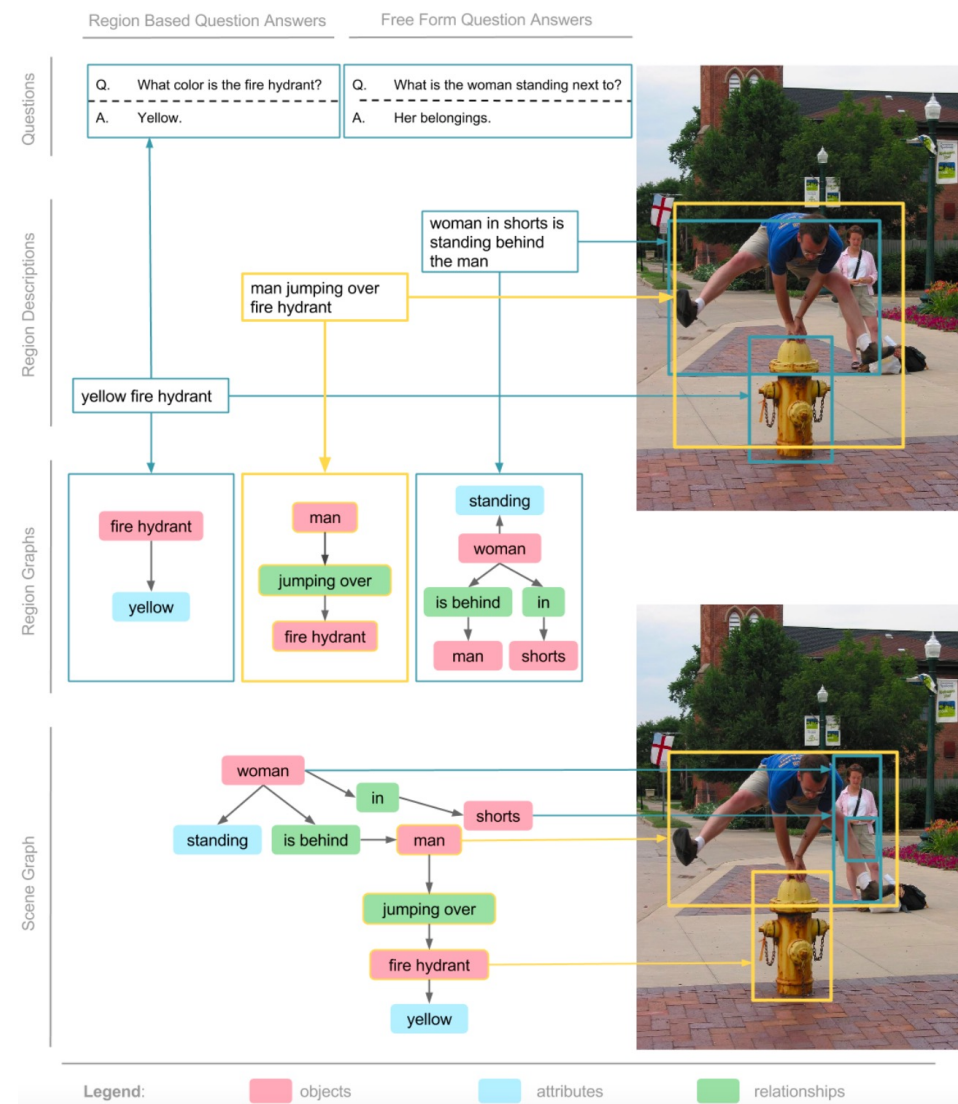


# Visual Genome [Krishna+ (Stanford Univ.), 17]

- ✖ 多くの高品質なラベルが付与された画像データセット
  - ✖ 10万8000枚以上
- ✖ 各画像に対し、平均50のregion description (1~16単語)
- ✖ 平均35のobjectで構成され、bboxで囲まれている
- ✖ 各objectはWordNetのsynset IDに正規化
  - ✖ man.n.03、person.n.01 など
  - ✖ personがmanの上位概念
- ✖ 平均26のattribute
  - ✖ objectには、0個以上のattribute
  - ✖ 色(例:yellow)、状態(例:standing) など
- ✖ 平均21のrelationship
  - ✖ jumping\_over(man, fire hydrant)と表現
- ✖ 各Regionの有向グラフ表現 (Region Graph)
  - ✖ ノードは object、attribute、relationship
- ✖ 全体を表す1つのScene Graph
  - ✖ Region Graphを組み合わせたもの
- ✖ 2種類のQAペア
  - ✖ 画像全体に基づく自由形式のQAと、画像のある領域に基づいたQA



# YOLOv3 [Redmon+ (Univ. of Washington), 18]

## ✖ Darknet-53

✖ ResNet-152と同程度の性能で2倍高速

## ✖ K-means法で学習用のbboxを作成

## ✖ 画像ごとに3つの3D tensorを予測

## ✖ ネットワークが $t_x, t_y, t_w, t_h$ を予測

✖  $c_x, c_y$ はオフセット、

$p_w, p_h$ は作成したbboxの幅と高さ

✖ 学習はground truthとの二乗誤差

✖ bboxの選択にはロジスティック分類を使用

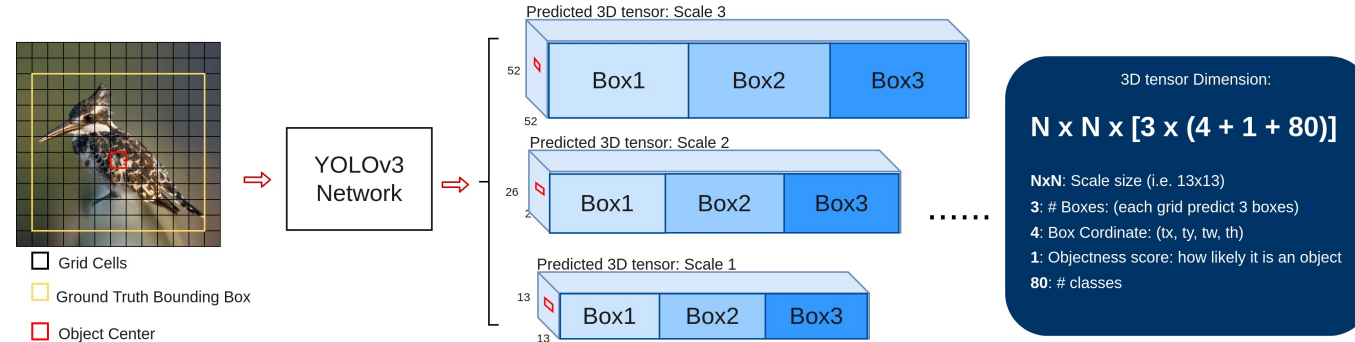
✖ 各bboxのobjectness scoreを予測

## ✖ クラスはマルチラベル分類

✖ ロジスティック分類を使用

✖ 学習にはbinary cross-entropy loss

## ✖ v2は19層モデルに対しv3は53層

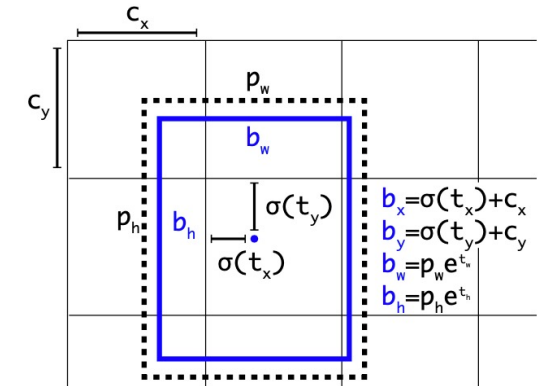


$$b_x = \sigma(t_x) + c_x$$

$$b_y = \sigma(t_y) + c_y$$

$$b_w = p_w e^{t_w}$$

$$b_h = p_h e^{t_h}$$



損失関数  $\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i,j}^{obj} [(t_x - \hat{t}_x)^2 + (t_y - \hat{t}_y)^2 + (t_w - \hat{t}_w)^2 + (t_h - \hat{t}_h)^2]$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i,j}^{obj} [-\log(\sigma(t_o)) + \sum_{k=1}^C BCE(\hat{y}_k, \sigma(s_k))]$$

$$+ \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i,j}^{noobj} [-\log(1 - \sigma(t_o))]$$

# Linear Transformers [Katharopoulos+ (Idiap Research Institute), ICML21]

✖ Attentionの計算に内積でなくカーネル関数を使用し、 $O(n)$ を達成

✖ 長い配列の自己回帰予測では最大4000倍の速度

✖  $i$ 行目のAttention  $V'_i$ を考える  $A_l(x) = V' = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$ .  $V'_i = \frac{\sum_{j=1}^N \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^N \text{sim}(Q_i, K_j)}$   $\text{sim}(q, k) = \exp\left(\frac{q^T k}{\sqrt{D}}\right)$

✖  $\text{sim}(q, k) = \phi(q)^T \phi(k)$  のように  $q, k$  を分解

$$V'_i = \frac{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j) V_j^T}{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j)} \quad i \text{に} \text{関与しないため一定} \quad \begin{aligned} \sum_{j=1}^n \phi(k_j) \cdot v_j &= A, \\ \sum_{j=1}^n \phi(k_j) &= B \end{aligned} \quad O(n) \quad \text{Attention}(Q, K, V)_i = \frac{\phi(q_i)^T A}{\phi(q_i)^T B} \quad O(1)$$

✖ Softmaxは近似が難しいので  $\phi(x) = \text{elu}(x) + 1$  を使用

✖ 自己回帰モデルの学習

✖  $i$ が $j$ から影響を受けるのは  $j \leq i$  の場合のみ(後ろから影響を受けない)

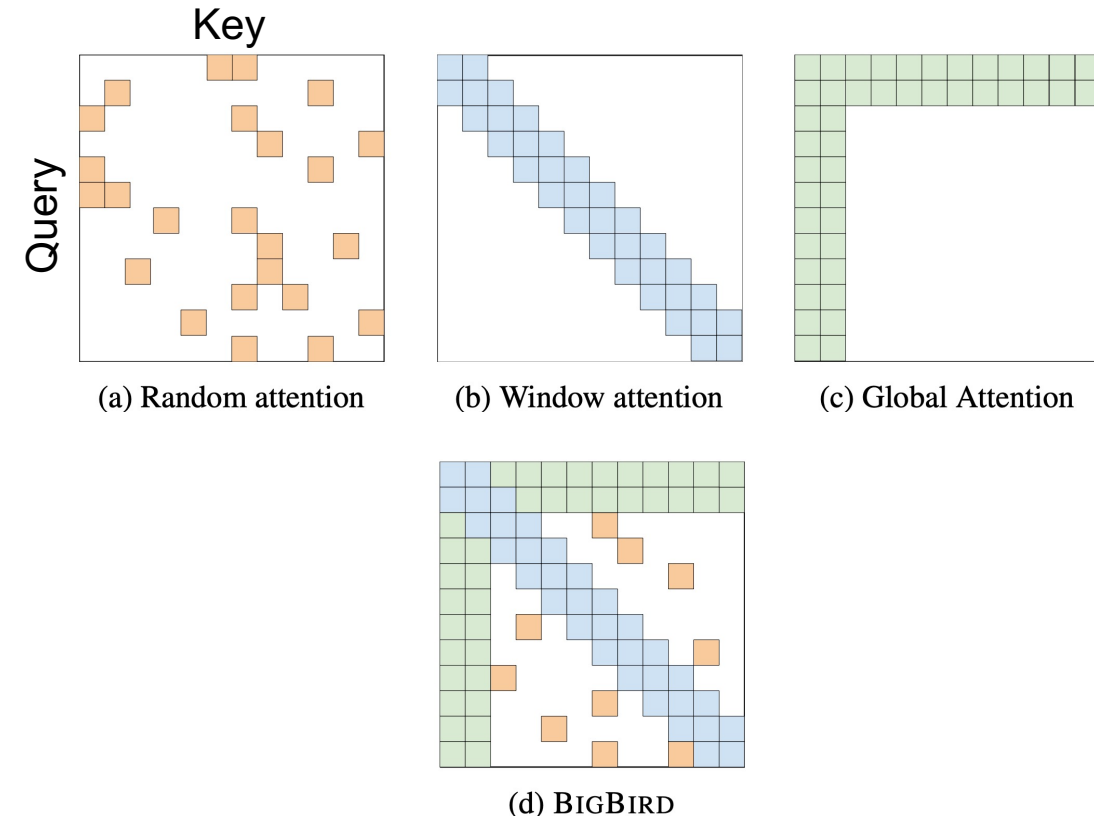
$$V'_i = \frac{\sum_{j=1}^i \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^i \text{sim}(Q_i, K_j)} \quad V'_i = \frac{\phi(Q_i)^T S_i}{\phi(Q_i)^T Z_i} \quad S_i = \sum_{j=1}^i \phi(K_j) V_j^T \quad Z_i = \sum_{j=1}^i \phi(K_j),$$

✖  $S, Z$ を漸化式とみなすことでメモリ消費量を削減可能

✖  $f_l$ はFFN  $s_0 = 0, \quad s_i = s_{i-1} + \phi(x_i W_K) (x_i W_V)^T, \quad y_i = f_l \left( \frac{\phi(x_i W_Q)^T s_i}{\phi(x_i W_Q)^T z_i} + x_i \right).$   
 $z_0 = 0, \quad z_i = z_{i-1} + \phi(x_i W_K),$

# BigBird [Zaheer+ (Google Research), NeurIPS20]

- ✂ Attentionにスパース性を導入し $O(n)$ を達成
  - ✂ 長い文書を必要とするタスクについて高性能、文書要約タスクでSoTA
- ✂ QueryとKeyの乗算 $O(n^2)$ を近似
- ✂ Random Attention
  - ✂ 要素をランダムに選びattentionを計算
    - ✂  $r = 2$ を選択、1行につき2要素をランダム
- ✂ Window Attention
  - ✂ 近傍のみattentionを計算
    - ✂  $w = 3$ を選択、自分自身とその隣2つ
- ✂ Global Attention
  - ✂ 選択した要素と他の全要素のattentionを計算
    - ✂  $g = 2$ を選択、行と列それぞれから2つ選択
- ✂ 上3つを足し合わせた箇所をAttention
  - ✂ 上は全て $O(n)$ なので、Big Birdも $O(n)$





# Mobile-Former [Chen+ (Microsoft), 21]

✂ ローカルな特徴量を扱うMobileNet+グローバルな特徴量を扱うTransformer

✂ 計算コストを節約しつつ、精度を高めることが可能

✂ 2つの入力

✂ 画像  $X_0 \in \mathbb{R}^{H \times W \times 3}$

✂ 学習可能なトークン  $Z_0 \in \mathbb{R}^{M \times d}$

✂ ランダムに初期化  $M$ : トークン数、 $M \leq 6$

✂ Mobile  $\rightarrow$  Former

✂  $z^{out} = z + \left[ \text{Attention}(z_h W_h^Q, x_h, x_h) \right]_{h=1:H} W^O$

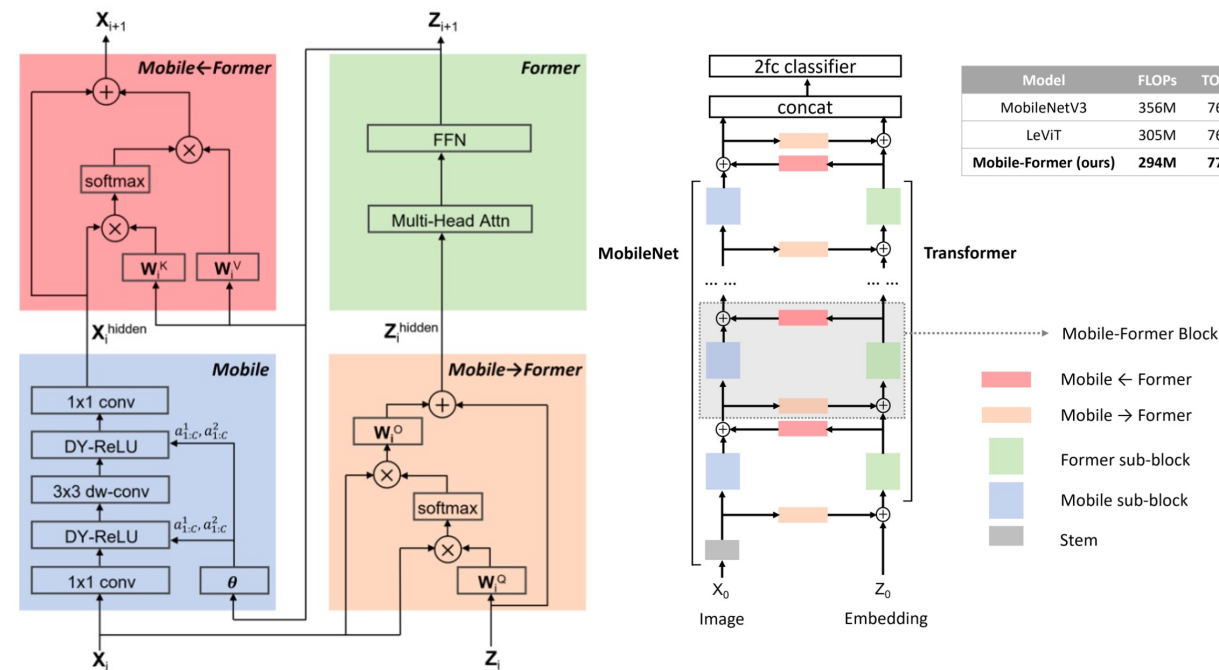
✂ Former  $\rightarrow$  Mobile

✂  $x^{out} = x + \left[ \text{Attention}(x_h, z_h W_h^K, z_h W_h^V) \right]_{h=1:H}$

✂  $H$  headのmulti-head attenのため $x, z$ をsplit  
 $x = [x_h]$  and  $z = [z_h]$  ( $1 \leq h \leq H$ )

✂ 最終的な $x, z$ をconcatし、線形層

✂ ImageNet等で学習



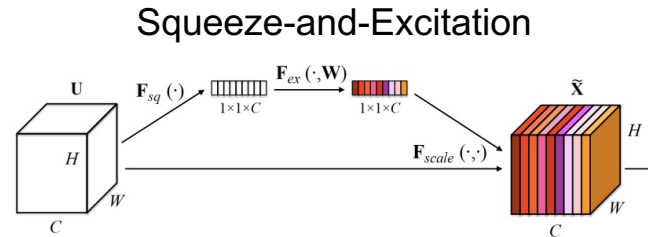
# MobileNetV3 [Howard+ (Google AI), ICCV19]

## ✖ 携帯電話のCPUに合わせて調整された軽量のCNNモデル

- ✖ object detectionやsemantic segmentationに適用
- ✖ MobileNetV3-Large
  - ✂ V2からlatencyを20%削減し、ImageNet分類精度を3.2%向上
- ✖ MobileNetV3-Small
  - ✂ V2と同じlatencyでImageNet分類精度を6.6%の向上

## ✖ Inverted Residual Block (V2)

- ✖ 3x3 Depthwise Conv.を1x1(Pointwise) Conv.で挟む



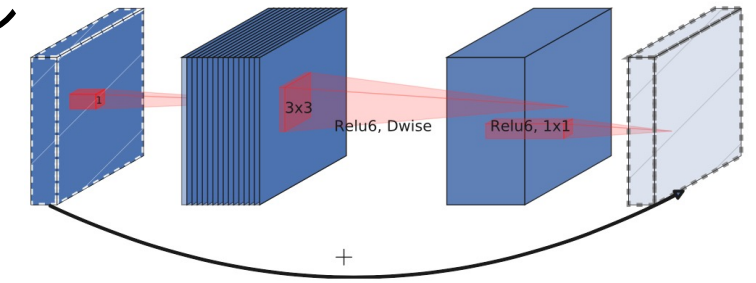
## ✖ V3ではSqueeze-and-Excitationを導入

- ✖ GAP→全結合層したものをつけ合わせる

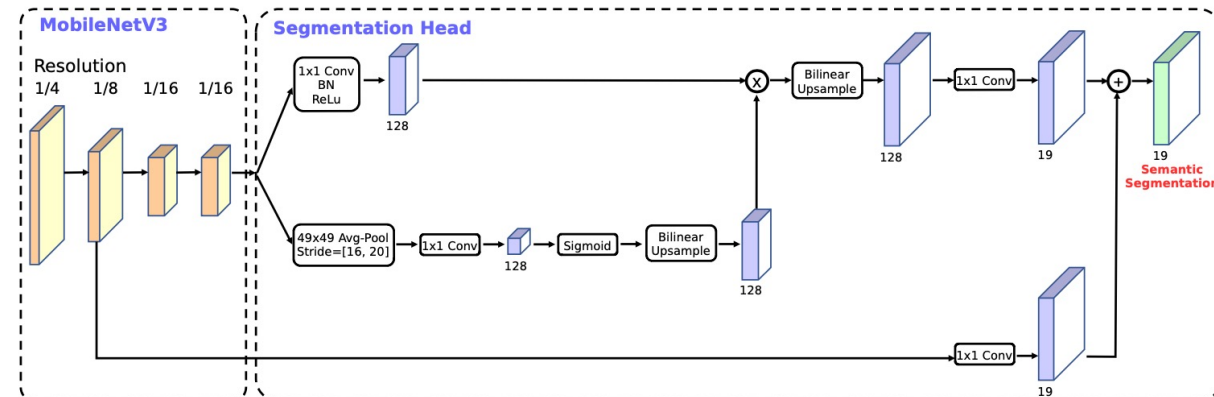
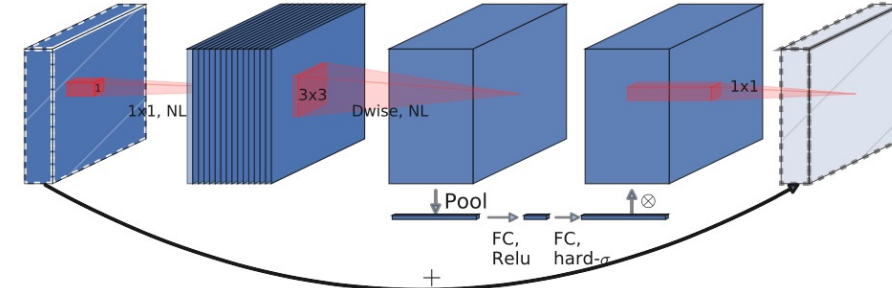
## ✖ Lite R-ASPP

- ✖ 軽量のSegmentation Head
  - ✂ R-ASPPよりわずかに高速で、性能も向上

MobileNet V2: bottleneck with residual



MobileNet V3 block



# DeepLab v3+ [Chen+ (Google Inc.), ECCV18]

## ✖ Semantic SegmentationのためのEncoder-Decoderモジュール

### ✖ Atrous Convolution

✖ マルチスケールの情報を捉える畳み込み

$$y[i] = \sum_k x[i + r \cdot k]w[k]$$

### ✖ Atrous Spatial Pyramid Pooling (ASPP)

✖ 異なるスケールのAtrous Convolutionを特徴マップに適用

### ✖ Decoder

✖ Low-level featuresを活用

✖ bilinearly upsamplingを16倍から4倍に改善

✖ 詳細なmapを作成

