

EdgeViTs: Competing Light-weight CNNs on Mobile Devices with Vision Transformers

Junting Pan^{1*}, Adrian Bulat², Fuwen Tan², Xiatian Zhu², Lukasz Dudziak²,
Hongsheng Li¹, Georgios Tzimiropoulos^{2,3} and Brais Martinez²

¹ The Chinese University of Hong Kong

² Samsung AI Cambridge

³ Queen Mary University of London

Abstract. Self-attention based models such as vision transformers (ViTs) have emerged as a very competitive architecture alternative to convolutional neural networks (CNNs) in computer vision. Despite increasingly stronger variants with ever higher recognition accuracies, due to the quadratic complexity of self-attention, existing ViTs are typically demanding in computation and model size. Although several successful design choices (e.g., the convolutions and hierarchical multi-stage structure) of prior CNNs have been reintroduced into recent ViTs, they are still not sufficient to meet the limited resource requirements of mobile devices. This motivates a very recent attempt to develop light ViTs based on the state-of-the-art MobileNet-v2, but still leaves a performance gap behind. In this work, pushing further along this under-studied direction we introduce **EdgeViTs**, a new family of light-weight ViTs that, for the first time, enable attention based vision models to compete with the best light-weight CNNs in the tradeoff between accuracy and on-device efficiency. This is realized by introducing a highly cost-effective *local-global-local* (LGL) information exchange bottleneck based on optimal integration of self-attention and convolutions. For device-dedicated evaluation, rather than relying on inaccurate proxies like the number of FLOPs or parameters, we adopt a practical approach of focusing directly on on-device latency and, for the first time, energy efficiency. Extensive experiments on image classification, object detection and semantic segmentation validate high efficiency of our EdgeViTs when compared to the state-of-the-art efficient CNNs and ViTs in terms of accuracy-efficiency tradeoff on mobile hardware. Specifically, we show that our models are Pareto-optimal when both accuracy-latency and accuracy-energy tradeoffs are considered, achieving strict dominance over other ViTs in almost all cases and competing with the most efficient CNNs. Code is available at <https://github.com/saic-fi/edgevit>.

1 Introduction

Vision transformers (ViTs) have rapidly superseded convolutional neural networks (CNNs) on a variety of visual recognition tasks [10,52], particularly when

* Work done during an internship at Samsung AI Cambridge.

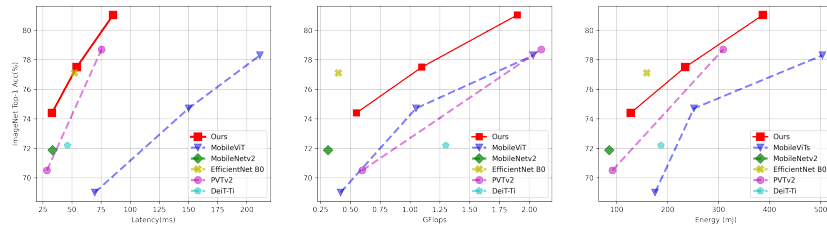


Fig. 1. Our EdgeViTs yield comparable or superior tradeoff between accuracy and efficiency (e.g., run speed, latency) on ImageNet-1K, in comparison to state-of-the-art efficient CNNs [43,48], representative generic ViTs [56,52] and latest MobileViTs [37]. Note that, EdgeViTs outperform MobileViTs, the only specially designed model for mobile device, across all the metrics by a large margin. Whilst our lightest EdgeViT consumes more GFLOPs than EfficientNet B0 [48], it runs faster (33.3 vs. 52.1 ms per image), with a reduced gap in *on-device energy consumption*. *Testing device:* Samsung Galaxy S21 (latency), Snapdragon 888 Hardware Development Kit (energy).

the priors and successful designs of previous CNNs are reintroduced for leveraging the induction bias of visual data such as local grid structures [55,31,7,16]. Due to the quadratic complexity of ViTs and the high-dimension property of visual data, it is indispensable that the computational cost needs to be taken into account in design [55]. Three representative designs to make computationally viable ViTs are (1) the use of a hierarchical architecture with the spatial resolution (i.e., the token sequence length) progressively down-sampled across the stages [12,55,7], (2) the locally-grouped self-attention mechanisms for controlling the length of input token sequences and parameter sharing [31,16], (3) the pooling attention schemes to subsample the **key** and **value** by a factor [54,35,60,55]. The general trend has been on designing more complicated and stronger ViTs to challenge the dominance of top-performing CNNs [20,40,48] in computer vision by achieving ever higher accuracies [53,66]. These advances however are still insufficient to satisfy the design requirements and constraints for mobile and edge platforms (e.g., smart phones, robotics, self-driving cars, AR/VR devices), where the vision tasks need to carry out in a timely manner under certain computational budgets. Prior efficient CNNs (e.g., MobileNets [23,43,22], ShuffleNets [36,64], EfficientNets [48,51,49], and etc.) remain the state-of-the-art network architectures for such platforms in the tradeoff between running latency and recognition accuracy (Fig. 1).

In this work, we focus on the development of largely under-studied *efficient ViTs* with the aim to surpass the CNN counterparts on mobile devices. We consider a collection of very practical design requirements for running a ViT model on a target real-world platform as follows: **(1)** *Inference efficiency* needs to be high (e.g., *low latency and energy consumption*) so that the running cost becomes generically affordable and more on-device applications can be supported. This is a direct metric that we really care about in practice. In contrast, the often-used efficiency metric, FLOPs (i.e., the number of multiply-adds), cannot

directly translate into the latency and energy consumption on a specific device, with several conditional factors including memory access cost, degree of parallelism, and the platform’s characteristics [36]. This is, *not* all operations of a model can be carried out at the same speed and energy cost on a device. Hence, FLOPs is merely an approximate and indirect metric of efficiency. **(2)** *Model size* (i.e., parameter number) is affordable for modern average devices. Given the availability of ever cheaper and larger storage spaces, this constraint has been relaxed significantly. For example, an average smart phone often comes with 32GB or more storage. As a consequence, using it as a threshold metric is no longer valid in most cases. **(3)** *Implementational friendliness* is also critical in real-world applications. For a wider range of deployment, it is necessary that a model can be implemented efficiently using the standard computing operations supported and optimized in the generic deep learning frameworks (e.g., ONNX, TensorRT, and TorchScript), without costly per-framework specialization. Otherwise, the on-device speed of a model might be unsatisfactory even with low FLOPs. For instance, the cyclic shift and its reverse operations introduced in Swin Transformers [31] are rarely supported by the mainstream frameworks, i.e., deployment unfriendly. In the literature, very recent MobileViTs [37] are the only series of ViTs designed for mobile devices. In architecture design, they are a straightforward combination of MobileNetv2 [43] and ViTs [10]. As a very initial attempt in this direction, MobileViTs still lag behind CNN counterparts. Further, its evaluation protocol takes the *model size* (i.e. the parameter number) as the competitor selection criteria (i.e., comparing the accuracy of models only with similar parameter numbers), which however is no longer a hard constraint with modern hardware as discussed above and is hence out of date.

We present a family of light-weight attention based vision models, dubbed as **EdgeViTs**, for the first time, enabling ViTs to compete with the best light-weight CNNs (e.g., MobileNetv2 [43] and EfficientNets [48]) in terms of accuracy-efficiency tradeoff on mobile devices. This sets a milestone in the landscape of light-weight ViTs *vs.* CNNs in the low resource regime. Our EdgeViTs are based on a novel factorization of the standard self-attention for more cost-effective information exchange within every individual layer. This is made possible by introducing a highly light-weight and easy-to-implement *local-global-local* (LGL) information exchange bottleneck characterized with three operations: **(i)** Local information aggregation from neighbor tokens (each corresponding to a specific patch) using efficient depth-wise convolutions; **(ii)** Forming a sparse set of evenly distributed *delegate tokens* for long-range information exchange by self-attention; **(iii)** Diffusing updated information from *delegate tokens* to the *non-delegate tokens* in local neighborhoods via transposed convolutions. As we show in experiments, this design presents a favorable hybrid of self-attention, convolutions, and transposed convolutions, achieving the best accuracy-efficiency tradeoff. It is efficient in that the self-attention is applied to a sparse set of delegate tokens. To support a variety of computational budgets, with our primitive module we establish a family of EdgeViT variants with three computational complexities: small (S), extra-small (XS), extra-extra-small (XXS).

We make the following **contributions**: **(1)** We investigate the design of light-weight ViTs from the practical on-device deployment and execution perspective. **(2)** For best scalability and deployment, we present a novel family of efficient ViTs, termed as EdgeViTs, designed based on an optimal decomposition of self-attention using standard primitive operations. **(3)** Regarding on-device performance, towards relevance for real-world deployment, we directly consider latency and energy consumption of different models rather than relying on high-level proxies like number of FLOPs or parameters. Our results experimentally verify efficiency of our models in a practical setting and refute some of the claims made in the existing literature. More specifically, extensive experiments on three visual tasks show that our EdgeViTs can match or surpass state-of-the-art light-weight CNNs, whilst consistently outperform the recent MobileViTs in accuracy-efficiency tradeoff, including largely ignored on-device energy evaluation. Importantly, EdgeViTs are consistently Pareto-optimal in terms of both latency and energy efficiency, achieving strict dominance over other ViTs in almost all cases and competing with the most efficient CNNs. On ImageNet classification our EdgeViT-XXS outperforms MobileNetv2 by 2.2% subject to the similar energy-aware efficiency.

2 Related Work

Efficient CNNs. Since the advent of modern CNN architectures [20,46], there has been a steady stream of works focusing on efficient architecture design for on-device deployment. The first widely adopted families bring depthwise separable convolutions in a ResNet-like structure, e.g., MobileNets [23,43], ShuffleNets [64,36]. These works define a space of well-performing efficient architectures, resulting in widespread usage. Successive works further exploit this design space by automating the architectural design choices [48,22,47,50]. As a parallel line of research, net pruning creates efficient architectures by removing spurious parts of a larger network with close-to-zero weights [17,57], or via first training a super-network that is further slimmed to meet a pre-specified computational budget [32,3]. Dynamic computing has also been explored, consisting of the mechanisms that condition the network parameters on the input data [62,6]. Finally, using low bit-width is a very critical technique that can offer different tradeoffs between the accuracy and efficiency [17,24,5].

Vision transformers. ViTs [10] quickly popularize transformer-based architectures for computer vision. A series of works followed instantly, offering large improvements to the original ViTs in terms of data efficiency [52,30] and architecture design [31,63,12,4]. Among these works, one of the main modifications is to introduce hierarchical designs in multiple stages from convolutional architectures [31,7,55,56]. Several works also focus on improving the positional encoding by using a relative positional embedding [45,44], making it learnable [14], or even replacing it by an attention bias element [15]. All these approaches mostly aim to improve the model performance.

Recently, more efforts have been made towards finding efficient alternatives to the multi-head self-attention (MHSA) module, which is typically the computational bottleneck in the ViT architectures. A particularly effective solution is to reduce the internal spatial dimensions within the MHSA. The MHSA involves projecting the input tensor into key, query and value tensors. Several recent works, e.g. [7,55,56], find that the key and value tensors could be downsampled with a limited loss in accuracy, leading to a better efficiency-accuracy tradeoff. Our work extends this idea by also downsampling the query tensors, which further improves the efficiency, as shown in Fig. 2.

There are also alternative approaches reducing the number of tokens dynamically [41,61,38,13]. That is, in the forward pass, tokens deemed to not contain the important information for the target task are pruned or pooled together, reducing the overall complexity thereafter. Finally, encouraged by their potential complementarity, many works have attempted to combine convolutional designs with self-attentions. This ranges from using convolutions at the stem [59], integrating convolutional operations into the MHSA block [26,58], or incorporating the MHSA block into ResNet-like architectures [45]. It is interesting to note that even the original ViTs explored similar tradeoffs. [10].

Vision transformers for mobile devices. Whilst the efficiency issue has been taken into account in designing the ViT variants discussed above, they are still not dedicated and satisfactory architectures for on-device applications. There is only one exception, MobileViTs [37], which are introduced very recently. However, compared to the current best light-weight CNNs such as MobileNets [43,22] and EfficientNets [48], these ViTs are still clearly inferior in terms of the on-device accuracy-efficiency tradeoff. In this work, we present the first family of efficient ViTs that can deliver comparable or even superior tradeoffs in comparison to the best CNNs and ViTs. We also extensively carry out the critical yet largely lacking on-device evaluations with energy consumption analysis.

3 EdgeViTs

3.1 Overview

For designing light-weight ViTs suitable for mobile/edge devices, we adopt a hierarchical pyramid network structure (Fig. 2(a)) used in recent ViT variants [55,8,7,56,12]. A pyramid transformer model typically reduces the spatial resolution but expands the channel dimension across different stages. Each stage consists of multiple transformer-based blocks processing tensors of the same shape, mimicking the ResNet-like networks. The transformer-based blocks heavily rely on the self-attention operations at a quadratic complexity w.r.t the spatial resolution of the visual features. By progressively aggregating the spatial tokens, pyramid vision transformers are potentially more efficient than isotropic models [10]. In this work, we dive deeper into the transformer-based block and introduce a cost-effective bottleneck, *Local-Global-Local* (LGL) (Fig. 2(b)). LGL further reduces the overhead of self-attention with a sparse attention module (Fig. 2(c)), achieving better accuracy-latency balancing.

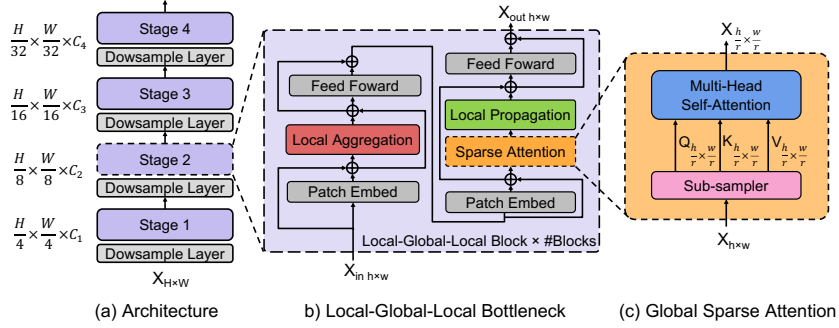


Fig. 2. (a) Schematic overview of our four stages EdgeViT architecture, with each stage consisting of a stack of (b) *Local-Global-Local* (LGL) blocks constructed with local aggregation module, sparse-self-attention and local propagation module, patch embedding (PE) and Feed Forward Network (FFN). In this example, h and w refer to input height and width of stage-2: $h = \frac{H}{8}$ and $w = \frac{W}{8}$. C_i refers to the number of channels for stage- i and r denotes the sub-sampling rate.

3.2 *Local-Global-Local* bottleneck

Self-attention has been shown to be very effective for learning the global context or long-range spatial dependency of an image, which is critical for visual recognition. On the other hand, as images have high spatial redundancy (e.g., nearby patches are semantically similar) [18], applying attention to all the tokens, even in a down-sampled feature map, is inefficient. There is hence an opportunity to reduce the scope of tokens whilst still preserving the underlying information flows that model the global and local contexts. In contrast to previous transformer blocks that perform self-attention at each spatial location, our LGL bottleneck only computes self-attention for a subset of the tokens but enables full spatial interactions, as in the standard multi-head self-attention (MHSA) [10].

To achieve this, we decompose the self-attention into consecutive modules that process the spatial tokens within different ranges (Fig. 2(b)). We introduce three efficient operations: i) *Local aggregation* that integrates signals only from locally proximate tokens; ii) *Global sparse attention* that model long-range relations among a set of delegate tokens where each of them is treated as a representative for a local window; iii) *Local propagation* that diffuses the global contextual information learned by the delegates to the non-delegate tokens with the same window. Combining these, our LGL bottleneck enables information exchanges between any pair of tokens in the same feature map at a low-compute cost. Each of these components is described in detail below:

- *Local aggregation*: for each token, we leverage depth-wise and point-wise convolutions to aggregate information in local windows with a size of $k \times k$ (Fig. 3(a)).
- *Global sparse attention*: we sample a sparse set of *delegate tokens* distributed evenly across the space, one *token* for each $r \times r$ window. Here, r denotes the

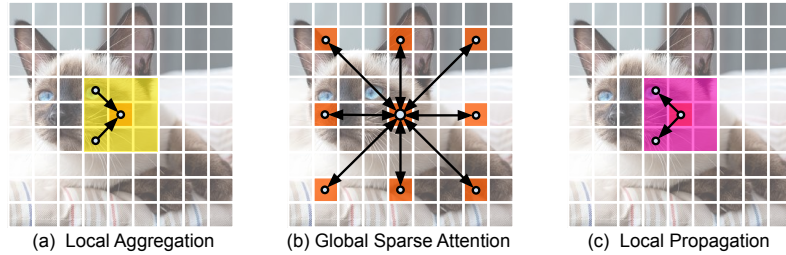


Fig. 3. Illustration of three key operations involved in the proposed *Local-Global-Local* (LGL) transformer block. In this example, we showcase how *the target token* (the orange square) at the center conducts information exchange with all the others in three sequential steps: (a) Local information from neighbor tokens within the yellow area is first aggregated to the target token. (b) Global sparse attention is then computed among the target token and other selected delegates in orange color. (c) Global context information encoded in the target token is last propagated to its neighbor *non-delegate* tokens within the pink area.

sub-sample rate. We then apply self-attention on these selected tokens only (Fig. 3(b)). This is distinct from all the existing ViTs [7,55,56] where all the spatial tokens are involved as queries in the self-attention computation.

- *Local propagation*: We propagate the global contextual information encoded in the delegate tokens to their neighbor tokens by transposed convolutions (Fig. 3(c)).

Formally, our LGL bottleneck can be formulated as:

$$\begin{aligned}
 X &= \text{LocalAgg}(\text{Norm}(X_{in})) + X_{in}, \\
 Y &= \text{FFN}(\text{Norm}(X)) + X, \\
 Z &= \text{LocalProp}(\text{GlobalSparseAttn}(\text{Norm}(Y))) + Y, \\
 X_{out} &= \text{FFN}(\text{Norm}(Z)) + Z.
 \end{aligned} \tag{1}$$

Here $X_{in} \in \mathcal{R}^{H \times W \times C}$ indicates the input tensors. **Norm** is the layer normalization operation [2]. **LocalAgg** represents the local aggregation operator, **FFN** is a two-layer perceptron, similar to the position-wise feed-forward network introduced in [10]. **GlobalSparseAttn** is the global sparse self-attention. **LocalProp** is the local propagation operator. For simplicity, positional encoding is omitted. Note that, all these operators can be implemented by commonly used and highly optimized operations in the standard deep learning platforms. Hence, our LGL bottleneck is implementation friendly.

Comparisons to existing designs. Our LGL bottleneck shares a similar goal with the recent PVTs [55,56] and Twins-SVTs [7] models that attempt to reduce the self-attention overhead. However, they differ in the core design. PVTs [55,56] perform self-attention where the number of **keys** and **values** are reduced by strided-convolutions, whilst the number of **queries** remains the same.

Model	#Channels	#Blocks	#Heads	FLOPs	#Param
EdgeViT-XXS	[36, 72, 144, 288]	[1,1,3,2]	[1,2,4,8]	0.56G	4.1M
EdgeViT-XS	[48, 96, 240, 384]	[1,1,2,2]	[1,2,4,8]	1.1G	6.7M
EdgeViT-S	[48, 96, 240, 384]	[1,2,3,2]	[1,2,4,8]	1.9G	11.1M

Table 1. Configuration of three EdgeViT variants. “#Channels”: number of channels per stage. “#Blocks”: number of LGL blocks per stage. “#Heads”: number of attention heads in MHSA. “#Param”: number of parameters.

In other words, PVTs still perform self-attention at each grid location. In this work, we question the necessity of positional-wise self-attention and explore to what extent the information exchange enabled by our LGL bottleneck could approximate the standard MHSA (see Section 4 for more details). Twins-SVTs [7] combine local-window self-attention [31] with global pooled attention from PVTs [55]. This is different from our hybrid design using both self-attention and convolution operations distributed in a series of *local-global-local* operations. As demonstrated in the experiments (Table 2 and 3), our design achieve a better tradeoff between the model performance and the computation overhead (e.g. latency, energy consumption, etc).

3.3 Architectures

We build a family of EdgeViTs with the proposed LGL bottleneck at different computational complexities (i.e. 0.5G, 1G, and 2G). The configurations are summarized in Table 1. Following the hierarchical ViTs [7,55,31,26], EdgeViTs consist of four stages with the spatial resolution (*i.e.*, the token sequence length) gradually reduced throughout, and their self-attention module replaced with our LGL bottleneck. For the stage-wise down-sampling, we use a conv-layer with a kernel size of 2×2 and stride 2, except for the first stage where we down-sample the input feature by $\times 4$, and use a 4×4 kernel and a stride of 4. We adopt the conditional positional encoding [8] that has been shown to be superior to the absolute positional encoding. This can be implemented using 2D depth-wise convolutions with a residual connection. In our model, we use 3×3 depth-wise convolutions with zero paddings. It is placed before the local aggregation and global sparse self-attention. The FFN consists of two linear layers with GeLU non-linearity [21] placed in-between. Our local aggregation operator is implemented as a stack of pointwise and depthwise convolutions. The global sparse attention is composed of a spatial uniform sampler with sample rates of (4, 2, 2, 1) for the four stages, and a standard MHSA. The local propagation is implemented with a depthwise separable transposed convolution with the kernel size and stride equal to the sample rate used in the global sparse attention. The exact architecture for the LGL bottleneck is described in the **supplementary material**.

4 Experiments

We benchmark EdgeViTs on visual recognition tasks. We pre-train EdgeViTs on the Imagenet1K recognition task [42], comparing the performances and computation overheads against alternative approaches. We also evaluate the generalization capacity of EdgeViTs on downstream dense prediction tasks: object detection and instance segmentation on the COCO benchmark [29], and semantic segmentation on the ADE20K Scene Parsing benchmark [65]. For on-device execution, we report execution time (latency) and energy consumption of all relevant models on ImageNet. We do not report on-device measurements on downstream tasks as they reuse ImageNet models.

4.1 Image Classification on ImageNet-1K

Training Settings. ImageNet-1K [42] provides 1.28 million training images and 50,000 validation images from 1000 categories. We follow the training recipe introduced in DeiT [52]. We optimize the models using AdamW [33] with a batch size of 1024, weight decay of 5×10^{-2} , and momentum of 0.9. The models are trained from *scratch* for 300 epochs with a linear warm-up during the first 5 epochs. Our base learning rate is set as 1×10^{-3} , and decay after the warm-up using a cosine schedule [34]. We apply the same data augmentations as in [52,7,56,26] which include random cropping, random horizontal flipping, mixup, random erasing and label-smoothing. During training, the images are randomly cropped to 224×224 . During testing, we use a single center crop of 224×224 . We report the top-1 accuracy on the validation set.

Benchmarking Settings. For latency measurements, we use a Samsung Galaxy S21 mobile phone equipped with a Snapdragon 888 chipset. All relevant models are benchmarked by running a forward pass 50 times using TorchScript lite interpreter via the Android benchmarking app provided by PyTorch [39]. We use CPU implementation, full precision and batch 1 to execute all operations. This choice comes from the fact that this is the only combination that was able to robustly execute all of the models from our paper. In general, more efficient implementations exist, such as those utilizing specialized hardware like Neural Processing Units (NPU). However, these put more restrictions on what can be executed and many models failed to run in our experiments when trying to use different hardware targets.

For energy measurements, we use a Monsoon High Voltage Power Monitor connected to a Snapdragon 888 Hardware Development Kit (HDK8350) to obtain accurate power readings over the course of running a forward pass of each test model 50 times. The same TorchScript runtime is used as in latency measurements. From the power signal reported by the monitor, we derive the average per-inference power and energy consumption by first subtracting background power consumption (i.e., power readings when not running any model) and then identifying 50 continuous regions of significantly higher power draw. Each region like that is considered a single inference and we calculate its total energy

Model	#Params	FLOPs	CPU(ms)	Acc Top-1 (%)
MobileNet-v2[43]	3.4M	0.3G	33.3 \pm 5.3	72.0
MobileNet-v3 0.75 [22]	4.0M	0.16G	23.0\pm3.7	73.3
EfficientNet-B0 [48]	5.3M	0.4G	52.1\pm7.4	77.1
MobileViT-XXS [37]	1.3M	0.4G	69.5 \pm 5.1	69.0
PVT-v2-B0[56]	3.4M	0.6G	26.0 \pm 6.9	70.5
Uniformer-Tiny*[26]	3.9M	0.6G	40.5 \pm 3.1	74.1
Twins-SVT-Tiny*[7]	4.1M	0.6G	36.9 \pm 2.3	71.2
EdgeViT-XXS	4.1M	0.6G	32.8\pm2.7	74.4
T2T-ViT-7 [63]	4.3M	1.1G	48.8 \pm 6.5	71.7
MobileViT-XS [37]	2.4M	1.1G	150.1 \pm 6.1	74.7
DeiT-Tiny [52]	5.7M	1.3G	46.2 \pm 13.6	72.2
TNT-Tiny [16]	6.1M	1.4G	86.4 \pm 6.0	73.9
EdgeViT-XS	6.7M	1.1G	54.1\pm2.2	77.5
T2T-ViT-12 [63]	6.9M	1.9G	69.9 \pm 5.6	76.5
PVT-v2-B1[56]	14M	2.1G	75.4 \pm 2.3	78.7
MobileViT-S [37]	5.6M	2.0G	221.3 \pm 9.3	78.3
LeViT-384\dagger [15]	39.1M	2.4G	71.3\pm2.3	79.5
EdgeViT-S	11.1M	1.9G	85.3\pm3.9	81.0

Table 2. Results on ImageNet-1K validation set. All models are tested on input scale of 224×224 , except for MobileViTs [37] that are tested with 256×256 according to their original implementation. * indicates down-scaled architectures beyond original definitions by authors to fit the mobile compute budget. LeViT-384 \dagger denotes the LeViT model retrained under the same setting as our EdgeViT.

as the integral over the individual power samples. Analogously we also calculate average power consumption by averaging over the same set of samples. After energy and power are calculated for each inference, the final statistics of a model are obtained by again averaging over the 50 identified runs. Our methodology follows what can be found in the literature [1].

Results. We compare EdgeViTs to a variety of baseline models, including the classic efficient CNNs, e.g. MobileNetV2 [43], MobileNetV3 [22], EfficientNet [48], and the state-of-the-art ViTs, e.g. MobileViT [37], PVT-v2 [56], DeiT [52], LeViT [15]. As the original LeViT [15] was optimized in a large-scale setting (i.e. 1000 epochs) with knowledge distillation, we perform a comparison by re-training LeViT under the same setting (300 epochs) as EdgeViT without knowledge distillation. We denote the retrained LeViTs as LeViT-384 \dagger . We select the baselines with a complexity of less than 2 GFLOPs as i) in real-world applications, the computational cost remains the top concern; ii) whilst FLOPs is an indirect metric for the latency, it is the most used cost metric in prior works. This selection criterion is different from [37] that instead uses the model size (i.e. the parameter number) which however has become a less restricted facet in mobile devices.

From Table 2, we can learn: i) EdgeViTs significantly outperform other lightweight *ViTs* at a similar level of GFLOPs complexity. Compared to the PVT-

Model	Top-1 (%)	CPU (ms)	Energy (mJ)	Power (W)	Efficiency (%/msW)
MobileNet-v2[43]	72.0	33.3	85.7 \pm 7.4	3.31 \pm 0.26	0.841
MobileNet-v3 0.75[22]	73.3	23.0	63.0\pm9.6	3.46\pm0.4	1.164
EfficientNet-B0[48]	77.1	52.1	159.0\pm26.2	3.62\pm0.45	0.485
PVT-v2-B0[56]	70.5	26.0	91.7 \pm 19.7	3.94 \pm 0.68	0.769
PVT-v2-B1[56]	78.7	75.4	309.0 \pm 65.8	4.63 \pm 0.71	0.255
Twins-SVT-Tiny*[7]	71.2	36.9	114.5 \pm 17.3	3.71 \pm 0.24	0.622
DeiT-Tiny [52]	72.2	46.2	187.2 \pm 7.6	4.77 \pm 0.21	0.386
Uniformer-Tiny*[26]	74.1	40.5	134.7 \pm 27.3	4.1 \pm 0.71	0.55
T2T-ViT-12 [63]	76.5	69.9	266.2 \pm 42.6	4.37 \pm 0.36	0.287
TNT-Tiny [16]	73.9	86.4	308.7 \pm 70.5	3.94 \pm 0.63	0.239
LeViT-384† [15]	79.5	71.3\pm2.2	455.2\pm125.8	6.18\pm0.74	0.173
MobileViT-XXS [37]	69.0	69.5	175.3 \pm 28.7	2.77 \pm 0.24	0.394
MobileViT-XS [37]	74.7	150.1	251.5 \pm 81.1	2.63 \pm 0.61	0.297
MobileViT-S [37]	78.3	221.3	503.6 \pm 117.0	2.76 \pm 0.21	0.155
EdgeViT-XXS	74.4	32.8	127.4\pm27.3	4.27\pm0.67	0.584
EdgeViT-XS	77.5	54.1	234.6\pm44.0	4.77\pm0.84	0.33
EdgeViT-S	81.0	85.3	386.7\pm43.5	4.8\pm0.26	0.209

Table 3. On-device energy evaluation on ImageNet-1K. All relevant metrics are reported as mean values per forward pass across 50 executions. For facilitating comparison, we define an energy-aware *efficiency* metric as the average gain in top-1 accuracy from each 1W run for 1ms (equivalent to consuming 1mJ of energy). (*Pareto-optimal models* are highlighted in bold in the last column).

v2 family [56], our EdgeViT-XXS/EdgeViT-S achieve 3.9%/2.3% improvements over PVT-v2-B0/PVT-v2-B1. Compared to MobileViTs, EdgeViTs achieve 5.4%, 2.8% and 2.7% gains in the three complexity settings. ii) **ViTs vs. CNNs:** Our EdgeViTs lift the performance of efficient ViTs to approach the level of well-established efficient CNNs. For example, the EdgeViT-XXS performs superior to MobileNet-v2 and MobileNet-v3-0.75 at a similar level of model size, but requires more GFLOPs. However, we observe that the efficient CNNs still surpass efficient ViTs in the accuracy-FLOPs tradeoff by a small margin.

On the other hand, as discussed early, numbers of FLOPs or parameters are merely indicative but do not fully reflect the on-device efficiency [36,11]. We further consider on-device latency and energy consumption directly. Other than the representative ViTs and CNNs, we also compare two recent ViT variants [26,7] with the number of channels and layers re-scaled to fit the complexity need. As presented in Table 2, EdgeViTs demonstrate strong performance with latencies comparable to MobileNets: EdgeViT-XXS achieves a gain of 2.4% over MobileNet-V2 while running slightly faster. EdgeViT-XXS also surpasses MobileNet-V3 by 1.1% but at the cost of being 9.8ms slower. EdgeViT-XS performs on par with the auto-searched EfficientNet-B0 model. We believe our models could also benefit from the automatic architecture search techniques as use

	CPU	Top1	CPU	Top1	CPU	Top1
w/o LA	33.9ms	72.7 max	34.8ms	74.3 w/o LP	32.4ms	73.9
LA(LSA)	36.1ms	74.0 avg	34.5ms	74.5 LP(Bilinear)	34.1ms	74.1
LA(Ours)	32.8ms	74.4 center	32.8ms	74.4 LP(Ours)	32.8ms	74.4

(a) **Local Aggregation** (b) **Global attention** (c) **Local Propagation****Table 4. Ablation on ImageNet-1K.** LA: the local aggregation operator. LP: the local propagation operator. LSA: the Locally-grouped Self-Attention used in [7].

in MobileNet-V3 and EfficientNets. Our models yield clear advantages over alternative ViT models. Compared to MobileViTs in the three GFLOPs settings, EdgeViTs excel by 5.4%, 2.8%, and 2.7% while being $\times 2$, $\times 2.7$, $\times 2.6$ faster.

Energy results are presented in Table 3. In addition to the raw energy and power numbers, for comparison simplicity, we define an energy-aware *efficiency* metric as the average gain in top-1 accuracy (in percentages) from each consumed 1mJ of energy. We observe that less accurate models tends to be more efficient. This is not a surprise in that improvements in accuracy scale sublinearly with model complexity. However, this also means that comparing efficiency of models with very different top-1 scores might be severely biased by the sole difficulty of achieving certain accuracy levels, which is independent from a model. Therefore, we limit our comparison to identifying *pareto-optimal models*, those upon which no other models can improve in either accuracy or energy efficiency without degrading other metrics. We can see that our EdgeViT family is able to dominate almost all other ViTs, with the only exception being LeViT-384 \dagger whose accuracy and efficiency fall between our EdgeViT-S and EdgeViT-XS. When compared to CNNs, our EdgeViTs compete with MobileNet-v3 and EfficientNet-B0 that are more efficient but also less accurate. MobileNet-v2 achieves decent results but is dominated by its newer version, MobileNet-v3. PVT-v2-B0, although high on the efficiency side, is rather inaccurate and hence is favored by highly efficient CNNs. Visibly at the end of the spectrum are the latest MobileViT models which turn out to be neither efficient nor accurate, when compared to the rest. Unlike them, our EdgeViT models, although not as efficient as best CNNs in the absolute sense, exhibit favourable trade-off between efficiency and accuracy by being rather highly accurate while not sacrificing efficiency.

Ablation study. We conduct detailed ablations to validate our design choices in the LGL bottleneck. We use EdgeViT-XXS as the base model and re-scale the alternative designs to ~ 0.5 GFLOPs for fair comparison.

Local aggregation. We compare our local aggregation (LA) operation to the Locally-grouped Self-Attention (LSA) used in [7,31]. It is shown in Table 4a that applying LA consistently improve the performance. Our convolutional LA mod-

Backbone	RetinaNet 1×							Mask R-CNN 1×						
	#Par.	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	#Par.	AP ^b	AP ₅₀ ^b	AP ₇₅ ^b	AP ^m	AP ₅₀ ^m	AP ₇₅ ^m
PVTv2-B0 [56]	13.0	37.2	57.2	39.5	23.1	40.4	49.7	23.5	38.2	60.5	40.7	36.2	57.8	38.6
EdgeViT-XXS	13.1	38.7	59.0	41.0	22.4	42.0	51.6	23.8	39.9	62.0	43.1	36.9	59.0	39.4
EdgeViT-XS	16.3	40.6	61.3	43.3	25.2	43.9	54.6	26.5	41.4	63.7	45.0	38.3	60.9	41.3
ResNet18 [20]	21.3	31.8	49.6	33.6	16.3	34.3	43.2	31.2	34.0	54.0	36.7	31.2	51.0	32.7
PVTv1-Tiny [55]	23.0	36.7	56.9	38.9	22.6	38.8	50.0	32.9	36.7	59.2	39.3	35.1	56.7	37.3
PVTv2-B1 [56]	23.8	41.2	61.9	43.9	25.4	44.5	54.3	33.7	41.8	64.3	45.9	38.8	61.2	41.6
EdgeViT-S	22.6	43.4	64.9	46.5	26.9	47.5	58.1	32.8	44.8	67.4	48.9	41.0	64.2	43.8

Table 5. Comparison to other visual backbones using RetinaNet and Mask-RCNN on COCO val2017 object detection and instance segmentation. “#Par.” refers to number of parameters in million. AP^b and AP^m indicate bounding box AP and mask AP.

ule performs better than the self-attention based operator (LSA). This validates our choice of using **depth-wise convolutions** in LA for local context learning.

Global sparse attention. We explore three options for delegate token sampling: **max**, **avg**, and **center**. All choices perform similarly in terms of accuracy, with our default design **center** being slightly faster.

Local propagation. We investigate two alternatives to the local propagation operator: i) **w/o LP**: We simply remove the local propagation. Note that EdgeViTsw/o LP has similar complexity to standard EdgeViTs. ii) **Bilinear**: we use the bilinear interpolation, instead of the transposed convolution, to up-sample the delegate tokens. Table 4c shows that adding LP improves the top-1 accuracy by 0.5%, with only 0.4ms overhead.

4.2 Dense Prediction

Following [55,56], we also evaluate the proposed EdgeViTs on COCO Objection Detection/Instance Segmentation [29] and ADE20K Scene Parsing [65]. Here, we use the EdgeViTs as the feature extractor for the main model and initialize it with the ImageNet1K-pretrained weights obtained in our previous experiments.

COCO Object Detection/Instance Segmentation. We demonstrate the performance of our model in main-stream object detection and instance segmentation frameworks: RetinaNet [28] for object detection, Mask R-CNN[19] with the FPN [27] for instance segmentation. Following the training protocol in [56,7,26], we resize the training images to have a shorter side of 800 pixels while keeping the longer side to be smaller than 1333 pixels. During testing, the images are re-scaled to have a shorter size of 800 pixels. The models are fine-tuned with 1× schedule (i.e. 12 epochs) by AdamW[33] using an initial learning rate of 1×10^{-4} and a batch size of 16. We train the models on the COCO 2017 training set and report the mAP@100 score on the COCO 2017 validation set.

Results. In Table 5, our EdgeViTs perform consistently better than other visual backbones on RetinaNet [28] and Mask R-CNN[19]. Our smallest vari-

Backbone	Semantic FPN		
	#Param (M)	GFLOPs	mIoU (%)
PVTv2-B0[56]	7.6	25.0	37.2
EdgeViT-XXS	7.9	24.4	39.7
EdgeViT-XS	10.6	27.7	41.4
ResNet18 [20]	15.5	32.2	32.9
PVTv1-Tiny [55]	17.0	33.2	35.7
PVTv2-B1[56]	17.8	34.2	42.5
EdgeViT-S	16.9	32.1	45.9

Table 6. Semantic segmentation results on the validation set of ADE20K. Segmentation model: Semantic FPN [25]. GFLOPs: Calculated at 512×512 input size.

ant EdgeViT-XXS, when used on RetinaNet [28], achieves 1.5 higher AP than PVTv2-B0. When used on Mask R-CNN [19], EdgeViT-XXS also surpasses PVTv2-B0 by 1.7 on the bounding box detection task (AP^b), and by 0.7 on the mask segmentation task (AP^m). For EdgeViT-S, we observe even larger gains when comparing to PVTv2-B1: +2.2 on RetinaNet [28], +3.0 AP^b and +1.2 AP^m on Mask R-CNN[19].

ADE20K Scene Parsing. We incorporate the pretrained EdgeViT in the Semantic FPN segmentation model [25]. As in [55,56], we create 512×512 random crops of the images during training and resize the images to have a shorter side of 512 pixels during inference. The models are finetuned by AdamW [33] using an initial learning rate of 1×10^{-4} and a batch size of 16. We train the models for 80K iterations on the ADE20K training set, and report the mean Intersection over Union (mIoU) score on the validation set.

Results In Table 6, we compare EdgeViTs to both CNN (ResNet-18 [20]) and ViT backbones (PVTs[55,56]) for FPN based Semantic Segmentation [25]. EdgeViTs achieves better performance than all counterparts at similar compute costs. Particularly, EdgeViT-XXS outperforms PVTv2-B0 by 2.5% in mIoU, EdgeViT-S surpasses PVTv2-B1 by a margin of 3.4%.

5 Conclusion

In this work, we investigate the design of efficient ViTs from the on-device deployment perspective. By introducing a novel decomposition of self-attention, we present a family of EdgeViTs that, for the first time, achieve comparable or even superior accuracy-efficiency tradeoff on generic visual recognition tasks, in comparison to a variety of state-of-the-art efficient CNNs and ViTs. We conduct extensive on-device experiments using practically critical and previously underestimated metrics (e.g., energy-aware efficiency) and reveal new insights and observations in the comparison of light-weight CNN and ViT models.

Acknowledgements. We thank Victor Escorcia, Yassine Ouali and Javier Fernandez for helpful discussions.

References

1. Almeida, M., Laskaridis, S., Mehrotra, A., Dudziak, L., Leontiadis, I., Lane, N.D.: Smart at what cost? characterising mobile deep neural networks in the wild. In: ACM Internet Measurement Conference (2021)
2. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
3. Berman, M., Pishchulin, L., Xu, N., Blaschko, M.B., Medioni, G.: AOWS: adaptive and optimal network width search with latency constraints. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
4. Bulat, A., Perez Rúa, J.M., Sudhakaran, S., Martinez, B., Tzimiropoulos, G.: Space-time mixing attention for video transformer. Advances on Neural Information Processing Systems (2021)
5. Bulat, A., Tzimiropoulos, G.: Bit-Mixer: Mixed-precision networks with runtime bit-width selection. In: IEEE International Conference on Computer Vision (2021)
6. Chen, Y., Dai, X., Liu, M., Chen, D., Yuan, L., Liu, Z.: Dynamic convolution: Attention over convolution kernels. In: IEEE Conference on Computer Vision and Pattern Recognition (2020)
7. Chu, X., Tian, Z., Wang, Y., Zhang, B., Ren, H., Wei, X., Xia, H., Shen, C.: Twins: Revisiting the design of spatial attention in vision transformers. Advances on Neural Information Processing Systems (2021)
8. Chu, X., Tian, Z., Zhang, B., Wang, X., Wei, X., Xia, H., Shen, C.: Conditional positional encodings for vision transformers. arXiv preprint arXiv:2102.10882 (2021)
9. Deb, K.: Multi-objective optimization. In: Search methodologies, pp. 403–449. Springer (2014)
10. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021)
11. Dudziak, L., Chau, T., Abdelfattah, M.S., Lee, R., Kim, H., Lane, N.D.: BRP-NAS: Prediction-based NAS using GCNs. In: Advances on Neural Information Processing Systems (2020)
12. Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C.: Multiscale vision transformers. In: IEEE International Conference on Computer Vision (2021)
13. Fayyaz, M., Koohpayegani, S.A., Jafari, F.R., Sommerlade, E., Joze, H.R.V., Pirsavash, H., Gall, J.: ATS: adaptive token sampling for efficient vision transformers. arXiv preprint arXiv:2111.15667 (2021)
14. Gehring, J., Auli, M., Grangier, D., Yarats, D., Dauphin, Y.N.: Convolutional sequence to sequence learning. In: International Conference on Machine Learning (2017)
15. Graham, B., El-Nouby, A., Touvron, H., Stock, P., Joulin, A., Jégou, H., Douze, M.: LeViT: a vision transformer in convnet’s clothing for faster inference. In: IEEE International Conference on Computer Vision (2021)
16. Han, K., Xiao, A., Wu, E., Guo, J., Xu, C., Wang, Y.: Transformer in transformer. Advances on Neural Information Processing Systems (2021)
17. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In: International Conference on Learning Representations (2016)

18. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: IEEE Conference on Computer Vision and Pattern Recognition (2022)
19. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: IEEE International Conference on Computer Vision (2017)
20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
21. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 (2016)
22. Howard, A., Pang, R., Adam, H., Le, Q.V., Sandler, M., Chen, B., Wang, W., Chen, L., Tan, M., Chu, G., Vasudevan, V., Zhu, Y.: Searching for MobileNetV3. In: IEEE International Conference on Computer Vision (2019)
23. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861 (2017)
24. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., Kalenichenko, D.: Quantization and training of neural networks for efficient integer-arithmetic-only inference. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
25. Kirillov, A., Girshick, R., He, K., Dollár, P.: Panoptic feature pyramid networks. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
26. Li, K., Wang, Y., Peng, G., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unified transformer for efficient spatial-temporal representation learning. In: International Conference on Learning Representations (2022)
27. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
28. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: IEEE International Conference on Computer Vision (2017)
29. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: European Conference on Computer Vision (2014)
30. Liu, Y., Sangineto, E., Bi, W., Sebe, N., Lepri, B., De Nadai, M.: Efficient training of visual transformers with small datasets. Advances on Neural Information Processing Systems (2021)
31. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. IEEE International Conference on Computer Vision (2021)
32. Liu, Z., Li, J., Shen, Z., Huang, G., Yan, S., Zhang, C.: Learning efficient convolutional networks through network slimming. In: IEEE International Conference on Computer Vision (2017)
33. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
34. Loshchilov, I., Hutter, F.: SGDR: stochastic gradient descent with warm restarts. In: International Conference on Learning Representations (2017)
35. Lu, J., Yao, J., Zhang, J., Zhu, X., Xu, H., Gao, W., Xu, C., Xiang, T., Zhang, L.: Soft: Softmax-free transformer with linear complexity. Advances on Neural Information Processing Systems (2021)
36. Ma, N., Zhang, X., Zheng, H.T., Sun, J.: ShuffleNetV2: Practical guidelines for efficient CNN architecture design. In: European Conference on Computer Vision (2018)

37. Mehta, S., Rastegari, M.: MobileViT: light-weight, general-purpose, and mobile-friendly vision transformer. *International Conference on Learning Representations (2022)*
38. Pan, B., Jiang, Y., Panda, R., Wang, Z., Feris, R., Oliva, A.: IA-RED²: Interpretability-aware redundancy reduction for vision transformers. In: *Advances on Neural Information Processing Systems (2021)*
39. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: PyTorch: An imperative style, high-performance deep learning library. In: *Advances on Neural Information Processing Systems (2019)*
40. Radosavovic, I., Kosaraju, R.P., Girshick, R., He, K., Dollár, P.: Designing network design spaces. In: *IEEE Conference on Computer Vision and Pattern Recognition (2020)*
41. Rao, Y., Zhao, W., Liu, B., Lu, J., Zhou, J., Hsieh, C.J.: DynamicViT: Efficient vision transformers with dynamic token sparsification. *Advances on Neural Information Processing Systems (2021)*
42. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International Journal on Computer Vision (2015)*
43. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: MobileNetV2: Inverted residuals and linear bottlenecks. In: *IEEE Conference on Computer Vision and Pattern Recognition (2018)*
44. Shaw, P., Uszkoreit, J., Vaswani, A.: Self-attention with relative position representations. In: *North American Chapter of the Association for Computational Linguistics (2018)*
45. Srinivas, A., Lin, T.Y., Parmar, N., Shlens, J., Abbeel, P., Vaswani, A.: Bottleneck transformers for visual recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition (2021)*
46. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: *IEEE Conference on Computer Vision and Pattern Recognition (2016)*
47. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: MnasNet: Platform-aware neural architecture search for mobile. In: *IEEE Conference on Computer Vision and Pattern Recognition (2019)*
48. Tan, M., Le, Q.: EfficientNet: Rethinking model scaling for convolutional neural networks. In: *International Conference on Machine Learning (2019)*
49. Tan, M., Le, Q.: EfficientNetV2: Smaller models and faster training. In: *International Conference on Machine Learning (2021)*
50. Tan, M., Le, Q.V.: Mixconv: Mixed depthwise convolutional kernels. In: *British Machine Vision Conference (2019)*
51. Tan, M., Pang, R., Le, Q.V.: Efficientdet: Scalable and efficient object detection. In: *IEEE Conference on Computer Vision and Pattern Recognition (2020)*
52. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: *International Conference on Machine Learning (2021)*
53. Touvron, H., Cord, M., Sablayrolles, A., Synnaeve, G., Jégou, H.: Going deeper with image transformers. *arXiv preprint arXiv:2103.17239 (2021)*
54. Wang, S., Li, B., Khabsa, M., Fang, H., Ma, H.: Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768 (2020)*

55. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: IEEE International Conference on Computer Vision (2021)
56. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: PVTv2: Improved baselines with pyramid vision transformer. Computational Visual Media (2022)
57. Wen, W., Wu, C., Wang, Y., Chen, Y., Li, H.: Learning structured sparsity in deep neural networks. In: Advances on Neural Information Processing Systems (2016)
58. Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., Zhang, L.: CvT: Introducing convolutions to vision transformers. IEEE International Conference on Computer Vision (2021)
59. Xiao, T., Singh, M., Mintun, E., Darrell, T., Dollár, P., Girshick, R.: Early convolutions help transformers see better. Advances on Neural Information Processing Systems (2021)
60. Xiong, Y., Zeng, Z., Chakraborty, R., Tan, M., Fung, G., Li, Y., Singh, V.: Nyströmformer: A nyström-based algorithm for approximating self-attention. In: AAAI Conference on Artificial Intelligence (2021)
61. Xu, Y., Zhang, Z., Zhang, M., Sheng, K., Li, K., Dong, W., Zhang, L., Xu, C., Sun, X.: Evo-vit: Slow-fast token evolution for dynamic vision transformer. In: AAAI Conference on Artificial Intelligence (2022)
62. Yang, B., Bender, G., Le, Q.V., Ngiam, J.: CondConv: Conditionally parameterized convolutions for efficient inference. In: Advances on Neural Information Processing Systems (2019)
63. Yuan, L., Chen, Y., Wang, T., Yu, W., Shi, Y., Jiang, Z., Tay, F.E., Feng, J., Yan, S.: Tokens-to-token ViT: Training vision transformers from scratch on ImageNet. In: IEEE International Conference on Computer Vision (2021)
64. Zhang, X., Zhou, X., Lin, M., Sun, J.: ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
65. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralla, A.: Scene parsing through ade20k dataset. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
66. Zhou, D., Kang, B., Jin, X., Yang, L., Lian, X., Jiang, Z., Hou, Q., Feng, J.: Deepvit: Towards deeper vision transformer. arXiv preprint arXiv:2103.11886 (2021)

A Appendix

A.1 Computing Complexity

We calculate the computational cost of spatial context modeling involved in our proposed LGL bottleneck. We omit point-wise operations for simplicity as the key difference is on the spatial modeling part. Let us assume an input $X \in \mathcal{R}^{h \times w \times c}$ where h , w , c denotes the height, the width, and the channel dimension, respectively. The cost of the local aggregation is $\mathcal{O}(k^2hwc)$, where k^2 is the local group size. By selecting one delegate out of r^2 tokens with r the sub-sampling rate, the complexity of our Sparse Global Self-Attention is then $\mathcal{O}(\frac{h^2w^2}{r^4}c)$. Finally, the local propagation step takes a cost of $\mathcal{O}(r^2hwc)$. Putting all these together we have a total cost of LGL is $\mathcal{O}(k^2hwc + \frac{h^2w^2}{r^4}c + r^2hwc)$. When comparing with the cost of a standard multi-head self-attention $\mathcal{O}(h^2w^2c)$, we can see that our LGL significantly reduces the computation overhead when $k \ll h, w$; and $r > 1$. In our experiments, for simplicity we set $k = 3$, and r to (4,2,2,1) for the four stages.

A.2 Implementation Details

All variants of EdgeViTs can be built upon these components according to the schematic overview (Fig. 2a of the main paper), and the model configuration parameters (Table 1a. of the main paper). For more details with the ablation studies in the main paper, we have replaced or removed one of these blocks with the details given below.

(1) In Table 4a, for the case of w/o LA, we aim to test the importance of separate local and global context modeling. Thus we remove both `LocalAgg` and `GlobalSparseAttn`, and instead use the Spatial-Reduced Self-Attention⁴ introduced in PVT [56], resulting in a single Self-Attention Block for both local and global context modeling. For the case of LA(LSA) we simply replace `LocalAgg` with Local-grouped Self-Attention⁵ introduced in [7].

(2) In Table 4b, we replace the default sampler (`Center`) with `Avg` and `Max` functions which can be implemented with `AvgPool2d()` and `MaxPool2d()` in Pytorch [39], respectively. Note, for both cases the kernel size is set to `sample_rate`.

(3) In Table 4c, in the case of w/o LP, we replace our `GlobalSparseAttn` with Spatial-Reduce Self-Attention from PVT [56], but different from w/o LA, we keep the `LocalAgg`. For the case of LP(Bilinear), the `LocalProp` is instantiated as a bilinear interpolation function (`Upsample(mode='bilinear')` in Pytorch).

Note, the number of layers for each of these variants is down-scaled to have 0.5GFLOPs for fair comparison.

⁴ https://github.com/whai362/PVT/blob/v2/classification/pvt_v2.py#L54-L126

⁵ <https://github.com/Meituan-AutoML/Twins/blob/main/gvt.py#L32-L71>

A.3 Accuracy-Speed Pareto-Optimal Models

In order to facilitate the Accuracy vs. Speed interpretation. We identify pareto optimal models when comparing trade-off between accuracy and latency [9]. In our context, the accuracy-latency pareto-optimal models are defined as those upon which no other models can improve in either accuracy or latency without degrading other metrics. As shown in Fig. 1, our EdgeViTs are well comparable with best efficient CNNs [43,22,48], whilst significantly dominating over all prior ViT counterparts. Specifically, EdgeViTs are all pareto-optimal in both trade-offs.

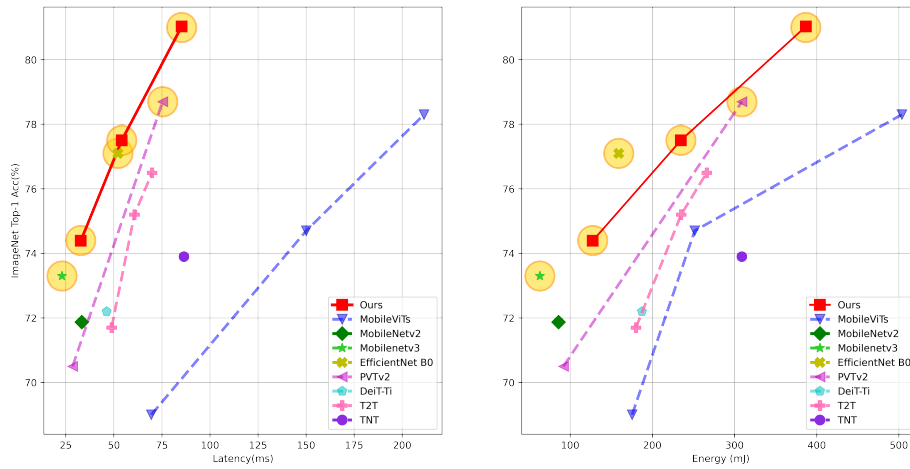


Fig. 4. Accuracy/Latency and Accuracy/Energy trade off on ImageNet-1K. Note that, all three variants of our EdgeViTs are pareto-optimal, which are highlighted with **amber circle**. *Testing device:* Samsung Galaxy S21 (latency), Snapdragon 888 Hardware Development Kit (energy).

A.4 Efficiency in detection/segmentation

This work proposes a genetic transformer-based network and demonstrate its efficacy when used as the backbone in detection/segmentation. We provide a evaluation by measuring only the inference time, energy and efficiency of the backbones for detection/segmentation. As shown in Tab. 7 and 8, EdgeViTs demonstrate higher efficiency compared to the baselines.

Model	AP (%)	CPU (s)	Energy (J)	Efficiency (%/msW)
PVTv2-B0	37.2	1.34 \pm 0.05	3.50 \pm 0.77	0.011
ResNet18	31.8	0.58 \pm 0.02	2.22 \pm 0.35	0.014
PVTv1-Tiny	36.7	1.91 \pm 0.18	4.40 \pm 0.94	0.008
PVTv2-B1	41.2	2.81 \pm 0.26	5.29 \pm 1.49	0.008
EdgeViT-XXS	38.7	0.59 \pm 0.02	2.02 \pm 0.58	0.019
EdgeViT-XS	40.6	0.90 \pm 0.03	2.89 \pm 0.66	0.014
EdgeViT-S	43.4	1.88 \pm 0.05	4.36 \pm 1.06	0.010

Table 7. Detection Efficiency. Input size: 800 \times 800

Model	mIoU (%)	CPU (s)	Energy (J)	Efficiency (%/msW)
PVTv2-B0	37.2	0.35 \pm 0.01	1.10 \pm 0.30	0.034
ResNet18	32.9	0.23 \pm 0.01	1.03 \pm 0.20	0.032
PVTv1-Tiny	35.7	0.49 \pm 0.02	1.63 \pm 0.32	0.022
PVTv2-B1	42.5	0.75 \pm 0.03	2.13 \pm 0.82	0.020
EdgeViT-XXS	39.7	0.19 \pm 0.01	0.71 \pm 0.11	0.056
EdgeViT-XS	41.4	0.31 \pm 0.01	1.11 \pm 0.28	0.037
EdgeViT-S	45.9	0.52 \pm 0.02	1.73 \pm 0.36	0.027

Table 8. Segmentation Efficiency. Input size: 512 \times 512