

ClawCraneNet [Liang+ (Univ. of Technology Sydney), 21]

✖ Text-based Video Segmentation taskに焦点

- ✖ 各フレームで全オブジェクトのマスクから言語情報で参照されるものを選択
- ✖ Referring Youtube - VOS challengeで1位

✖ Object features $v^j = \text{MLP}(\text{Max}(F_v \odot o^j))$

✖ j は候補物体(N_v 個)、 F_v は画像全体の特徴量

✖ Positional Relation Module

$$p_i = (x_{min}^i, y_{min}^i, x_{max}^i, y_{max}^i, x_c^i, y_c^i, w_i, h_i, r_i^x, r_i^y) \quad \mathcal{V}_i = v_i + W_p(p_i)$$

✖ o^i の最小外接boxの左上、右下、中心の座標、幅、高さ

✖ r_i^x と r_i^y は、x, y軸の相対位置のインデックス

✖ Temporal Relation Module $\mathcal{S}_s(x_i, x_j) = \mathcal{S}_c(\mathcal{V}_i^r, \mathcal{V}_j^r) + \alpha * U(x_i, x_j)$

✖ 隣接フレームからの2オブジェクト x_i, x_j

✖ \mathcal{S}_c : コサイン類似度、 U : IoUの重なり

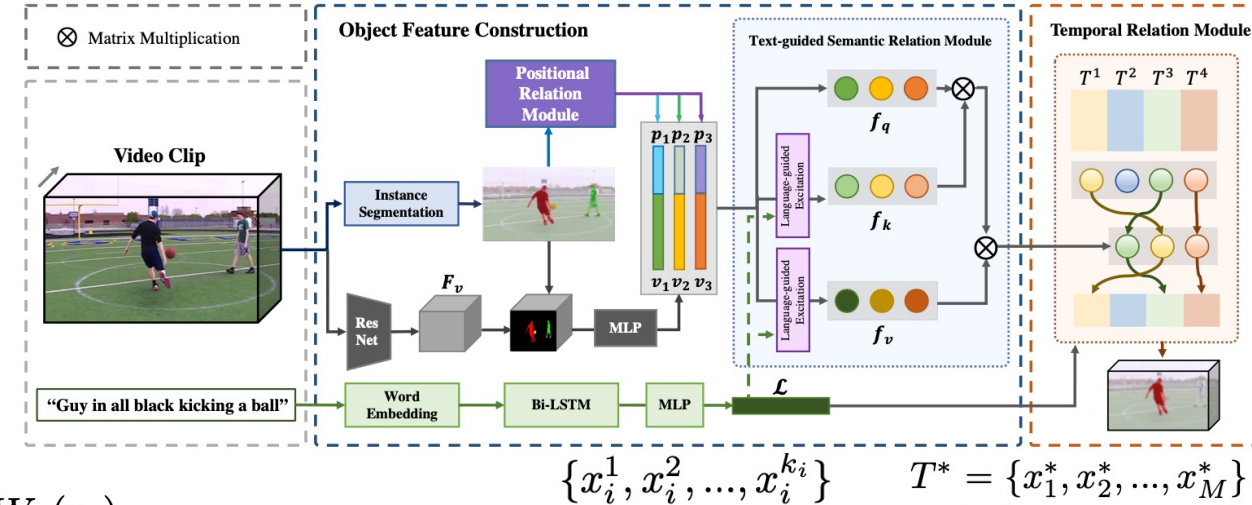
✖ 値 γ よりも大きく、 β 回のマッチングで更新されない場合に適応

✖ 言語特徴量 L とのコサイン類似度が最大の \mathcal{V}_i^r を採用

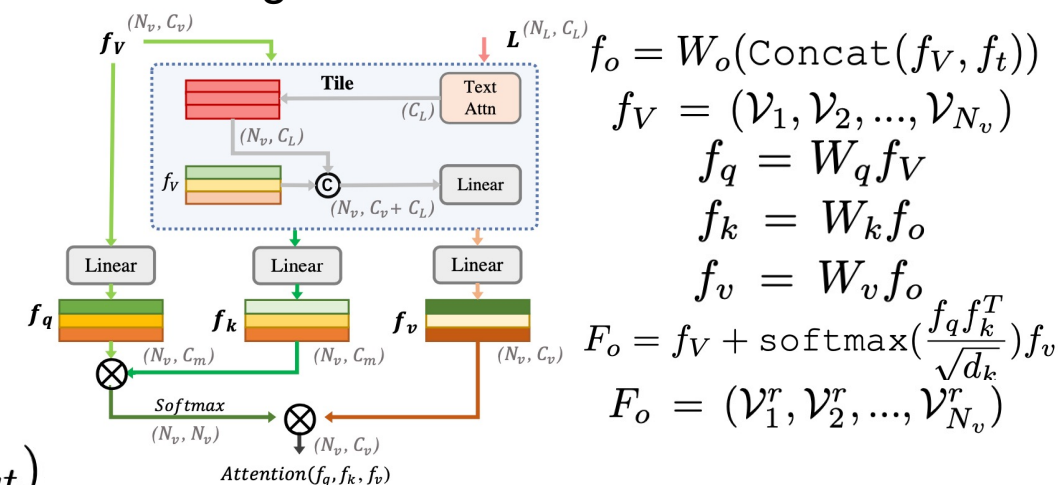
✖ 損失関数

$$s_i = \frac{\exp(\mathcal{V}_i^{r^T} \mathcal{L} / \tau)}{\sum_{j=1}^{N_v} \exp(\mathcal{V}_j^{r^T} \mathcal{L} / \tau)} \quad \text{loss} = -\log(s_{gt})$$

✖ gt : ground truth



Text-guided Semantic Relation Module



$$\begin{aligned} f_o &= W_o(\text{Concat}(f_v, f_t)) \\ f_v &= (v_1, v_2, \dots, v_{N_v}) \\ f_q &= W_q f_v \\ f_k &= W_k f_o \\ f_v &= W_v f_o \\ F_o &= f_v + \text{softmax}\left(\frac{f_q f_k^T}{\sqrt{d_k}}\right) f_v \\ F_o &= (\mathcal{V}_1^r, \mathcal{V}_2^r, \dots, \mathcal{V}_{N_v}^r) \end{aligned}$$

CMF [Yang+ (Beijing Univ.), ICIP21]

✖ Referring Image SegmentationにおいてG-Ref、UNC、UNC+でSoTA

✖ Attention Matrix $A_{i,t} = \text{softmax}[(W_v v_i)^T (W_h h_t)]$, $l_i = \sum_{t=1}^T A_{i,t} h_t$

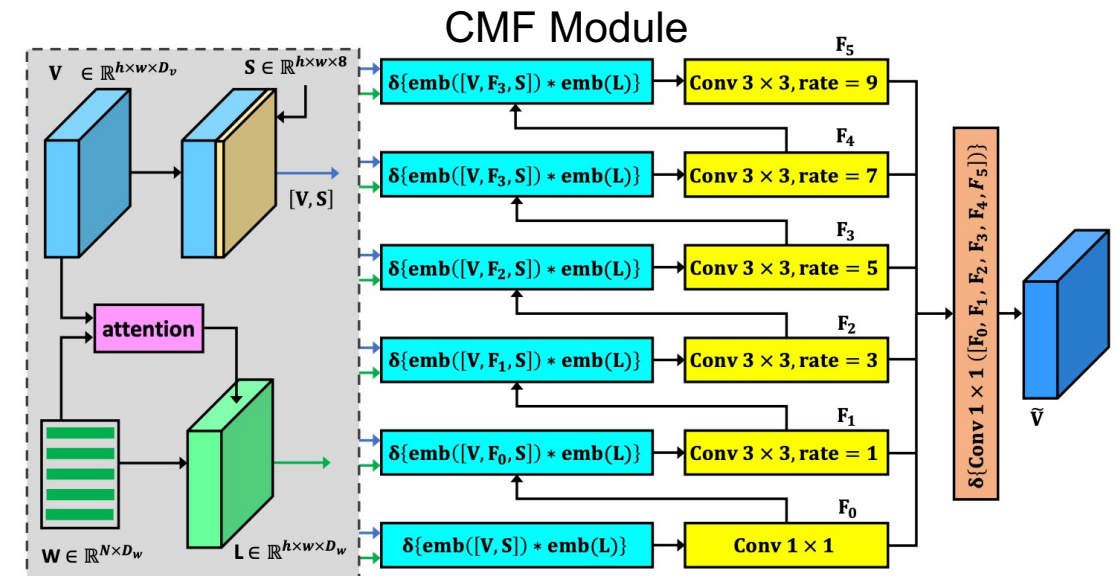
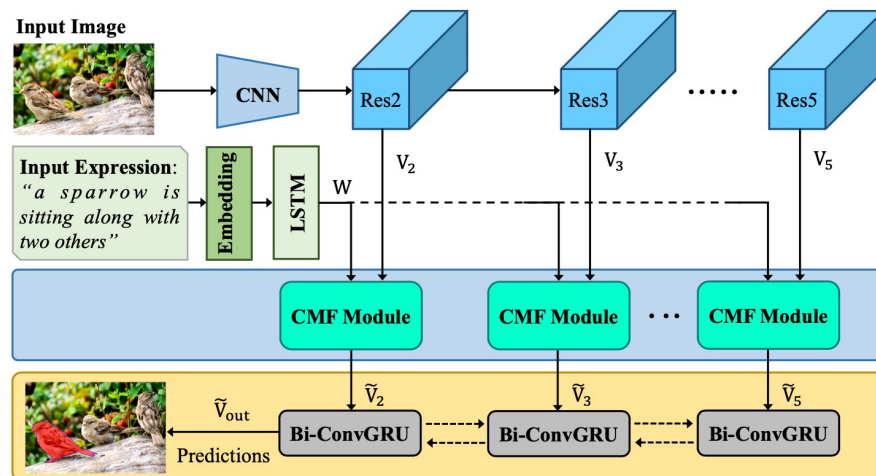
✖ i は視覚特徴のindex、 t は t 番目の単語の言語特徴

✖ Cascaded Multi-modal Fusion(CMF)

✖ $\text{emb}(\cdot)$: 1x1 conv、 δ : ReLU、 $*$: アダマール積

✖ bi-directionally convolutional GRU $\tilde{V}_{out} = \text{ReLU}(W_p^{\vec{H}} \vec{H}_l + W_p^{\overleftarrow{H}} \overleftarrow{H}_l + b)$

✖ マルチモーダルな特徴を
ボトムアップとトップダウンの2つの方法で統合



BUSNet [Yang+ (ShanghaiTech Univ.), CVPR21]

✖ 推論の途中経過が可視化できるRISモデル

✖ UNC、UNC+、G-RefでSoTA

✖ Bottom-Up Shift (BUS)

✖ グラフを生成

✖ ノード o_n / 有向エッジ e_k

✖ 名詞句 / 前置詞、動詞

✖ マルチモーダル特徴 X_n を生成、 X'_n にアップデート

✖ V : 画像特徴、 P : 位置、 \bar{h}_n : o_n の言語特徴、 \mathcal{E}_n : o_n に繋がるエッジ

✖ $*$ は要素ごとの乗算、 $;$ はConcatenate

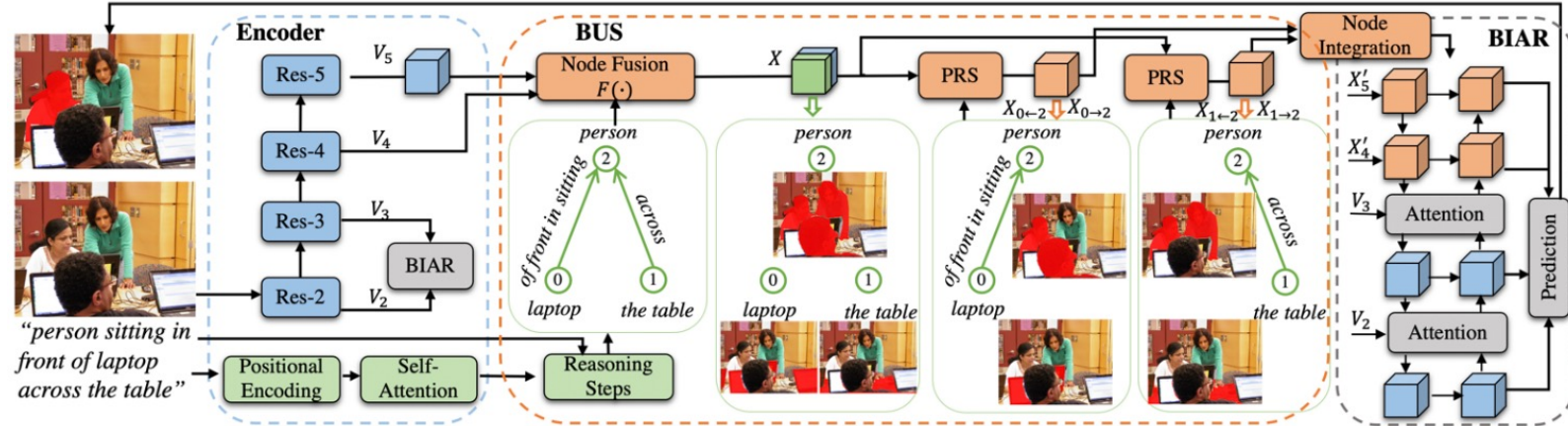
✖ Pairwise Relation Shift (PRS)

✖ \odot はピクセルごとの乗算、 γ はtanh

✖ Bidirectional Attentive Refinement (BIAR)

✖ $\{V_2, V_3, V_4, X_4', X_5'\}$ を改良し、 $\{G_1, G_2, G_3, G_4, G_5\}$ を得る

✖ アップサンプリングして合計し予測



$$X_n = \text{Conv}_v([V; P]) * \text{Tile}(\mathbf{W}_{\bar{h}} \bar{h}_n) \quad X_n = F(V, P, o_n)$$

$$X_{n \leftarrow m}, X_{n \rightarrow m} = \text{PRS}^{(3)}(X_n, e_k^{(r)}, X'_m),$$

$$X'_n = \frac{\sum_{o_m \in e_k^{(o)} \& e_k \in \mathcal{E}_n} X_{n \leftarrow m} + X_n}{|\mathcal{E}_n| + 1}$$

$$\mathbf{A}_{s \leftarrow o} = \gamma(\text{Conv}_r^{-1}(\mathbf{X}_o)), \mathbf{X}_{s \leftarrow o} = F(\mathbf{A}_{s \leftarrow o} \odot \mathbf{V}, \mathbf{P}, e^{(s)}),$$

$$\mathbf{A}_{s \rightarrow o} = \gamma(\text{Conv}_r(\mathbf{X}_s)), \mathbf{X}_{s \rightarrow o} = F(\mathbf{A}_{s \rightarrow o} \odot \mathbf{V}, \mathbf{P}, e^{(o)}),$$

$$\mathbf{A}_i^{td} = \sigma(\text{Conv}_c(\text{Conv}_a(\mathbf{G}_i) + \text{Conv}_b(\text{Up}(\mathbf{G}_{i+1}^{td}))))$$

$$\mathbf{G}_i^{td} = \begin{cases} \text{Conv}_i(\mathbf{G}_i), & \text{if } i \in \{5\} \\ \text{Conv}_i(\mathbf{G}_i + \text{Up}(\mathbf{G}_{i+1}^{td})), & \text{if } i \in \{4\} \\ \text{Conv}_i(\mathbf{A}_i^{td} \odot \mathbf{G}_i + \text{Up}(\mathbf{G}_{i+1}^{td})), & \text{if } i \in \{1, 2, 3\} \end{cases}$$