

StyleCLIP [Patashnik+ (Adobe Research), CVPR21]

✂ StyleGANとCLIPを組み合わせた、テキストによる画像操作を行うモデル

✂ Latent Optimization : $\arg \min D_{\text{CLIP}}(G(w), t) + \lambda_{L2} \|w - w_s\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$

✂ D_{CLIP} : CLIP埋め込み空間における、生成画像とテキストのコサイン距離

✂ $L2$: 元の w との距離、入力画像との乖離を防ぐ

✂ L_{ID} : 入力・生成画像に対して生成した埋め込みのコサイン類似度 $\mathcal{L}_{\text{ID}}(w) = 1 - \langle R(G(w_s)), R(G(w)) \rangle$

✂ R はArcFace、人物の乖離を制御

✂ Latent Mapper

✂ w は3つのmapperに入力

✂ coarse, medium, fine

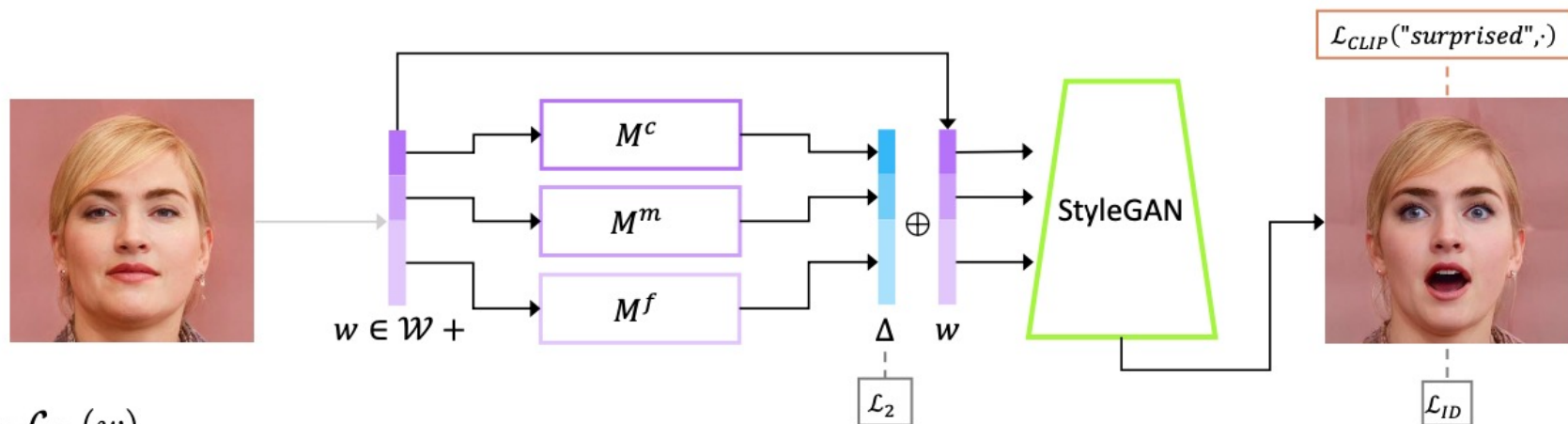
$$w = (w_c, w_m, w_f)$$

$$M_t(w) = (M_t^c(w_c), M_t^m(w_m), M_t^f(w_f))$$

✂ 損失関数

$$\mathcal{L}(w) = \mathcal{L}_{\text{CLIP}}(w) + \lambda_{L2} \|M_t(w)\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$$

$$\mathcal{L}_{\text{CLIP}}(w) = D_{\text{CLIP}}(G(w + M_t(w)), t)$$



✖ 未来の出来事にキャプションを付けるタスク

✖ 次のイベントの畳み込み特徴量を予測し、その特徴量に基づいてキャプションを生成

✖ Temporal Feature Predictor

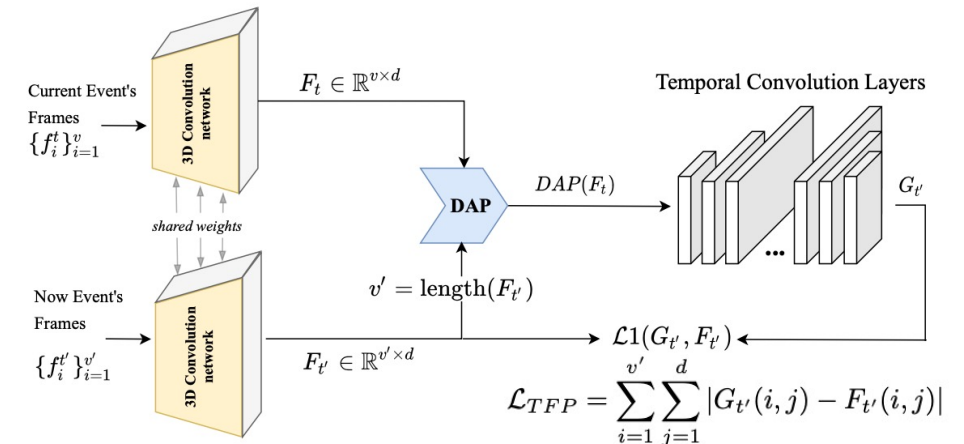
✖ t番目を元にt'番目のイベントの画像特徴量を予測

✖ $G_{t'} = \text{Conv}(\text{DAP}(F_t))$, where $G \in \mathbb{R}^{v' \times d}$

✖ v: フレーム数 d: 次元数

$G_{t'}^{\text{final}} = \lambda \cdot \text{AVG}(F_t) \oplus (1 - \lambda) \cdot \text{AVG}(G_{t'})$

✖ AVG: 時間次元でのAvg Pooling \oplus : 要素ごとの和

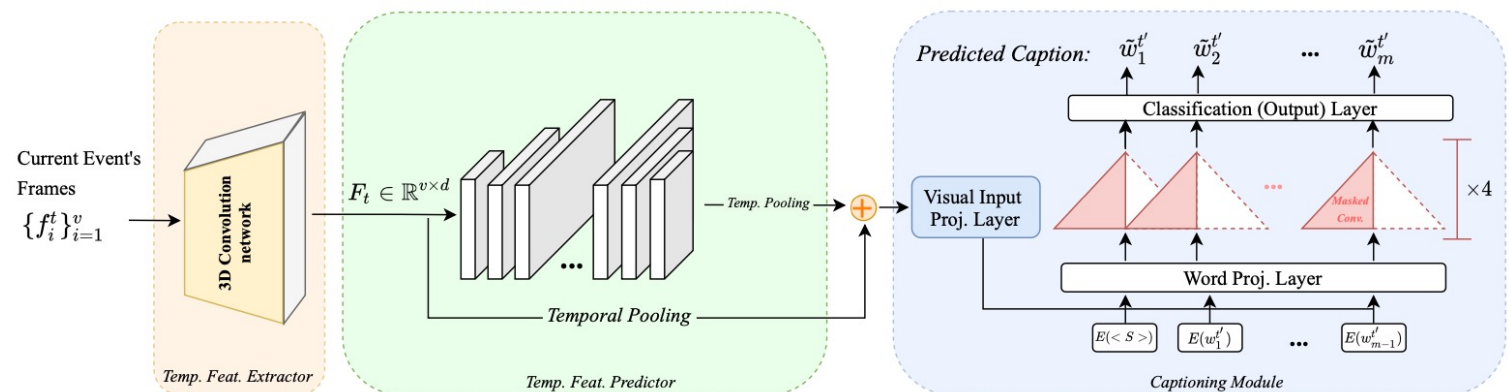


✖ Captioning Module

✖ m-1番目までの単語とG_{t'}からm番目の単語を予測

✖ classification layerの上にsoftmaxを使用

✖ Lossはcross-entropy



HAMLET [Islam+ (Univ. of Virginia), IROS20]

- ✖ human activity recognition (HAR)において、テストした全データセットでベースラインを上回る
 - ✖ ユニモーダルなデータから特徴量を抽出し、それらを分離・融合するマルチモーダルなAttentionメカニズム
 - ✖ HAMLET: Hierarchical Multimodal Self-attention based HAR

- ✖ Input: X^m ($B \times S^m \times E^m$)

- ✖ B : バッチサイズ、 S^m : モダリティ m のセグメント数、 E^m : 特徴次元 ($channel(C^m) \times height(H^m) \times width(W^m)$)

- ✖ Unimodal Self-Attention (UAT)

- ✖ H^m : Temporal Feature Encoderの出力

$$\begin{aligned} Q_i^m &= H^m W_i^{Q,m} & head_i^m &= Attn(Q_i^m, K_i^m, V_i^m) & F_m &= \sum_{s \in S^m} F_{m,s}^a \\ K_i^m &= H^m W_i^{K,m} & F_m^a &= [head_1^m; \dots; head_h^m] W^{O,m} \\ V_i^m &= H^m W_i^{V,m} \end{aligned}$$

- ✖ Multimodal Attention based Feature Fusion (MAT)

- ✖ F^{Gu} にmulti-head self-attention、出力は F^{Ga}

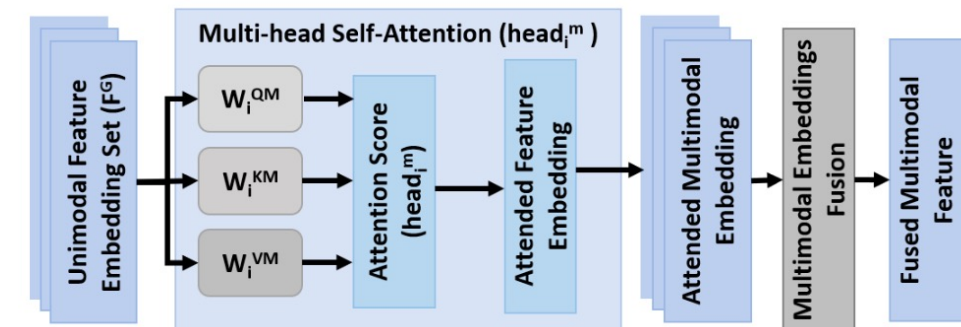
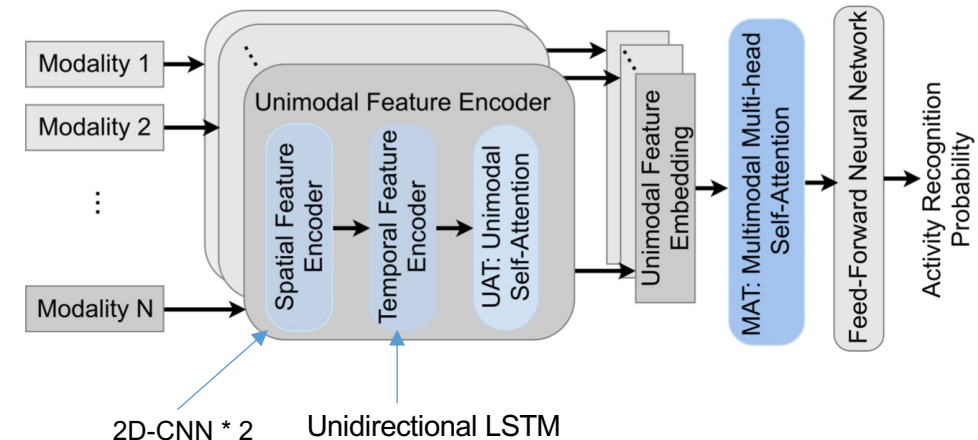
$$\text{✖ } F^{Gu} = (F_1, F_2, \dots, F_M) \text{ } M \text{は全モダリティ数、順不同}$$

- ✖ F^{Ga} から2つの方法で F^G を算出 (CONCATの方が高性能)

- ✖ MAT-SUM
$$F^G = \sum_{m=1}^M F_m^{Ga}$$

- ✖ MAT-CONCAT
$$F^G = [F_1^{Ga}; F_2^{Ga}; \dots; F_M^{Ga}]$$

- ✖ F^G は最終的に全結合層へ $loss(y, \hat{y}) = \frac{1}{B} \sum_{i=1}^B y_i \log \hat{y}_i$
 - ✖ Lossはcross-entropy



VLT [Ding+ (Nanyang Technological Univ), ICCV21]

✖ Attention networkによるReferring Segmentation

✖ 参照画像の助けを借りて参照文を理解するQuery Generation/Balance Moduleを提案

✖ Query Generation Module (QGM)

✖ 画像特徴量 F_{vq} 、言語特徴量 F_t を線形投射

$$\otimes f_{vqn} \in R^{1 \times (HW)}, n = 1, 2, \dots, N_q$$

$$\otimes f_{ti} \in R^{1 \times C}, i = 1, 2, \dots, N_l$$

$$\otimes a_{ni} = \sigma(f_{vqn} W_v) \sigma(f_{ti} W_a)^T$$

$$F_{qn} = A_n \sigma(F_t W_t)$$

✖ Query Balance Module (QBM)

✖ 各 C_{qn} は、クエリ F_{qn} が予測されたコンテキストにどれだけ適合するかを示す

✖ Mask Decoder

✖ 3つの 3×3 convの後、 1×1 convでマスクを出力

✖ 損失関数はマスクのBinary Cross Entropy

