PAUMER: Patch Pausing Transformer for Semantic Segmentation

Evann Courdier^{*,a,b}, Prabhu Teja S^{*,a,b} François Fleuret^{a, b, c}

*Equal Contribution, aldiap Research Institute, ^bEPFL, ^cUniversity of Geneva, Switzerland, Contact: evann.courdier@idiap.ch, prabhu.teja@idiap.ch







One training to pause them all!: Training to handle various run-time requirements

- For each batch, sample a layer $l \sim \mathcal{U}\{L^{\mathsf{lo}}, L^{\mathsf{hi}}\}$ and a patch pausing proportion $au_l \sim \mathcal{U}[\tau_l^{\mathrm{lo}}, \tau_l^{\mathrm{hi}}].$
- We employ an auxiliary decoder after the operations of layer l, and pause τ_l patches with the lowest entropy.
- Training both main and aux decoder with cross entropy.

Reduce computation for 'easier' parts of the input image.

Not all tokens are equal



(a) Input image.

(b) The first layer at which the patch is correctly classified. Darker is lower.

Figure: How long do we have to process each part of the image for correct predictions with a ViT-Ti backbone? It is apparent that not all patches of the images need a homogenous amount of processing i.e., not all patches need the full network and can be reliably decoded after a few layers itself. We exploit this to build patch pausing for efficient transformers for segmentation.

Pause processing tokens which have been processed 'enough'

Pausing patches at multiple layers









(b.2) Entropy Computed



(c.2) Entropy Computed







Figure: Our method progressively stops processing patches after they reach a low enough prediction entropy. (a) is the input image. (b.1, c.1) show the patches that are stopped from being processed after the 3rd and 5th layers. (b.2) and (c.2) show entropy computed from the auxiliary decoders that is used to pause patches. The network automatically pauses easy parts of the image while allocating more computation to the parts that correspond to boundaries, and to smaller and rarer classes in (c.3).



Figure: We modify Segmenter[1] to enable pausing of patches, and feed them directly to the decoder. We add a simple auxiliary 1 imes 1 convolution decoder, and use the predicted posterior entropy $\mathbb{H}[\hat{Y}|X]$ of each component of X to reorder the feature representation X. A portion τ of this feature representation is paused and fed to the decoder directly. The rest of the features (of size N' < N) are processed further.

Posterior entropy correlates with correctness



Traversing the pareto front of throughput-mloU tradeoff



Figure: Comparing PAUMER to baseline and random pausing on Cityscapes validation set.

Figure: Comparing PAUMER's speed-accuracy tradeoff with other architectures.

[1] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In ICCV, 2021.