

# EfficientNet [Le+ (Google Research), ICML19]

✂ 深さ、幅、解像度のすべての次元を一様にスケーリング

✂ 既存の最高のGPipeの精度を上回り、パラメータは8.4倍少なく、推論速度は6.1倍

✂ compound scaling method

✂  $\phi$  で広さ、深さ、解像度を右式で決定

✂  $\alpha = 1.2, \beta = 1.1, \gamma = 1.15$

$$\text{depth: } d = \alpha^\phi$$

$$\text{width: } w = \beta^\phi$$

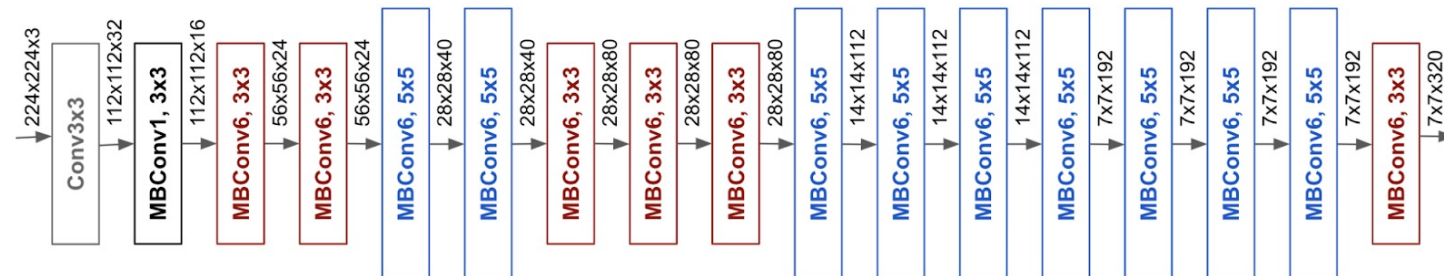
$$\text{resolution: } r = \gamma^\phi$$

$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

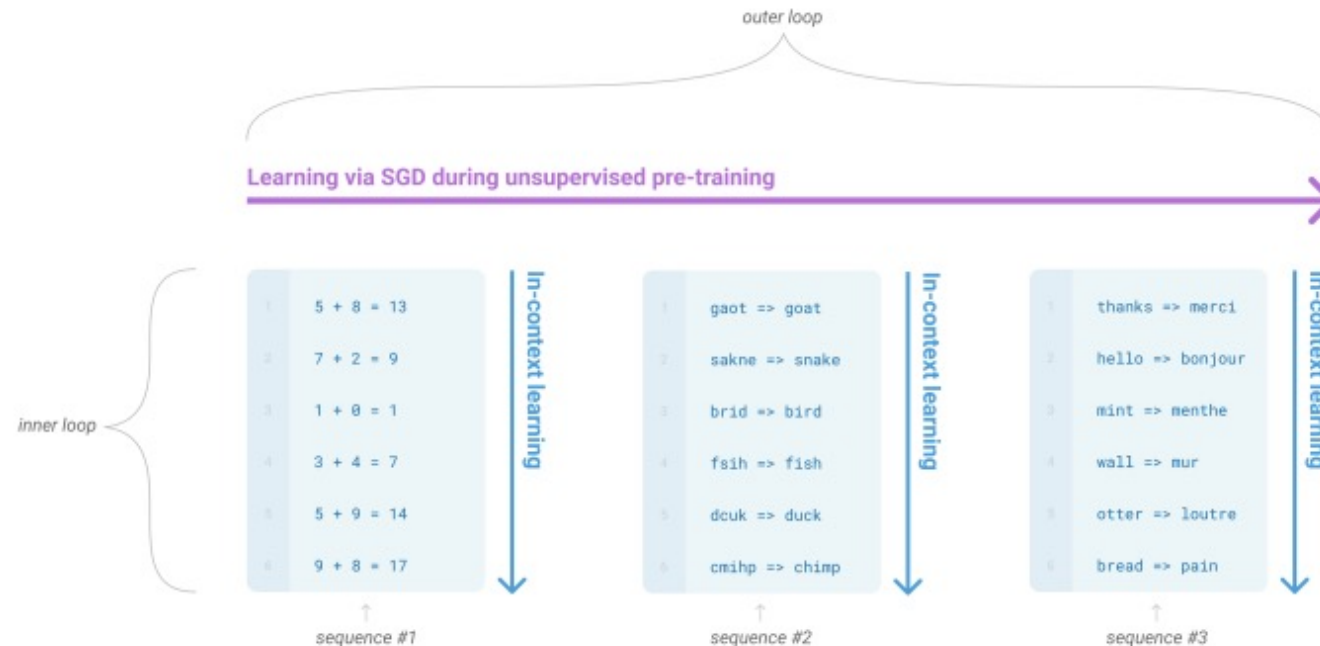
✂ BaselineはMnasNetを使用

✂ MBConv: Mobile Inverted BottleneckにSEモジュールを追加



# GPT-3 [Brown+ (OpenAI), NeurIPS20]

- ✖ Fine-tuningなしで様々なタスクを解くモデル
  - ✖ タスクに応じたデータセットが不要
  - ✖ 偏ったデータでfine-tuningすることによる汎化性能の低下を防ぐ
- ✖ 約45TBの大規模なテキストデータを約1750億個のパラメータを使用して学習
  - ✖ GPT-2と比較して使用データは1100倍以上、パラメータ数は117倍以上
- ✖ 各タスクにおいて、Fine-tuningしたモデルと同等の性能を達成



# T5 [Raffel+ (Google Research), JMLR20]

✖ すべてのNLPタスクを“text-to-text”の形式として扱う

✖ テキストを入力し、生成したテキストを出力

✖ 入力テキストの接頭辞はハイパーパラメータ

✖ 要約、質疑応答、テキスト分類など多くのタスクでSOTA

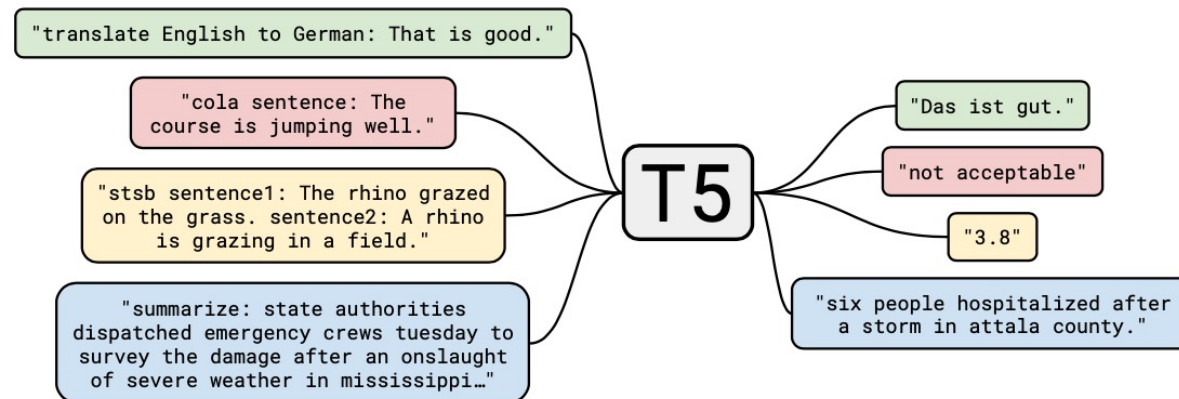
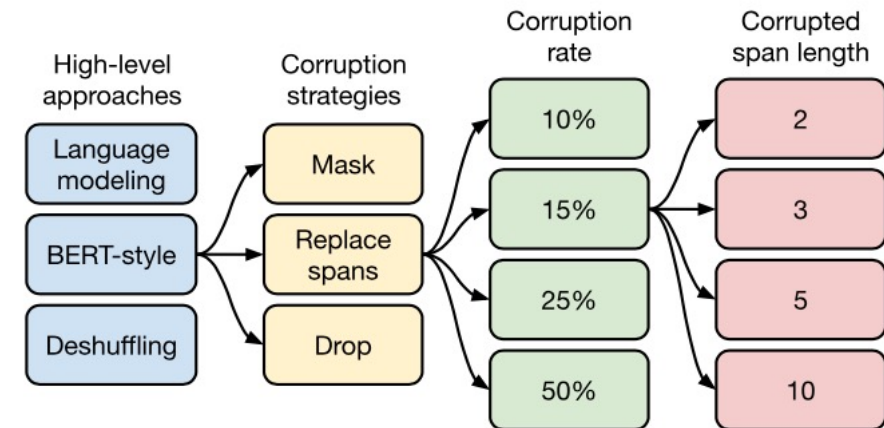
✖ Dataset: Colossal Clean Crawled Corpus(C4)

✖ Common Crawl をクリーンアップ

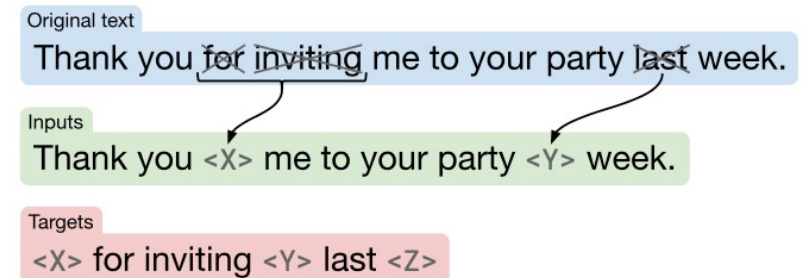
✖ そのまま使用する場合より精度が向上

✖ 事前学習について性能比較(右)

✖ Corrupted Span lengthは大差なく、Baselineをそのまま使用



BERT-style + Replace spans



# Habitat [Savva+ (Facebook AI Research), ICCV19]

## ✂ Embodied AIのためのプラットフォーム

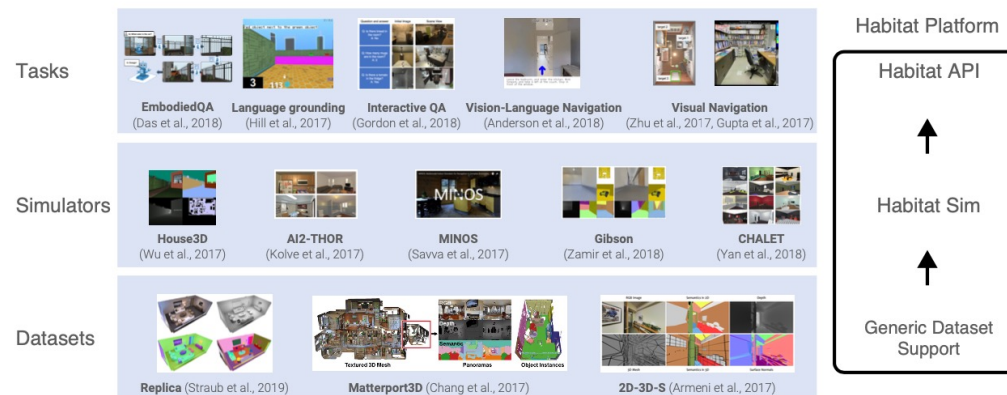
- ✂ エージェントを高効率でフォトリアリスティックな3Dシミュレーションで育成

## ✂ Habitat-Sim

- ✂ 柔軟で高性能な3次元シミュレータ
- ✂ エージェント、センサー、汎用的な3Dデータセットを扱う
  - ✂ シングルスレッドで数千fps、シングルGPUではマルチプロセスで10,000fps以上を達成

## ✂ Habitat-API

- ✂ 体現型AIアルゴリズムをエンドツーエンドで開発するためのモジュール式高レベルライブラリ
  - ✂ タスク(ナビゲーション、質問応答など)の定義、エージェントの設定とトレーニング等を行う



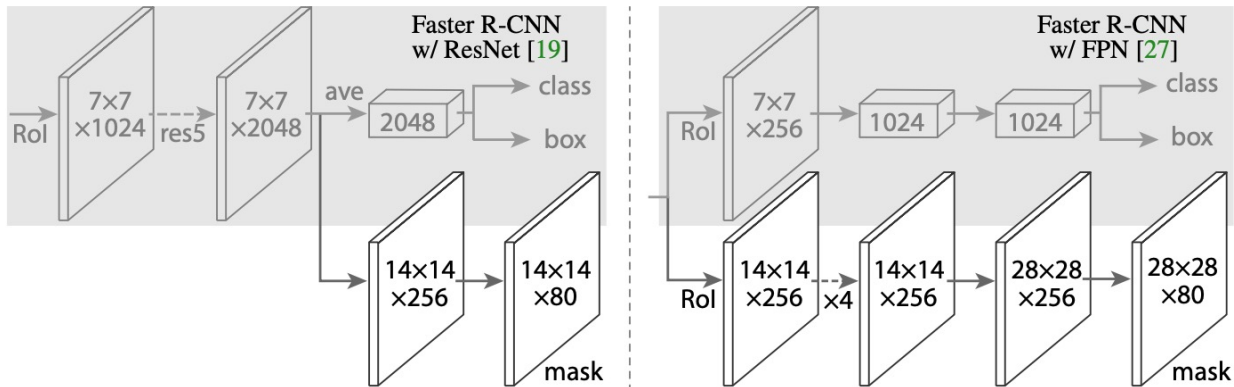
Savva, M., Kadian, A., Maksymets, O., Zhao, Y., Wijmans, E., Jain, B., ... & Batra, D. (2019).

Habitat: A platform for embodied ai research. In Proceedings of the IEEE/CVF International Conference on Computer Vision (pp. 9339-9347).



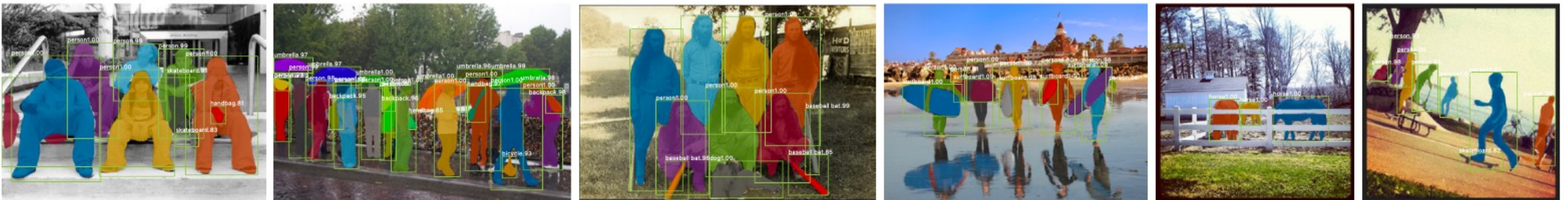
# Mask R-CNN [He+ (Facebook AI Research), ICCV17]

- ✂ Object Detectionに加え、Instance Segmentationに対応したモデル
  - ✂ Instance Segmentation タスクにおいてSoTAとなる37.1mAP
- ✂ Faster R-CNNにbranchを追加
  - ✂ それぞれのRoI(Region of Interest)に対してSegmentation Maskを予測



$$L = L_{cls} + L_{box} + L_{mask}$$

## Binary Cross Entropy



# MAttNet [Yu+ (University of North Carolina at Chapel Hill), CVPR18]

✖ Referring expression comprehension においてSoTA

✖ 文章表現を3つに分解

- ✖ subject appearance: 全体像
- ✖ Location: 位置を示す
- ✖ Relationship: 物体同士の関係性

✖ FC層で3つのモジュールの重みに変換

$$e_t = \text{embedding}(u_t)$$

$$\vec{h}_t = \text{LSTM}(e_t, \vec{h}_{t-1})$$

$$\tilde{h}_t = \text{LSTM}(e_t, \tilde{h}_{t+1})$$

$$h_t = [\vec{h}_t, \tilde{h}_t].$$

$$[w_{subj}, w_{loc}, w_{rel}] = \text{softmax}(W_m^T [h_0, h_T] + b_m)$$

✖ Loss関数

✖  $o_i$ は候補オブジェクト、 $r_i$ はexpression

$$S(o_i|r) = w_{subj}S(o_i|q^{subj}) + w_{loc}S(o_i|q^{loc}) + w_{rel}S(o_i|q^{rel})$$

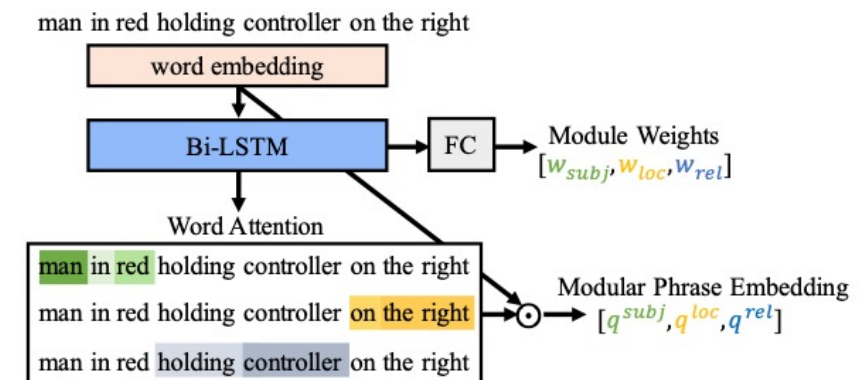
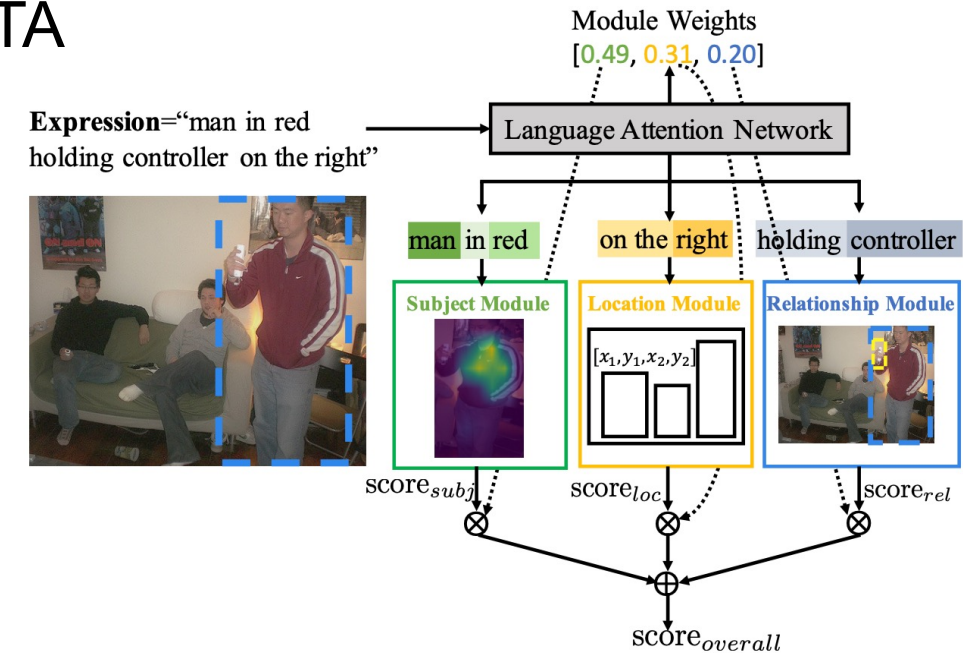
$$L_{rank} = \sum_i [\lambda_1 \max(0, \Delta + S(o_i|r_j) - S(o_i|r_i))$$

$$L = L_{subj}^{attr} + L_{rank}$$

$$+ \lambda_2 \max(0, \Delta + S(o_k|r_i) - S(o_i|r_i))]$$

Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., & Berg, T. L. (2018).

Mattnet: Modular attention network for referring expression comprehension. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1307-1315).



# Causal Attention [Yang+ (Nanyang Technological University), CVPR21]

## ✖ 事前学習によるAttentionの悪化 (右図)

- ✖ トレーニングセットのバイアスによる

## ✖ 因果関係を考慮したCausal Attention(CATT)

- ✖ 他のサンプルを用いてAttentionを算出
- ✖ LXMERTに組み込むことで  
サイズの大きいUNITERと同等の精度

## ✖ In-Sample attention (IS-ATT)

- ✖  $K_I, V_I$ は現在の入力サンプル

**Input:**  $Q_I, K_I, V_I$ ,

**Prob:**  $A_I = \text{Softmax}(Q_I^T K_I)$

**Ouput:**  $\hat{Z} = V_I A_I$

## ✖ Cross-Sample attention (CS-ATT)

- ✖  $K_C, V_C$ はトレーニングセットの他のサンプル

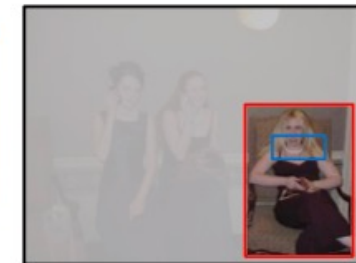
**Input:**  $Q_C, K_C, V_C$ ,

**Prob:**  $A_C = \text{Softmax}(Q_C^T K_C)$

**Ouput:**  $\hat{X} = V_C A_C$



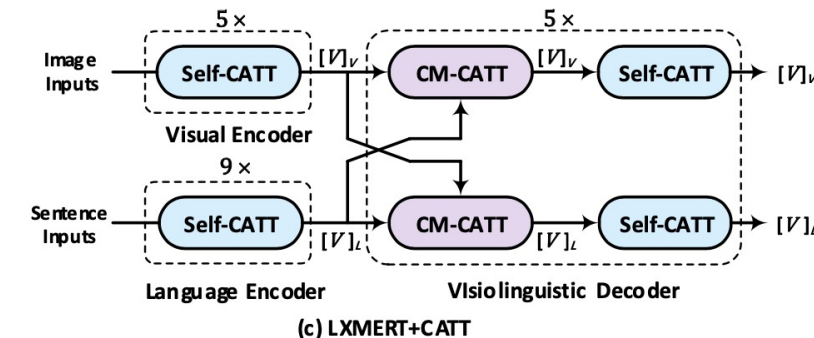
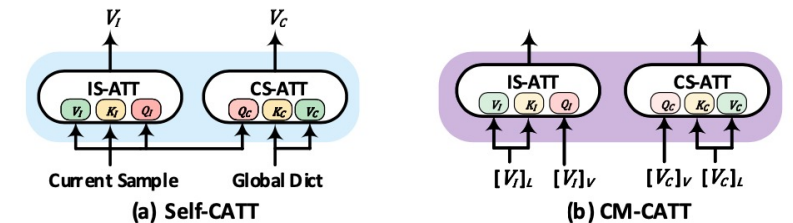
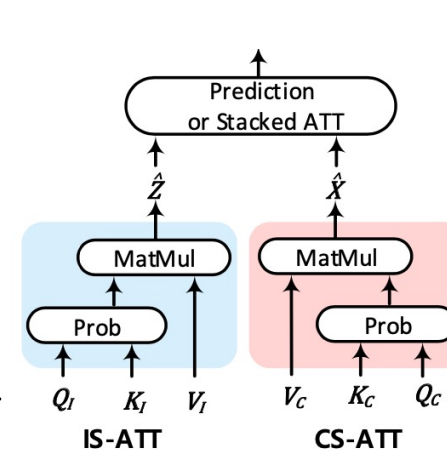
#"Sport+Man" / #"Sport+Screen"=213  
Q: What sport is being shown on the screen?  
A: Dancing (Bowling)



#"Color+Girl" / #"Color+Necklace"=54  
Q: What color is the girl's necklace?  
A: Black (White)



#"Board+Man" / #"Board+Woman"=20  
Q: What gender is the person standing up?  
A: Male (Female)





# Transform and Tell [Tran+ (Australian National University), CVPR20]

✖ ニュース記事と埋め込まれた画像からキャプションを生成

✖ BLEU-4: 0.89→6.05、CIDEr: 13.1→53.8

✖ Image Encoder: ResNet-152のプーリング層の前の出力

✖ Face Encoder: MTCNNでbboxを検出しFaceNetに通す、MFは顔の数

✖ Object Encoder: YOLOv3でbboxを検出しResNet-152に通す、MOはObject数

✖ Article Encoder: RoBERTaを使用、すべての層の加重和

✖ Decoder: 3つの入力からキャプションを生成

✖ 1: 前ステップで生成したトークンのembedding

$z_{0t} \in \mathbb{R}^{D^E}$   $D^E$ はhidden size

✖ 2: 前に生成された全トークンのembedding

$Z_{0<t} = \{z_{00}, z_{01}, \dots, z_{0t-1}\}$

✖ 3: Encoderからの $X^I$ ,  $X^A$ ,  $X^F$ ,  $X^O$

✖ 入力はL個のtransformer blockに送られる

$z_{1t} = \text{Block}_1(z_{0t} | Z_{0<t}, X^I, X^A, X^F, X^O)$

$z_{2t} = \text{Block}_2(z_{1t} | Z_{1<t}, X^I, X^A, X^F, X^O)$

...

$z_{Lt} = \text{Block}_L(z_{L-1t} | Z_{L-1<t}, X^I, X^A, X^F, X^O)$

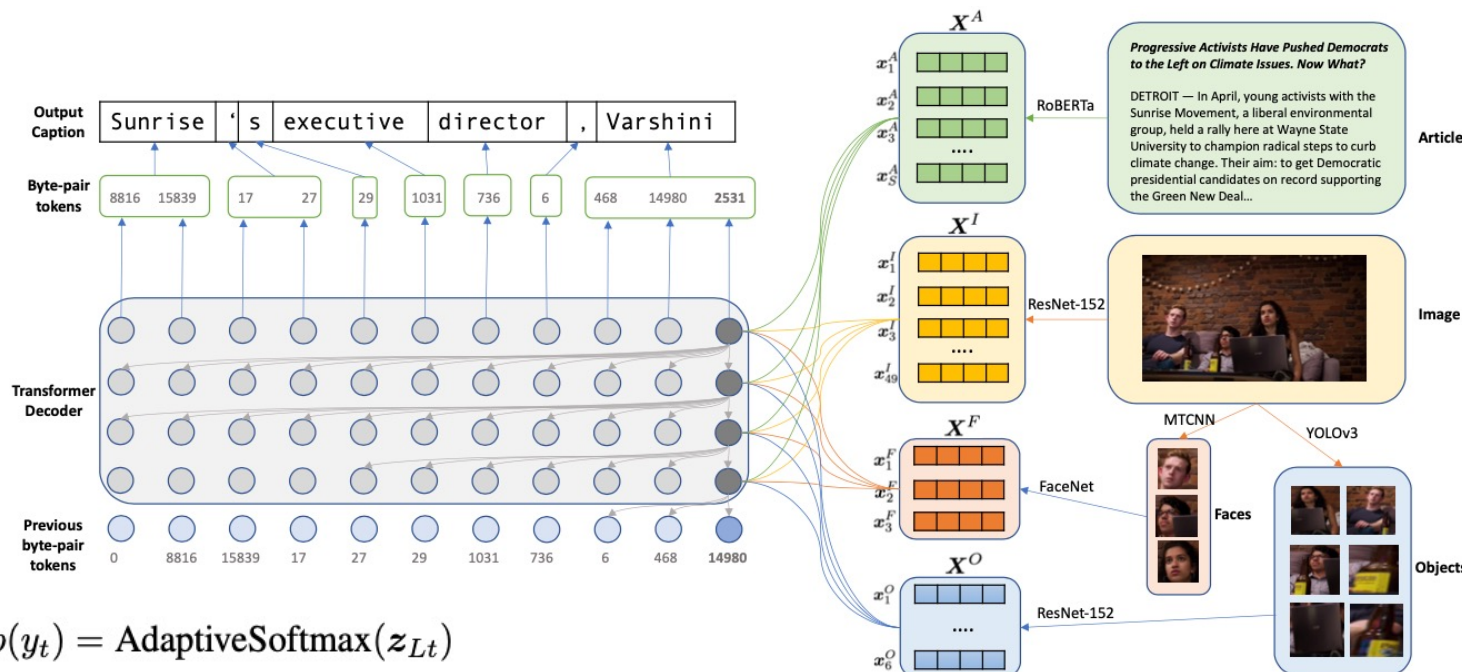
✖  $Z_{Lt}$ からt番目トークンの生成確率 $p(y_t)$ を推定:  $p(y_t) = \text{AdaptiveSoftmax}(z_{Lt})$

$X^I = \{x_i^I \in \mathbb{R}^{D^I}\}_{i=1}^{M^I}$ , where  $D^I = 2048$  and  $M^I = 49$

$X^F = \{x_i^F \in \mathbb{R}^{D^F}\}_{i=1}^{M^F}$ , where  $D^F = 512$

$X^O = \{x_i^O \in \mathbb{R}^{D^O}\}_{i=1}^{M^O}$ , where  $D^O = 2048$

$X^A = \{x_i^A \in \mathbb{R}^{D^T}\}_{i=1}^{M^T}$ , where  $D^T = 1024$





# Oscar [Li+ (Microsoft Corporation), ECCV20]

## ✂ 画像-テキスト表現を学習するVLP手法

✂ 複数のV+LベンチマークでSoTA

## ✂ アンカーポイントとしてオブジェクトタグを使用

## ✂ Word Tokens (w)

✂ テキストの単語埋め込みベクトル

## ✂ Object Tag (q)

✂ 画像から検出されたオブジェクトタグの単語埋め込みベクトル

## ✂ Region Features (v)

✂ 画像の領域ベクトルの集合

## ✂ 事前学習タスク

✂ Masked Token Loss  $\mathcal{L}_{MTL} = -\mathbb{E}_{(v,h) \sim \mathcal{D}} \log p(h_i | h_{\setminus i}, v)$

✂ w, qをマスクし、復元するよう学習

✂ Contrastive Loss  $\mathcal{L}_C = -\mathbb{E}_{(h',w) \sim \mathcal{D}} \log p(y | f(h', w))$

✂ qを50%の確率で異なるタグに置き換え、適切なタグかどうか予測

