# What Makes for Hierarchical Vision Transformer?

Yuxin Fang[1*], Xinggang Wang[1(✉)], Rui Wu[2], Jianwei Niu[2], Wenyu Liu[1]
[1]School of EIC, Huazhong University of Science & Technology    [2] Horizon Robotics

{yxf, xgwang}@hust.edu.cn

## Abstract

*Recent studies show that hierarchical Vision Transformer with interleaved non-overlapped intra window self-attention & shifted window self-attention is able to achieve state-of-the-art performance in various visual recognition tasks and challenges CNN's dense sliding window paradigm. Most follow-up works try to replace shifted window operation with other kinds of cross window communication while treating self-attention as the de-facto standard for intra window information aggregation. In this short preprint, we question whether self-attention is the only choice for hierarchical Vision Transformer to attain strong performance, and what makes for hierarchical Vision Transformer? We replace self-attention layers in Swin Transformer and Shuffle Transformer with simple linear mapping and keep other components unchanged. The resulting architecture with 25.4M parameters and 4.2G FLOPs achieves 80.5% Top-1 accuracy, compared to 81.3% for Swin Transformer with 28.3M parameters and 4.5G FLOPs. We also experiment with other alternatives to self-attention for context aggregation inside each non-overlapped window, which all give similar competitive results under the same architecture. Our study reveals that the **macro architecture** of Swin model families (i.e., interleaved intra window & cross window communications), other than specific aggregation layers or specific means of cross window communication, may be more responsible for its strong performance and is the real challenger to CNN's dense sliding window paradigm.*

## 1. Introduction

Recently, the impregnable position of convolutional neural networks (CNNs) in computer vision seems to be weakened by the emerging hierarchical Vision Transformer families. As one representative, Swin Transformer [8] and its variants (*e.g.*, [2, 5]) with an interleaved non-overlapped in-

| | MHSA | Linear | DW Linear | MLP |
|---|---|---|---|---|
| Shifted Window [8] | 80.5 | 79.7 | 79.8 | 79.9 |
| Shuffle [5] | 80.4 | 79.6 | 79.6 | 79.8 |

Table 1: Study of two representative hierarchical Vision Transformer with different information aggregation layers on ImageNet-1$k$ [10]. All models are trained with 200 epochs using the optimization scheme in [8, 13]. The model using MHSA layers keeps the same configurations as [5, 8]. Other models are tuned and selected with our best effort to consume $\sim$ 4.5G FLOPs budgets and parameters ranging from 24M to 29M with only model depth & width adaptations. We remove all densely slid conv-layers in the network stem and each block of [5] for a clearer study of different aggregation layers in non-overlapped windows. *The results motivate us to focus more on the macro architecture design than specific aggregation layers or specific means of cross window communication.*

tra window self-attention & shifted window self-attention paradigm are able to achieve strong performance in image recognition and demonstrate excellent transferability on various downstream computer vision tasks. Swin Transformer partitions feature maps to a series of non-overlapped local windows, and uses multi-head self-attention (MHSA) layers to aggregate information in each window individually. Instead of using a dense sliding window paradigm for cross window token mixing like CNNs, Swin Transformer propose to shift windows between consecutive layers. Most follow-up works of Swin Transformer focusing on replacing shift window operation with other kinds of cross window communication, while the use of window-based MHSA is usually taken for granted.

In this short preprint, we question whether MHSA is the only choice to aggregate information for Swin model families. MHSA is good at capture dense and long-range contextual information [6, 17] in visual recognition, and what makes for hierarchical Vision Transformer. Intuitively, it is somewhat too aggressive to use MHSA to model contextual relation in a local window with only $7 \times 7 = 49$ tokens. This motivates us to replace MHSA layers with linear map-

---

| | Conv | MHSA | Linear & MLP |
|---|---|---|---|
| Columnar | Missing Piece | [1, 13] | [4, 7, 9, 11, 12] |
| Hierarchical | Too Many | [2, 3, 5, 8, 14] | **Missing Piece** |

Table 2: A cursory summary of macro architectures (*Column*) and specific aggregation layers (*Row*). There are two "missing pieces" remain.

ping, one of the most common & simplest components in neural architecture design, in two representative hierarchical Vision Transformer instances, *i.e.*, Swin Transformer [8] and Shuffle Transformer [5]. The resulting model is termed as LINEARMAPPER. We find that LINEARMAPPER with simple linear mappings is sufficient for local context aggregation, and is able to achieve very competitive performance in ImageNet-1$k$ using strong data augmentations and regularization methods [8, 13].

Furthermore, we experiment with other variants for information aggregation inside each local window (*e.g.*, depth-wise linear mapping & MLP). We find they all give similar competitive results under the same architecture, as shown in Tab. 1.

Based on the available evidence, we hypothesize that the high-level design methodology of Swin model families (*i.e.*, interleaved non-overlapped intra window token mixing & cross window communications), other than specific aggregation layers such as MHSA or specific means of cross window communication such as shifted window & spatial shuffle, may be more responsible for the strong performance and is the real challenger to CNN's dense sliding window paradigm. This short preprint is **not** an attempt to show that linear or MLP layers are superior to MHSA. On the contrary, we find MHSA is better than linear & MLP in terms of accuracy with fewer parameters & FLOPs budgets (see Tab. 5). Our purpose is to abstract away from specific aggregation layers and highlight the importance of macro architecture. We hope this short preprint can encourage the community to rethink the role of attention in neural architecture design and shed a little light on future studies of general visual representation learning.

## 2. Background

Highly mature and robust training recipes [13] enable standard Transformer architecture [1, 15] inherit from NLP to attain excellent performance in the image recognition task even with limited data. To efficiently apply Vision Transformers to other downstream tasks in computer vision, two key issues need to be solved: (1) involving hierarchical architecture to establish multi-scale feature representations, and (2) reduce memory & computation costs from global attention. [16] processes Transformer features under multi-resolution stages instead of in a columnar man-

**Algorithm 1** - `LinearMapper` Pseudocode.

```
# B: num_windows, C: channel
# ws: window size, gs: group size
# t(dim1, dim2): # transpose dim1 & dim2

lin_map_h = Linear(ws*gs, ws*gs)
lin_map_w = Linear(ws*gs, ws*gs)
proj = PointWiseConv(C, C)

# LinearMapper
# x: input features with shape of (B, C, ws*ws)
def LinearMapper(x):

    # Height dim linear mapping for each window
    hf = x.view(B, C//gs, gs*ws, ws)
    hf = lin_map_h(hf.t(-1, -2)).t(-1, -2)
    hf = hf.view(B, C, ws*ws)

    # Width dim linear mapping for each window
    wf = x.view(B, C//gs, ws, gs*ws)
    wf = lin_map_w(wf)
    wf = wf.view(B, C, ws*ws)

return proj(hf + wf)
```

ner. [8, 14] propose to compute attention in a local window. Many follow up hierarchical window-based Vision Transformers emerge (*e.g.*, [2, 5]) and challenge the hegemonic position of CNN in computer vision. Recently, another series of work explores the performance of architectures based exclusively on columnar structured MLPs (*e.g.*, [11, 12]) in image recognition tasks.

As cursorily summarized in Tab. 2, there are two "missing pieces" remain, *i.e.*, a CNN with columnar architectures, and an MLP with hierarchical architectures. This preprint gives a very brief study to the latter one: a straightforward, incremental, yet must-know model in computer vision. We argue it is inevitable to investigate the potential of hierarchical linear & MLP structures in computer vision, which encourages and leads us to rethink the role between macro model design methodology and specific network building blocks.

## 3. Method

We attempt to use linear mapping, one of the simplest components in neural architecture design, as a touchstone to reveal the macro architecture (*i.e.*, non-overlapped intra window & cross window communication in an alternating fashion) seems to be more responsible for Swin model families' strong performance other than specific aggregation layers such as MHSA.

We choose two representative and publicly available instantiations of Swin model families, *i.e.*, Swin Transformer [8] and Shuffle Transformer [5]. We directly replace their `WindowAttention` modules with `LinearMapper` described in Algorithm 1, other components and configurations keep unchanged. The linear mappings are performed on tokens' height and width dimensions separately similar

| Method | Model Width | Model Depth | #Params. (M) | FLOPs (G) | Throughput (Img/s) | Top-1 Acc. |
|---|---|---|---|---|---|---|
| Swin Transformer [8] | 96 | $\{2, 2, 6, 2\}$ | 28.3 | 4.5 | 378 | 81.3 |
| Swin LINEARMAPPER (Ours) | 64 | $\{2, 4, 22, 4\}$ | 25.4 | 4.2 | 320 | 80.5 |
| Shuffle Transformer [5] | 96 | $\{2, 2, 6, 2\}$ | 28.3 | 4.5 | 476 | 81.4 |
| Shuffle LINEARMAPPER (Ours) | 64 | $\{2, 4, 22, 4\}$ | 25.4 | 4.2 | 445 | 80.4 |

Table 3: Comparisons with two hierarchical Vision Transformers on ImageNet-$1k$ with $300$ epochs training.

| Method | #Params. (M) | FLOPs (G) | Throughput (Img/s) | Top-1 Acc. |
|---|---|---|---|---|
| MLP Mixer-B/16 [7, 11] | 59 | 12.7 | 227 | 77.3 |
| ResMLP-24 [12] | 30 | 6.0 | 497 | 79.4 |
| ResMLP-36 [12] | 45 | 8.9 | 327 | 79.7 |
| gMLP-S [7] | 20 | 4.5 | 419 | 79.6 |
| Swin LINEARMAPPER (Ours) | 25 | 4.2 | 320 | 80.5 |
| Shuffle LINEARMAPPER (Ours) | 25 | 4.2 | 445 | 80.4 |

Table 4: Comparison with some global & columnar MLP variants on ImageNet-$1k$ with $300$ epochs training.

to the merit of [6].

In our default instantiation, the weights of `lin_map_h` and `lin_map_w` are *shared* across different groups. Linear mappings with separate parameters for different groups are denoted as depth-wise linear mapping (DW Linear in Tab. 1 and Tab. 6) in the context of this short preprint.

As a (very important) by-product of our study, the proposed LINEARMAPPER along with its variants enables a fully linear & MLP architecture to *directly process input images with arbitrary shapes*. This property allows MLP architectures to be easily transferred to dense prediction tasks such as object detection and scene parsing. We will study the transferability of LINEARMAPPER to other downstream computer vision tasks in the future and focus on the image recognition task in this preprint.

## 4. Experiments

### 4.1. Setup

The experiments are conducted on the public available codebase of [5, 8] and `timm` library [18]. All models are trained and evaluated on ImageNet-$1k$ [10] following the setup in [8, 13]. We train models with 300 epochs for main results in Tab. 3, and 200 epoch for other tables and analysis. The input resolution is $224 \times 224$ and the window size is $7 \times 7$ for all experiments. For a clearer study of different aggregation layers in non-overlapped windows, we remove all densely slid conv-layers in the network stem and each block of [5] in this preprint. Model throughput data[1] during inference are measured using a single Titan Xp GPU with batch size $64$.

### 4.2. Main Results

**Comparisons with Hierarchical Vision Transformers.** As shown in Tab. 3, LINEARMAPPER is flexible and able to achieve competitive performance on ImageNet-$1k$ image recognition benchmark with two different cross window communication paradigm (*i.e.*, shifted window and spatial shuffle). Along with the results in Tab. 1, it is interesting to note that (1) shifted window and spatial shuffle give similar results under the same aggregation layer, and (2) different aggregation layers are all quite competitive under the same cross window token mixing approach. These results support our proposal: the macro architecture of Swin model families (*i.e.*, interleaved intra window & cross window token mixing), other than specific aggregation layers or specific means of cross window communication, may be more responsible for its strong performance.

**Comparison with Global & Columnar MLP Variants.** The window partitioning operation and hierarchical architecture introduce 2D locality bias and invariance to linear & MLP architectures. It is not a surprise that LINEARMAPPER is more efficient and competitive than global & columnar MLP structures by leveraging these CNN design priors, as shown in Tab. 4.

### 4.3. Analysis

In this section, we study the impact of model width & depth configurations, weight sharing property of linear mapping layers, as well as the number of groups. Overall, we conclude that LINEARMAPPER is quite robust to different model choices and configurations thanks to the macro model architecture inherit from [8] and highly mature training recipes provide by [13].

---

[1] To our knowledge, the inference speed gap between Shuffle Transformer (w/o dense conv-layers) and Swin Transformer shown in Tab. 3 mainly comes from different normalization layers (*i.e.*, `BatchNorm` for Shuffle Transformer *v.s.* `LayerNorm` for Swin Transformer).

| Method | Model Width | Model Depth | #Params. (M) | FLOPs (G) | Throughput (Img/s) | Top-1 |
|---|---|---|---|---|---|---|
| Swin Transformer [8] | 96 | {2, 2, 6, 2} | 28.3 | 4.5 | 378 | 80.5 |
| Swin LINEARMAPPER (Baseline) | 96 | {2, 2, 6, 2} | 22.0 | 3.6 | 413 | 78.8 |
| Swin LINEARMAPPER-Wider | 112 | {2, 2, 6, 2} | 30.6 | 4.8 | 343 | 79.8 |
| Swin LINEARMAPPER-Deeper | 64 | {2, 4, 22, 4} | 25.4 | 4.2 | 320 | 79.7 |

Table 5: Wider *v.s.* deeper macro structure for LINEARMAPPER (200 epochs training on ImageNet-1$k$).

| #Group | Linear | DW Linear | #Params. | FLOPs | Top-1 |
|---|---|---|---|---|---|
| 32 | ✓ |  | 22 | 3.6 | 78.8 |
| 32 |  | ✓ | 28 | 3.6 | 78.9 |
| 48 | ✓ |  | 22 | 3.6 | 78.7 |
| 48 |  | ✓ | 26 | 3.6 | 78.9 |

Table 6: Linear *v.s.* depth-wise Linear in LINEARMAPPER (width: 96, depth: {2, 2, 6, 2}, 200 epochs training on ImageNet-1$k$).

| #Group | gs in Alg. 1 | #Params. | FLOPs | Throughput | Top-1 |
|---|---|---|---|---|---|
| 96 | 1 | 22 | 3.6 | 413 | 78.4 |
| 48 | 2 | 22 | 3.6 | 416 | 78.7 |
| 32 | 3 | 22 | 3.6 | 413 | 78.8 |
| 16 | 6 | 22 | 3.6 | 410 | 78.8 |
| 8 | 12 | 25 | 3.6 | 420 | 78.6 |

Table 7: Number of groups in LINEARMAPPER (width: 96, depth: {2, 2, 6, 2}, 200 epochs training on ImageNet-1$k$).

**Going Wider or Deeper?** We adjust the model width and depth of LINEARMAPPER to align with the FLOPs budgets of Swin-T Transformer. As shown in Tab. 5, LINEARMAPPER-Wider seems to be more speed-friendly while LINEARMAPPER-Deeper is a lot more parameter efficient. We choose LINEARMAPPER-Deeper's model width and height as our default instantiation for the main results.

**Linear or Depth-wise Linear?** Depth-wise linear (DW Linear) layers refer to linear mapping with separate parameters for each group in the context of this preprint. Therefore the model parameters increase while theoretical FLOPs keep unchanged when using DW linear layers instead of shared linear weights. DW linear layers bring no significant improvement as shown in Tab. 6. Therefore we choose to share linear weights across different groups as our default instantiation for better efficiency.

**Number of Groups (#Group) for Linear Layers.** In Tab. 7, we study the impact of different #Groups in linear layers. The weights of linear layers are shared across groups. We find setting #Groups too large or too small is harmful to performance, while other choices yield similar results. In this preprint, #Groups = 32 (gs = 3 in Alg. 1) is chosen as the default instantiation.

## 5. Conclusion and Future Work

In this short preprint, we raise a crucial question: "What makes for hierarchical Vision Transformer?", and attempt to give an answer: the macro architecture design methodology may be more important than specific network layers and components. The proposed LINEARMAPPER along with its variants also enables a fully linear & MLP architecture to directly process input images with arbitrary shapes, which makes it possible for investigating the transferability of linear & MLP architectures on various downstream tasks other than image recognition.

However, the available evidence is insufficient to some extent for we only conduct study on two specific hierarchical Vision Transformer instances with a limited model-size range. Moreover, the transferability of LINEARMAPPER on other challenging downstream tasks such as object detection and scene parsing is not evaluated in this preprint. We leave them as future work.

## References

[1] 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. *arXiv preprint arXiv:2010.11929*, 2020.

[2] Jiemin Fang, Lingxi Xie, Xinggang Wang, Xiaopeng Zhang, Wenyu Liu, and Qi Tian. Msg-transformer: Exchanging local spatial information by manipulating messenger tokens. *arXiv preprint arXiv:2105.15168*, 2021.

[3] Peng Gao, Jiasen Lu, Hongsheng Li, Roozbeh Mottaghi, and Aniruddha Kembhavi. Container: Context aggregation network. *arXiv preprint arXiv:2106.01401*, 2021.

[4] Meng-Hao Guo, Zheng-Ning Liu, Tai-Jiang Mu, and Shi-Min Hu. Beyond self-attention: External attention using two linear layers for visual tasks. *arXiv preprint arXiv:2105.02358*, 2021.

[5] Zilong Huang, Youcheng Ben, Guozhong Luo, Pei Cheng, Gang Yu, and Bin Fu. Shuffle transformer: Rethink-

ing spatial shuffle for vision transformer. *arXiv preprint arXiv:2106.03650*, 2021.

[6] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *ICCV*, 2019.

[7] Hanxiao Liu, Zihang Dai, David R So, and Quoc V Le. Pay attention to mlps. *arXiv preprint arXiv:2105.08050*, 2021.

[8] 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. *arXiv preprint arXiv:2103.14030*, 2021.

[9] Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly well on imagenet. *arXiv preprint arXiv:2105.02723*, 2021.

[10] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015.

[11] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021.

[12] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, and Hervé Jégou. Resmlp: Feedforward networks for image classification with data-efficient training. *arXiv preprint arXiv:2105.03404*, 2021.

[13] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.

[14] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12894–12904, 2021.

[15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.

[16] 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. *arXiv preprint arXiv:2102.12122*, 2021.

[17] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.

[18] Ross Wightman. Pytorch image models. https://git.io/fjVdB, 2019.