



Token Turing Machines are Efficient Vision Models

Purvish Jajal¹, Nick John Eliopoulos¹, Benjamin Shiue-Hal Chou¹, George K. Thiruvathukal², James C. Davis¹, Yung-Hsiang Lu¹

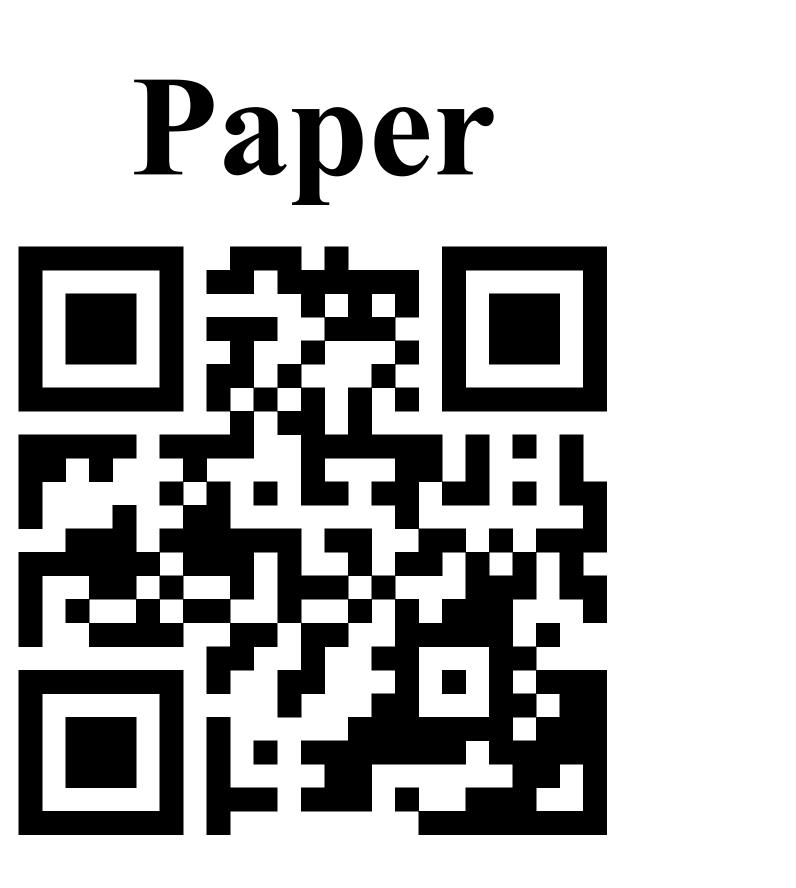


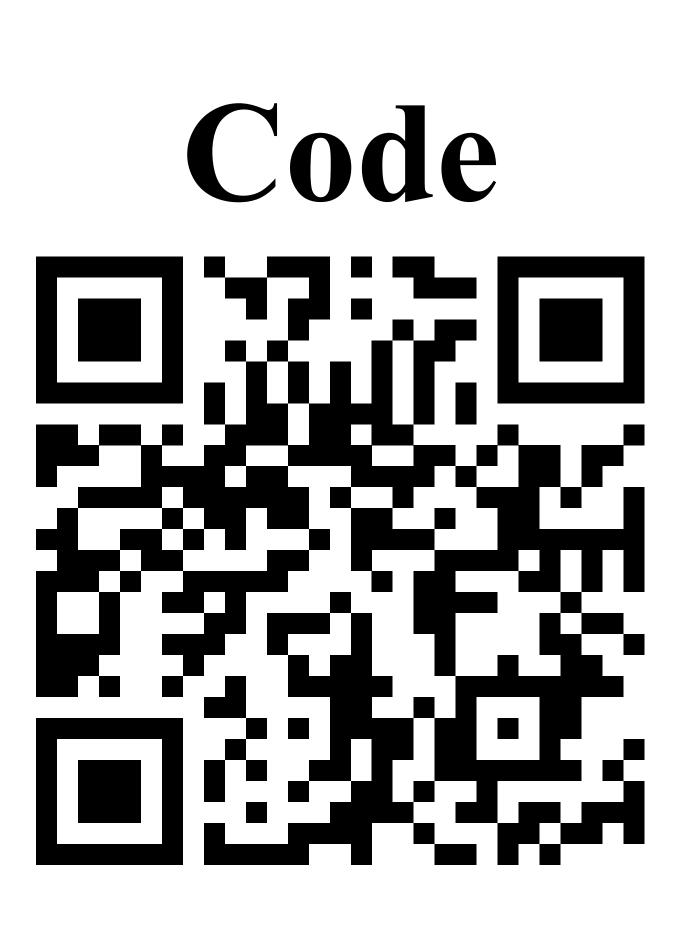
Background

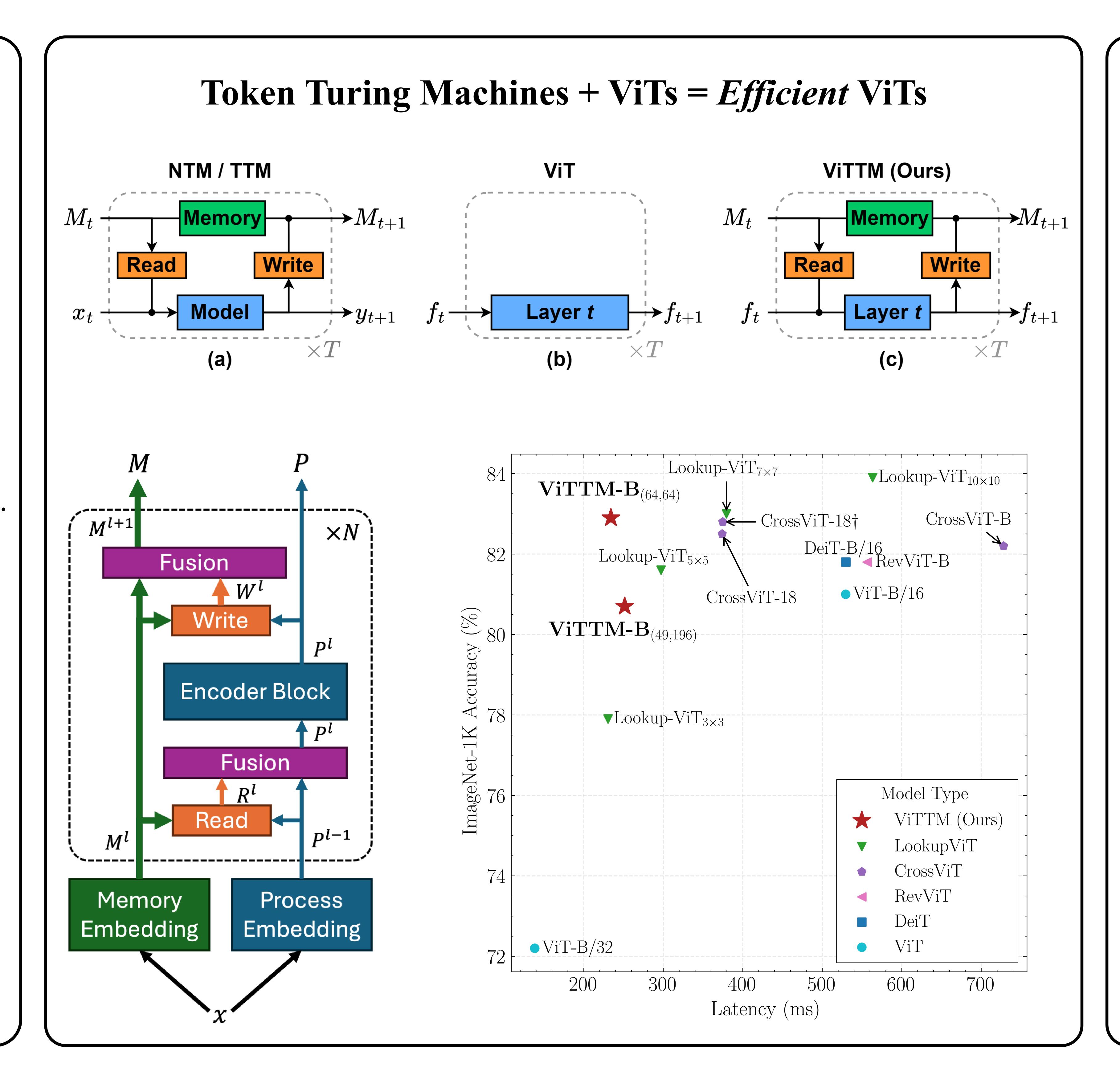
- Vision Transformers (ViTs) have quadratic complexity with respect to input size.
- Computational costs can be reduced by processing fewer tokens but this compromises accuracy.

Contributions

- A novel efficient memory-augmented ViT (ViTTM) with competitive accuracy-latency trade-offs for classification and segmentation.
- A comprehensive ablation exploring various designs and identify key decisions that impact the quality of the architecture.







ImageNet 1K, 56% faster

| Model Class | Model | Params (M) | GFLOPs \downarrow | Latency (ms) | Top-1 (%) |
|-------------------|-----------------------------|------------|----------------------------|---|------------------|
| | ViT-S/16 | 22 | 4.25 | 149.5 | 74.2 |
| | DeiT-S/16 | 22 | 4.25 | 152.0 | 79.8 |
| ViT/DeiT | ViT-B/32 | 88 | 4.37 | 138.3 | 72.2 |
| | ViT-B/16 | 87 | 16.87 | 529.5 | 81.0 |
| | DeiT-B/16 | 87 | 16.87 | 529.7 | 81.8 |
| | CrossViT-S | 27 | 5.63 | 149.5 152.0 138.3 529.5 | 81.0 |
| | CrossViT-15 | 28 | 5.81 | 249.1 | 82.3 |
| | CrossViT-15† | 28 | 6.13 | 252.3 | 81.5 |
| | CrossViT-B | 105 | 21.22 | 728.1 | 82.2 |
| | CrossViT-18 | 43 | 9.05 | 374.1 | 82.5 |
| Two-Stream | CrossViT-18† | 44 | 9.50 | 149.5 152.0 138.3 529.5 529.7 235.7 249.1 252.3 728.1 374.1 378.2 556.5 230.5 297.2 379.5 563.4 172.7 161.6 138.2 125.9 115.2 103.3 94.0 551.1 515.7 440.6 402.0 367.0 330.2 296.4 | 82.8 |
| Two-Stream | Rev-ViT-B | 86 | 17.49 | 556.5 | 81.8 |
| | LookupViT $_{3\times3}$ | 90 | 5.26 | 230.5 | 77.9 |
| | LookupViT $_{5\times5}$ | 90 | 6.94 | 297.2 | 81.6 |
| | LookupViT _{7×7} | 90 | 9.45 | 379.5 | 83.0 |
| | $LookupViT_{10\times 10}$ | 90 | 14.80 | 563.4 | 83.9 |
| | $ViT-S/16_{(r=2)}$ | 22 | 4.31 | 172.7 | 74.0 |
| | $ViT-S/16_{(r=4)}$ | 22 | 4.02 | 161.6 | 73.8 |
| | $ViT-S/16_{(r=8)}$ | 22 | 3.41 | 138.2 | 73.1 |
| | $ViT-S/16_{(r=10)}$ | 22 | 3.14 | 125.9 | 72.5 |
| | $ViT-S/16_{(r=12)}$ | 22 | 2.85 | 115.2 | 71.6 |
| | $ViT-S/16_{(r=14)}$ | 22 | 2.57 | 149.5 152.0 138.3 529.5 529.7 235.7 249.1 252.3 728.1 374.1 378.2 556.5 230.5 297.2 379.5 563.4 172.7 161.6 138.2 125.9 115.2 103.3 94.0 551.1 515.7 440.6 402.0 367.0 330.2 296.4 | 70.4 |
| Token Merging [1] | $ViT-S/16_{(r=16)}$ | 22 | 2.30 | | 68.1 |
| Token Merging [1] | $ViT-B/16_{(r=2)}$ | 86 | 16.46 | 551.1 | 81.0 |
| | $ViT-B/16_{(r=4)}$ | 86 | 15.34 | 515.7 | 80.9 |
| | $ViT-B/16_{(r=8)}$ | 86 | 13.12 | 440.6 | 80.4 |
| | $ViT-B/16_{(r=10)}$ | 86 | 12.02 | 402.0 | 80.1 |
| | $ViT-B/16_{(r=12)}$ | 86 | 10.93 | 367.0 | 79.6 |
| | $ViT-B/16_{(r=14)}$ | 86 | 9.84 | 330.2 | 78.9 |
| | $ViT-B/16_{(r=16)}$ | 86 | 8.78 | 296.4 | 77.6 |
| Ours | ViTTM-S _(64,64) | 33 | 1.84 | 77.7 | 79.2 |
| | $ViTTM-B_{(64,64)}$ | 127 | 7.08 | 234.1 | 82.9 |
| | ViTTM-B _(49,196) | 125 | 7.10 | 251.5 | 80.9 |

Semantic Segmentation, +94% FPS

| Model (M) | FPS (A100) † | FPS (A30) † | mIoU † |
|------------------------------|--------------|-------------|--------|
| ViT-B/16 ₃₈₄ [35] | 23.8 | 13.8 | 48.06 |
| ViT-B/16 ₂₂₄ | 23.8 | 13.8 | 45.65 |
| $ViTTM-B_{(64,64)}$ | (+37%) 32.5 | (+94%) 26.8 | 45.17 |
| ViTTM-B _(49,196) | (+38%) 32.8 | (+94%) 26.7 | 43.60 |

Measured on the ADE20K dataset using Segmenter.