

SwinIR

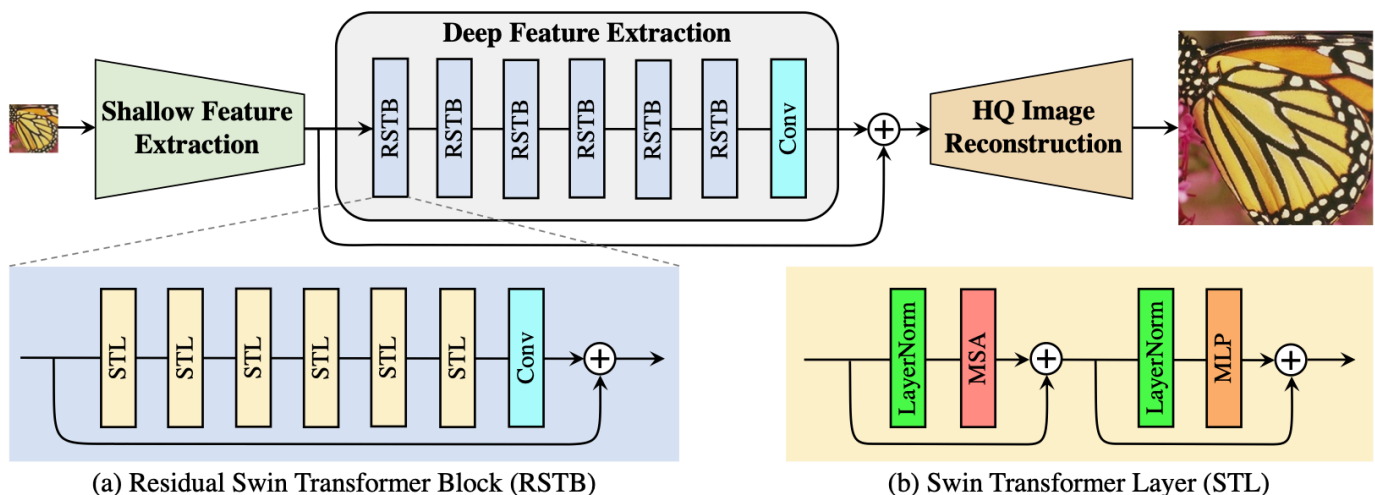
論文ソース

- [SwinIR: Image Restoration Using Swin Transformer](#)

概要

- SOTA達成(当時のSOTAはCNNベース)
- Super ResolutionをCNNベースでやってるけどSwinTransformer使ったほうが精度いいよというお話
- 3層からなるアーキテクチャ
 - shallow feature extraction
 - deep feature extraction
 - high-quality(HQ) image reconstruction

SwinIRの構造



Shallow Feature Extraction

- 1層の 3×3 convoluntinal layerで畳み込みしてチャンネル数変更してるだけ
- 入力: low-quality input $I_{LQ} \in \mathbb{R}^{H \times W \times C_{in}}$
- 出力: shallow feature $F_0 \in \mathbb{R}^{H \times W \times C}$

- $F_0 = H_{SF}(I_{LQ})$

Deep Feature Extraction

- 入力 F_0 で出力 $F_{DF} \in \mathbb{R}^{H \times W \times C}$ で形状不変
 - $F_{DF} = H_{DF}(F_0)$
- H_{DF} : deep feature extraction module
 - K 個の residual Swin Transformer blocks(RSTB)と最後に1層のconvolutionから成る
 - 数式で書くと以下
 - $F_i = H_{RSTB_i}(F_{i-1}), \quad i = 1, 2, \dots, K,$
 - $F_{DF} = H_{CONV}(F_K)$

RSTB

- L 個の Swin Transformer layers(STL)と最後の1層のconvolutionから成る
 - 最初と最後で残差接続
- $F_{i,j} = H_{STL_{i,j}}(F_{i,j-1}), \quad j = 1, 2, \dots, L,$
- $F_{i,out} = H_{CONV_i}(F_{i,L}) + F_{i,0}$

STL

- Swin Transformer layerと同じ
- $M \times M$ のウィンドウに分割するので入力が $H \times W \times C$ から $\frac{HW}{M^2} \times M^2 \times C$ にreshape
- $\frac{HW}{M^2}$ 個のwindowごとにMulti head Self Attention
 - $Q = XP_Q, K = XP_K, V = XP_V$
 - $\text{Attention}(Q, K, V) = \text{SoftMax}\left(\frac{QK^T}{\sqrt{d}} + B\right)V$
- 全体としてはSwin Transformerと同じくLayer Normalization -> Multi head Self Attention -> residual -> LN -> MLP(2層GELU) -> residual
- MSAはSW-MSA->W-MSAを交互に繰り返す
 - ソースコード見てないけどL偶数になっているはず

HQ Image Reconstruction

- shallow featureとdeep featureの残差接続が入力
 - $I_{RHQ} = H_{REC}(F_0 + F_{DF})$
- sub-pixel convolution layerを使って実装
 - ESPCN
 - <https://arxiv.org/pdf/1609.05158.pdf>