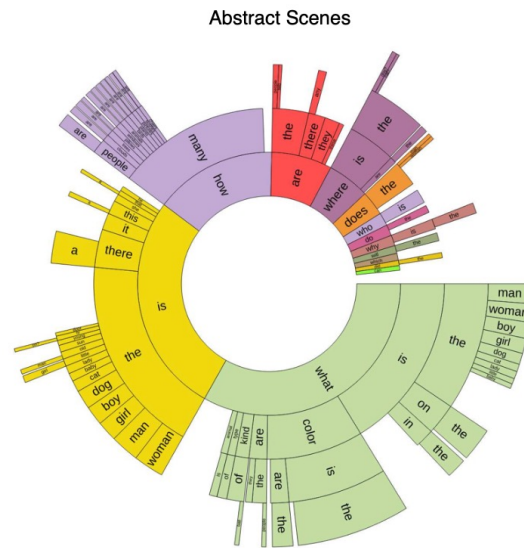
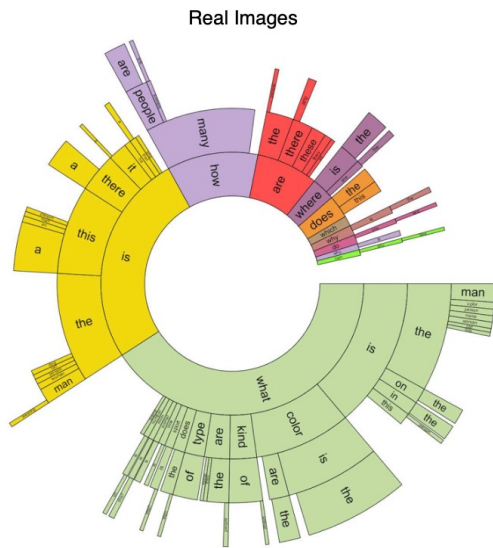


VQA [Agrawal+ (Virginia Tech), ICCV15]

- ✖ 画像と自然言語の質問を入力とし、自然言語で答えを生成するタスク
 - ✖ 250K以上の画像、760Kの質問、および約10Mの回答を含むデータセット
- ✖ real image と abstract scene が存在
- ✖ 多くのanswerは数語のみ、89.32%が1単語
 - ✖ 自動評価にも適している



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

VQA 2.0 [Goyal+ (Virginia Tech), CVPR17]

✖ 言語バイアスを大幅に削減し、バランスの取れたVQAデータセット(サイズは約2倍)

✖ 「What sport is」から始まる質問のうち41%で「tennis」が正解だった

✖ Do you see a ... という文に対してはyesと答えるだけで87%正解だった

✖ 被験者は、VQAから(Image, Question, Answer)の組 (I, Q, A)が与えられ、Iに似ているが、Qに対する答えがA'になる画像I'を特定

✖ 言語のみのモデルでは(Q, I)と(Q, I')を区別する根拠がない

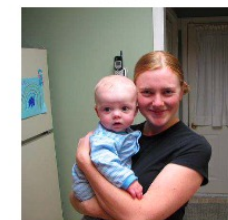
✖ 画像理解の進歩をより正確に反映

✖ 既存のVQAで学習したモデルはVQA 2.0ではパフォーマンスが低下

Who is wearing glasses?
man woman



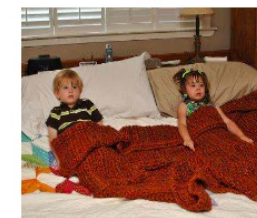
Where is the child sitting?
fridge arms



Is the umbrella upside down?
yes no



How many children are in the bed?
2 1



ORT [Herdade+ (Yahoo Research), NeurIPS19]

✂ Image Captioningに特化したTransformer

- ✂ 検出されたオブジェクト間の位置とサイズの関係性をエンコード

✂ Relation box

- ✂ Transformerにおける Ω_A (右) の要素 ω_A に位置情報の関係を組み込む

$$\Omega_A = \frac{QK^T}{\sqrt{d_k}}$$

- ✂ (x_m, y_m, w_m, h_m) はbboxである m の中心座標、幅、高さ

$$\lambda(m, n) = \left(\log \left(\frac{|x_m - x_n|}{w_m} \right), \log \left(\frac{|y_m - y_n|}{h_m} \right), \log \left(\frac{w_n}{w_m} \right), \log \left(\frac{h_n}{h_m} \right) \right)$$

$$\omega_G^{mn} = \text{ReLU}(\text{Emb}(\lambda)W_G)$$

$$\omega^{mn} = \frac{\omega_G^{mn} \exp(\omega_A^{mn})}{\sum_{l=1}^N \omega_G^{ml} \exp(\omega_A^{ml})}$$

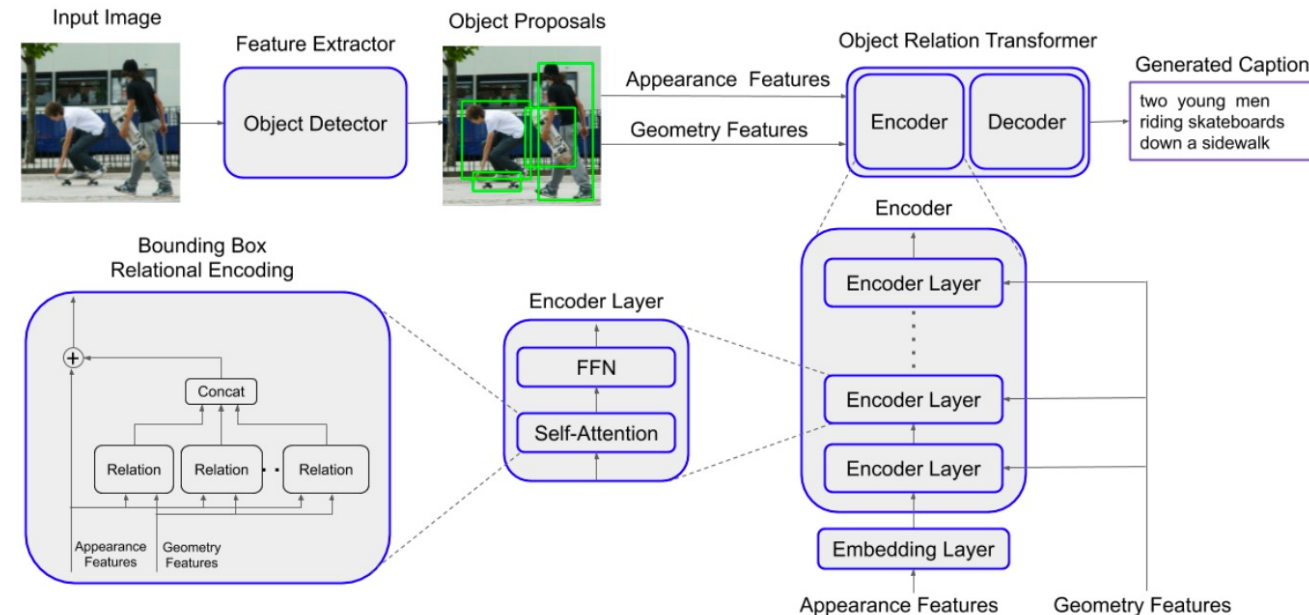
✂ Encoder Output

- ✂ Ω は ω^{mn} で構成される $N \times N$ 行列

$$\text{head}(X) = \text{self-attention}(Q, K, V) = \Omega V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O$$

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



NetVLAD++ [Giancola+ (KAUST, Saudi Arabia), CVPR21]

✖ Action Spotting: タイムスタンプにおける瞬間的なイベント(=Action)を特定するタスク

✖ NetVLAD++: 新しい特徴プーリング手法

✖ 過去 $[-T_b, 0]$ と未来 $[0, T_a]$ の特徴量を考慮

✖ $V = \square(V_b, V_a)$

✖ \square は V_b と V_a の aggregation

✖ V_a は $[0, T_a]$ において NetVLAD pool された特徴量

✖ NetVLAD

✖ VLAD に自由度を与えたもの

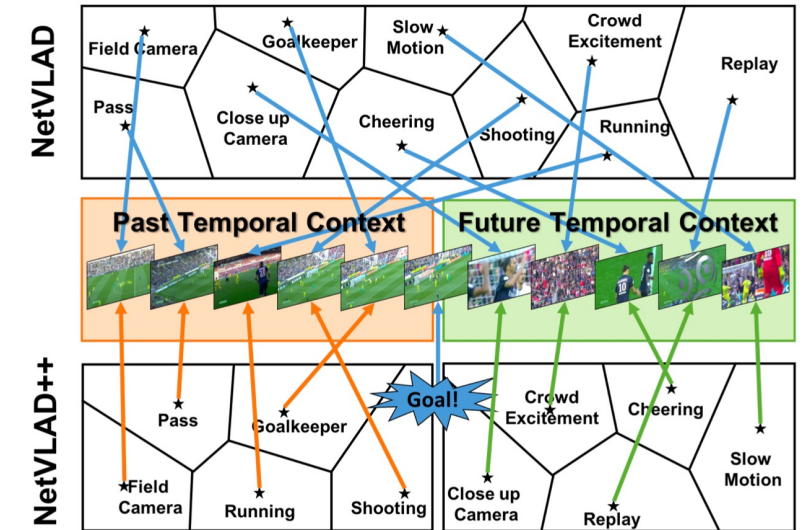
✖ W_k, b_k, c_k を最適化

$$V(j, k) = \sum_{i=1}^N \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \mathbf{x}_i + b_{k'}}} (\mathbf{x}_i(j) - \mathbf{c}_k(j))$$

✖ Average-mAP: 53.4%

✖ Dataset: SoccerNet-v2

✖ 当時の SoTA と比較して 12.7% の向上



	SoccerNet-v2	visible	unshown	Ball out	Throw-in	Foul	Ind. free-kick	Clearance	Shots on tar.	Shots off tar.	Corner	Substitution	Kick-off	Yellow card	Offside	Dir. free-kick	Goal	Penalty	Yel. → Red	Red card
MaxPool [19]	18.6	21.5	15.0	38.7	34.7	26.8	17.9	14.9	14.0	13.1	26.5	40.0	30.3	11.8	2.6	13.5	24.2	6.2	0.0	0.9
NetVLAD [19]	31.4	34.3	23.3	47.4	42.4	32.0	16.7	32.7	21.3	19.7	55.1	51.7	45.7	33.2	14.6	33.6	54.9	32.3	0.0	0.0
AudioVid [41]	39.9	43.0	23.3	54.3	50.0	55.5	22.7	46.7	26.5	21.4	66.0	54.0	52.9	35.2	24.3	46.7	69.7	52.1	0.0	0.0
CALF [7]	40.7	42.1	29.0	63.9	56.4	53.0	41.5	51.6	26.6	27.3	71.8	47.3	37.2	41.7	25.7	43.5	72.2	30.6	0.7	0.7
NetVLAD++	53.4	59.4	34.8	70.3	69.0	64.2	44.4	57.0	39.3	41.0	79.7	68.7	62.1	56.7	39.3	57.8	71.6	79.3	3.7	4.0

- ✖ チームスポーツの選手を所属チームに応じて教師なしで分類するタスク
 - ✖ ジャージの色やデザインは事前に分からないとする
 - ✖ 選手の位置に関するヒートマップを正確に計算可能
- ✖ Mask R-CNNで物体検出
- ✖ ラベル付きデータに基づいて学習したCNNで審判を分類
- ✖ 選手画像の特徴ベクトルに対し、k-meansで2チームのcluster centerを推定
- ✖ cluster centerに基づきプレイヤーをチームに帰属
- ✖ 1フレームの教師なし学習で94%の精度
 - ✖ 500フレーム(17秒)以内に97%の精度

