

Pay Attention to MLPs

1. Introduction

Motivation

- It remains an open question whether the inductive bias in self-attention is essential to the remarkable effectiveness of Transformers.
- Study the necessity of self-attention modules in key language and vision applications

MLP	Self-Attention
MLPs with static parameterization can represent arbitrary functions	The attention mechanism introduces the inductive bias that the model can be dynamically parameterized based on the input representations

- Propose gMLP, and show experiments [image classification, Masked Language Model]
 - both pretraining and finetuning metrics for gMLPs improve as quickly as for Transformers
- Transformers can be more practically advantageous over gMLPs on tasks that require cross-sentence alignment (e.g., by 1.8% on MNLI), even with similar capacity and pretraining perplexity.
- **Overall, our results suggest that self-attention is not a necessary ingredient for scaling up machine learning models**

Inductive bias

“Inductive bias 란 학습자가 지금까지는 만나보지 않았던 상황에서 정확한 예측을 하기 위해 사용하는 추가적인 가정 (additional assumptions) 을 의미한다. ... 성공적인 학습 이후에, 학습자는 훈련동안에는 보이지 않던 예들 까지도 정확한 출력에 가까워지도록 추측한다. 그런 경우에 어떤 추가적인 가정이 없이는 (보이지 않는 상황이 가상의 출력값을 가지기 때문에) 해결될 수 없다. Target function 의 성질에 대해 필요한 가정과 같은 것이 inductive bias 라는 말에 포함된다”[2]

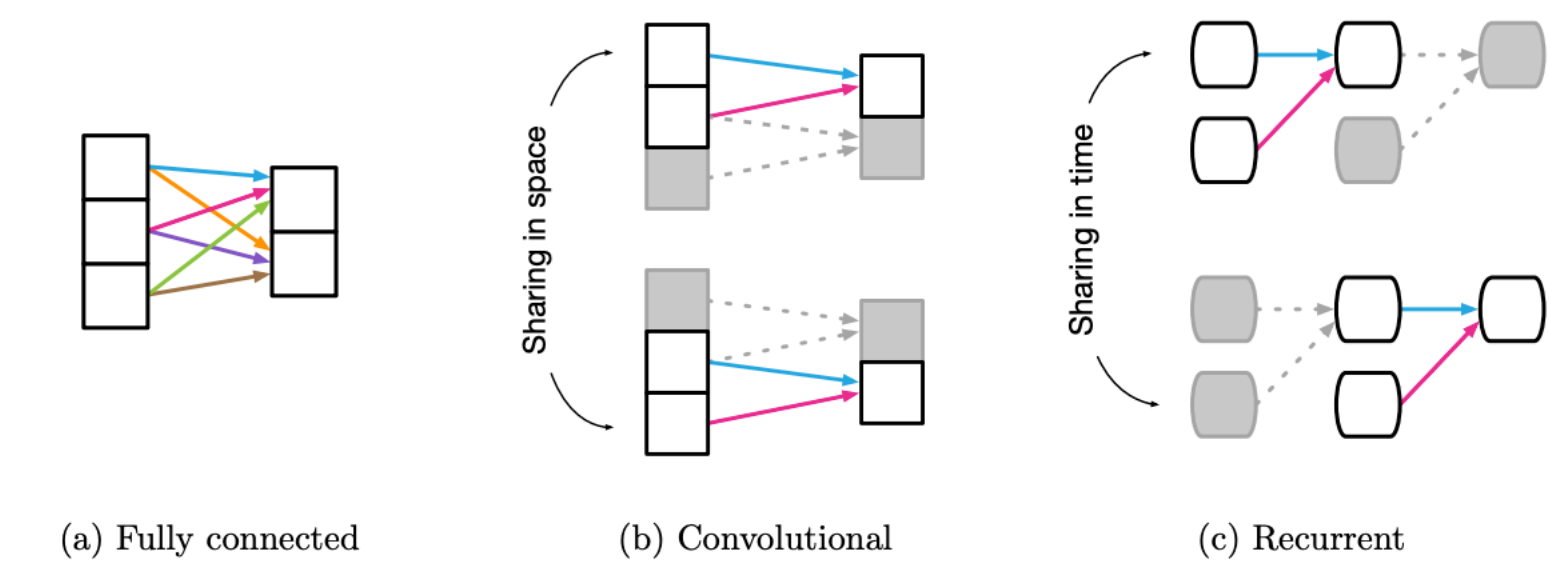


Figure 1: Reuse and sharing in common deep learning building blocks. (a) Fully connected layer, in which all weights are independent, and there is no sharing. (b) Convolutional layer, in which a local kernel function is reused multiple times across the input. Shared weights are indicated by arrows with the same color. (c) Recurrent layer, in which the same function is reused across different processing steps.

Examples

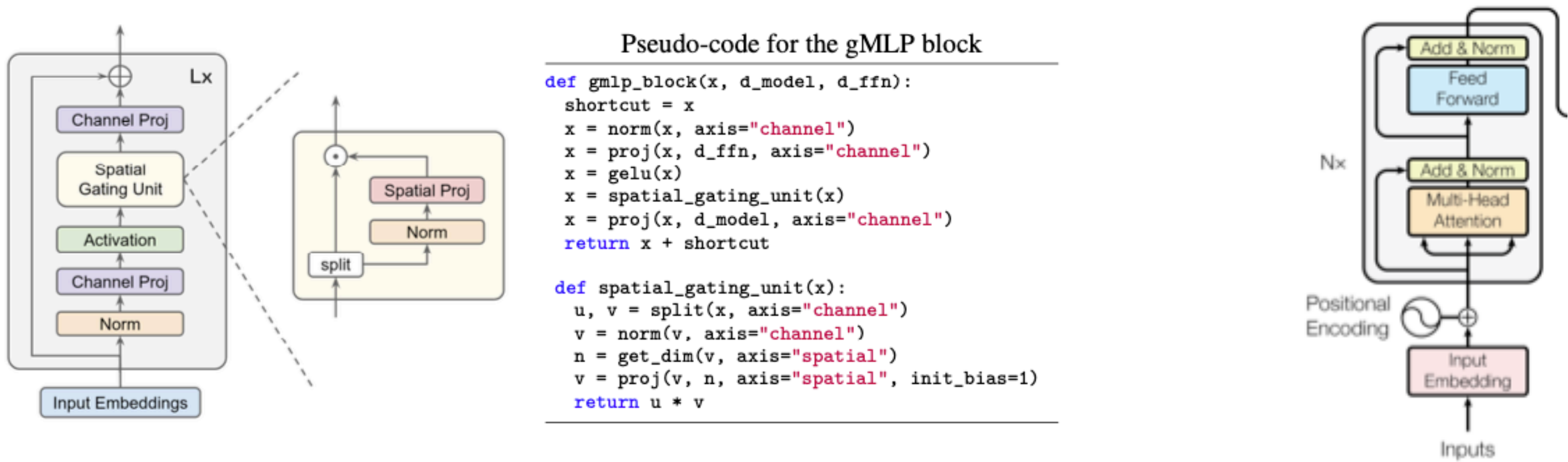
[1] Fully connected layers

- 모든 가중치가 독립적이고, 공유되지 않음 → The implicit relational inductive bias in a fully connected layer is very weak
- 모든 입력 unit은 어떤 출력 유닛의 값도 결정하는 데 interact 될 수 있다.

[2] Convolutional layers

- Conv layer는 Fully connected layer와 다르게, 어떤 중요한 관계적 inductive biases를 부여 함: **Locality** and **Translation invariance**
 - **Locality** reflects that the arguments to the relational rule are those entities in close proximity with one another in the input signal’s coordinate space, isolated from distal entities.
 - **Translation invariance** reflects reuse of the same rule across localities in the input.

2. Model



- Unlike Transformers, gMLPs do not require positional encodings, nor is it necessary to mask out the paddings during NLP finetuning.

<div>$Z = \sigma(XU)$$\tilde{Z} = s(Z)$$Y = \tilde{Z}V$</div>	<ul style="list-style-type: none">• Spatial Gating Unit: a layer which captures spatial interactions<ul style="list-style-type: none">→ our major focuses is therefore to design a good s capable of capturing complex spatial interactions across tokens→ our model <i>does not require position embeddings</i> because such information will be captured in $s(\cdot)$

2.1 Spatial Gating Unit

Equations	Plot
$Z = \sigma(XU)$ $\tilde{Z} = s(Z)$ $Y = \tilde{Z}V$ $f_{W,b}(Z) = WZ + b$ $s(Z) = Z_1 \odot f_{W,b}(Z_2)$	<p>The diagram illustrates the architecture of the Spatial Gating Unit. At the bottom, a blue grid representing input X is multiplied by a purple trapezoidal matrix U to produce a green grid $Z = \sigma(XU)$. A code snippet <code>u, v = x.chunk(2, dim=-1)</code> is shown. The green grid Z is then split into two paths. The left path produces a green grid Z_1. The right path produces a yellow grid Z_2, which is multiplied by a yellow rectangular matrix W to produce a green grid $f_{W,b}(Z_2)$. Finally, Z_1 and $f_{W,b}(Z_2)$ are combined via element-wise multiplication (indicated by a circle with a dot) to produce the final output $s(Z)$, shown as a small green grid.</p>

- For training stability, we find it critical to **initialize W** as near-zero values and b as ones, meaning that $s(\cdot)$ is approximately an identity mapping at the beginning of training.
- the magnitude for each element in Z can be rapidly tuned according to the gating function $f_{w,b}(\cdot)$

3. Image Classification

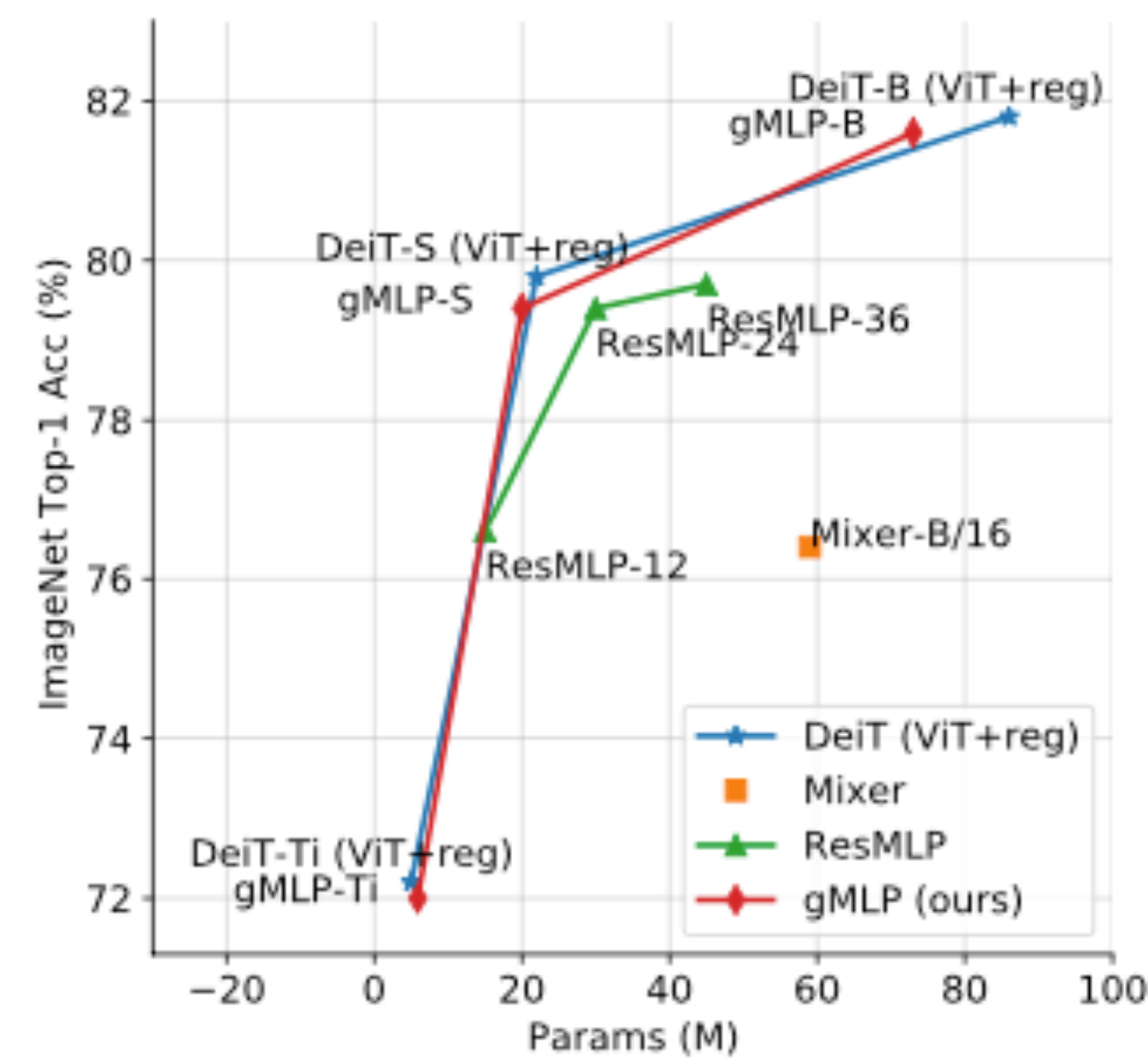


Figure 2: ImageNet accuracy vs model capacity.

- We compare our attention-free models with recent attentive models based on vanilla Transformers, including Vision Transformer (ViT) [7], DeiT [8] (ViT with improved regularization), and several other representative convolutional networks.
- The accuracy-parameter/FLOPs tradeoff of gMLPs surpasses all concurrently proposed MLP-like architectures , which we attribute to the effectiveness of our Spatial Gating Unit

3. Image Classification

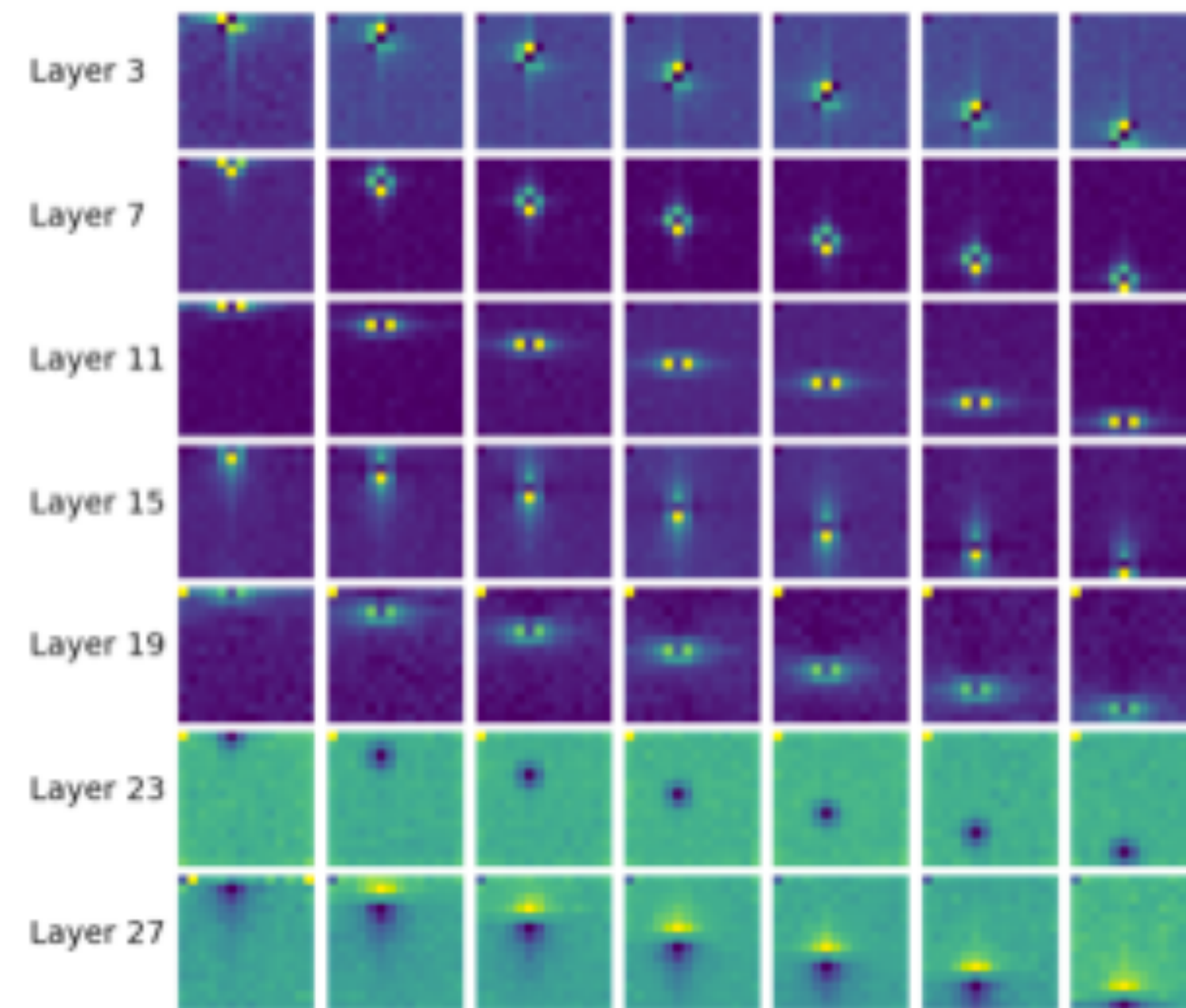


Figure 3: Spatial projection weights in gMLP-B. Each row shows the filters (reshaped into 2D) for a selected set of tokens in the same layer.

- The spatial weights after learning exhibit both locality and spatial invariance. In other words, each spatial projection matrix effectively learns to perform convolution with a data-driven, irregular (non-square) kernel shape.

4. Masked Language Modeling with BERT

- We do not use positional encodings.
- We also find it unnecessary to mask out <pad> tokens in gMLP blocks during finetuning as the model can quickly learn to ignore them.

4.1 Ablation: The Importance of Gating in gMLP for BERT’s Pretraining

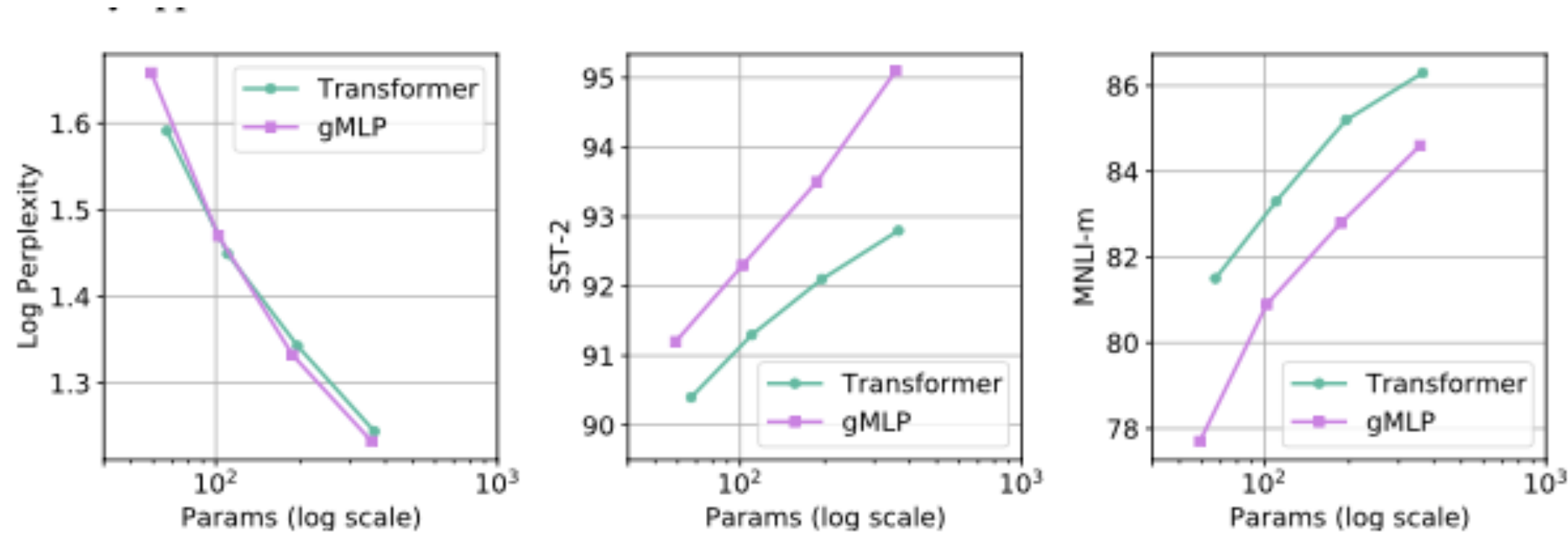
Model	Perplexity	Params (M)
BERT _{base} <small>BERT with a Transformer architecture and learnable absolute position embeddings.</small>	4.37	110
BERT _{base} + rel pos <small>BERT with a Transformer architecture and T5-style learnable relative position biases</small>	4.26	110
BERT _{base} + rel pos - attn	5.64	96
Linear gMLP, $s(Z) = f(Z)$	5.14	92
Additive gMLP, $s(Z) = Z + f(Z)$	4.97	92
Multiplicative gMLP, $s(Z) = Z \odot f(Z)$	4.53	92
Multiplicative, Split gMLP, $s(Z) = Z_1 \odot f(Z_2), Z = Z_1 Z_2$	4.35	102

1. SGU outperforms other variants in perplexity
2. gMLP with SGU also achieves perplexity comparable to Transformer.

4.2 Case Study: The Behavior of gMLP as Model Size Increases

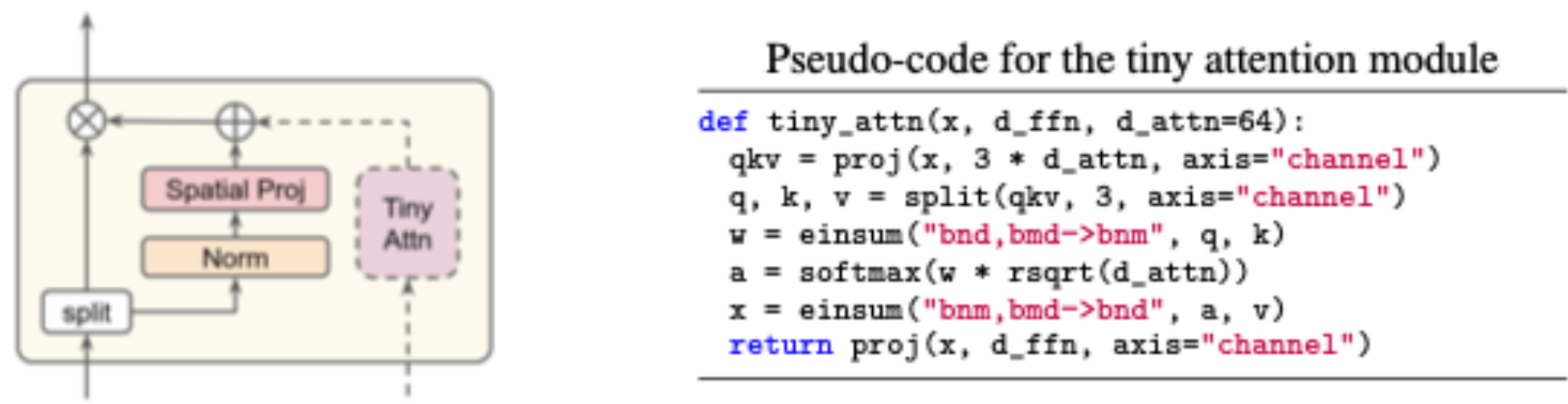
Model	#L	Params (M)	Perplexity	SST-2	MNLI-m
Transformer	6+6	67	4.91	90.4	81.5
gMLP	18	59	5.25	91.2	77.7
Transformer	12+12	110	4.26	91.3	83.3
gMLP	36	102	4.35	92.3	80.9
Transformer	24+24	195	3.83	92.1	85.2
gMLP	72	187	3.79	93.5	82.8
Transformer	48+48	365	3.47	92.8	86.3
gMLP	144	357	3.43	95.1	84.6

- The results above show that a deep enough gMLP is able to match and even **outperform** the perplexity of Transformers with comparable capacity

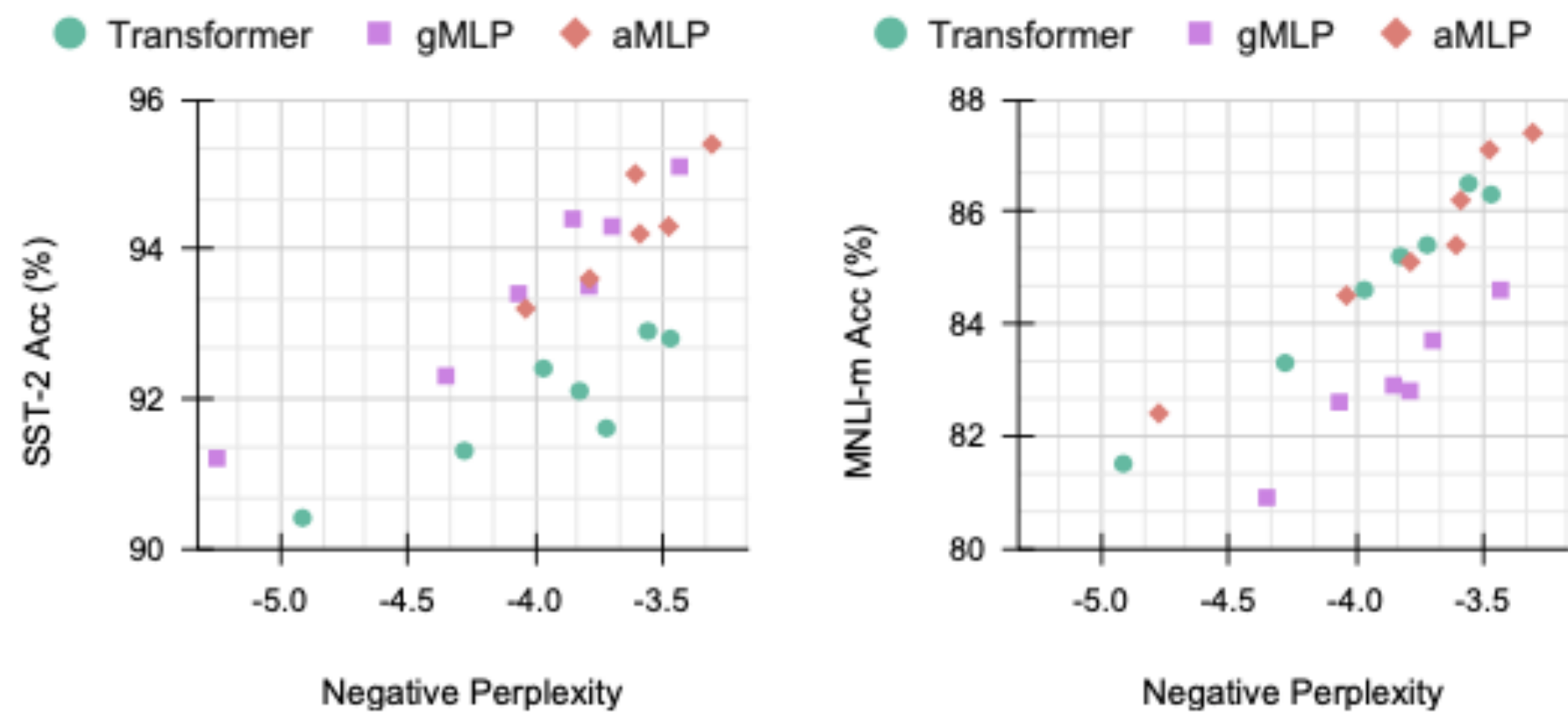


- **Our attention-free model is advantageous** on SST-2 but **worse** on MNLI is particularly informative—the former is a single-sentence task 이진 분류(감성 분석) whereas the latter involves sentence pairs (premise and hypothesis) {전제, 결론} 쌍으로 이루어진 자연어 추리

4.3 Ablation: The Usefulness of Tiny Attention in BERT’s Finetuning



- To isolate the effect of attention, we experiment with a hybrid model where a tiny self-attention block is attached to the gating function of gMLP (Figure 6)



4.4 Main Results for MLM in the BERT Setup

	Perplexity	SST-2	MNLI (m/mm)	SQuAD		Attn Size	Params (M)
				v1.1	v2.0		
BERT _{base} [2]	–	92.7	84.4/-	88.5	76.3	768 (64 × 12)	110
BERT _{base} (ours)	4.17	93.8	85.6/85.7	90.2	78.6	768 (64 × 12)	110
gMLP _{base}	4.28	94.2	83.7/84.1	86.7	70.1	–	130
aMLP _{base}	3.95	93.4	85.9/85.8	90.7	80.9	64	109
BERT _{large} [2]	–	93.7	86.6/-	90.9	81.8	1024 (64 × 16)	336
BERT _{large} (ours)	3.35	94.3	87.0/87.4	92.0	81.0	1024 (64 × 16)	336
gMLP _{large}	3.32	94.8	86.2/86.5	89.5	78.3	–	365
aMLP _{large}	3.19	94.8	88.4/88.4	92.2	85.4	128	316

Table 6: Pretraining perplexities and dev-set results for finetuning. “ours” indicates models trained using our setup. We report accuracies for SST-2 and MNLI, and F1 scores for SQuAD v1.1/2.0.

- our gMLP_{large} achieves 89.5% F1 on SQuAD-v1.1 without any attention or dynamic parameterization mechanism
- our hybrid model aMLP_{large} achieves 4.4% higher F1 than Transformers on the more difficult SQuAD-v2.0 task.

5 Conclusion

- We show that gMLPs, a simple **variant of MLPs with gating**, can be competitive with Transformers in terms of BERT’s pretraining perplexity and ViT’s accuracy.

References

- [1] Liu, Hanxiao, et al. "Pay Attention to MLPs." *arXiv preprint arXiv:2105.08050* (2021).
- [2] <http://www.aistudy.co.kr/neural/bias.htm>
- [3] Battaglia, Peter W., et al. "Relational inductive biases, deep learning, and graph networks." *arXiv preprint arXiv:1806.01261* (2018).