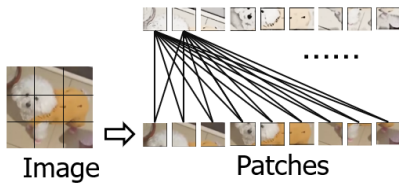


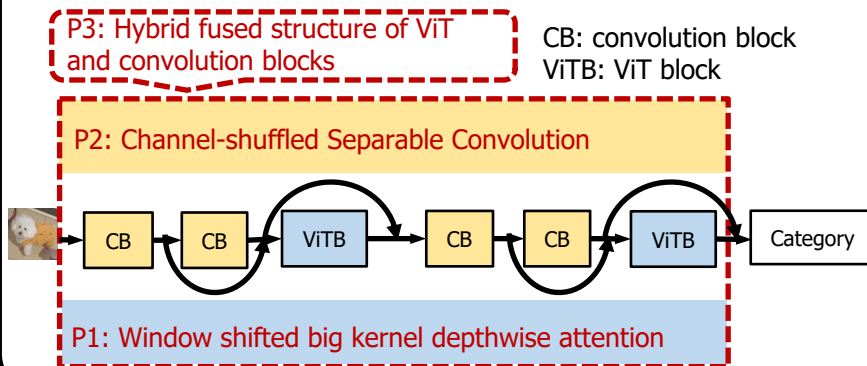
Background



Speeding up MobileViT is needed

Proposed Methods

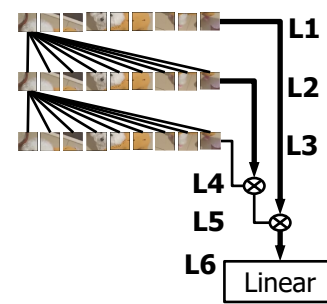
Overall framework



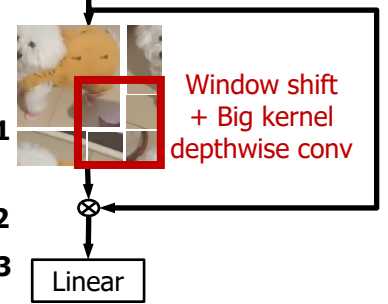
CB: convolution block
ViTB: ViT block

P1: Window Shifted Big kernel Depthwise Attention

Conventional attention

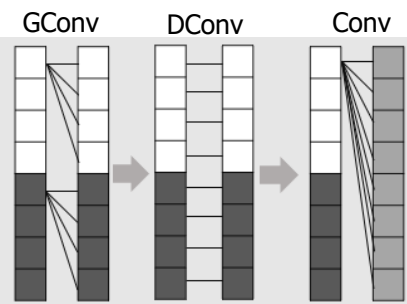


Proposed attention

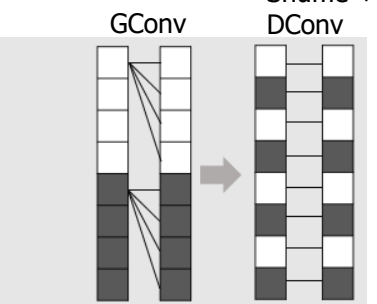


P2: Channel-shuffled Separable Convolution

Conventional CB



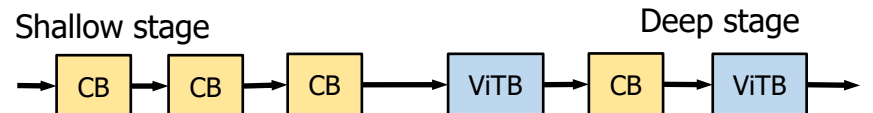
Proposed CB



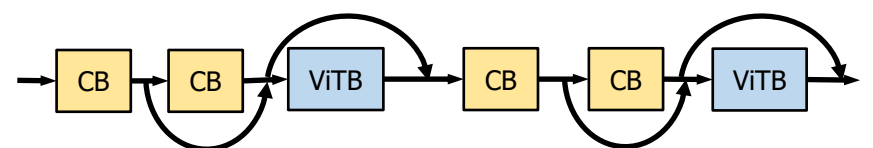
DConv: Depthwise Convolution GConv: Group Convolution

P3: Hybrid fused structure of ViT and convolution blocks

Conventional structure:



Hybrid fused structure:



Experimental Results

FLOPs: floating point operations

Model	FLOPs (M)	Layers	Accuracy (%)	Inference Time on GPU 1080Ti (ms)	Inference Time on CPU I7 6700k (ms)
Original MobileViT	64	63	70.18	5.51	5.67
Light MobileViT	51	59	70.35	5.15	5.19
Conv2Former	45	55	70.16	4.73	4.76
+P1	40	45	69.13	3.98	4.28
+P1+P2	36	38	67.95	3.40	3.64
+P1+P2+P3	45	38	70.92	3.42	3.81

Conclusion

- This work aims to speed up MobileViT while keeping accuracy
- Proposed MobileViT speeds up 38%, while accuracy is 0.74% higher than original MobileViT

