

# Bilinear Attention Networks

2018 VQA Challenge runner-up  
(1st single model)



*Jin-Hwa Kim, Jaehyun Jun, Byoung-Tak Zhang*

Biointelligence Lab.  
Seoul National University

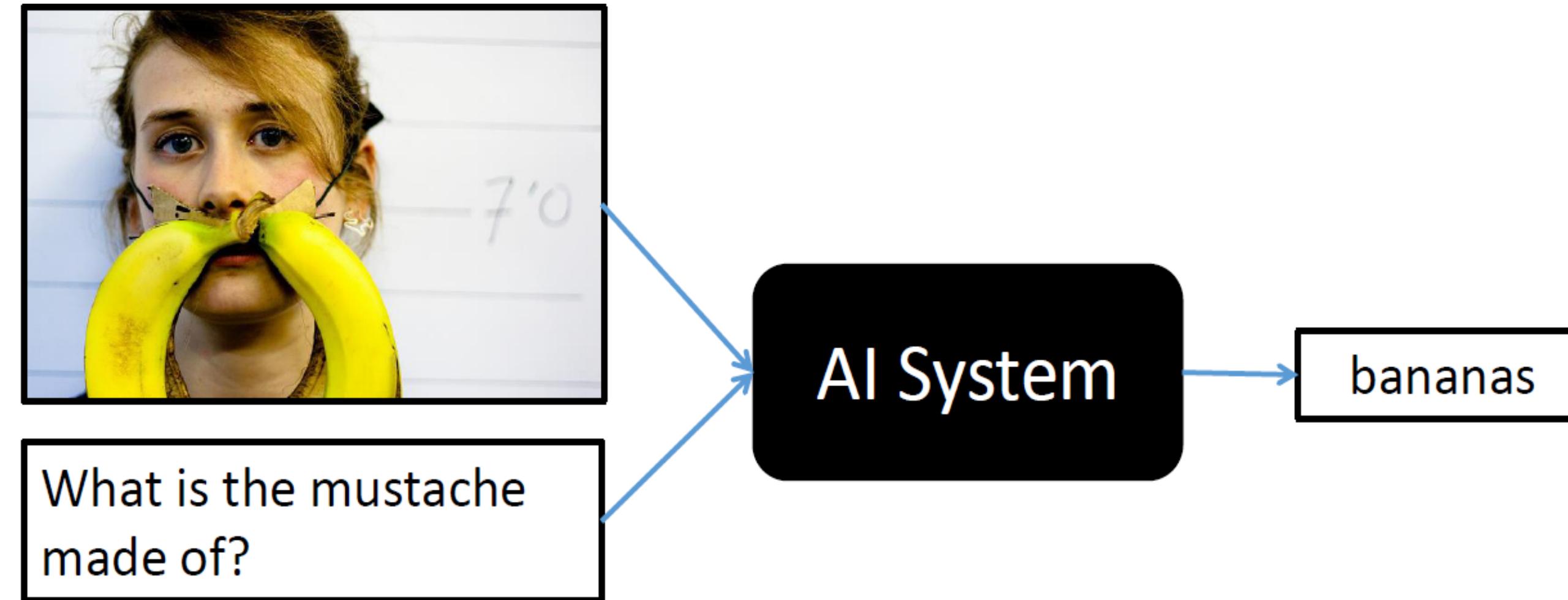


SEOUL  
NATIONAL  
UNIVERSITY



# Visual Question Answering

- Learning joint representation of multiple modalities
- Linguistic and visual information



# VQA 2.0 Dataset

- 204K MS COCO images
- 760K questions (10 answers for each question from unique AMT workers)
- train, val, test-dev (remote validation), test-standard (publications)
- and test-challenge (competitions), test-reserve (check overfitting)

$$accuracy = \min\left(\frac{\# \text{ humans that provided that answer}}{3}, 1\right)$$

- VQA Score is the avg. of 10 choose 9 accuracies →

<https://github.com/GT-Vision-Lab/VQA/issues/1>

<https://github.com/hengyuan-hu/bottom-up-attention-vqa/pull/18>

	Images	Questions	Answers
Train	80K	443K	4.4M
Val	40K	214K	2.1M
Test	80K	447K	unknown

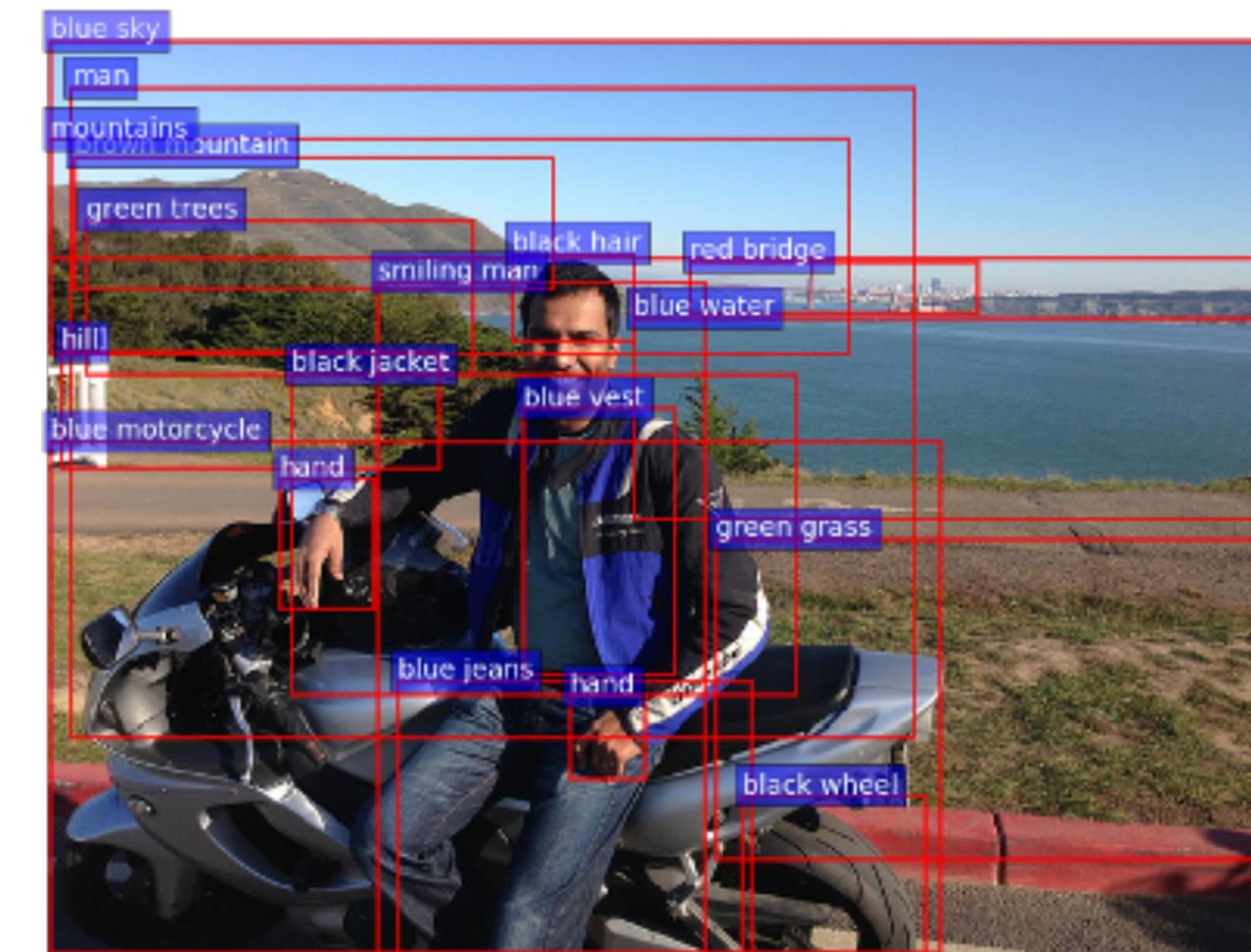
# annotations	VQA Score
0	0.0
1	0.3
2	0.6
3	0.9
>3	1.0

# Objective

- Introducing bilinear attention
  - *Interactions between words and visual concepts* are meaningful
  - Proposing an efficient method (with the same time complexity) on top of low-rank bilinear pooling
- Residual learning of attention
  - Residual learning with attention mechanism for incremental inference
- Integration of counting module (Zhang et al., 2018)

# Preliminary

- Question embedding (fine-tuning)
  - Use the **all outputs of GRU** (every time steps)
  - $X \in \mathbb{R}^{N \times p}$  where N is hidden dim., and  $p=|\{x_i\}|$  is # of tokens
- Image embedding (fixed **bottom-up-attention**)
  - Select 10-100 detected objects (rectangles) using pre-trained Faster RCNN, to extract *rich* features for each object (1600 classes, **400 attributes**)
  - $Y \in \mathbb{R}^{M \times \phi}$  where M is feature dim., and  $\phi=|\{y_j\}|$  is # of objects



# Low-rank Bilinear Pooling

- Bilinear model and its approximation (Wolf et al., 2007, Pirsiaavash et al., 2009)

$$f_i = \mathbf{x}^T \mathbf{W}_i \mathbf{y} \approx \mathbf{x}^T \mathbf{U}_i \mathbf{V}_i^T \mathbf{y} = \mathbf{1}^T (\mathbf{U}_i^T \mathbf{x} \circ \mathbf{V}_i^T \mathbf{y})$$

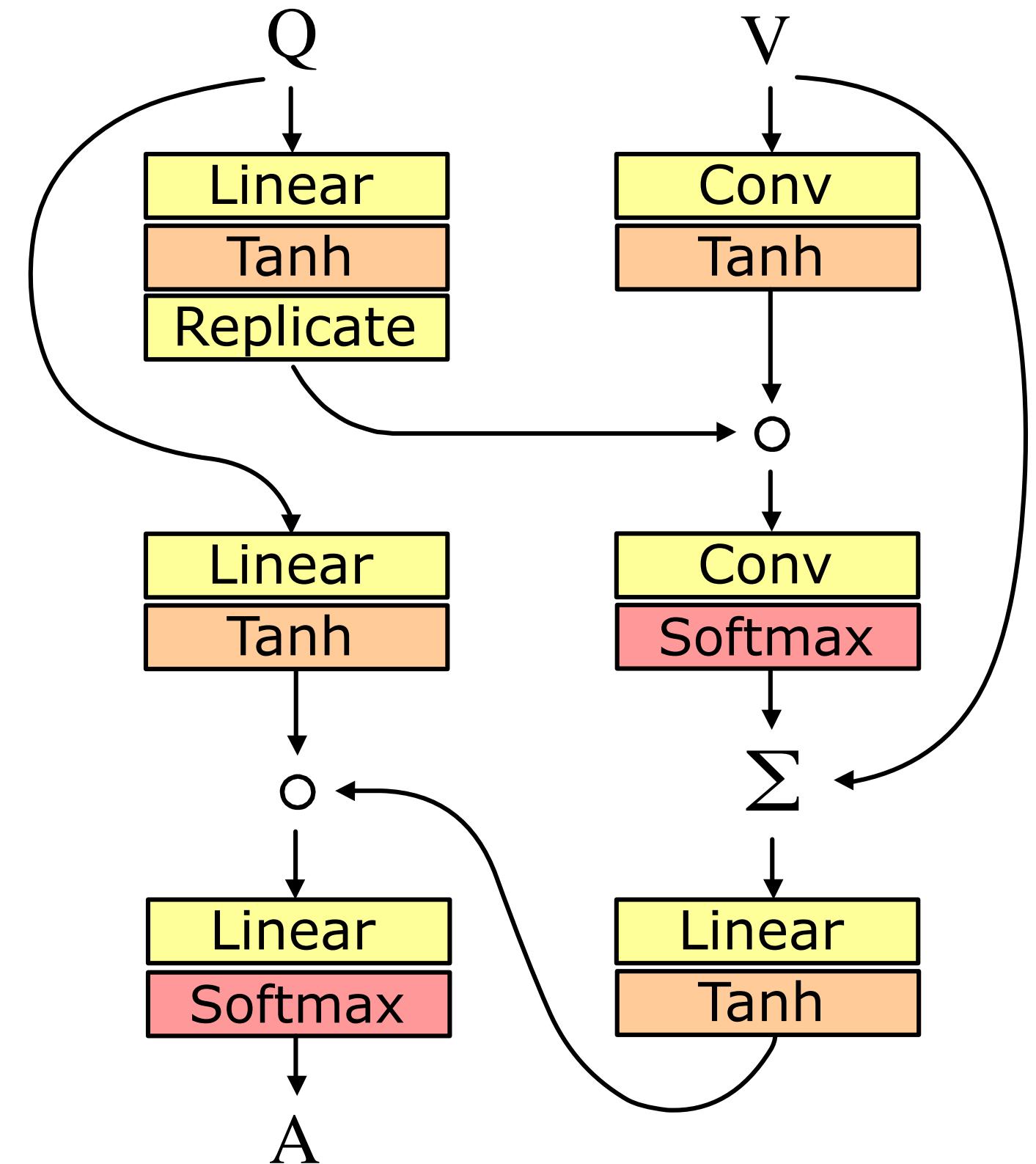
- Low-rank bilinear pooling (Kim et al., 2017)

$$\mathbf{f} = \mathbf{P}^T (\mathbf{U}^T \mathbf{x} \circ \mathbf{V}^T \mathbf{y})$$

For vector output, instead of using three-dimensional tensors  $\mathbf{U}$  and  $\mathbf{V}$ , replace the vector of ones with a pooling matrix  $\mathbf{P}$  (use three two-dimensional tensors).

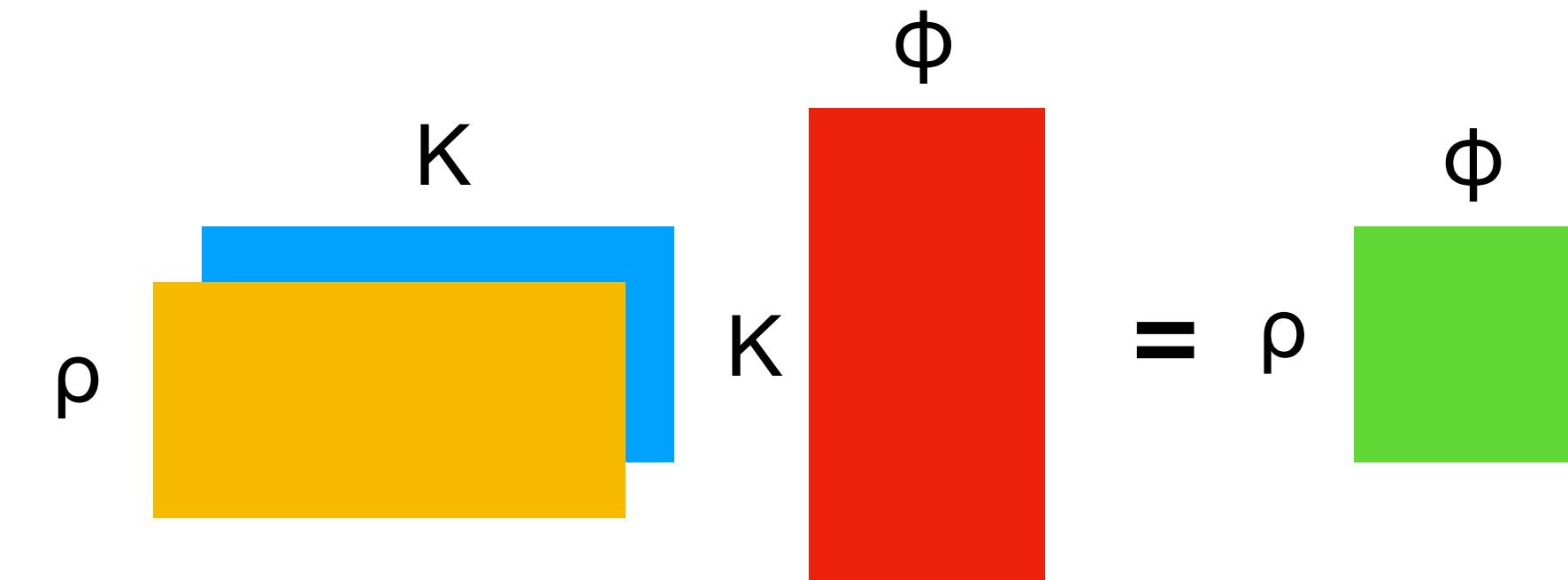
# Unitary Attention

- This pooling is used to get attention weights with a question embedding (single-channel) vector and visual feature vectors (multi-channel) as the two inputs.
- We call it *unitary attention* since a question embedding vector queried the feature vectors, *unidirectionally*.



# Bilinear Attention Maps

- $\mathbf{U}$  and  $\mathbf{V}$  are linear embeddings
- $\mathbf{p}$  is a learnable projection vector

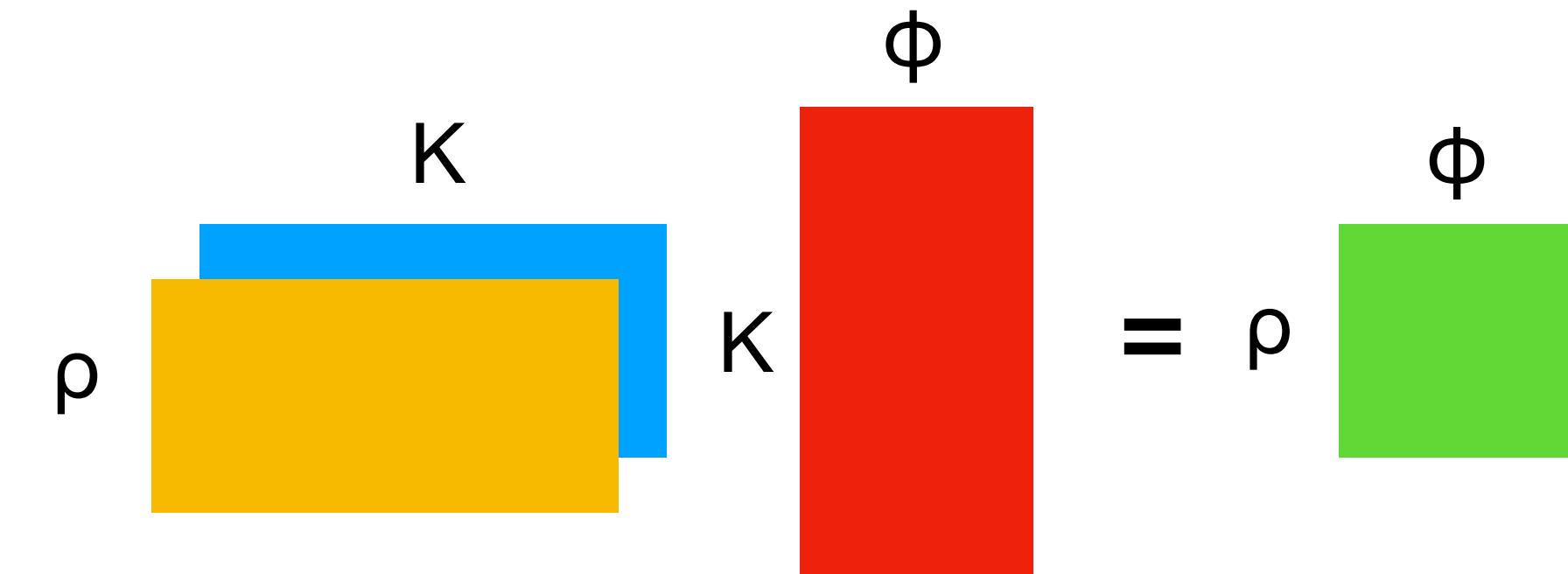


$$\mathcal{A} := \text{softmax} \left( \left( (1 \cdot \mathbf{p}^T) \circ \mathbf{X}^T \mathbf{U} \right) \mathbf{V}^T \mathbf{Y} \right)$$

element-wise multiplication  
↓  
 $\rho \times \Phi$        $\rho \times K$        $\rho \times N$        $N \times K$        $K \times M$        $M \times \Phi$   
 $\rho \times K$        $K \times \Phi$

# Bilinear Attention Maps

- Exactly the same approach with low-rank bilinear pooling



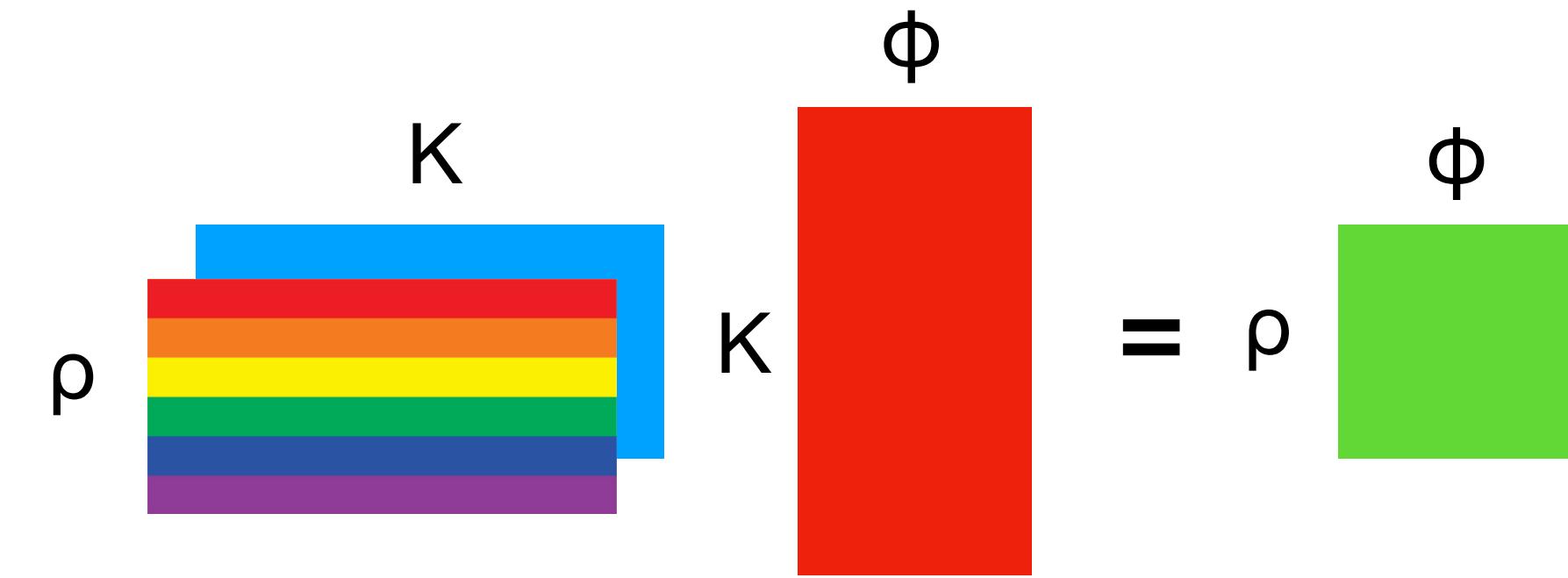
$$\mathcal{A} := \text{softmax}\left( \left( (\mathbf{1} \cdot \mathbf{p}^T) \circ \mathbf{X}^T \mathbf{U} \right) \mathbf{V}^T \mathbf{Y} \right)$$

element-wise multiplication  
↓

$$A_{i,j} = \mathbf{p}^T \left( (\mathbf{U}^T \mathbf{X}_i) \circ (\mathbf{V}^T \mathbf{Y}_j) \right).$$

# Bilinear Attention Maps

- Multiple bilinear attention maps are acquired by different projection vectors  $\mathbf{p}_g$  as:



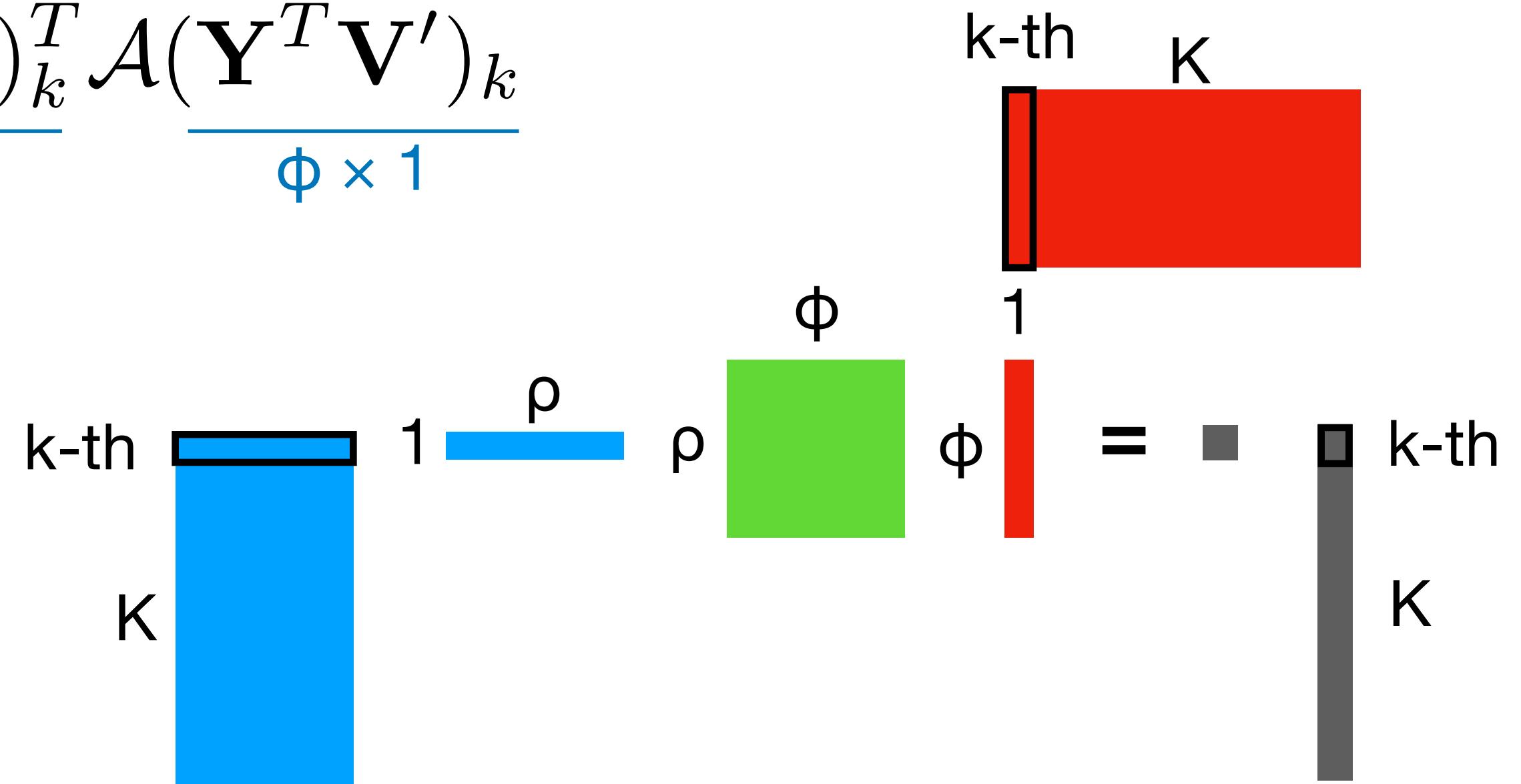
$$\mathcal{A}_g := \text{softmax} \left( \left( (\mathbf{1} \cdot \mathbf{p}_g^T) \circ \mathbf{X}^T \mathbf{U} \right) \mathbf{V}^T \mathbf{Y} \right)$$

↑

# Bilinear Attention Networks

- Each multimodal joint feature is filled with following equation ( $k$  is the index of  $K$ ; *broadcasting* in PyTorch let you avoid for-loop for this):

$$\mathbf{f}'_k = \frac{(\mathbf{X}^T \mathbf{U}')_k^T}{1 \times \rho} \mathcal{A} \frac{(\mathbf{Y}^T \mathbf{V}')_k}{\phi \times 1}$$



\* *broadcasting: automatically repeat tensor operations in api-level supported by Numpy, Tensorflow, Pytorch*

# Bilinear Attention Networks

- One can show that this is equivalent to a bilinear attention model where each feature is pooled by low-rank bilinear approximation

$$\mathbf{f}'_k = (\mathbf{X}^T \mathbf{U}')_k^T \mathcal{A} (\mathbf{Y}^T \mathbf{V}')_k$$

$$\begin{aligned}\mathbf{f}'_k &= \sum_{i=1}^{|\{\mathbf{x}_i\}|} \sum_{j=1}^{|\{\mathbf{y}_j\}|} \mathcal{A}_{i,j} (\mathbf{X}_i^T \mathbf{U}'_k) (\mathbf{V}'^T \mathbf{Y}_j) \\ &= \sum_{i=1}^{|\{\mathbf{x}_i\}|} \sum_{j=1}^{|\{\mathbf{y}_j\}|} \mathcal{A}_{i,j} \underline{\mathbf{X}_i^T (\mathbf{U}'_k \mathbf{V}'^T) \mathbf{Y}_j} \\ &\quad \text{low-rank bilinear pooling}\end{aligned}$$

# Bilinear Attention Networks

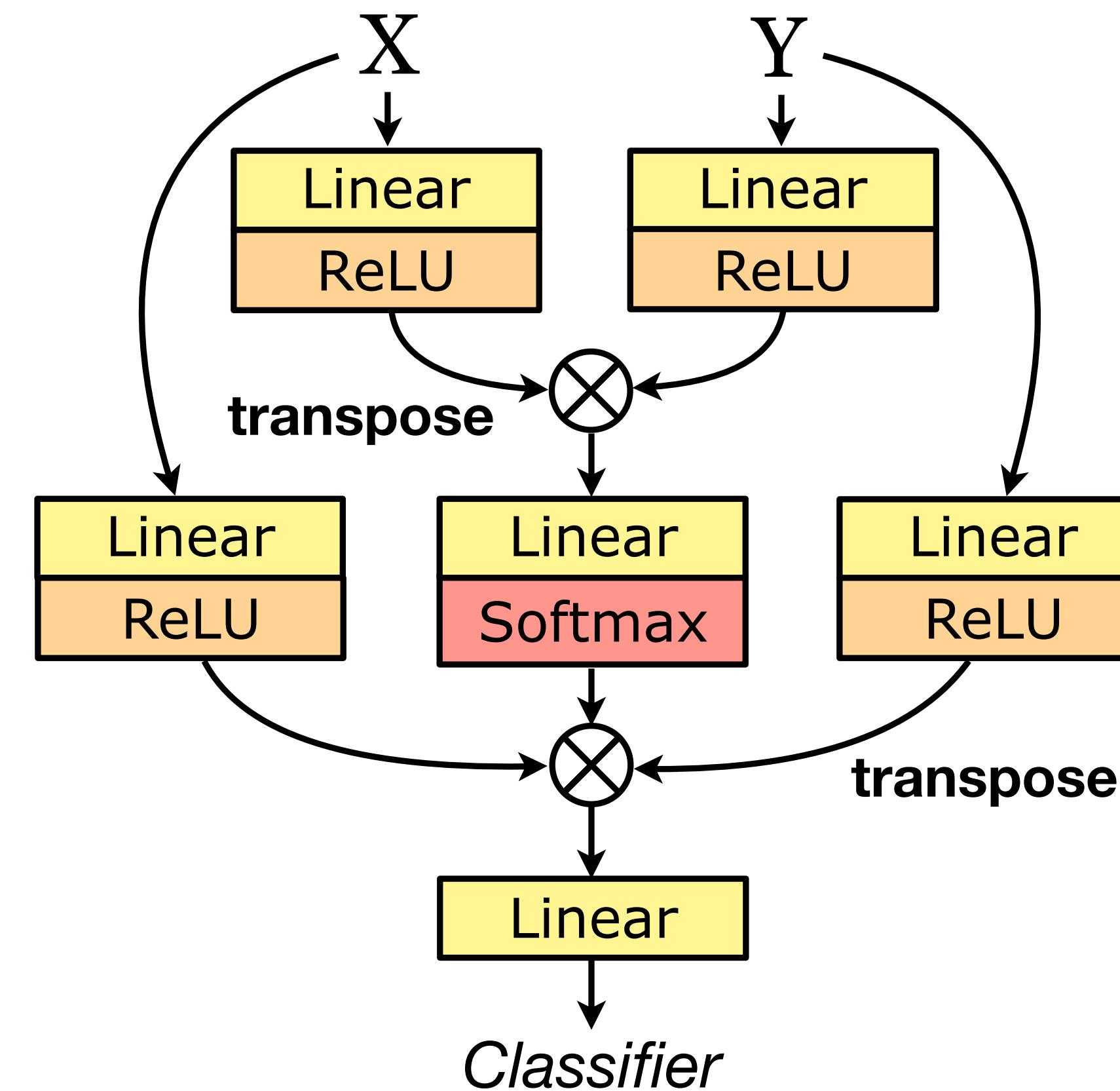
- One can show that this is equivalent to a bilinear attention model where each feature is pooled by low-rank bilinear approximation
- Low-rank bilinear feature learning ***inside*** bilinear attention

$$\begin{aligned}\mathbf{f}'_k &= \sum_{i=1}^{|\{\mathbf{x}_i\}|} \sum_{j=1}^{|\{\mathbf{y}_j\}|} \mathcal{A}_{i,j} (\mathbf{X}_i^T \mathbf{U}'_k) (\mathbf{V}'^T_k \mathbf{Y}_j) \\ &= \sum_{i=1}^{|\{\mathbf{x}_i\}|} \sum_{j=1}^{|\{\mathbf{y}_j\}|} \mathcal{A}_{i,j} \underline{\mathbf{X}_i^T (\mathbf{U}'_k \mathbf{V}'^T_k) \mathbf{Y}_j} \\ &\quad \text{low-rank bilinear pooling}\end{aligned}$$

# Bilinear Attention Networks

- One can show that this is equivalent to a bilinear attention model where each feature is pooled by low-rank bilinear approximation
- Low-rank bilinear feature learning *inside* bilinear attention
- Similarly to MLB (Kim et al., ICLR 2017), activation functions can be applied

# Bilinear Attention Networks



# Time Complexity

- Assuming that  $M \geq N > K > \phi \geq \rho$ , the time complexity of bilinear attention networks is  $O(KM\phi)$  where  $K$  denotes hidden size, since BAN consists of **matrix chain multiplication**
- Empirically, BAN takes 284s/epoch while unitary attention control takes 190s/epoch
- Largely due to the increment of input size for Softmax function,  $\phi$  to  $\phi \times \rho$

# Residual Learning of Attention

- Residual learning exploits multiple attention maps ( $\mathbf{f}_0$  is  $\mathbf{X}$  and  $\{\mathbf{a}_i\}$  is fixed to ones):

$$\mathbf{f}_{i+1} = \text{BAN}_i(\mathbf{f}_i, \mathbf{Y}; \mathcal{A}_i) \cdot \mathbf{1}^T + \underbrace{\alpha_i \mathbf{f}_i}_{\text{shortcut}}$$

bilinear attention map

↑

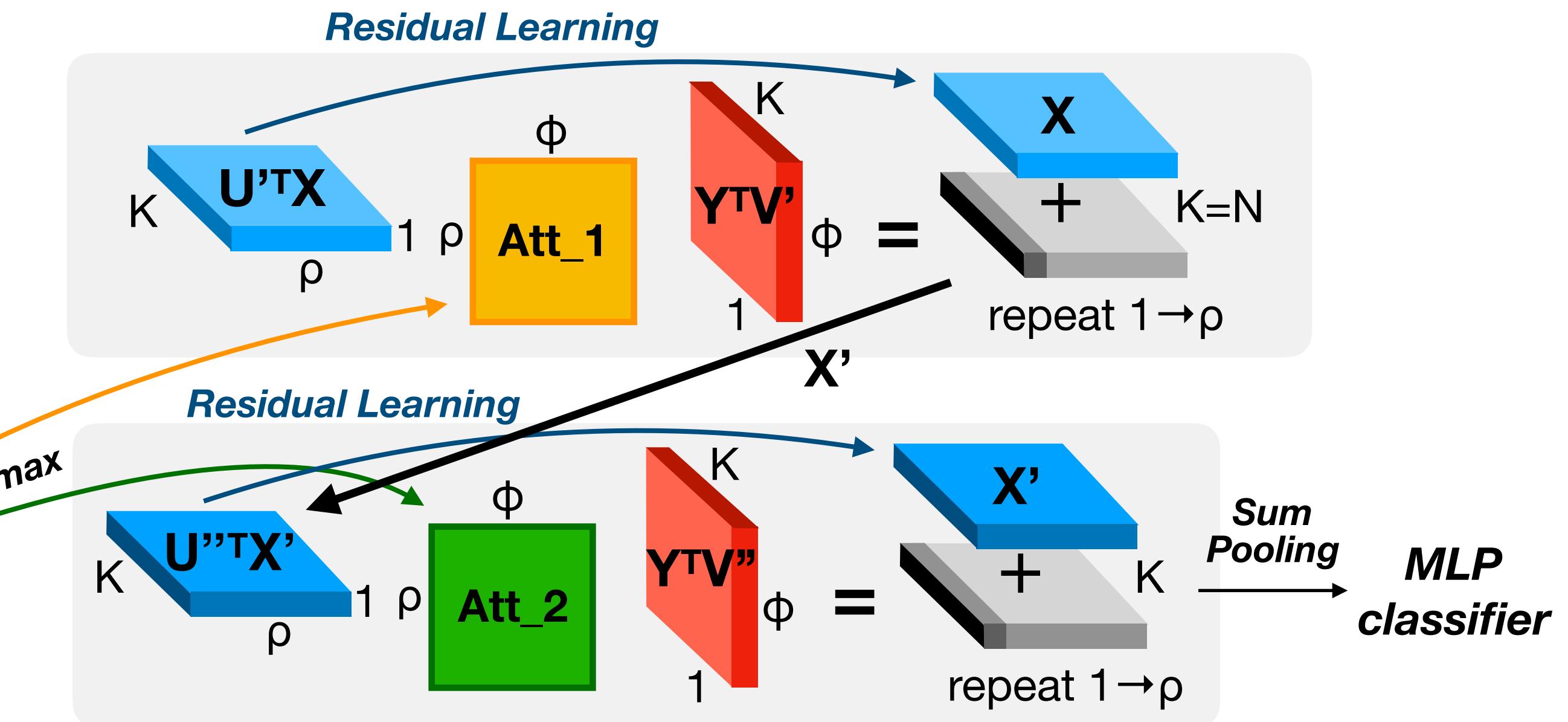
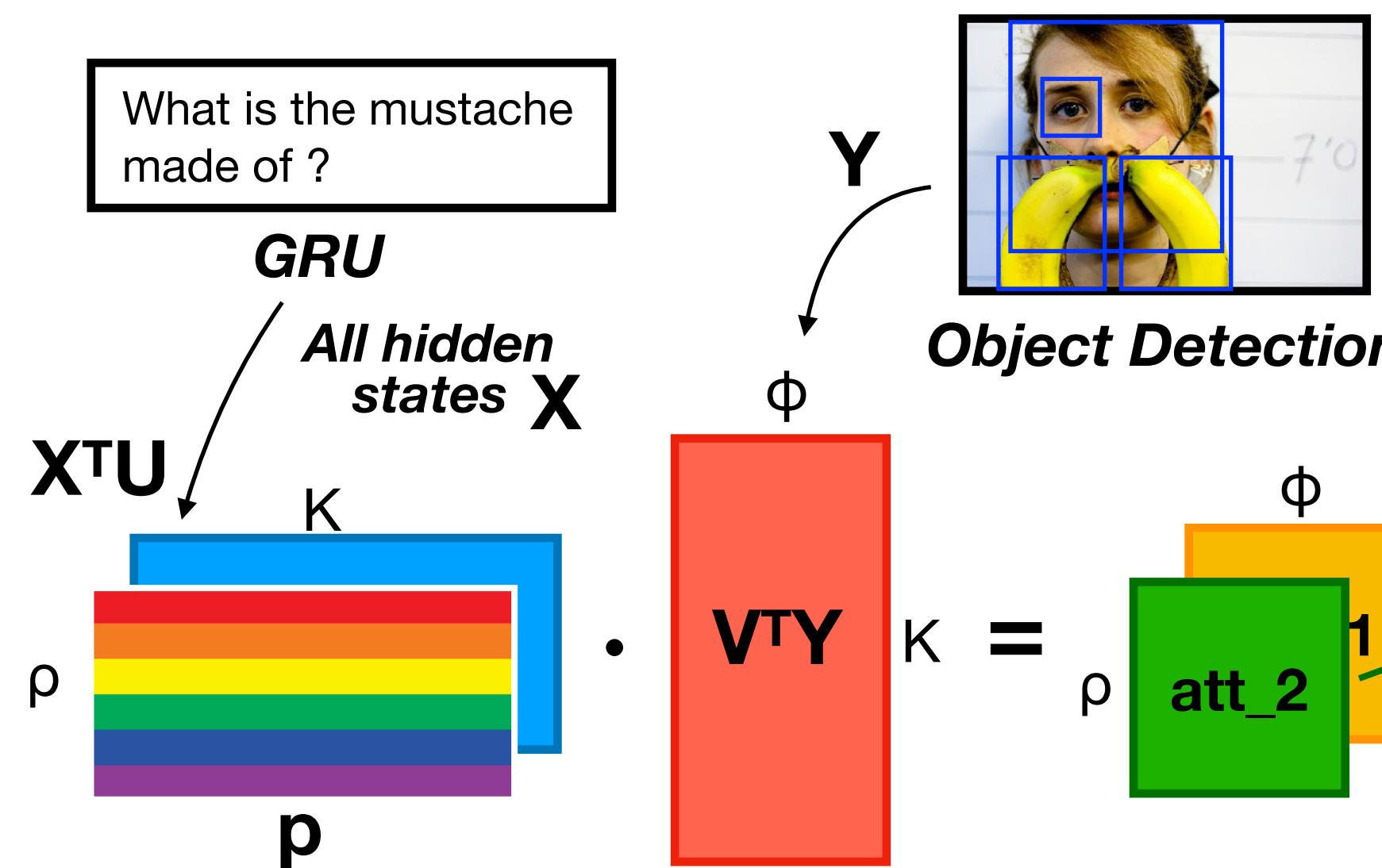
bilinear attention networks

↑

repeat {# of tokens} times

# Overview

- After getting bilinear attention maps, we can stack multiple BANs.



**Step 1. Bilinear Attention Maps**

**Step 2. Bilinear Attention Networks**

# Multiple Attention Maps

- Single model on validation score

	Validation VQA 2.0 Score	+%
<b>Bottom-Up (Teney et al., 2017)</b>	63.37	
<b>BAN-1</b>	65.36	+1.99
<b>BAN-2</b>	65.61	+0.25
<b>BAN-4</b>	65.81	+0.2
<b>BAN-8</b>	66.00	+0.19
<b>BAN-12</b>	66.04	+0.04

# Residual Learning

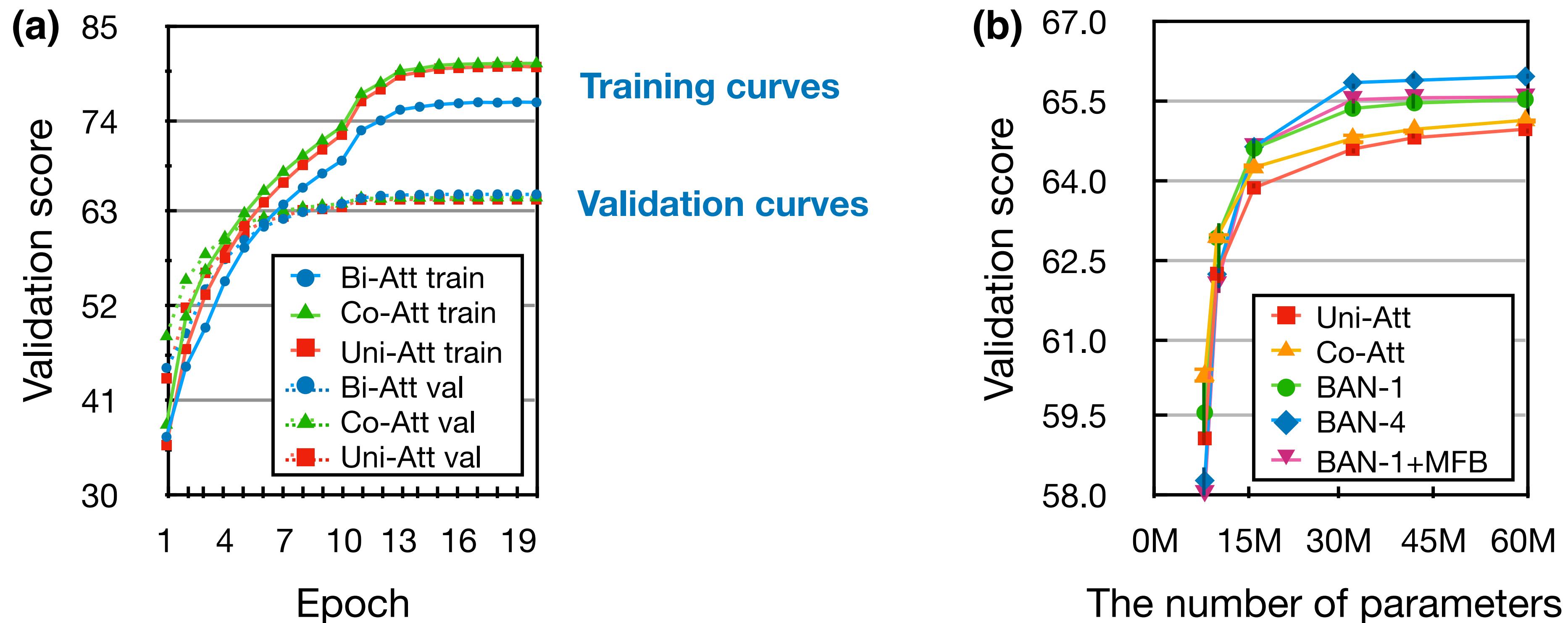
		Validation VQA 2.0 Score	+/-
$\sum_i \text{BAN}_i(\mathbf{X}, \mathbf{Y}; A_i)$	<b>BAN-4 (Residual)</b>	<b>65.81 ±0.09</b>	
$\left\  \sum_i \text{BAN}_i(\mathbf{X}, \mathbf{Y}; A_i) \right\ $	<b>BAN-4 (Sum)</b>	64.78 ±0.08	-1.03
$\left\  \sum_i \text{BAN}_i(\mathbf{X}, \mathbf{Y}; A_i) \right\ $	<b>BAN-4 (Concat)</b>	64.71 ±0.21	-0.07

# Comparison with Co-attention

BAN-1		Validation VQA 2.0 Score	+/-
		65.36 ±0.14	
		64.79 ±0.06	
	Co-Attention	64.59 ±0.04	-0.57
	Attention	64.59 ±0.04	-0.20

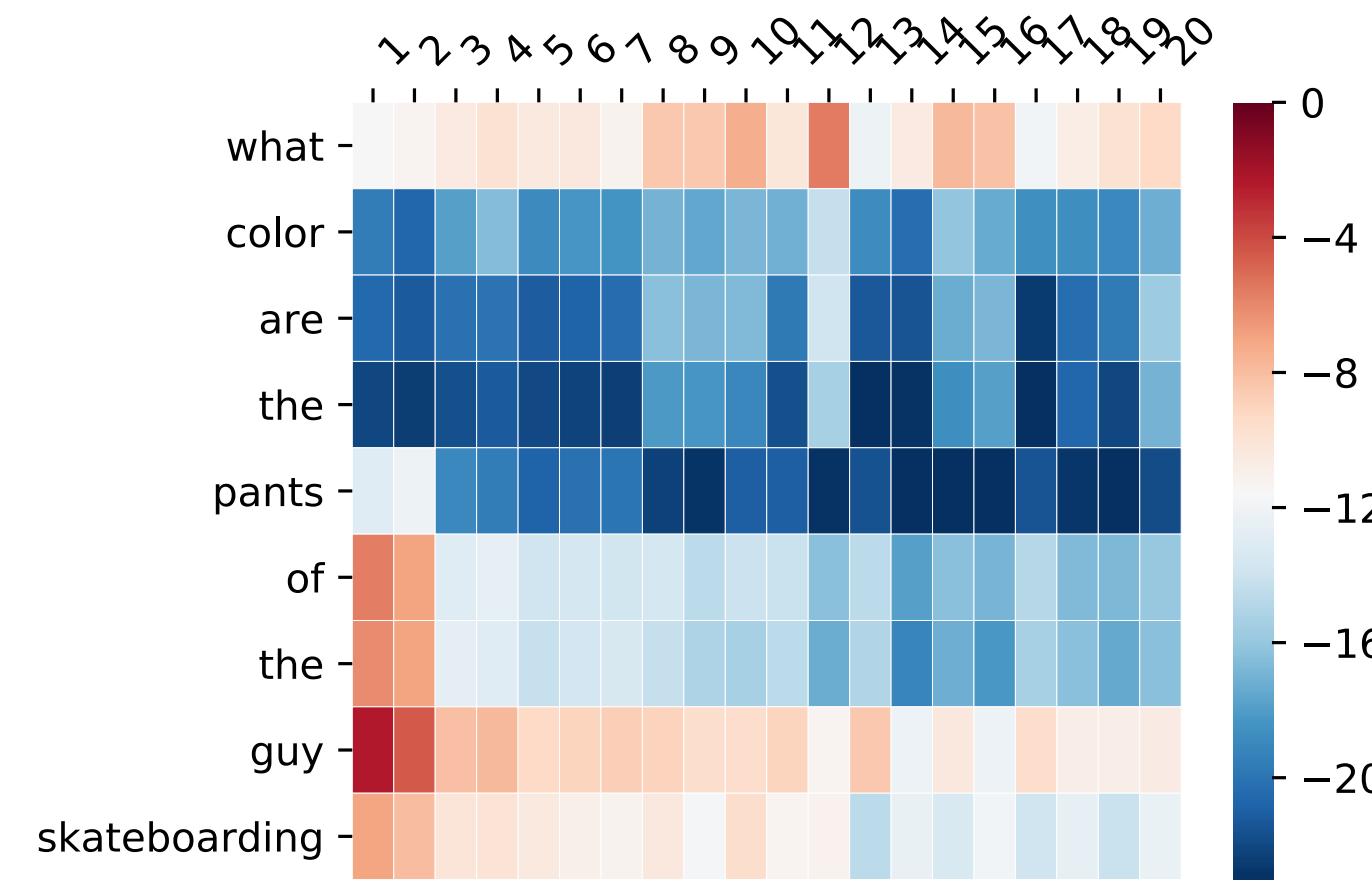
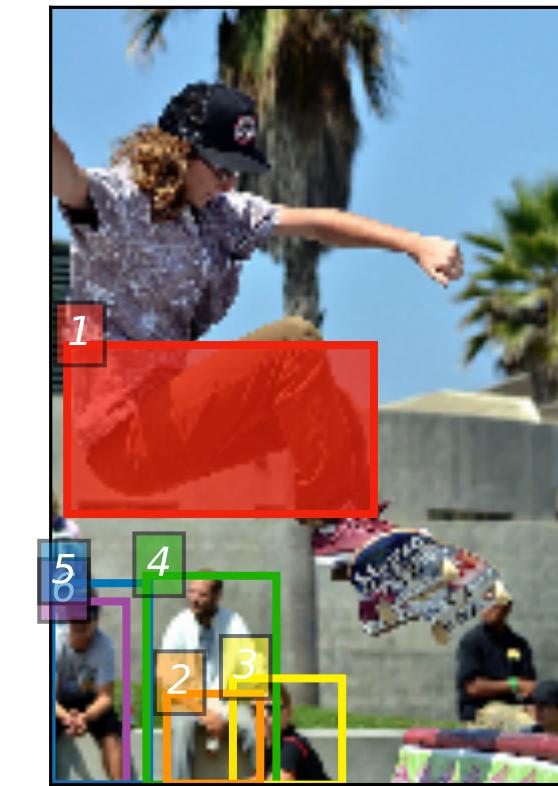
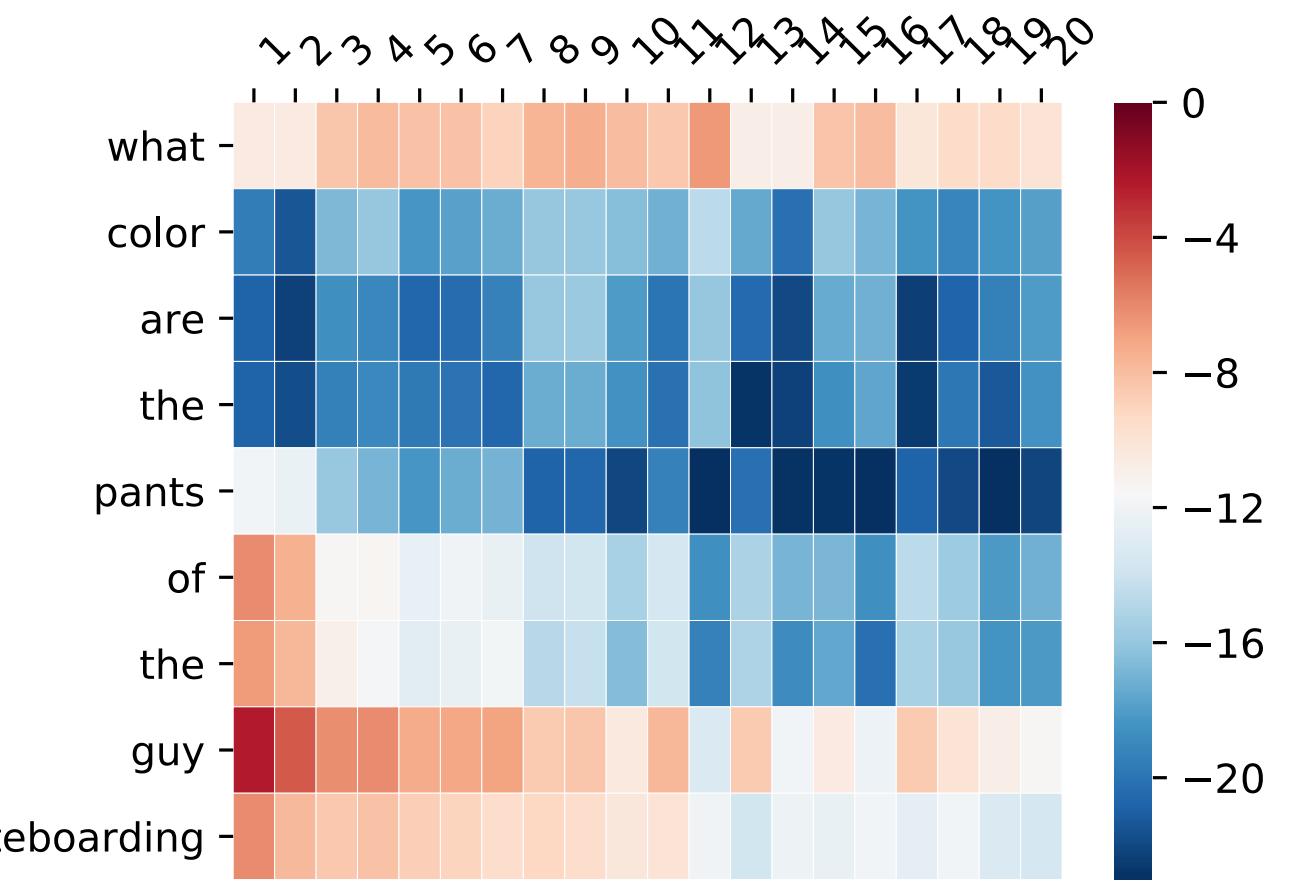
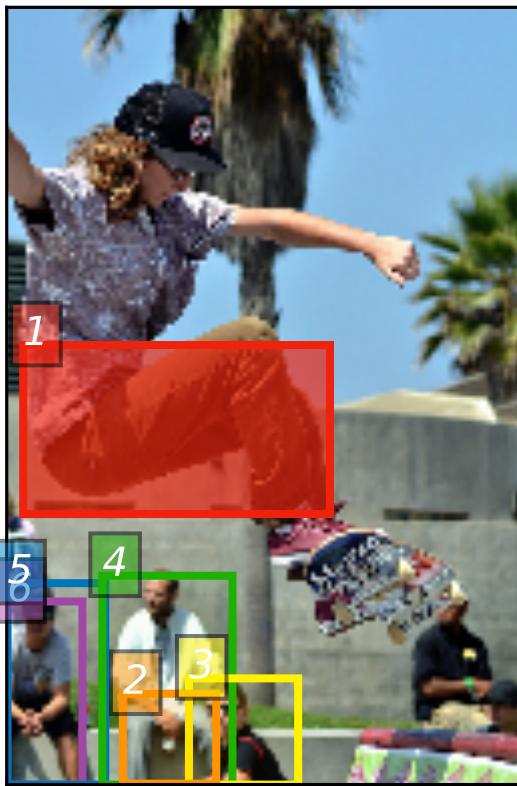
\* The number of parameters is controlled (all comparison models have 32M parameters).

# Comparison with Co-attention



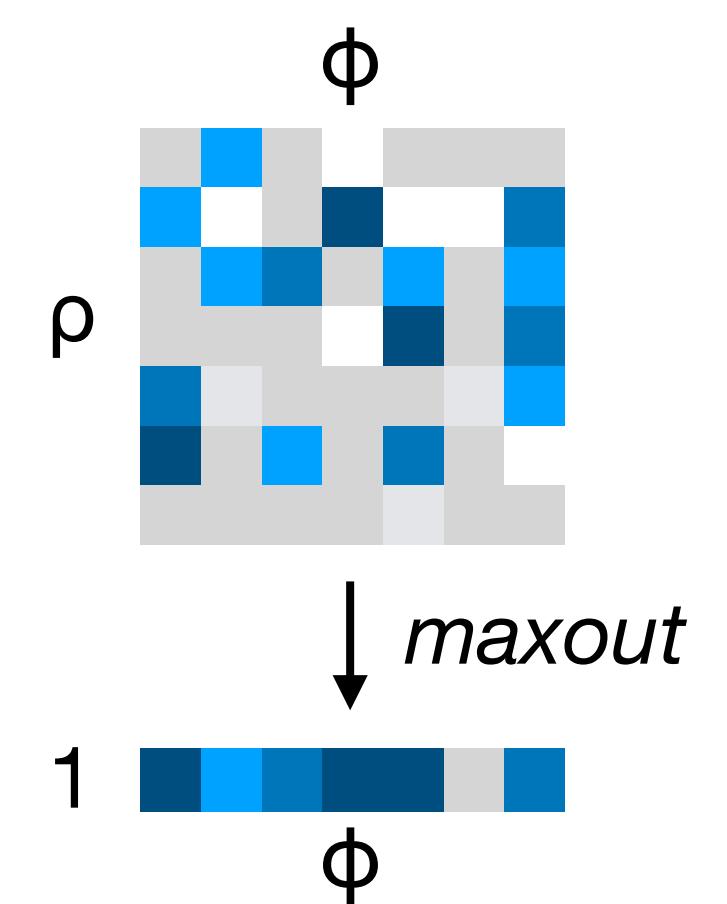
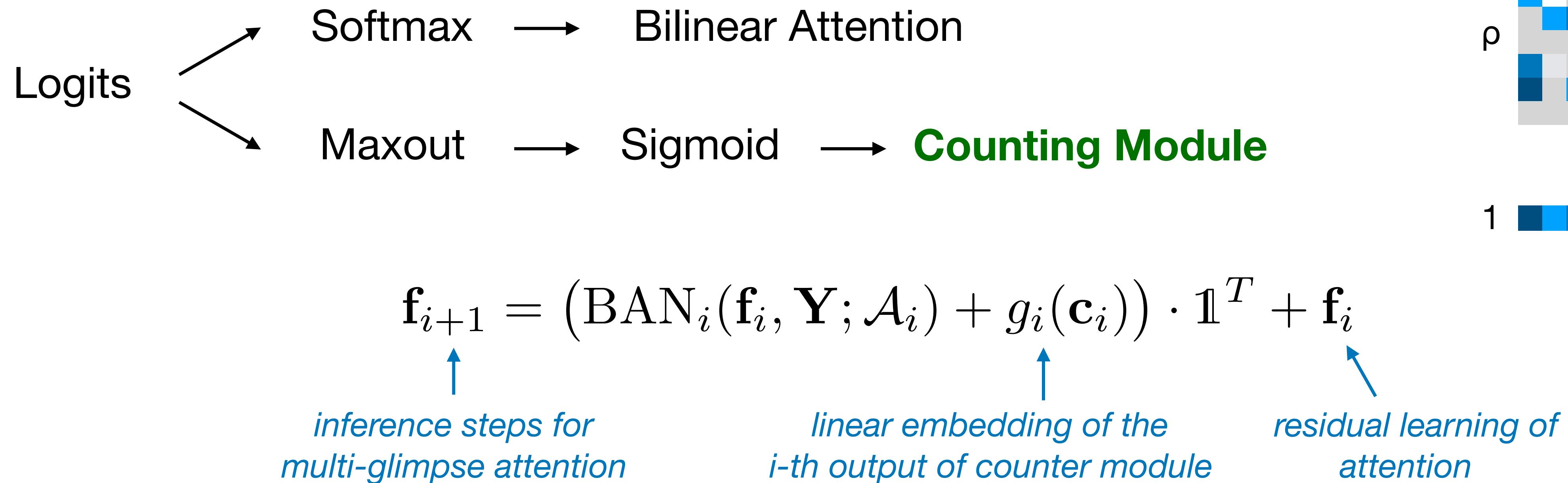
\* The number of parameters is controlled (all comparison models have 32M parameters).

# Visualization



# Integration of Counting module with BAN

- Counter (Zhang et al., 2018) gets a multinoulli distribution and detected box info
- **Our method** – *maxout* from bilinear attention distribution to unitary attention distribution, easy to adapt multitask learning



# Comparison with State-of-the-arts

	test-dev	Numbers
Zhang et al. (2018)		51.62
Ours		54.04

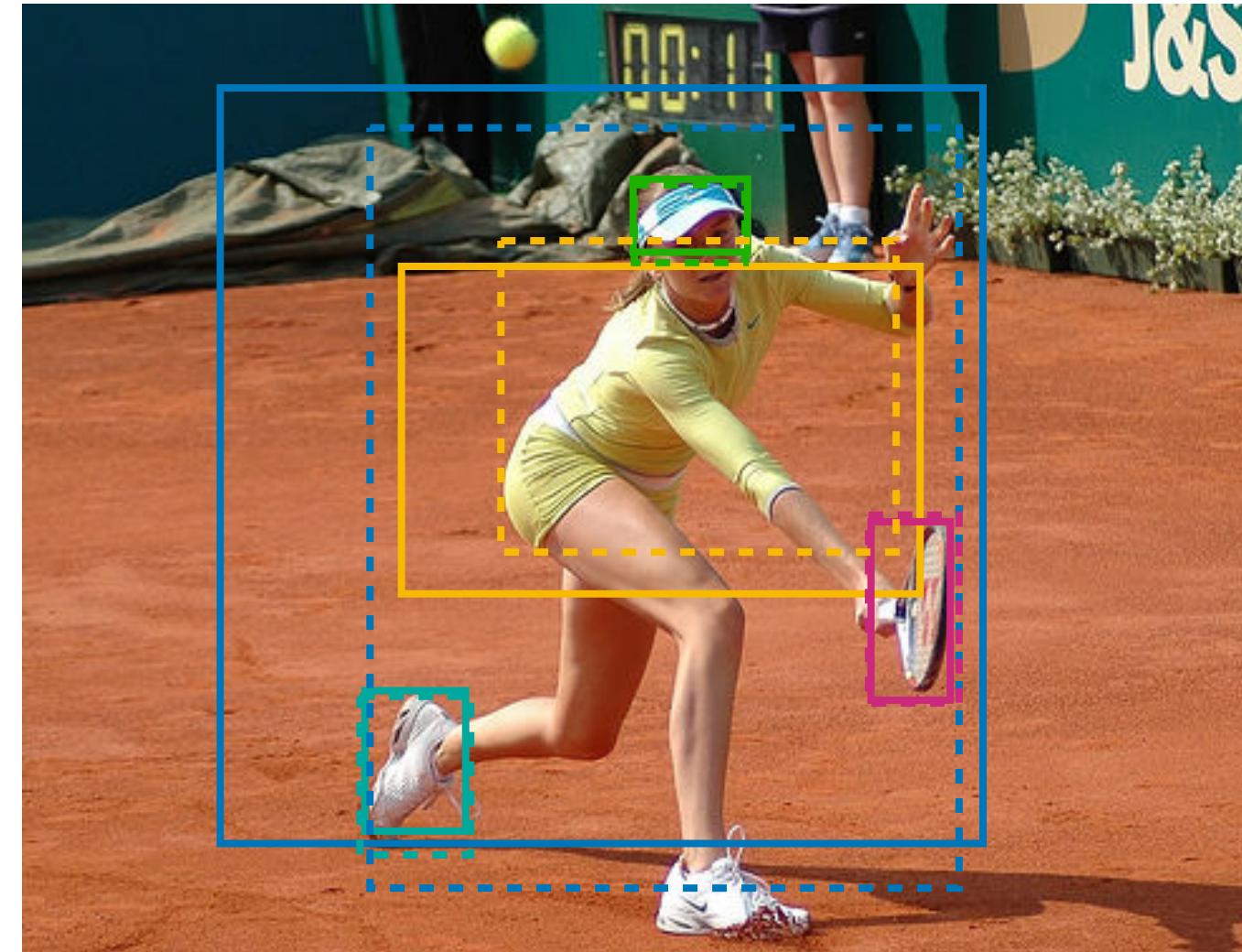
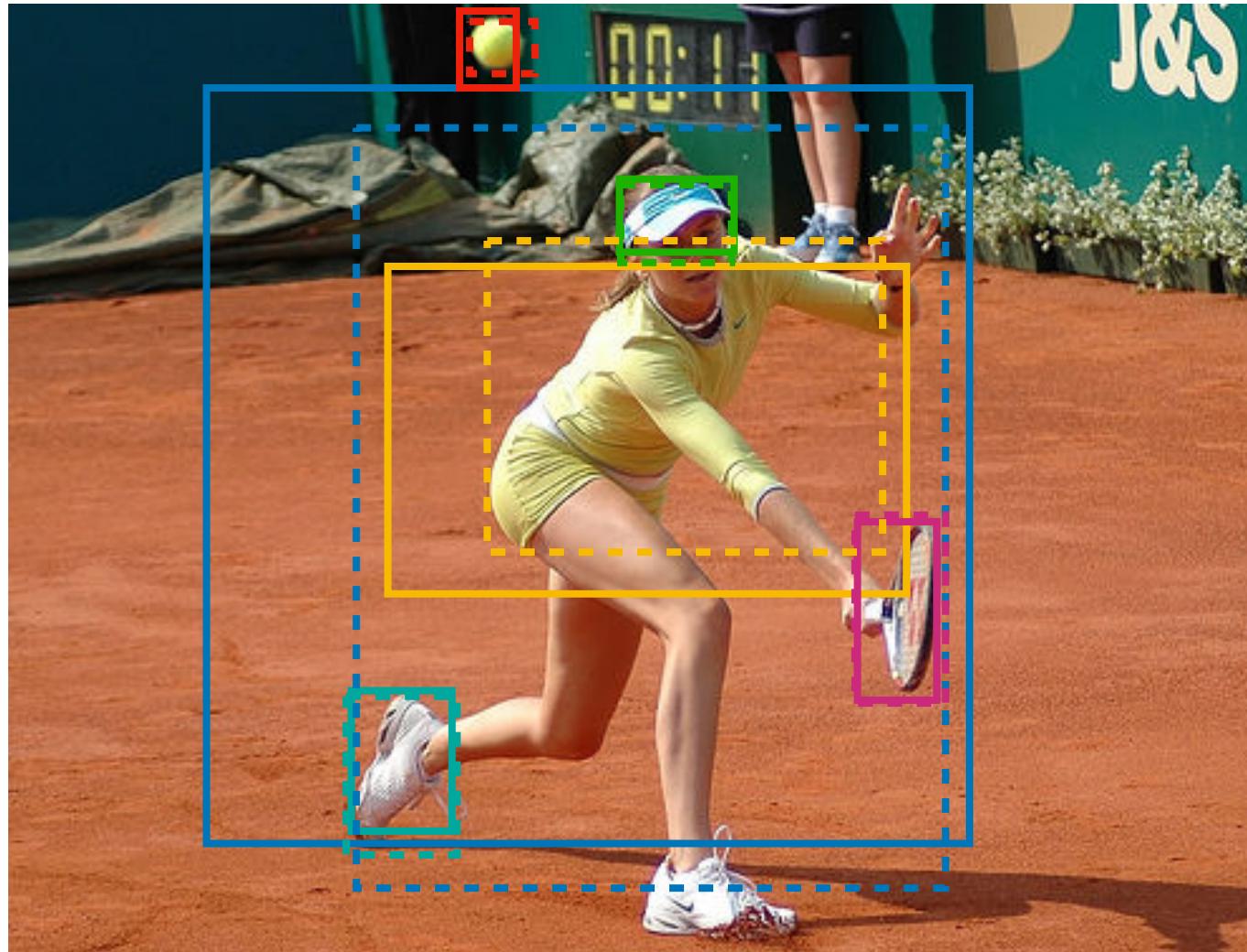
- Single model on test-dev score

	Test-dev VQA 2.0 Score	+%	
2016 winner	Prior	25.70	
	Language-Only	44.22	+18.52%
	MCB (ResNet)	61.96	+17.74%
2017 winner	Bottom-Up (FRCNN)	65.32	+3.36%
	MFH (ResNet)	65.80	+0.48%
2017 runner-up	MFH (FRCNN)	68.76	+2.96%
	BAN (Ours; FRCNN)	69.52	+0.76%
	BAN-Glove (Ours; FRCNN)	69.66	+0.14%
	BAN-Glove-Counter (Ours; FRCNN)	70.04	+0.38%

The diagram illustrates the architecture of the BAN model. It consists of three stacked components: 'image feature' (blue arrow), 'attention model' (red arrow), and 'counting feature' (green arrow). The 'image feature' component is at the top, followed by the 'attention model' in the middle, and the 'counting feature' at the bottom.

# Flickr30k Entities

- Visual grounding task – mapping entity phrases to regions in an image



[/EN#40120/people A girl] in [/EN#40122/clothing a yellow tennis suit] , [/EN#40125/other green visor] and [/EN#40128/clothing white tennis shoes] holding [/EN#40124/other a tennis racket] in a position where she is going to hit [/EN#40121/other the tennis ball] .

# Flickr30k Entities

- Visual grounding task – mapping entity phrases to regions in an image



[\[/EN#38656/people A male conductor\]](#) wearing [\[/EN#38657/clothing all black\]](#)  
[leading \[/EN#38653/people a orchestra\]](#) and [\[/EN#38658/people choir\]](#) on [\[/](#)  
[EN#38659/scene a brown stage\]](#) playing and singing [\[/EN#38664/other a musical](#)  
[number\]](#) .

# Flickr30k Entities Recall@1,5,10

	R@1	R@5	R@10
Zhang et al. (2016)	27.0	49.9	57.7
SCRC (2016)	27.8	-	62.9
DSPE (2016)	43.89	64.46	68.66
GroundeR (2016)	47.81	-	-
MCB (2016)	48.69	-	-
CCA (2017)	50.89	71.09	75.73
Yeh et al. (2017)	53.97	-	-
Hinami & Satoh (2017)	65.21		
BAN (ours)	<b>66.69</b>	<b>84.22</b>	<b>86.35</b>

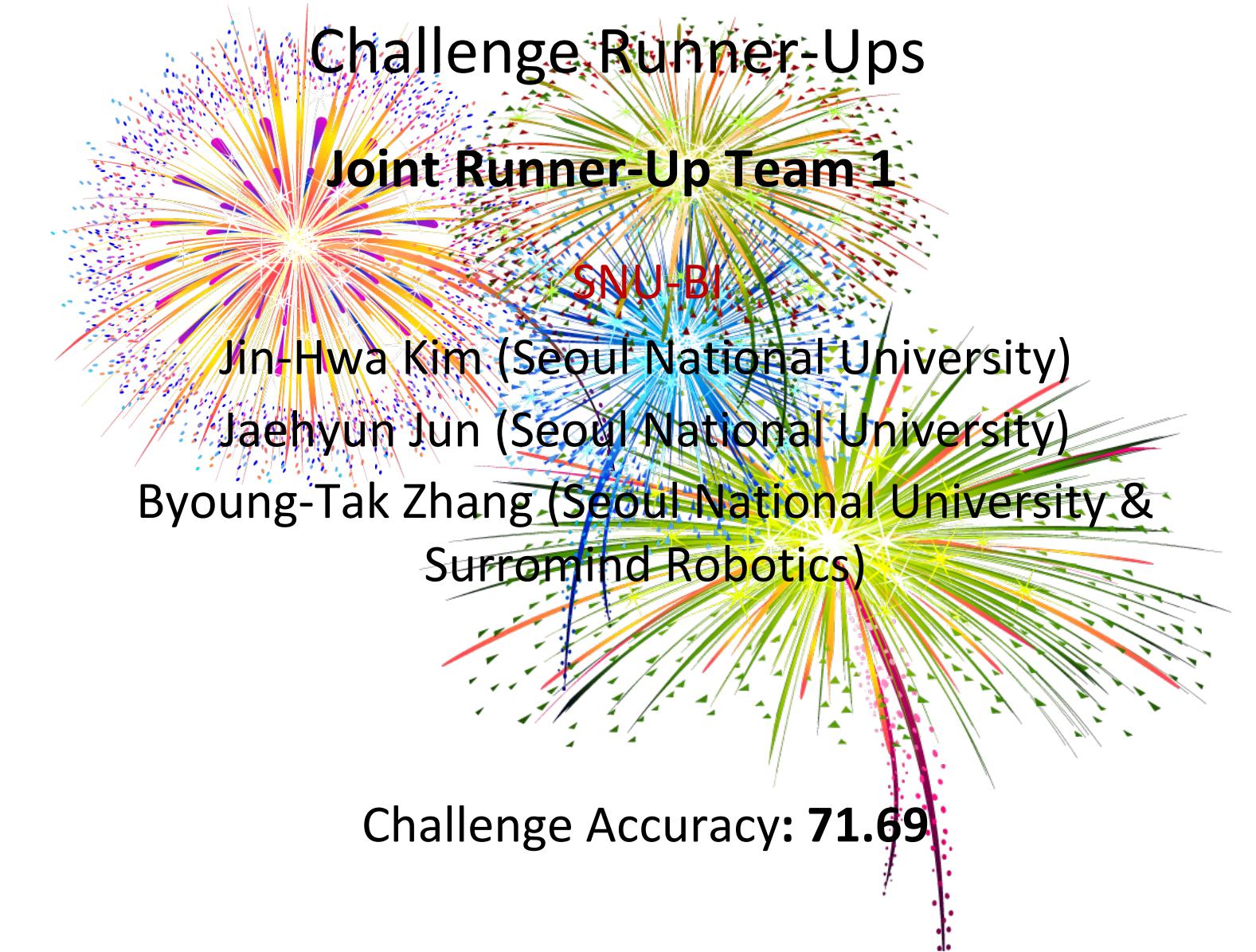
# Flickr30k Entities Recall@1,5,10

	<b>people</b>	<b>clothing</b>	<b>bodyparts</b>	<b>animals</b>	<b>vehicles</b>	<b>instruments</b>	<b>scene</b>	<b>other</b>
<b>#Instances</b>	5,656	2,306	523	518	400	162	1,619	3,374
<b>GroundeR (2016)</b>	61.00	38.12	10.33	62.55	68.75	36.42	58.18	29.08
<b>CCA (2017)</b>	64.73	46.88	17.21	65.83	68.75	37.65	51.39	31.77
<b>Yeh et al. (2017)</b>	68.71	46.83	19.50	70.07	73.75	39.50	60.38	32.45
<b>Hinami &amp; Satoh (2017)</b>	78.17	61.99	35.25	74.41	76.16	<b>56.69</b>	68.07	47.42
<b>BAN (ours)</b>	<b>79.90</b>	<b>74.95</b>	<b>47.23</b>	<b>81.85</b>	<b>76.92</b>	43.00	<b>68.69</b>	<b>51.33</b>

# Conclusions

- Bilinear attention networks gracefully extends unitary attention networks,
- as low-rank bilinear pooling inside bilinear attention.
- Furthermore, residual learning of attention efficiently uses multiple attention maps.

2018 VQA Challenge runners-up (shared 2nd place)  
1st single model (70.35 on test-standard)



**arXiv:1805.07932**

Code & pretrained model in PyTorch:

**[github.com/jnhwkim/ban-vqa](https://github.com/jnhwkim/ban-vqa)**

# Thank you!

## Any question?

We would like to thank Kyoung-Woon On, Bohyung Han, and Hyeonwoo Noh for helpful comments and discussion.

This work was supported by 2017 Google Ph.D. Fellowship in Machine Learning and Ph.D. Completion Scholarship from College of Humanities, Seoul National University, and the Korea government (IITP-2017-0-01772-VTT, IITP-R0126-16-1072-SW.StarLab, 2018-0-00622-RMI, KEIT-10044009-HRI.MESSI, KEIT-10060086-RISF).

The part of computing resources used in this study was generously shared by Standigm Inc.