

# Learning Visual Context by Comparison

Minchul Kim al., Lunit.

ECCV 2020 spotlight paper

2020-10-26

MI2RL

Kyuri Kim

# Introduction

## - Contributions

- We present Attend-and Compare Module (ACM) for capturing the difference between an object of interest and its corresponding context. **(The necessity of comparison between related regions in an image)**

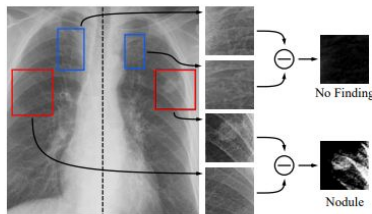


Fig. 1: An example of a comparison procedure for radiologists. Little differences indicate no disease (blue), the significant difference is likely to be a lesion (red).

- (1) We propose a novel context module called ACM that explicitly compares different regions, following [the way radiologists read](#) chest X-rays.
- (2) The proposed ACM captures [multiple comparative self-attentions](#) whose difference is beneficial to recognition tasks.
- (3) We demonstrate the effectiveness of ACM on three chest X-ray datasets and COCO detection & segmentation dataset with various architectures.

# Related Work

## - Context Modeling

In the visual recognition domain recent self-attention mechanisms,

- SENet(Squeeze and excitation networks) - CVPR 2018
- A Style-based Re-calibration Module (SRM) (Lunit)
- Convolutional block attention module (CBAM) - ECCV 2018 (Lunit)

Works that explicitly tackle the problem of using context stem from using pixel-level pairwise relationships,

- Non-local neural networks (NL) - CVPR 2018 (Kaiming He)
- Global-Context network (GC) - 2019 (Microsoft)
- Criss-cross attention (CC) - IEEE TPAMI 2020 & ICCV 2019

# Related Work

## - Context Modeling

- SENet(Squeeze and excitation networks)
  - “Our goal is to improve the representational power of a network by explicitly modelling the interdependencies between the channels of its convolutional features.”
  - Learns to model channel-wise attention using the spatially averaged feature.
- Non-local neural networks (NL)
  - Calculate pixel-level pairwise relationship weights and aggregate (weighted average) the features from all locations according to the weights.

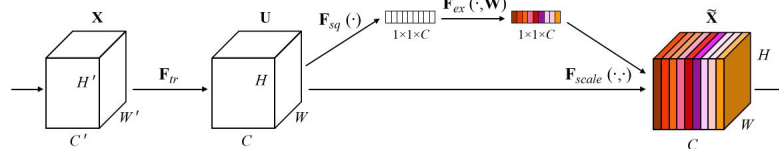


Figure 1: A Squeeze-and-Excitation block.

Weighted sum of **All** pixels with **similarity**

$$y_i = \frac{1}{C(x)} \sum_j f(x_i, x_j) g(x_j)$$

$$g(x_j) = W_g x_j$$

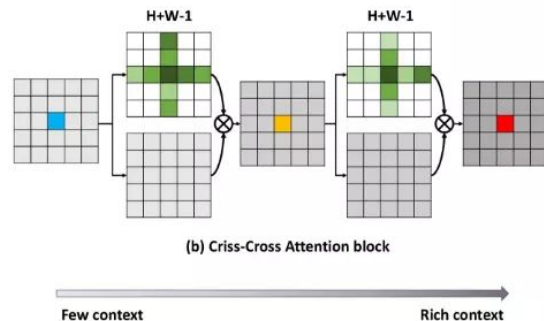
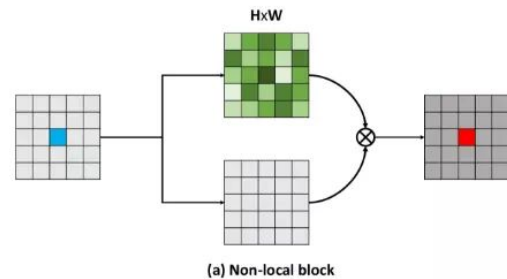
weight\*input pixel

# Related Work

## - Context Modeling

- Criss-cross attention (CC)

- For semantic segmentation reduces the computation cost of NL by replacing the pairwise relationship attention maps with criss-cross attention block which considers only horizontal and vertical directions separately.
- NL and CC explicitly model the pairwise relationship between regions with affinity metrics, but the qualitative results in demonstrate a tendency to aggregate features only among foreground objects or among pixels with similar semantics.



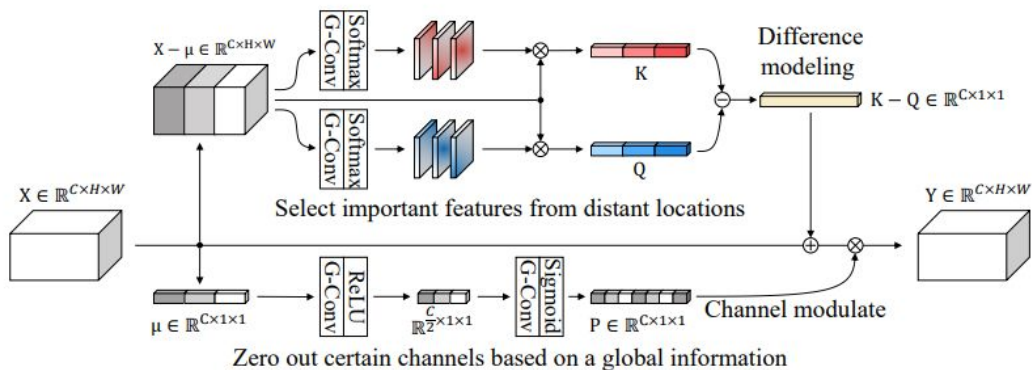
# Method

## - Attend-and-Compare Module

### 1. Overview

- Attend-and-Compare Module (ACM) extracts an object of interest and the corresponding context to compare, and enhances the original image feature with the comparison result.

$$Y = f_{\text{ACM}}(X) = P(X + (K - Q)), \quad (1)$$



# Method

## - Attend-and-Compare Module

### 2. Components of ACM

- Object of Interest and Corresponding Context

$$K = \sum_{i,j \in H,W} \frac{\exp(W_K X_{i,j})}{\sum_{H,W} \exp(W_K X_{h,w})} X_{i,j}, \quad (2)$$

① 1x1 conv for single ch.

② softmax for Normalize

③ Apply weighted avg.

- Channel Re-calibration

$$P = \sigma \circ \text{conv}_2^{1 \times 1} \circ \text{ReLU} \circ \text{conv}_1^{1 \times 1}(\mu), \quad (3)$$

# Method

## - Attend-and-Compare Module

### 2. Components of ACM

- Group Operation

$$K^g = \sum_{i,j \in H,W} \frac{\exp(W_K^g X_{i,j}^g)}{\sum_{H,W} \exp(W_K^g X_{h,w}^g)} X_{i,j}^g, \quad (4)$$

- group conv. (#G)
- $K = [K_1, \dots, K_G]$

- Loss Function

$$\ell_{\text{orth}}(K, Q) = \frac{K \cdot Q}{C}, \quad (5)$$

- C = channel num.

$$\ell_{\text{task}} + \lambda \sum_m^M \ell_{\text{orth}}(K_m, Q_m), \quad (6)$$

- Placement of ACMs

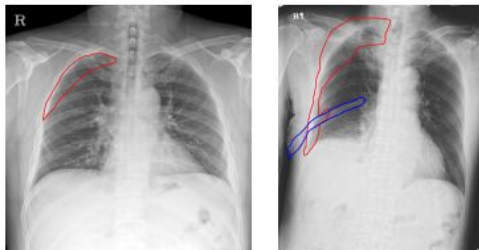


# Experiments

## - Experiment Dataset

- 1) Emergency-Pneumothorax (Em-Ptx) and Nodule (Ndl) datasets for lesion localization in chest X-rays
- 2) Chest X-ray 14 dataset for multi-label classification
- 3) COCO 2017 dataset for object detection and instance segmentation.

### <Emergency-Pneumothorax (Em-Ptx) Datasets>



(a) Emergency (b) Non-emergency

Table 1: Results on Em-Ptx dataset. Average of 5 random runs are reported for each setting with standard deviation. RN stands for ResNet [14].

Method	AUC-ROC	JAFROC	Method	AUC-ROC	JAFROC
RN-50	86.78 $\pm$ 0.58	81.84 $\pm$ 0.64	RN-101	89.75 $\pm$ 0.49	85.36 $\pm$ 0.44
RN-50 + SE [15]	93.05 $\pm$ 3.63	89.19 $\pm$ 4.38	RN-101 + SE [15]	90.36 $\pm$ 0.83	85.54 $\pm$ 0.85
RN-50 + NL [36]	94.63 $\pm$ 0.39	91.93 $\pm$ 0.68	RN-101 + NL [36]	94.24 $\pm$ 0.34	91.70 $\pm$ 0.83
RN-50 + CC [17]	87.73 $\pm$ 8.66	83.32 $\pm$ 10.36	RN-101 + CC [17]	92.57 $\pm$ 0.89	89.75 $\pm$ 0.89
RN-50 + ACM	<b>95.35<math>\pm</math>0.12</b>	<b>94.16<math>\pm</math>0.21</b>	RN-101 + ACM	<b>95.43<math>\pm</math>0.14</b>	<b>94.47<math>\pm</math>0.10</b>

# Results

Table 2: Performance with respect to varying module architectures and hyper-parameters on Em-Ptx dataset. All the experiments are based on ResNet-50.

Module	AUC-ROC	JAFROC
None	86.78 $\pm$ 0.58	81.84 $\pm$ 0.64
$X + (K - Q)$	94.25 $\pm$ 0.31	92.94 $\pm$ 0.36
$PX$	87.16 $\pm$ 0.42	82.05 $\pm$ 0.30
$P(X + K)$	94.96 $\pm$ 0.15	93.59 $\pm$ 0.24
$P(X + (K - Q))$	<b>95.35<math>\pm</math>0.12</b>	<b>94.16<math>\pm</math>0.21</b>

(a) Ablations on  $K, Q$  and  $P$ .

#groups	AUC-ROC	JAFROC
8	90.96 $\pm$ 1.88	88.79 $\pm$ 2.23
32	<b>95.35<math>\pm</math>0.12</b>	<b>94.16<math>\pm</math>0.21</b>
64	95.08 $\pm$ 0.25	93.73 $\pm$ 0.31
128	94.89 $\pm$ 0.53	92.88 $\pm$ 0.53

(b) Ablations on number of groups.

$\lambda$	AUC-ROC	JAFROC
0.00	95.11 $\pm$ 0.20	93.87 $\pm$ 0.20
0.01	95.29 $\pm$ 0.34	94.09 $\pm$ 0.41
0.10	<b>95.35<math>\pm</math>0.12</b>	<b>94.16<math>\pm</math>0.21</b>
1.00	95.30 $\pm$ 0.17	94.04 $\pm$ 0.11

(c) Ablations on orthogonal loss weight  $\lambda$ .

## <Nodule (Ndl) Datasets>

Table 3: Results on Ndl dataset. Average of 5 random runs are reported for each setting with standard deviation.

Method	AUC-ROC	JAFROC
ResNet-50	87.34 $\pm$ 0.34	77.35 $\pm$ 0.50
ResNet-50 + SE [15]	87.66 $\pm$ 0.40	77.57 $\pm$ 0.44
ResNet-50 + NL [36]	88.35 $\pm$ 0.35	80.51 $\pm$ 0.56
ResNet-50 + CC [17]	87.72 $\pm$ 0.18	78.63 $\pm$ 0.40
ResNet-50 + ACM	<b>88.60<math>\pm</math>0.23</b>	<b>83.03<math>\pm</math>0.24</b>

# Results

## <Chest X-ray 14 Datasets>

Table 5: Performance in average AUC of various methods on CXR14 dataset. The numbers in the bracket after model names are the input sizes.

Modules	DenseNet121(448)	ResNet-50(448)
None	(CheXNet [29]) 84.54	84.19
SE [15]	84.95	84.53
NL [36]	84.49	85.08
CC [17]	84.43	85.11
ACM	<b>85.03</b>	<b>85.39</b>

## <COCO Datasets>

Table 6: Results on COCO dataset. All experiments are based on Mask-RCNN [13].

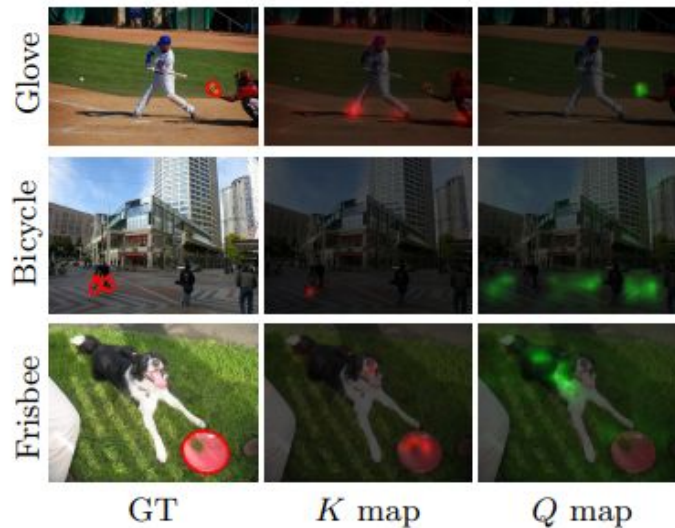
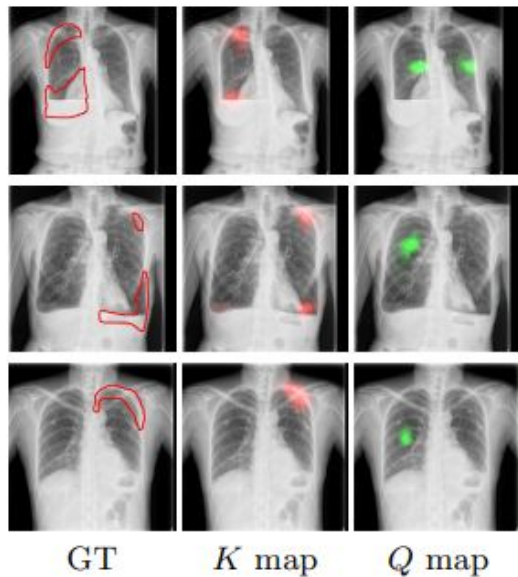
Method	AP <sup>bbox</sup>	AP <sup>bbox</sup> <sub>50</sub>	AP <sup>bbox</sup> <sub>75</sub>	AP <sup>mask</sup>	AP <sup>mask</sup> <sub>50</sub>	AP <sup>mask</sup> <sub>75</sub>
ResNet-50	38.59	59.36	42.23	35.24	56.24	37.66
ResNet-50+SE [15]	39.10	60.32	42.59	35.72	57.16	38.20
ResNet-50+NL [36]	39.40	60.60	43.02	35.85	57.63	38.15
ResNet-50+CC [17]	39.82	60.97	42.88	36.05	57.82	38.37
ResNet-50+ACM	<b>39.94</b>	<b>61.58</b>	<b>43.30</b>	<b>36.40</b>	<b>58.40</b>	<b>38.63</b>
ResNet-101	40.77	61.67	44.53	36.86	58.35	39.59
ResNet-101+SE [15]	41.30	62.36	45.26	37.38	59.34	40.00
ResNet-101+NL [36]	41.57	62.75	45.39	37.39	59.50	40.01
ResNet-101+CC [17]	<b>42.09</b>	63.21	<b>45.79</b>	<b>37.77</b>	59.98	<b>40.29</b>
ResNet-101+ACM	41.76	<b>63.38</b>	45.16	37.68	<b>60.16</b>	40.19
ResNeXt-101	43.23	64.42	47.47	39.02	61.10	42.11
ResNeXt-101+SE [15]	43.44	64.91	47.66	39.20	61.92	42.17
ResNeXt-101+NL [36]	43.93	65.44	48.20	39.45	61.99	42.33
ResNeXt-101+CC [17]	43.86	65.28	47.74	39.26	62.06	41.97
ResNeXt-101+ACM	<b>44.07</b>	<b>65.92</b>	<b>48.33</b>	<b>39.54</b>	<b>62.53</b>	<b>42.44</b>

# Results

## - Qualitative Results

(L) Utilize pneumothorax regions as objects of interest and normal lung regions as the corresponding context.

(R) Utilize the object of interest and the corresponding context in the natural image domain.



# Conclusions

- Key idea is to extract an object of interest and a corresponding context and explicitly compare them to make the image representation more distinguishable.
- The qualitative analysis shows that ACM automatically learns dynamic relationships. The objects of interest and corresponding contexts are different yet contain useful information for the given task.