

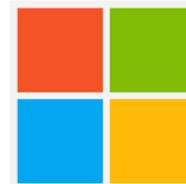
A Simple VQA Model with a Few Tricks and Image Features from Bottom-up Attention

Damien Teney¹, Peter Anderson^{2*}, David Golub^{4*}, Po-Sen Huang³,
Lei Zhang³, Xiaodong He³, Anton van den Hengel¹

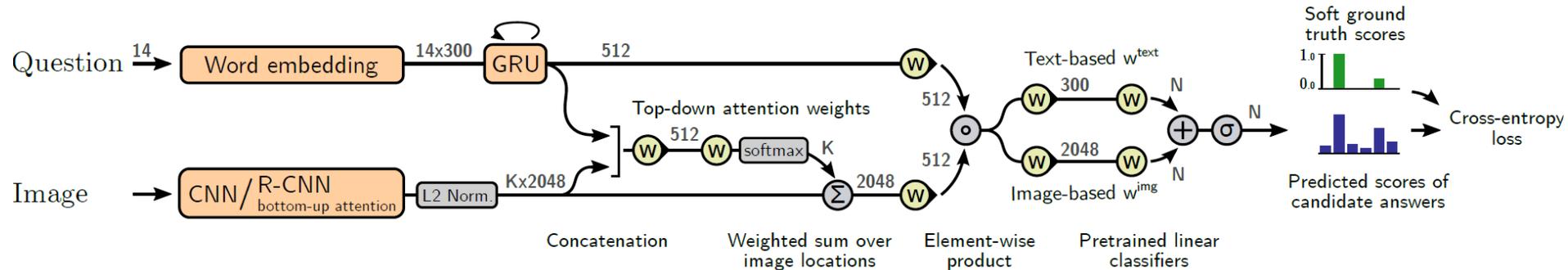
¹University of Adelaide ²Australian National University

³Microsoft Research ⁴Stanford University

**Work performed while interning at MSR*



Proposed model



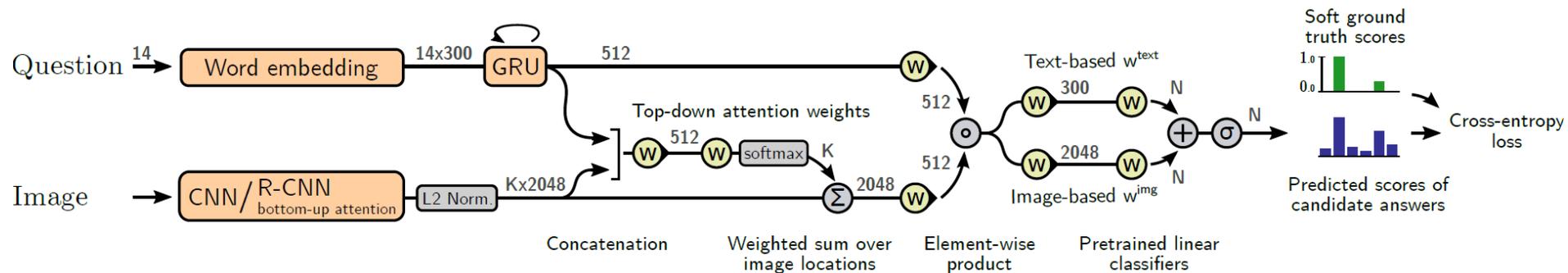
Straightforward architecture

- **Joint embedding** of question/image
- Single-head, question-guided **attention over image**
- Element-wise product

The devil is in the details

- Image features from **Faster R-CNN**
- **Gated tanh** activations
- Output as **regression** of answer scores, **soft scores as target**
- Output classifiers initialized with **pretrained representations of answers**

Gated layers



Non-linear layers: **gated hyperbolic tangent** activations

- Defined as: input x , output y

$$\tilde{y} = \tanh(Wx + b) \quad \text{intermediate activation}$$

$$g = \sigma(W'x + b') \quad \text{gate}$$

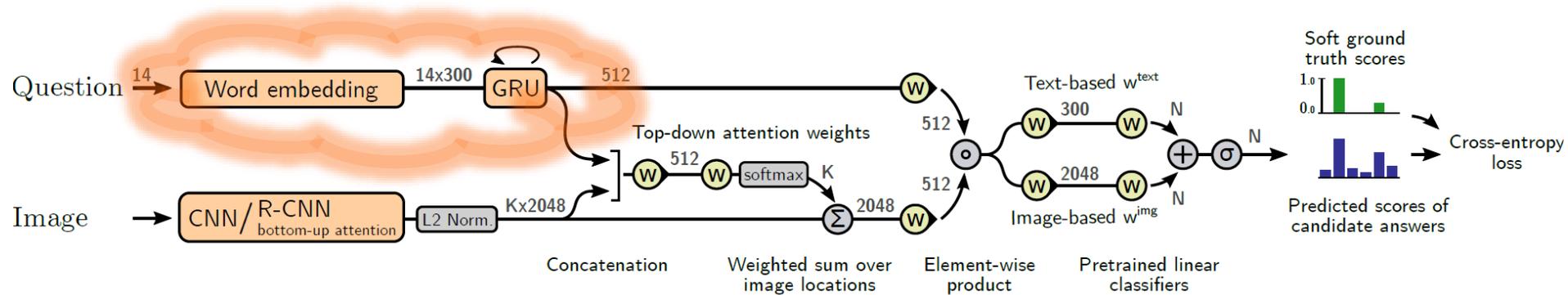
$$y = \tilde{y} \circ g \quad \text{combine with element-wise product}$$

- Inspired by gating in LSTMs/GRUs
- Empirically better than ReLU, tanh, gated ReLU, residual connections, etc.
- Special case of highway networks; used before in:

[1] Dauphin et al. Language modeling with gated convolutional networks, 2016.

[2] Teney et al. Graph-structured representations for visual question answering, 2017.

Question encoding



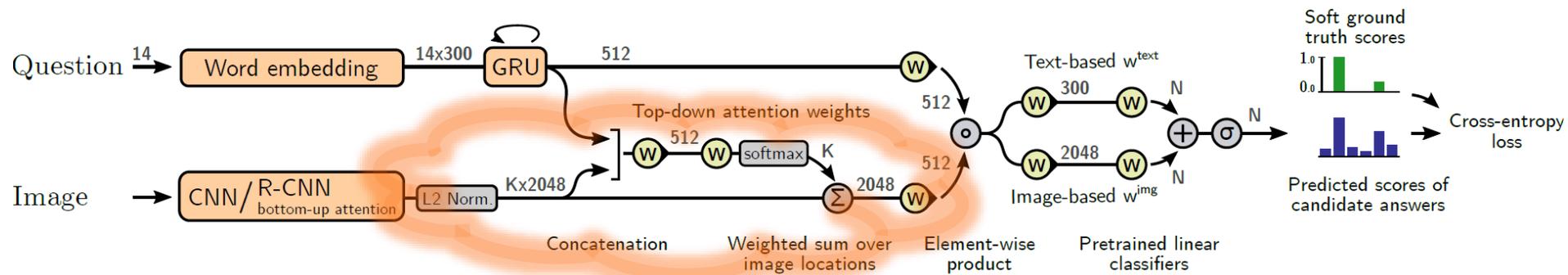
Chosen implementation

- Pretrained **GloVe** embeddings, $d=300$
- **GRU** encoder

Better than....

- Word embeddings learned from scratch
- GloVe of dimension 100, 200
- Bag-of-words (sum/average of embeddings)
- GRU backwards
- GRU bidirectional
- 2-layer GRU

Classical “top-down” attention on image features



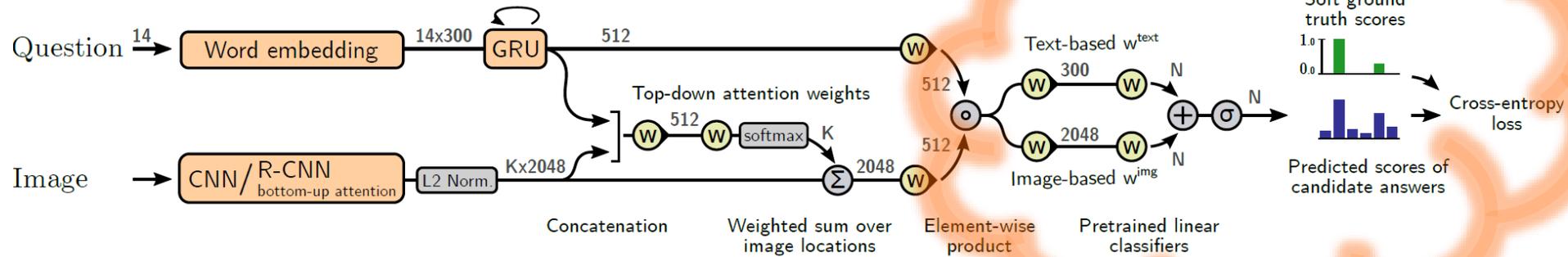
Chosen implementation

- Simple attention **on image feature maps**
- **One** head
- **Softmax** normalization of weights

Better than....

- No L2 normalization
- Multiple heads
- Sigmoid on weights

Output



Chosen implementation

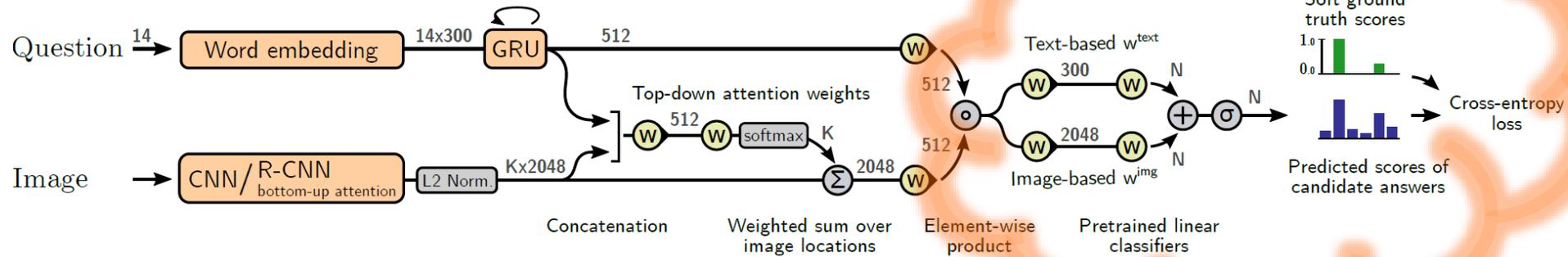
- **Sigmoid** output (regression) of answer scores:
allows **multiple answers** per question
- **Soft targets** in $[0,1]$
allows **uncertain** answers
- **Initialize classifiers** with representations of answers

$$y = \sigma(Wx) \quad W \text{ of dimensions } nAnswers \times d$$

Better than....

- Softmax classifier
- Binary targets $\{0,1\}$
- Classifiers learned from scratch

Output



Chosen implementation

- **Sigmoid** output (regression) of answer scores:
allows **multiple answers** per question
- **Soft targets** in $[0,1]$
allows **uncertain** answers
- **Initialize classifiers** with representations of answers

$y = \sigma(W^{\text{text}}x^{\text{text}} + W^{\text{img}}x^{\text{img}})$ Initialize W^{text} with GloVe **word embeddings**

Initialize W^{img} with **Google Images** (global ResNet features)

Training and implementation

- Additional training data from Visual Genome: questions with **matching answers** and **matching images** (about 30% of Visual Genome, *i.e.* ~485,000 questions)
- Keep **all questions**, even those with no answer in candidates, and with $0 < \text{score} < 1$
- Shuffle training data but keep **balanced pairs in same mini-batches**
- Large mini-batches of **512** QAs; sweet spot in {64, 128, 256, 384, 512, 768, 1024}
- 30-Network **ensemble**: different random seeds, sum predicted scores

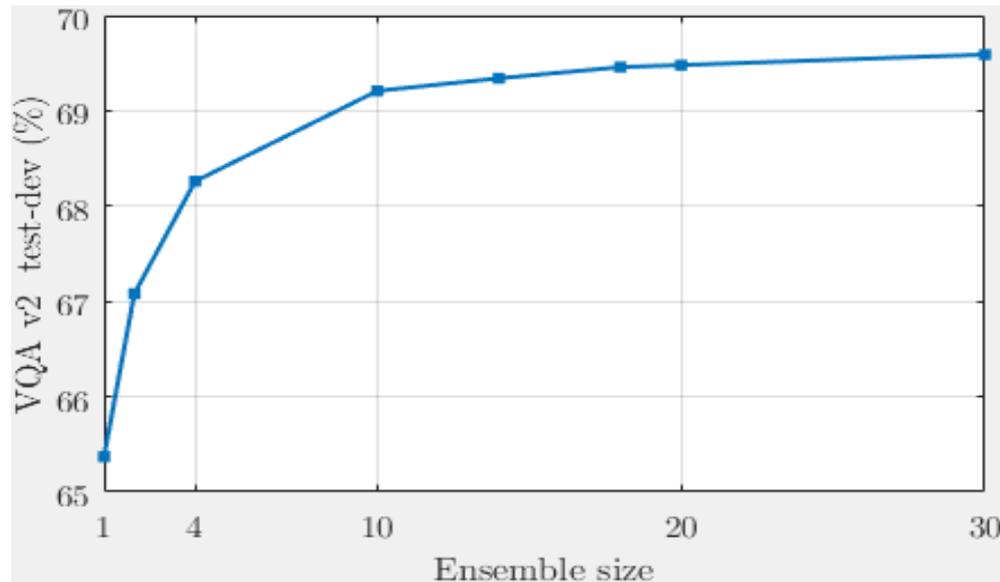
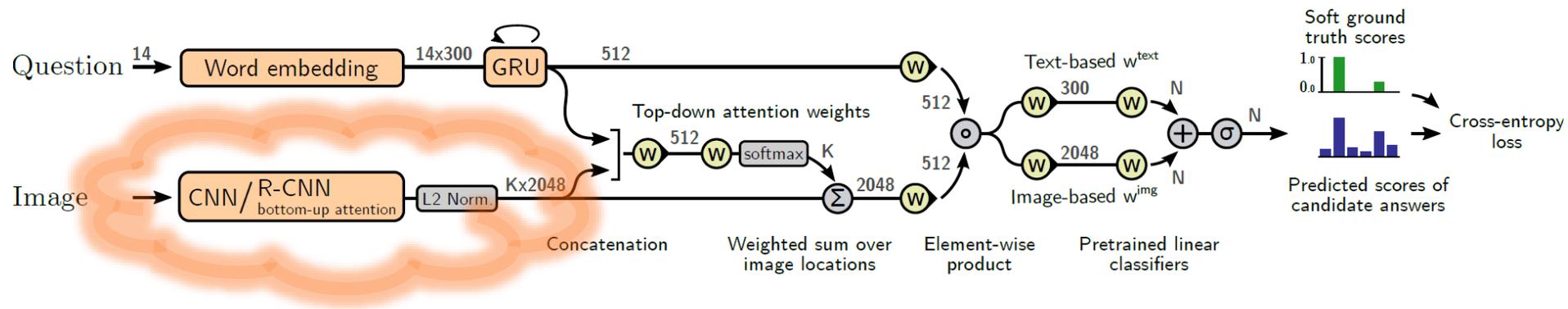
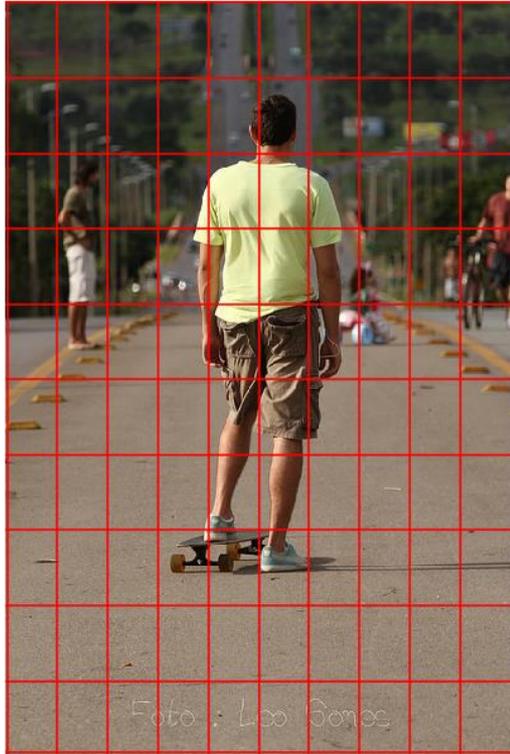


Image features from bottom-up attention

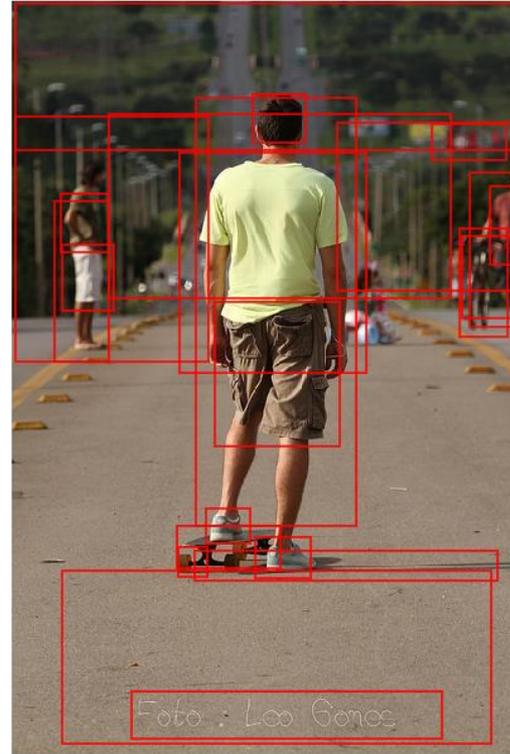


- Equally applicable to **VQA and image captioning**
- Significant relative **improvements: 6 – 8 %** (VQA / CIDEr / SPICE)
- Intuitive and interpretable (natural approach)

Bottom-up image attention

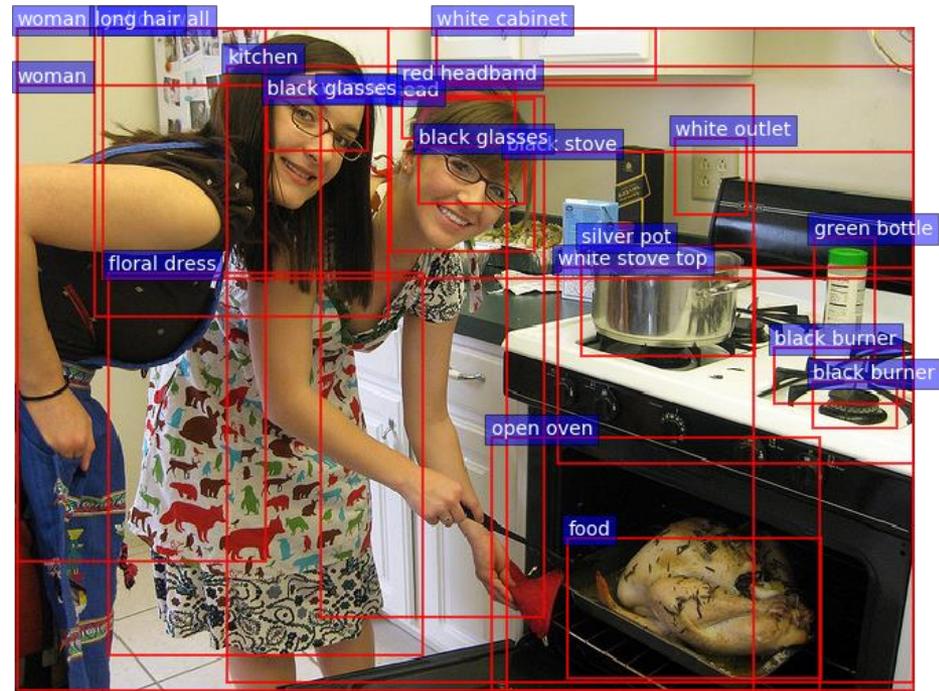
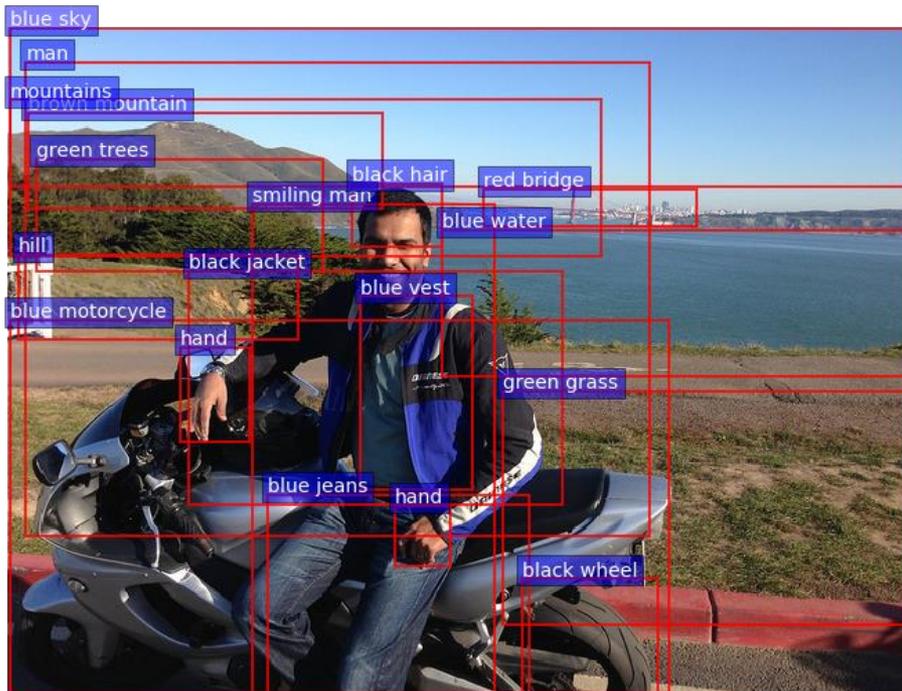


Typically, attention models operate on the spatial output of a CNN



We calculate attention at the level of **objects and other salient image regions**

Can be implemented with Faster R-CNN¹

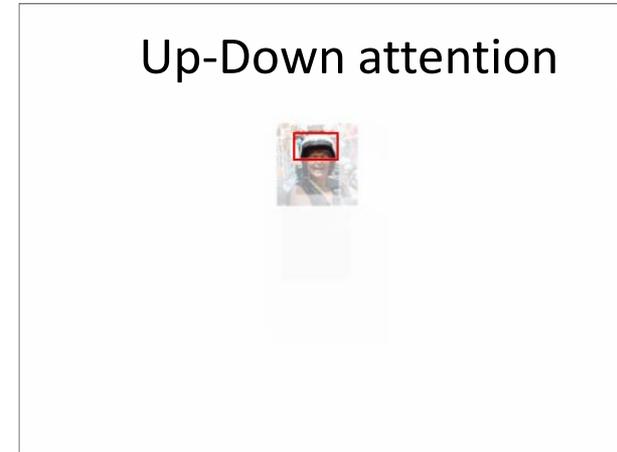


- Pre-train on 1600 objects and 400 attributes from **Visual Genome**²
- Select salient regions based on object detection confidence scores
- Take the mean-pooled ResNet-101³ feature from each region

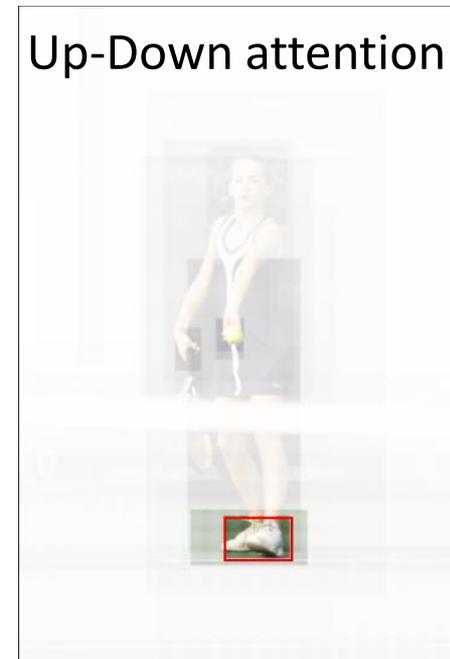
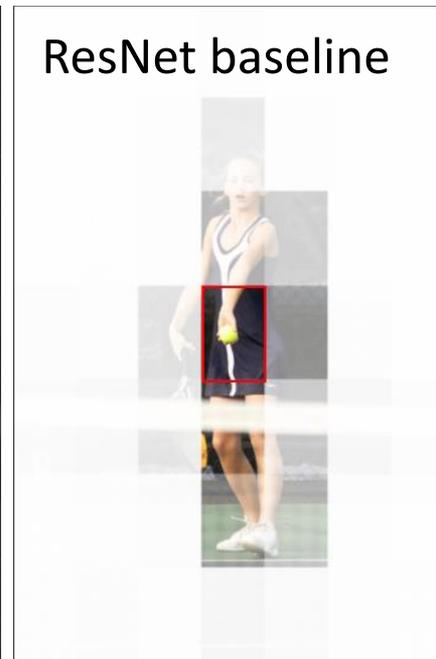
¹NIPS 2015, ²<http://visualgenome.org>, ³CVPR 2016

Qualitative differences in attention methods

Q: Is the person wearing a **helmet**?

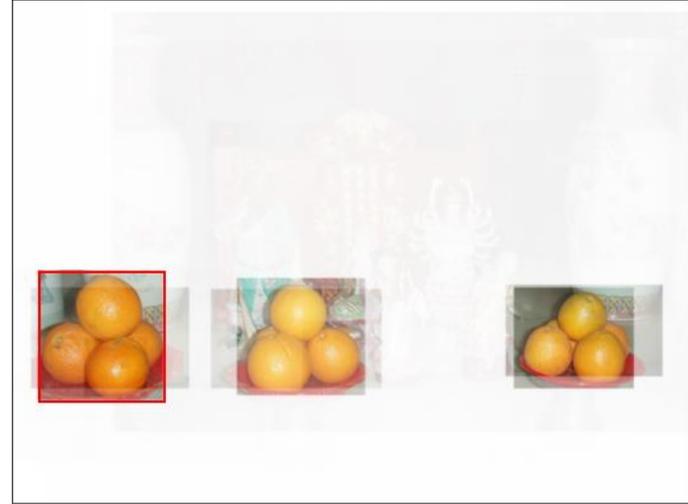


Q: What **foot** is in front of the other **foot**?

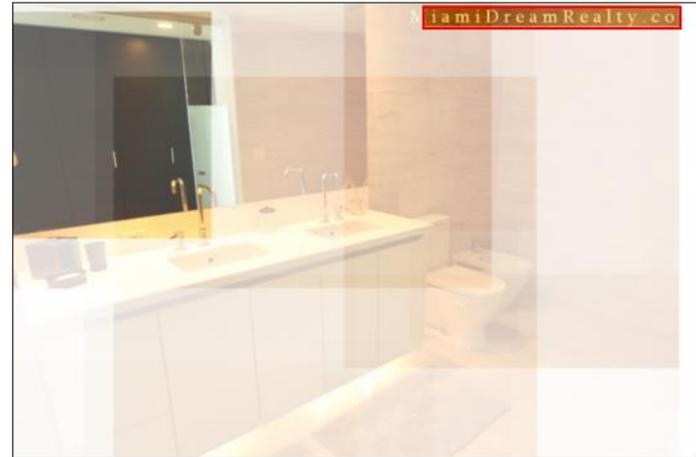


VQA failure cases: counting, reading

Q: **How many** oranges are sitting on pedestals?

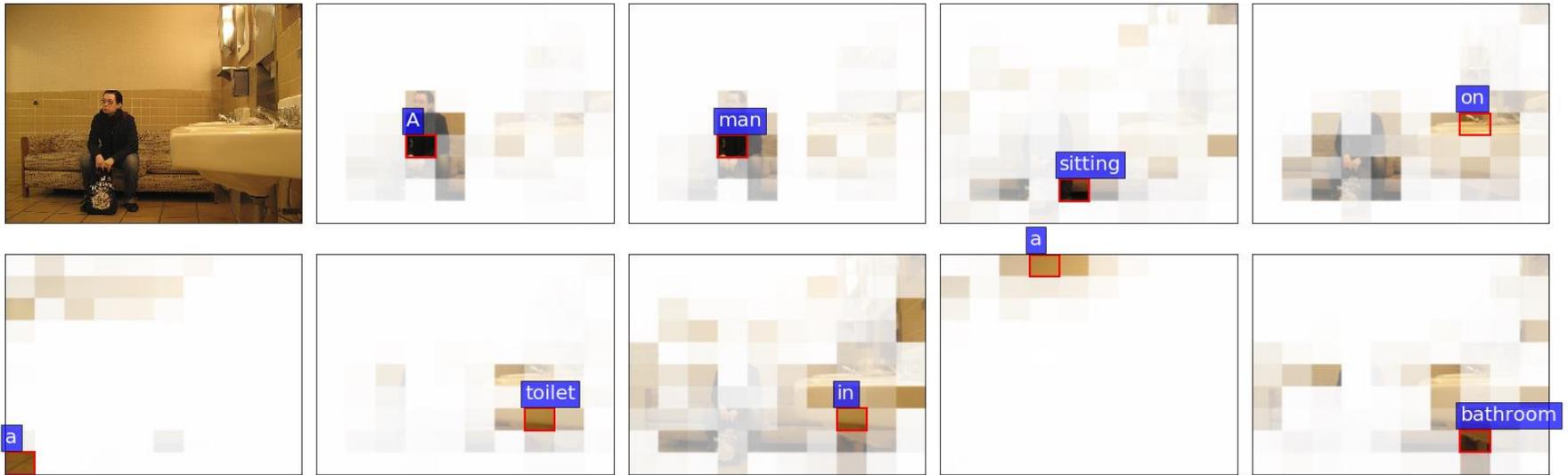


Q: What is the **name** of the realtor?



Equally applicable to Image Captioning

ResNet baseline: A man sitting on a **toilet** in a bathroom.



Up-Down attention: A man sitting on a **couch** in a bathroom.



MS COCO Image Captioning Leaderboard

Results															
#	User	BLEU-1		BLEU-2		BLEU-3		BLEU-4		METEOR		ROUGE-L		CIDEr-D	
		c5 ▲	c40 ▲	c5 ▲	c40 ▲	c5 ▲	c40 ▲								
1	panderson_msr	0.802	0.952	0.641	0.888	0.491	0.794	0.369	0.685	0.276	0.367	0.571	0.724	1.179	1.205
		(1)	(1)	(1)	(1)	(1)	(1)	(1)	(1)	(1)	(2)	(4)	(1)	(1)	(2)
2	TencentAI.v1	0.793	0.943	0.633	0.880	0.483	0.781	0.363	0.669	0.275	0.363	0.572	0.719	1.178	1.207
		(2)	(2)	(2)	(2)	(4)	(2)	(4)	(3)	(2)	(5)	(3)	(3)	(2)	(1)
3	xxzhu	0.786	0.935	0.629	0.871	0.485	0.778	0.368	0.670	0.275	0.364	0.572	0.721	1.173	1.194
		(5)	(4)	(3)	(3)	(2)	(3)	(2)	(2)	(3)	(4)	(1)	(2)	(3)	(3)
4	CASIA_IVA	0.786	0.934	0.629	0.870	0.484	0.776	0.368	0.669	0.274	0.362	0.572	0.719	1.170	1.188
		(4)	(5)	(4)	(4)	(3)	(4)	(3)	(4)	(4)	(7)	(2)	(4)	(4)	(4)
5	Anonymous	0.787	0.937	0.627	0.867	0.476	0.765	0.356	0.652	0.270	0.354	0.564	0.705	1.160	1.180
		(3)	(3)	(5)	(5)	(5)	(5)	(5)	(5)	(6)	(14)	(5)	(10)	(5)	(5)
6	etiennem	0.781	0.931	0.619	0.860	0.470	0.759	0.352	0.645	0.270	0.355	0.563	0.707	1.147	1.167
		(6)	(6)	(6)	(6)	(6)	(6)	(6)	(7)	(5)	(13)	(6)	(8)	(6)	(6)

- Bottom-up attention adds 6 – 8% improvement on SPICE and CIDEr metrics (see arXiv: [Bottom-Up and Top-Down Attention for Image Captioning and VQA](#))
- First place on almost all MS COCO leaderboard metrics

VQA experiments

- Current best results Ensemble, trained on tr+va+VG, eval. on test-std
Yes/no: 86.52 Number: 48.48 Other: 60.95 **Overall: 70.19**
- Bottom-up attention** adds **6%** relative improvement
(even though the baseline ResNet has twice as many layers)

Single-network, trained on tr+VG, eval. on va

	Yes/No	Number	Other	Overall
Ours: ResNet (1×1)	76.0	36.5	46.8	56.3
Ours: ResNet (14×14)	76.6	36.2	49.5	57.9
Ours: ResNet (7×7)	77.6	37.7	51.5	59.4
Ours: Up-Down	80.3	42.8	55.8	63.2
Relative Improvement	3%	14%	8%	6%

Take-aways and conclusions

- Difficult to predict effects of architecture, hyperparameters, ...
Engineering effort: good intuitions are valuable, then need **fast experiments**
Performance \approx (# Ideas) * (# GPUs) / (**Training time**)
- Beware of experiments with reduced training data
- Non-cumulative gains, performance saturates
Fancy tweaks may just add more capacity to network
May be redundant with other improvements
- Calculating attention at the level of **objects and other salient image regions**
(bottom-up attention) significantly improves performance
Replace pretrained CNN features with pretrained bottom-up attention features

Questions ?

arXiv:1708.02711: [Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge](#)

arXiv:1707.07998: [Bottom-Up and Top-Down Attention for Image Captioning and VQA](#)



Damien Teney, Peter Anderson, David Golub, Po-Sen Huang,
Lei Zhang, Xiaodong He, Anton van den Hengel

