DeFormer: Decomposing Pre-trained Transformers for Faster Question Answering











Qingqing Cao



Niranjan Balasubramanian

Resource Challenges in Large Transformers



How to Reduce Runtime Latency and Memory?



Why are Transformers Slow for QA?



How Can We Reduce the Context Bottleneck?



How critical is question dependence?



High question dependence

How critical is question dependence?



Less question dependence

Question Dependence in Different Layers



Decomposing the Dependence in Lower Layers



Fine-tuning DeFormer with Auxiliary Supervision



Evaluation

• How effective is DeFormer in latency, memory and accuracy?

• What is the impact of auxiliary supervision?

• What is the efficiency tradeoff at different layers?

DeFormer is Faster and Memory Efficient



DeFormer Gains Speedup on Diverse Devices



DeFormer (Mostly) Retains Original Performance



BERT-base DeFormer-BERT-base

A Large DeFormer is Faster and More Accurate than a Smaller Transformer



Evaluation

• How effective is DeFormer in latency, memory and accuracy?

• What is the impact of auxiliary supervision?

• What is the efficiency tradeoff at different layers?

LRS and KD both Provide Benefits in DeFormer



Evaluation

• How effective is DeFormer in latency, memory and accuracy?

• What is the impact of auxiliary supervision?

• What is the efficiency tradeoff at different layers?

Which Layer to Decompose?



Takeaways



DeFormer uses a simple decomposition technique that significantly speeds up inference, reduces memory while retaining most of the model's accuracy



DeFormer leverages distillation along with auxiliary supervision to reduce performance gap to the original Transformer



DeFormer remains largely the same as the Transformer, allowing reusing original model weights without repeating pre-training