

Linformer [Wang+, 20]

✂ TransformerのSelf-Attention

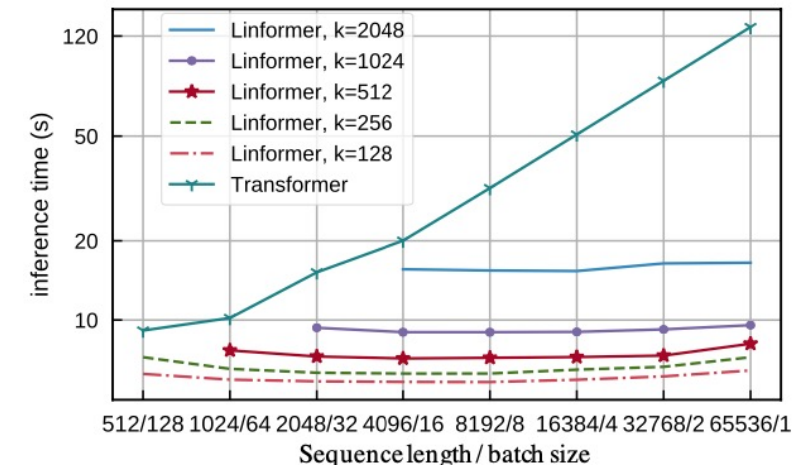
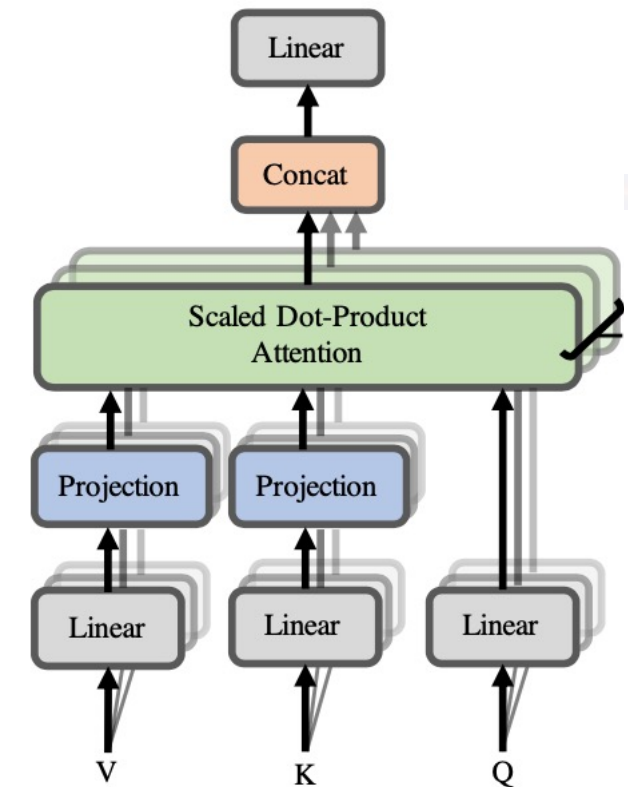
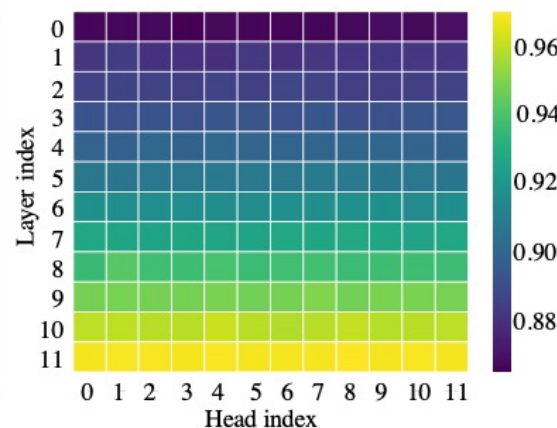
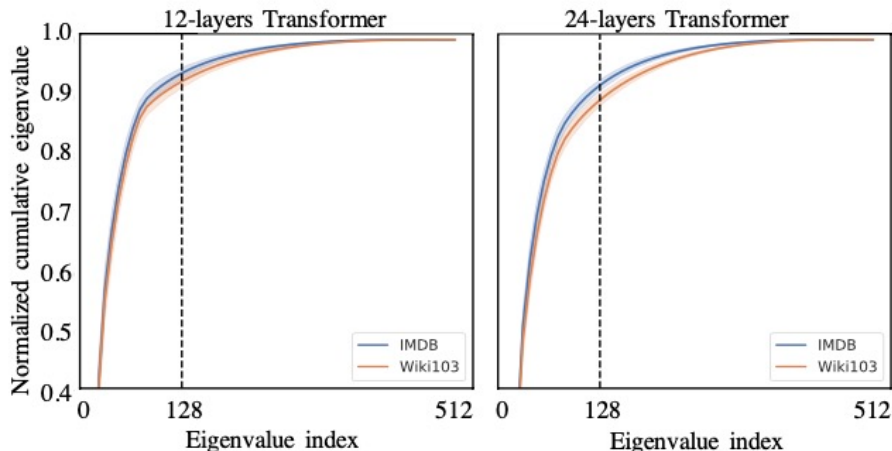
✂ 計算量 $O(n^2)$

✂ KeyとValueに対しprojection layerを導入

✂ 特異値分解を用いて低ランク行列に近似

✂ 行列の情報の殆どが
最初の数個の最大特異値から回収可

✂ 計算量 $O(n)$



UNiT [Hu(Facebook)+, 21]

✂ 複数のタスクを単一モデルで学習するTransformer

✂ 物体検知、自然言語理解、VQA等

✂ Image Encoder

✂ CNNとTransformer Encoder

✂ Text Encoder

✂ BERT

✂ Decoder

✂ DETR

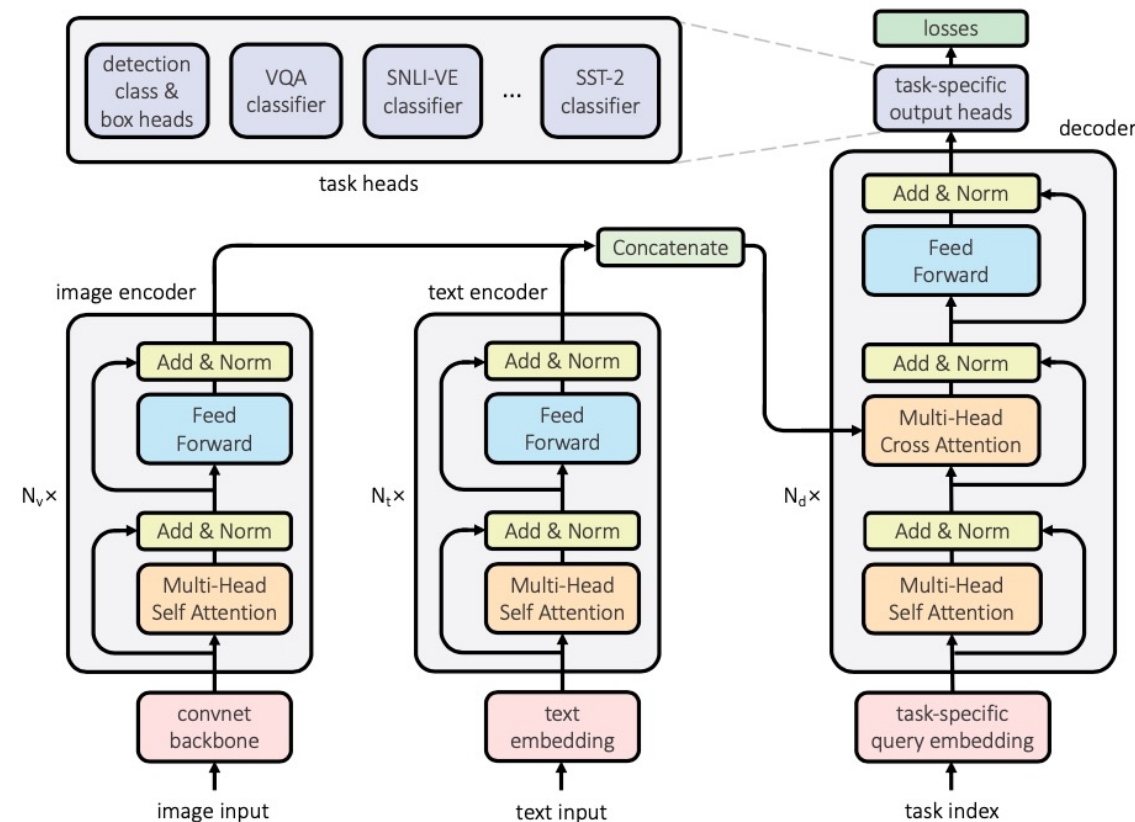
✂ タスクによって必要なEncoderを使用

✂ V&Lタスクでは各Encoderの出力をconcat

✂ Task-specific output heads

✂ 各タスクがそれぞれprediction headを所持

✂ すべてのタスクを c_t クラス間の分類タスクとしてモデル化



Perceiver [Jaegle+, 21]

✂ Transformerを何十万もの入力に対応

- ✂ ImageNet: 畳み込みせず50000ピクセルを処理

✂ Cross Attention

- ✂ 高次元のByte array(pixel arrayなど) を低次元のLatent array(サイズはハイパーパラメータ)へ投影してから複数のTransformerで処理

✂ Latent Transformerの入力が低次元なため、より深いネットワーク構成が可能

- ✂ 重みを共有し、パラメータ数を減少
- ✂ 複数のbyte-attend層
 - ✂ Latent arrayが必要に応じて入力画像から情報を反復的に抽出可能

