

An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

Koorye

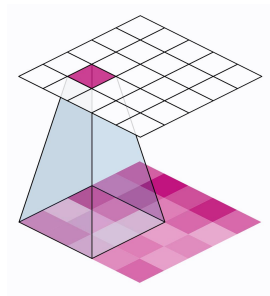
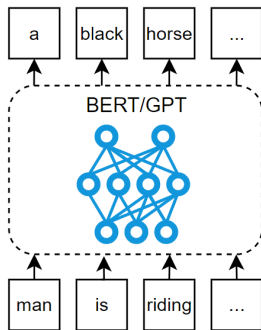
University of Electronic and Science Technology of China

13th May 2024

Outline

- 1 Research Aim
- 2 Method
- 3 Results
- 4 Conclusion & Discussion

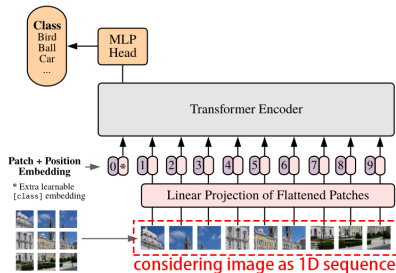
The GAP between Language Models and Vision Models:



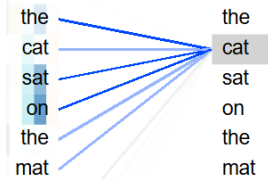
(a) Language Model receives 1D sequence (b) Vision Model receives 2D image pixels

Could the approach used in Language Models inspire a new era for visual tasks?

Method



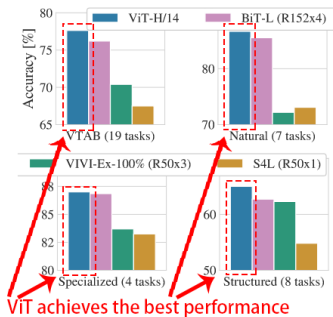
(c) Vision Transformer Model



(d) Attention Module

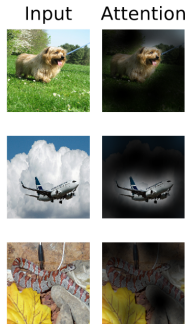
- Considering images as **1D sequences** of patches, like text sentences.
- Model the **relationship** between patches, like the process of human understanding.

Results



ViT achieves the best performance

(e) Model Accuracy Comparison



(f) Attention Visualization

- The **accuracy** of model reaches a new level with smaller model **size**.
- Attention modules capture **visually rich areas** of information.

Conclusion:

- Inspired by Language Models, processing images into sequences similar to text.
- Surpassing traditional models in accuracy and being more compact.

Discussion:

- Need more training data and larger models to maximize its potential.
- Unknown performances in more visual tasks.