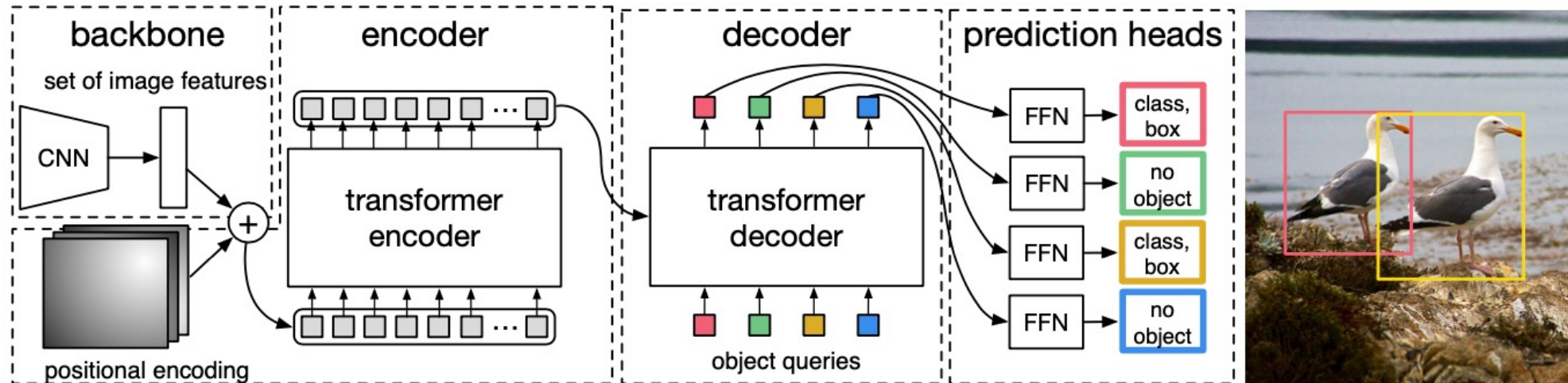


# DETR [Carion+ (Facebook AI), ECCV20]

- ✖ 一度の処理で画像内のすべての物体を検出しクラスを予測
- ✖ 画像内の物体の関係性を利用
- ✖ Transformerを用いることで並列にデコード可能
  - ✖ 既存研究は殆どが自己回帰モデル



- ✖ 小さい物体の検出や学習時間が課題

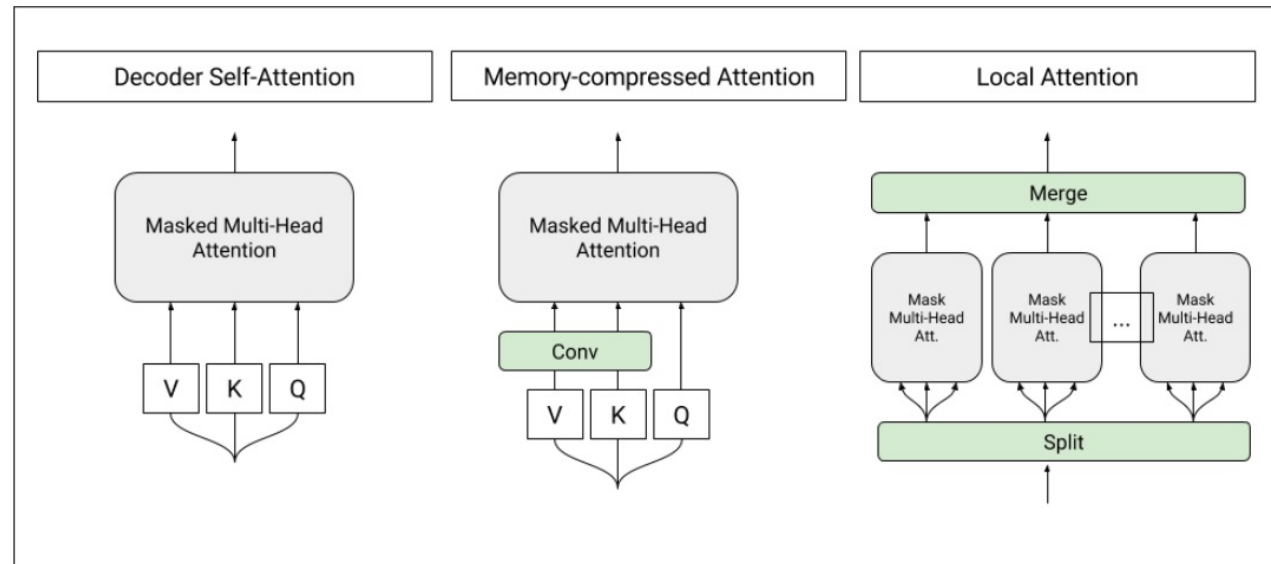
# RoBERTa [Liu+ (Facebook AI), ICLR20]

- ✖ BERTの事前学習を改良し、性能を大幅に向上
- ✖ 大きなバッチで学習
  - ✖ 256(BERT)→2k
- ✖ 多くのデータで、長い時間(10万～50万ステップ)で学習
  - ✖ Dataset: 16GB(BERT)→160GB
- ✖ 学習の度にマスクの箇所を変更
  - ✖ BERTでは最初に一度決定したマスクを使用
- ✖ より長い文章を入力として学習
  - ✖ BERTでは短いセグメントを多く使用
- ✖ Next Sentence Prediction(NSP)は効果がない
  - ✖ BERTで用いられていた事前学習タスク
  - ✖ 除くことで精度向上

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7
XLNet <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6

# Memory Compressed Transformer [J. Liu+ (Google Brain), ICLR18]

- ✖ 文章の要約タスクに関して長シーケンスをTransformerで扱う
- ✖ Memory-compressed Attention
  - ✖ Value、Keyに対して畳み込みを行い数を減らす
- ✖ Local Attention
  - ✖ シーケンスを同じような長さのブロックに分割し、各ブロックで独立してattentionを行う



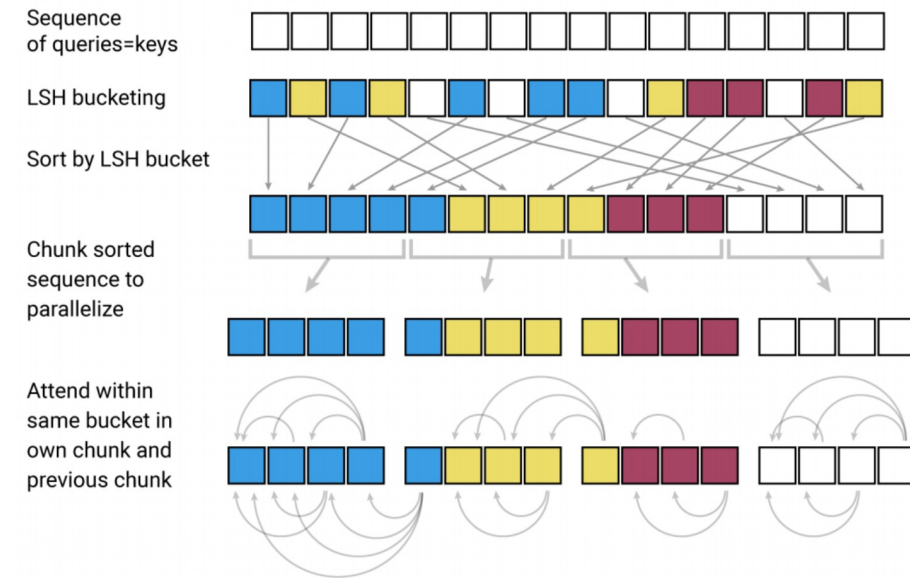
# Reformer [Nikita Kitaev+ (U.C. Berkeley & Google Research), 20]

✂ 長いシーケンスLで高速に動作し、メモリ効率が高いTransformer

✂ locality-sensitive hashing (LSH) により $O(L^2)$  から $O(L \log L)$ に削減

✂ LSH Attention

✂ 元々の位置より手前側、同一のバケット内、1つ前までのチャンク、という条件を満たす場合のみAttention

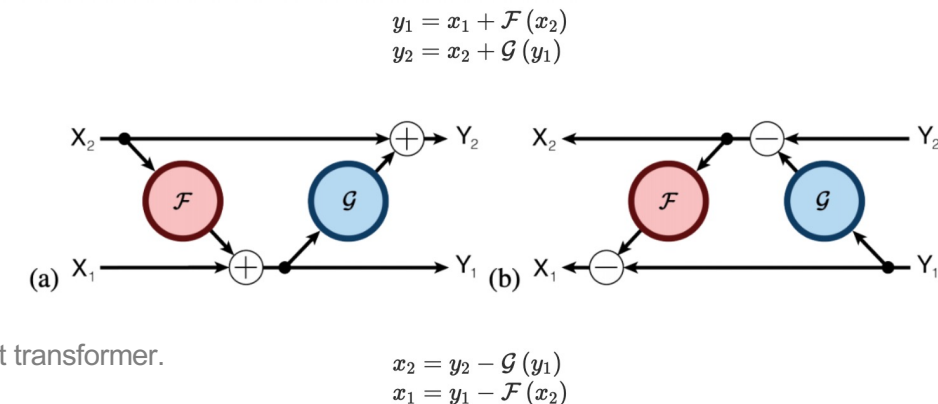


✂ Reversible layers

✂ ResNetに類似した構造

✂ 出力のactivationから入力のactivationを簡単な計算で復元

✂ レイヤー数Nに比例して増えるactivationを保存しておくメモリを削減



# CLIP [Radford+ (OpenAI), 21]

- ✖ インターネット上の画像とテキストのペアで事前学習することでラベル付けのコストを削減
  - ✖ 4億組の学習データ

## ✖ Contrastive Pre-Training

- ✖ 正解ラベルが複数の候補のうちどれかを選択する事前学習

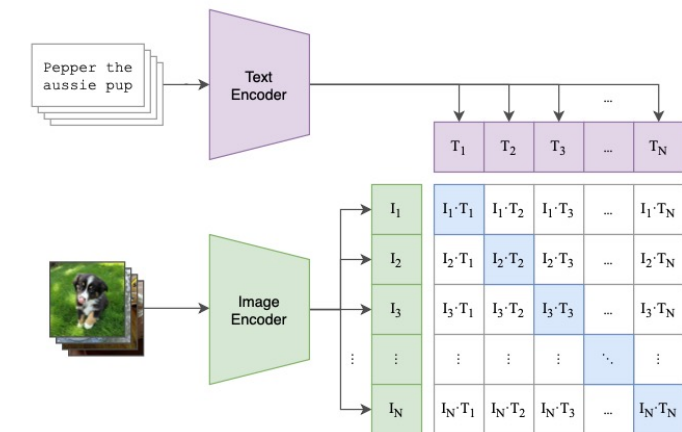
## ✖ Prompt Engineering

- ✖ 推論時、正解ラベルをA photo of {label}とすることで精度を向上
  - ✖ 事前学習時のラベルと近づく

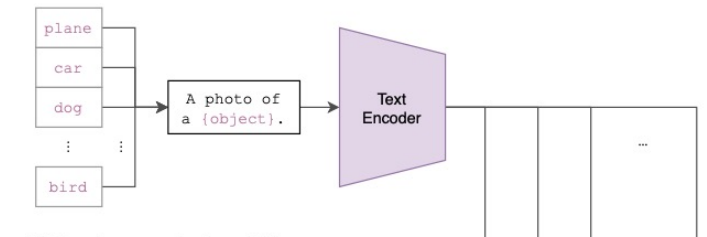
## ✖ Zero-shotで高精度の分類が可能

- ✖ 教師あり学習を行ったResNet-50と互角の性能
- ✖ 衛星画像や腫瘍検出等の複雑なタスクでは精度が下がる

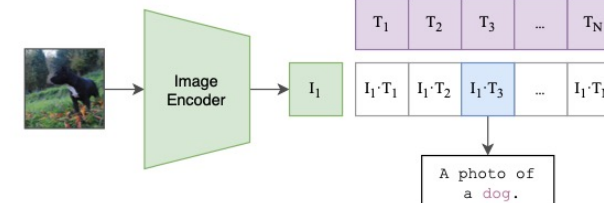
(1) Contrastive pre-training



(2) Create dataset classifier from label text

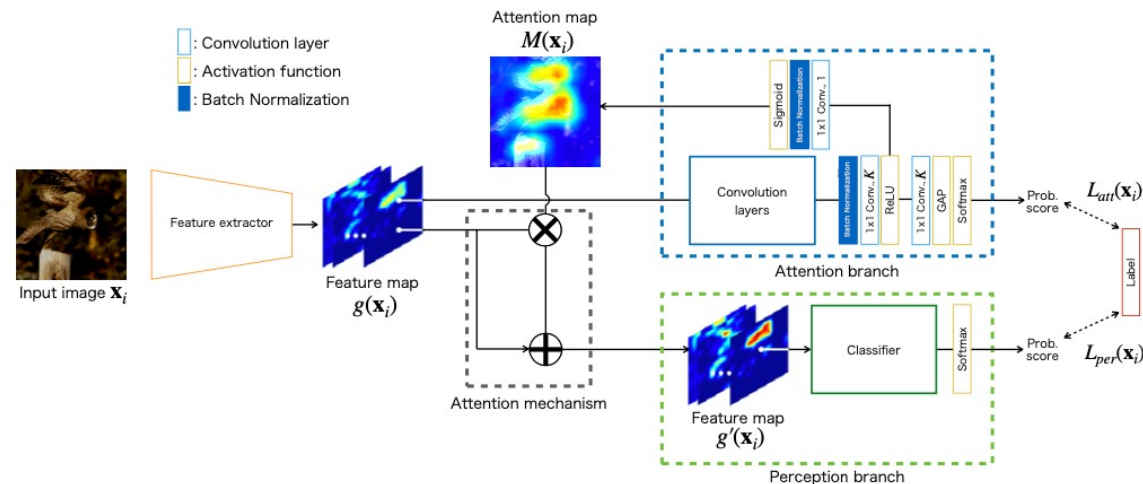


(3) Use for zero-shot prediction



# Attention Branch Network [Fukui+ (Chubu University), CVPR21]

- ✂ Attentionの可視化と精度向上を同時に実現
- ✂ CNNをFeature extractorとPerception branchに分割
- ✂ Feature extractorからAttention branchへ分岐
- ✂ Attention mapとFeature extractorの出力を合わせてPerception branchに入力
  - ✂ 両方のブランチでの損失を使用して学習可能



- ✂ 今後は、学習過程にラベルを含まない強化学習にもABNを適用する予定



# Iterative Shrinking with RL [Sun+ (University of Liverpool), CVPR21]

## ✂ 強化学習ベースの物体検出

✂ Proposal-freeのreferring expression grounding task

## ✂ 予測のboxを縮小する操作を強化学習のactionとして定義

✂ エージェントが縮小方向を決定

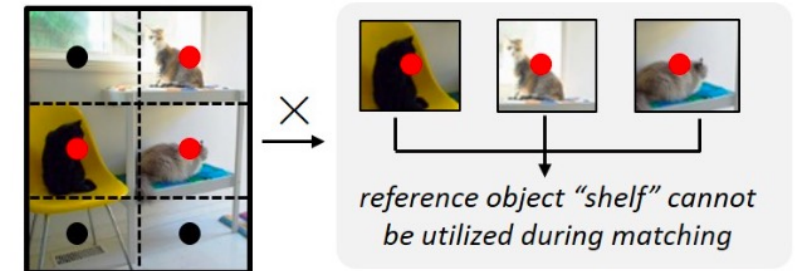
## ✂ 報酬

$$r_i = \begin{cases} 0 & IoU < 0.3 \text{ or } \Delta IoU \leq 0 \\ 1 & 0.3 \leq IoU < 0.5 \text{ and } 0 < \Delta IoU \\ 10 & 0.5 \leq IoU \text{ and } 0 < \Delta IoU \end{cases}$$

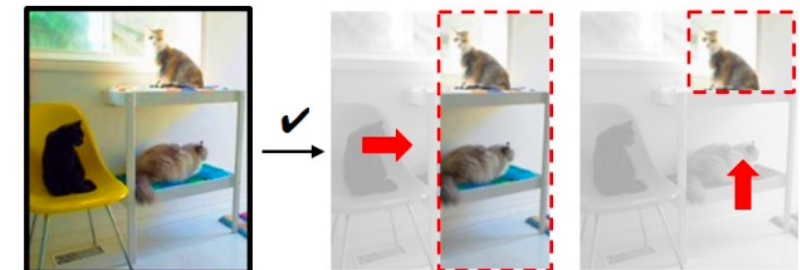
## ✂ 画像内のすべてのオブジェクトを総合的に判断

✂ 参照表現の理解に優位

「the cat above the shelf」



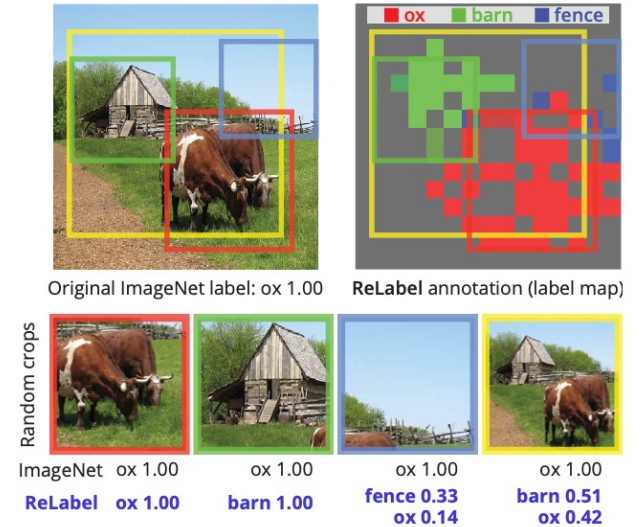
(a). feature-point level matching method



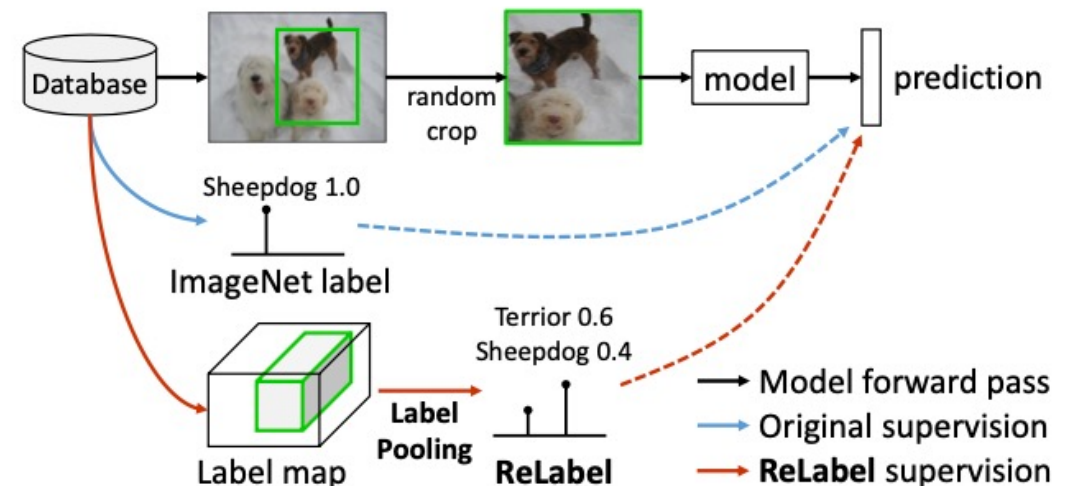
(b). proposed iterative shrinking method

# ReLabel [Yun+ (NAVER AI Lab), CVPR21]

- ✖ Image Net に多くのラベルノイズ
  - ✖ 多くのサンプルに複数のクラスが含まれている
  - ✖ Random crop時には正解が画像に含まれていないことも
- ✖ マルチラベルタスクとして画像ごとにアノテーション
- ✖ 膨大なデータセットで学習した識別器でlabel mapを作成



- ✖ LabelPooling
  - ✖ ランダムに切り抜かれた箇所に関して対応したラベルを学習に使用





# Copy-Paste [Golnaz Ghiasi+ (Google Research), CVPR21]

- ✖ 画像のスケーリング、コピー＆ペーストによる非常にシンプルなData Augmentation
- ✖ 片方の画像からランダムにオブジェクトのサブセットを選択し、もう片方の画像に貼り付け
- ✖ Instance SegmentationやObject Detectionで従来のデータ拡張より優れた結果



# SOHO [Huang+ (University of Science and Technology Beijing), CVPR21]

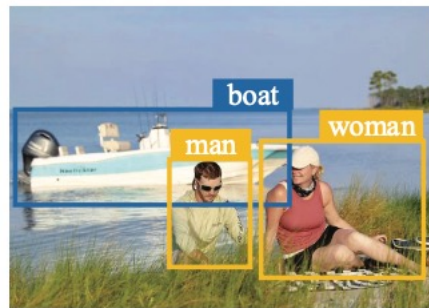
## ✖ Seeing Out of the bOx

## ✖ VL事前学習において自然言語のペアとして画像全体を入力

- ✖ bounding boxを抽出しないため推論を10倍以上高速化
- ✖ 検出した物体と周りの領域との関係を学習可能

## ✖ Visual Dictionary

- ✖ 視覚特徴をトークン化するために、類似した視覚特徴量を同じ特徴量に集約



### Task I: TR

**Baseline:** A couple **sit in a**

**boat** on the sea. ✖

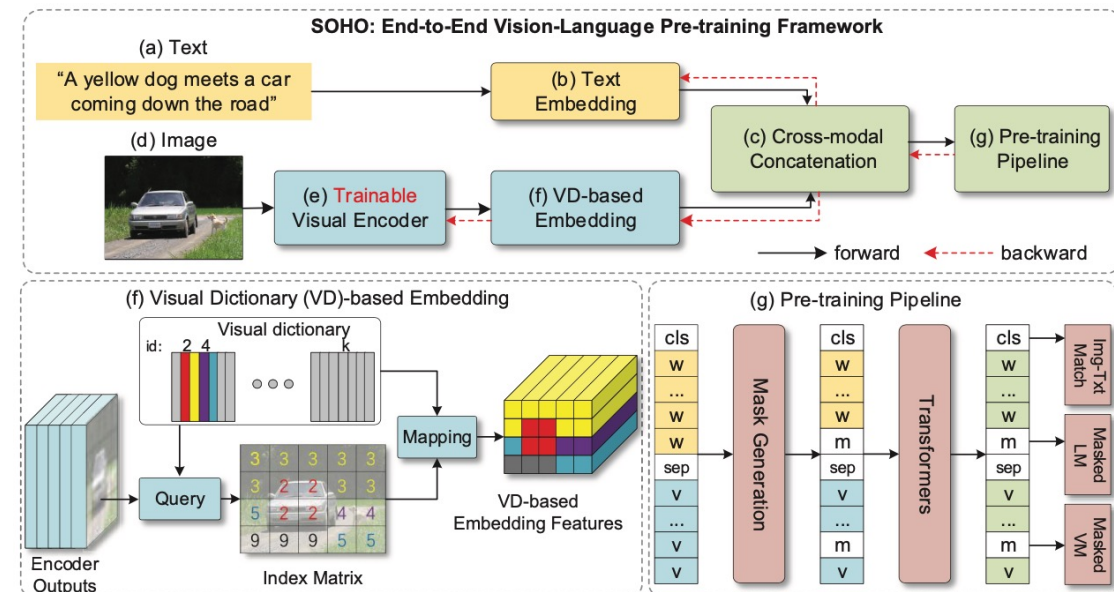
**Ours:** A couple **sit on the** **shore next to a boat** on the sea. ✔

### Task II: VQA

**Q:** What are the people doing?

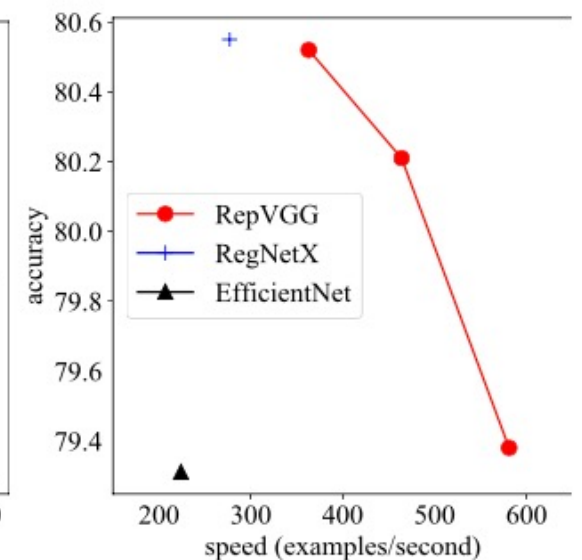
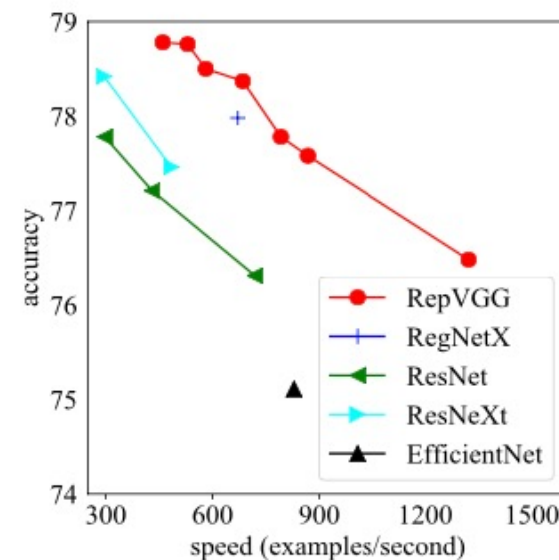
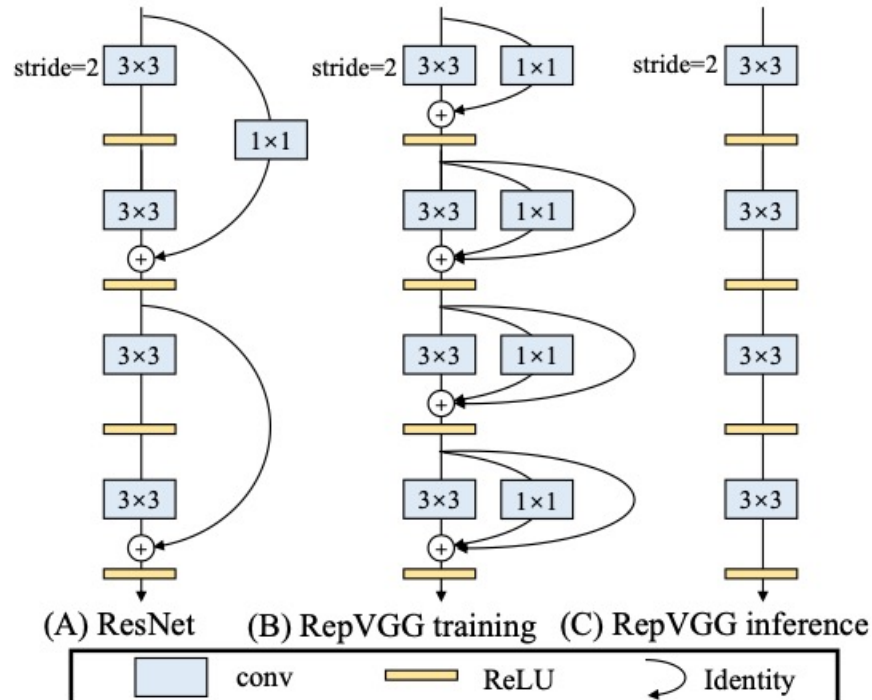
**Baseline:** Boating. ✖

**Ours:** Chatting. ✔



# RepVGG [Ding+ (BNRist), CVPR21]

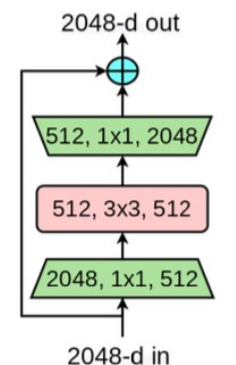
- ✖ 学習時は3x3conv、1x1conv、skip-connectionの3つのブランチを使用
- ✖ 推論時は分岐を使用しないシンプルなネットワーク
- ✖ ResNet-50よりも83%、ResNet-101よりも101%高速に動作し、高精度
  - ✖ EfficientNetのような最先端のモデルと比較しても、精度と速度の良好なトレードオフを実現



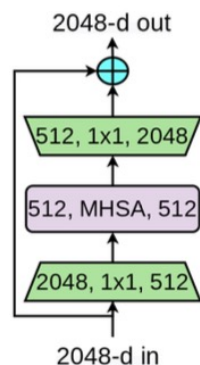


# BoTNet [Tsung-Yi Lin+ (Google Research), CVPR21]

- ✖ ResNetの最後の3つの $3 \times 3$ 畳み込みをmulti-head self-attentionに置き換え
  - ✖ パラメータを削減し、レイテンシのオーバーヘッドを最小限に抑える
  - ✖ ハイパーパラメータ等の変更が不要
- ✖ Instance segmentation、Object detectionでベースラインを大幅に改善
- ✖ EfficientNetと同等の精度を達成し、また1.64倍高速に動作



ResNet Bottleneck



Bottleneck Transformer

stage	output	ResNet-50	BoTNet-50
c1	$512 \times 512$	$7 \times 7, 64, \text{stride } 2$	$7 \times 7, 64, \text{stride } 2$
c2	$256 \times 256$	$3 \times 3 \text{ max pool, stride } 2$	$3 \times 3 \text{ max pool, stride } 2$
c3	$128 \times 128$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
c4	$64 \times 64$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
c5	$32 \times 32$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
c6	$16 \times 16$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ \text{MHSA}, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
# params.		$25.5 \times 10^6$	$20.8 \times 10^6$
M.Adds		$85.4 \times 10^9$	$102.98 \times 10^9$
TPU steptime		786.5 ms	1032.66 ms

