbone of our works in parentheses.

Backbone	Method	#Params (M)	FLOPs (G)	mIoU (%)
ResNet18	SemanticFPN [59]	15.5	128.8	32.9
PVT-Tiny [92]	SemanticFPN [59]	17.0	132.8	35.7
DeiT-Tiny [85]	UperNet [103]	10.7	142.8	37.8
DependencyViT-Lite-T	UperNet [103]	11.1	130.2 (7.8)	36.1
DependencyViT-T	UperNet [103]	11.1	145.1 (22.7)	40.3
ResNet50	SemanticFPN [59]	28.5	729.6	36.7
PVT-Small [92]	SemanticFPN [59]	28.2	712.0	39.8
DeiT-Small [85]	UperNet [103]	41.3	566.8	43.0
Swin-Tiny [63]	UperNet [103]	60.0	945.0	44.5
DependencyViT-Lite-S	UperNet [103]	43.1	515.2 (29.6)	41.2
DependencyViT-S	UperNet [103]	43.1	574.4 (88.8)	45.7

Note that to use the pretrained model provided by TimesFormer [6], we only apply our dynamic pooling scheme on TimesFormer without the message controller. We perform dynamic pruning in the 2_{th} , 5_{th} , 8_{th} , 11_{th} layers, with 20% tokens pruned each time on both the temporal and spatial dimension. We can see under three different settings, the lite models still maintain a good performance while the FLOPs are reduced to 25%.

As shown in Figure 8, we show DependencyViT-Lite can learn the temporal dependency from videos. The sampled 8 frames are parsed into three subtrees (in gray boxes). And we use black lines to show the dependencies between two subtrees. We see that the root subtree contains keyframes and the root frame is the most informative frame.

C.5. Related Work in NLPs

Unsupervised dependency parsing is also a long-standing task in NLP. This task aims to induce dependency trees from raw corpora that do not have human-annotated tree structures. Traditional dependency grammar induction methods [2, 33, 81] are based on Dependency Model with Valence (DMV) [47]. DMV-based methods induce dependency from the statistical relation between tokens and their Part-of-Speech Tagging. Despite being very successful in

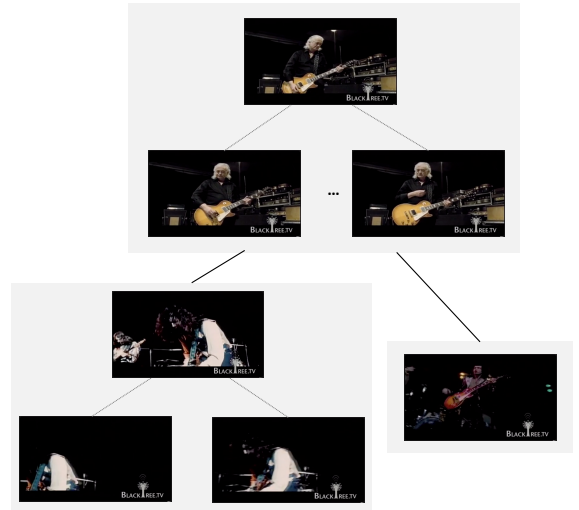


Figure 8. We show DependencyViT-Lite can learn the temporal dependency from videos. The sampled 8 frames are parsed into three subtrees (in gray boxes). And we use black lines to show the dependencies between two subtrees. We see that the root subtree contains keyframes and the root frame is the most informative frame. A few frames are enough for video recognition.

Table 7. COCO object detection and segmentation results with Mask R-CNN [36]. All models are trained with $1 \times$ schedule and multi-scale inputs. FLOPs are measured by 800×640 . The GFLOPs of the backbone of our DependencyViT and DependencyViT-Lite are marked in parentheses. The first three metrics are for object detection, while the last three for instance segmentation.

Backbone	#Params (M)	FLOPs (G)	Mask R-CNN 1x					
			AP^b	AP_{50}^b	AP_{75}^b	AP^m	AP_{50}^m	AP_{75}^m
ResNet18 [36]	31.2	190.0	34.0	54.0	36.7	31.2	51.0	32.7
PVT-Tiny [92]	32.9	195.0	36.7	59.2	39.3	35.1	56.7	37.3
DeiT-Tiny [85]	27.3	244.6	30.6	46.8	32.8	27.4	44.7	28.9
DependencyViT-Lite-T	27.8	238.1 (3.5)	35.2	58.8	38.6	34.1	56.2	36.1
DependencyViT-T	27.8	245.6 (11.0)	37.8	62.1	41.4	36.0	59.3	38.6
ResNet50 [36]	44.2	260.0	38.0	58.6	41.4	34.4	55.1	36.7
PVT-Small [92]	44.1	245.0	40.4	62.9	43.8	37.8	60.1	40.3
DeiT-Small [85]	44.9	276.2	36.9	55.1	39.7	32.7	52.3	34.5
DependencyViT-Lite-S	46.85	249.9 (13.2)	38.1	62.5	41.8	36.2	59.4	38.4
DependencyViT-S	46.85	280.0 (43.3)	42.4	66.5	46.4	38.5	62.7	41.9

Table 8. Comparison of image classification on ImageNet-1K when different number of tokens are pruned.

Model	kept tokens	#Params (M)	FLOPs (G)	Top-1 (%)
DependencyViT-Lite-32	32	6.2	0.6	72.4
DependencyViT-Lite-64	64	6.2	0.8	73.7
DependencyViT-Lite-128	128	6.2	1.0	74.9
DependencyViT	196	6.2	1.3	75.4

the natural language domain, similar methods can not be directly applied to visual dependency induction due to two reasons: 1) DMV-based methods require discrete tokens as input, whereas visual inputs are continuous values; 2) they also heavily rely on the sequential order of input tokens, whereas visual inputs have at least two dimensions. In recent years, researchers proposed several transformer-based unsupervised dependency parsing methods, including Structformer [76] and UDGNet [75]. However, unsupervised vision dependency parsing using transformers is still very challenging because images are composed of pixels that contain no significant semantic or syntactic meaning. In contrast, natural language is composed of words expressing abstract concepts and belonging to specific syntactic roles. To overcome the challenge, DependencyViT adapts a progressive parsing schema that gradually composes low-level representations to high-level representations and makes progressive parsing decisions alongside the level of abstractness.

D. Training Details

D.1. Details of Model Configuration

In this work, we simply follow the design strategy suggested by the standard ViT (DeiT) [26, 85]. The non-overlapping patch embedding layer is implemented by stride convolution. The convolutional kernel and stride value are 16 and 16, respectively. We stack our dependency

blocks with the resolution and feature dimension kept the same. We set the number of attention heads $H = 12$ and the number of dependency blocks $L = 12$ for all models. We set token dimensions $C = 192$ for the tiny model and $C = 384$ for the small model. In the head selector, we introduce a temperature hyper-parameter for the softmax function, which is set to 0.1 for all models.

For DependencyViT-Lite, similar to current hierarchical models that divide the entire architecture into four stages, we perform dynamic pruning in the 2_{th} , 5_{th} , 8_{th} , 11_{th} layers with a token kept number as 160, 128, 96, and 64, respectively. For dense prediction tasks, the tree architecture is still maintained by recording relationships (probability distributions) between the pruned nodes and their parents to form a complete tree. After the end of the network, we retrieve those pruned nodes by a soft aggregation from their parents, preserving the model capability and generating a dense representation. As a result, the proposed architecture can conveniently replace the backbone networks in existing methods for various vision tasks.

D.2. Image Classification on ImageNet

The ILSVRC 2012 classification dataset (ImageNet-1K) [20] consists of 1,000 classes, with a number of 1.2 million training images and 50,000 validation images.

We compare different methods on ImageNet-1K [20]. We implement our DependencyViT on the timm framework [96]. Following [24, 60, 63, 97, 114], we use the same

Table 9. Video-level accuracy on the Kinetics-400 validation set.

Method	Top-1 (%)	Top-5 (%)	FLOPs (G)	Frames	Resolution
TimeSformer	76.9	92.7	0.20	8	224
TimeSformer-Lite	70.6	89.3	0.08	8	224
TimeSformer-HR	78.1	93.3	1.70	16	448
TimeSformer-HR-Lite	73.1	90.4	0.67	16	448
TimeSformer-L	79.8	94.1	2.38	96	224
TimeSformer-L-Lite	74.1	91.3	0.61	96	224

set of data augmentation and regularization strategies used in [85] after excluding repeated augmentation [5, 40] and exponential moving average (EMA) [69]. We train all the models for 300 epochs with a batch size 2048 and use AdamW [65] as the optimizer. The weight decay is set to 0.05 and the maximal gradient norm is clipped to 1.0. We use a simple triangular learning rate schedule [80] as in [86]. The stochastic depth drop rates are set to 0.1 and 0.2 for our tiny and small models, respectively. During training, we crop images randomly to 224×224 , while a center crop is used during evaluation on the validation set. For fair comparisons, neither token labeling [44] nor distillation [85] is used in all experiments.

D.3. Object Detection on COCO

The COCO dataset [61] contains over 200,000 images labeled with object detection bounding boxes and instance segmentation masks. We evaluate our approach on the val2017, containing 5000 images.

We benchmark our models on object detection with COCO 2017 [61]. The pre-trained models are used as visual backbones and then plugged into two representative pipelines, RetinaNet [60] and Mask R-CNN [35]. All models are trained on the 118k training images and results reported on the 5K validation set. We follow the standard to use two training schedules, $1 \times$ schedule with 12 epochs and $3 \times$ schedule with 36 epochs. The same multi-scale training strategy as in [63] by randomly resizing the shorter side of the image to the range of [480, 800] is used. During training, we use AdamW [65] for optimization with initial learning rate 10^{-4} and weight decay 0.05. We use 0.1 and 0.2 stochastic depth drop rates to regularize the training for our tiny and small models, respectively.

D.4. Semantic Segmentation on ADE20k

Besides the instance segmentation results above, we further evaluate our model on semantic segmentation, a task that usually requires high-resolution input and long-range interactions. ADE20K [115] is a scene-centric containing 20 thousands images annotated with 150 object categories.

We benchmark our method on ADE20K [115]. Specifically, we use UperNet [103] as the segmentation method

and our DependencyViT as the backbone. For all models, we use a standard recipe by setting the input size to 512×512 and train the model for 160k iterations with batch size 16.