

•Ú{Ç^ºe‡0 DynamicViT: Efficient Vision Transformers with Dynamic Token Sparsification0 cĐQúN†N yÍ
Transformerŷ ViTŷ N-šØeHW0Rjg•Q—OYNärL0 NâN f/^ºe‡v„N;‰•Q...[1Tœ•!s.ŷ

1. h8_Ã` `ó

- **‰•Â[ß**ŷ ‰•Æ‰•ÉTransformerv„g ~È~„mKNÂO••VNŽN \ •èR Oá`o`İg Y'v„NärLŷ tokensŷ ŷ QvNÖQ—
- **e!İÖ**ŷ cĐQúR`` NärLz u•S hFgŹŷ • •Ç•{“~š~„mKj!WWŷ prediction moduleŷ • \BRjg•Q—OYNärLŷ N

2. Qs•.b€g/

- **R`` NärLRjg•**ŷ
 - ~„mKj!WWWúNŽ_SRMry_•O0ç;kİN*NärLv„'Í‰•`'R epŷ u b Nœ•ŪR6Q³{Vc©x ŷ binary decision mask
 - NärLRjg•f/\Bk!S v„ŷ hierarchicalŷ ŷ • keXŽR Rjg•kÔOç0
- **zİR0ziO S **ŷ
 - O•u(Gumbel-Softmax‰•ãQ³—^Sİ_@v„'Çh7•î~0
 - cĐQúlèa Rç©x {Vueŷ attention maskingŷ ŷ • •Ç—;e-^«Rjg•NärLN QvNÖNärLv„N¤N'ŷ [žs°Sİ^v^Lç;{—
- **ç-~Âvİh **ŷ
 - ~ÓT R {c_Y1ŷ cross-entropyŷ 0 „„™•c_Y1ŷ self-distillationŷ 0 KLec^!c_Y1ŷ KL divergenceŷ TœkÔO

3. [žšœ~Ógœ

- **`'€ý**ŷ
 - W(ImageNetN ŷ Rjg•66%v„NärLT ŷ FLOPsQİ\ 31%~37%ŷ T T 'İcĐSG40%NâN ŷ |¾^!N –MN •...•Ç
 - W(DeiT0 LV-ViT{lŷ!WçN šœçÂN†g eH`'0
- **O Rç**ŷ
 - vøkÔ~Óg„S Rjg•ŷ Y,CNNv„l`S ŷ ŷ R`` Rjg•fôpum;N xINöSĚY}0
 - Sİ‰•ÆS f>y:j!Wç€ý• ke€Zq&NŽVpPİN-v„Qs•.S:Wß0

4. R e°p¹

- **R`` `!***ŷ Rjg•Q³{VWúNŽ•“QeQ...[1ŷ € —^Vú[šj!_ 0
- **šØeH`!***ŷ • •Çlèa Rç©x [žs°xINöSĚY}v„z u•ç;{—0
- **• u(`!***ŷ Sİ^"u(NŽY yÍ‰•Æ‰•ÉTransformergŹg„ŷ Y,ViT0 DeiT0 LV-ViTŷ 0

5. ^"u(N \Ug

- • u(NŽVpPİR {NûRjŷ g*geSİbİ\U•ó‰•Æ~'R {Tœ[Æ–Æ~„mKNûRj0
- Năx _ n•ŷ <https://github.com/raoyongming/DynamicViT0>

`;~Ó

DynamicViT• •ÇR`` Rjg•Q—OYNärLŷ f>„WcĐSGN†‰•Æ‰•ÉTransformerv„eHs‡ŷ T eöOÝc N†šØ|¾^!ŷ