



Vision and Language: from Perception to Creation



Tao Mei, Ph.D.

AI Research, JD.COM
tmei@jd.com

image captioning



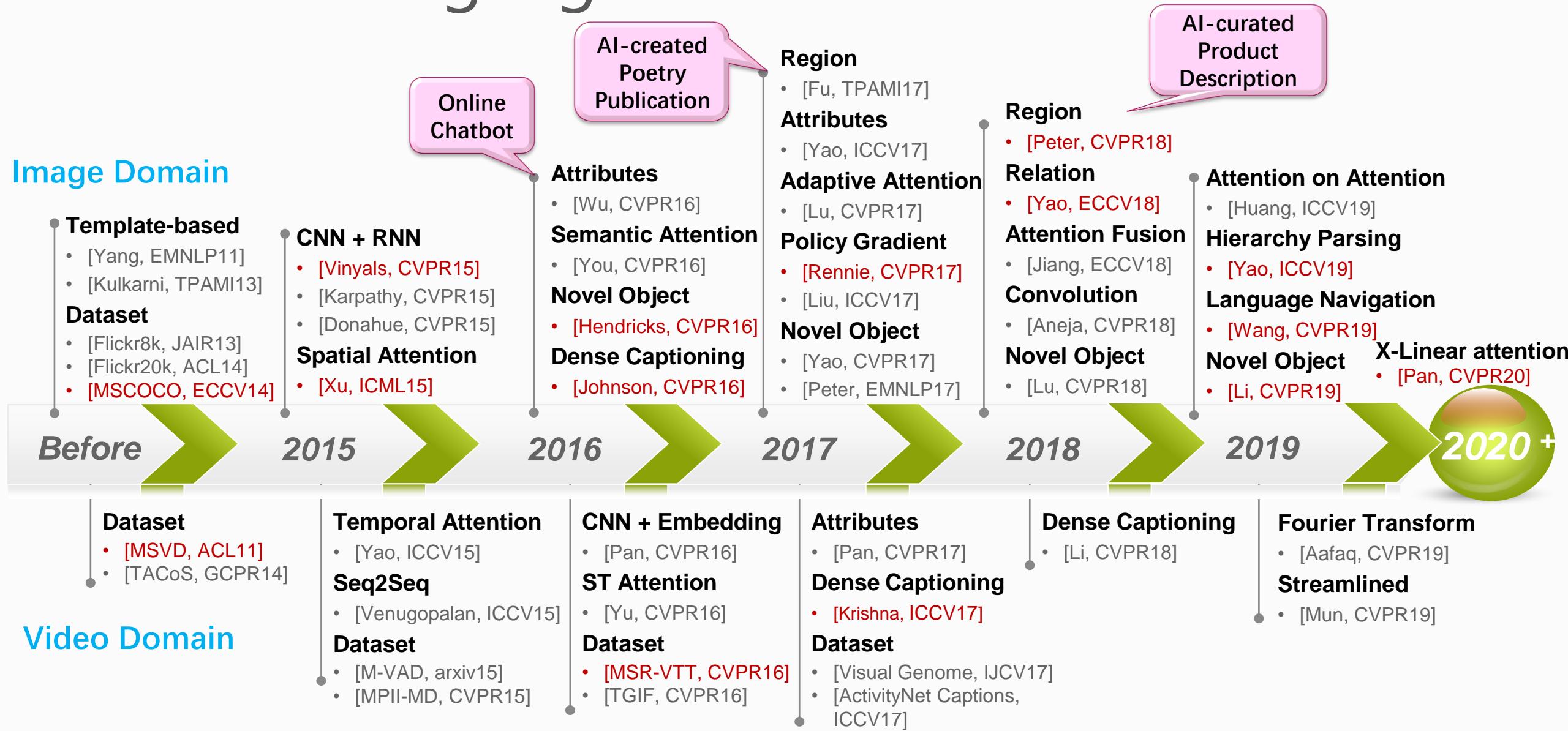
a group of zebras standing in a field [Vinyals, CVPR15]
a group of zebras grazing on grass [You, CVPR16]
a group of zebras grazing in a field [Yao, ICCV17]
a group of zebras and a rainbow in the sky [Peter, CVPR18]
a group of zebras grazing in a field with a rainbow in the sky [Yao, ECCV18]
a group of zebras and other animals grazing in a field with a rainbow in the sky [Pan, CVPR20]

video captioning



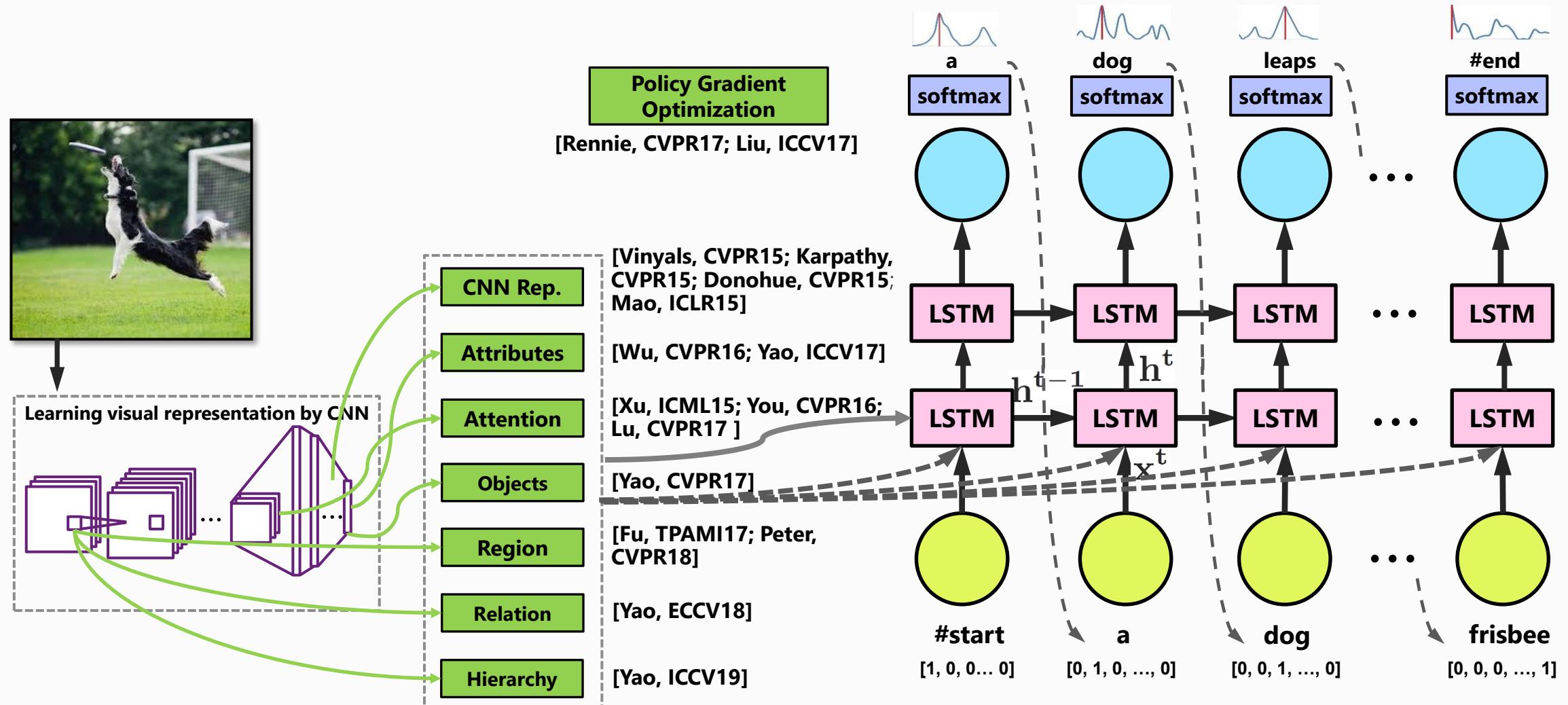
a man is singing a song [Venugopalan, ICCV15]
a man is dancing [Pan, CVPR16]
a woman is dancing [Pan, CVPR17]
a woman is dancing in the rain [Chen, AAAI19]

Vision to Language



Mainstream: CNN Encoder + LSTM Decoder

[Google15, Stanford15, Berkeley15, Baidu/UCLA15, UdeM15, Rochester16, UAdelaide16, Virginia Tech17, THU17, MSR17&18, IBM17, U of Oxford & Google17, JD AI18&19]



* Note that this figure only shows prediction process.

Technology Trends

- Aim for a thorough image/video understanding for captioning
 - > Direction I: enhance image encoder with X
- Explore the (1st , 2nd , ...) interaction across multi-modal inputs (hidden states of sentence decoder & the encoded image features)
 - > Direction II: integrate encoder/decoder via X attention
- Learn a universal encoder-decoder structure for VL tasks
 - > Direction III: vision-language pre-training

Direction I: enhance image encoder with X

X = visual attributes

[You, CVPR16; Wu,
CVPR16; Yao, ICCV17]



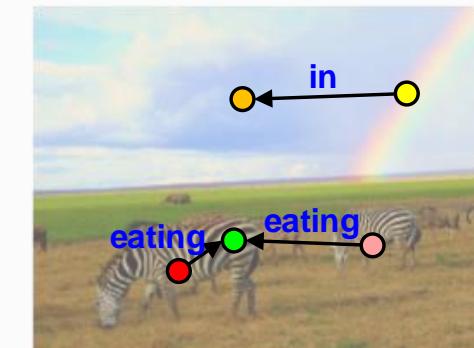
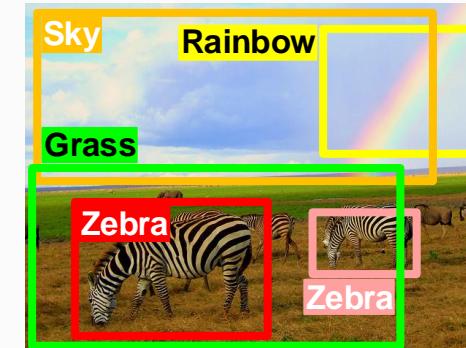
X = object/entity

[Yao, CVPR17; Li, CVPR19]



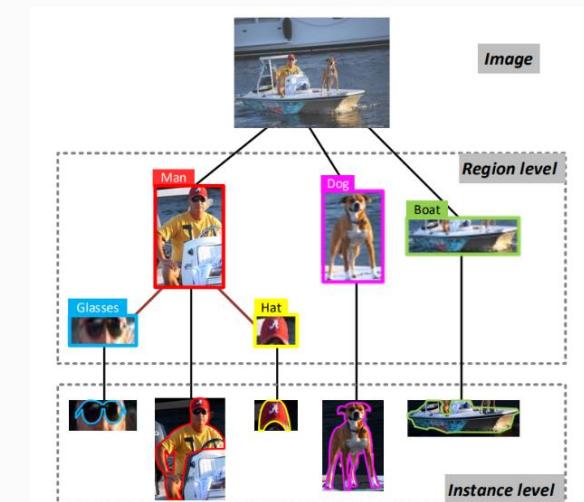
X = region/relation

[Peter, CVPR18; Yao, ECCV18]



X = instance/hierarchy

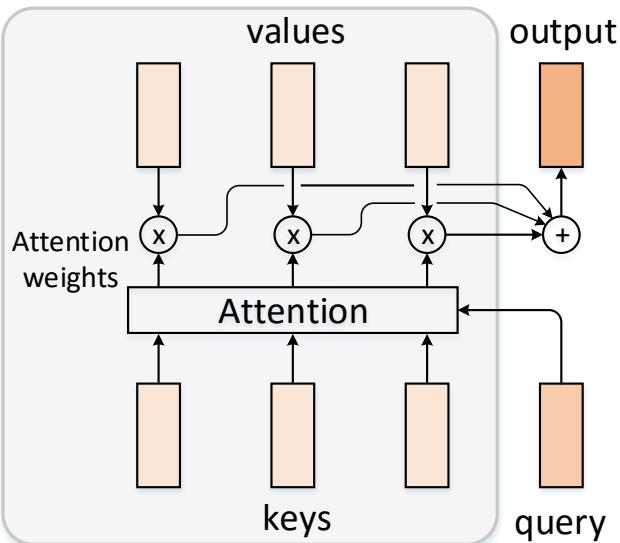
[Yao, ICCV19]



Direction II: integrate encoder/decoder via X attention

basic concept

Memory (key-value pairs)



Query (Q): hidden state from language decoder

Keys (K) = Values (V): region-level representations from image encoder

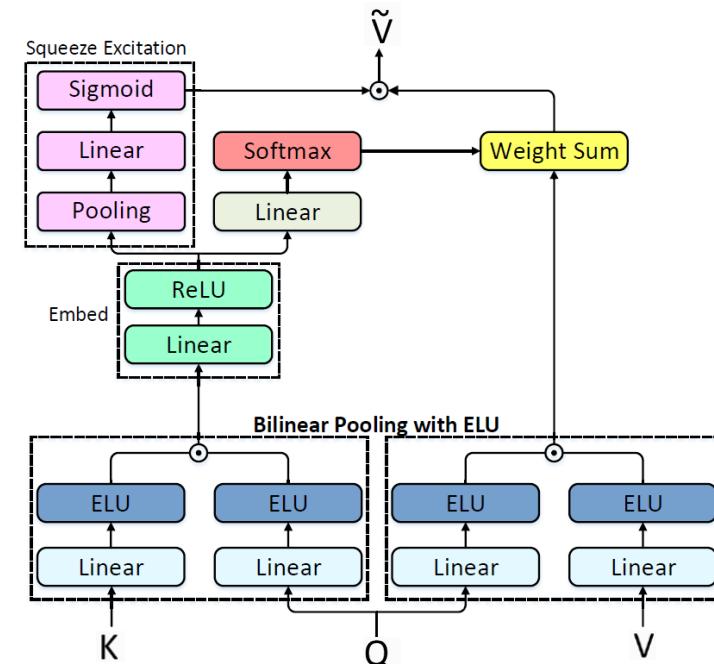
X = visual/top-down attention [Xu, ICML15; Peter, CVPR18]

[Sharma, ACL18]

X = multi-head attention [Pan, CVPR20]

[Sharma, ACL18]

X = x-linear attention [Pan, CVPR20]



Evaluations on COCO test server [March, 2020]

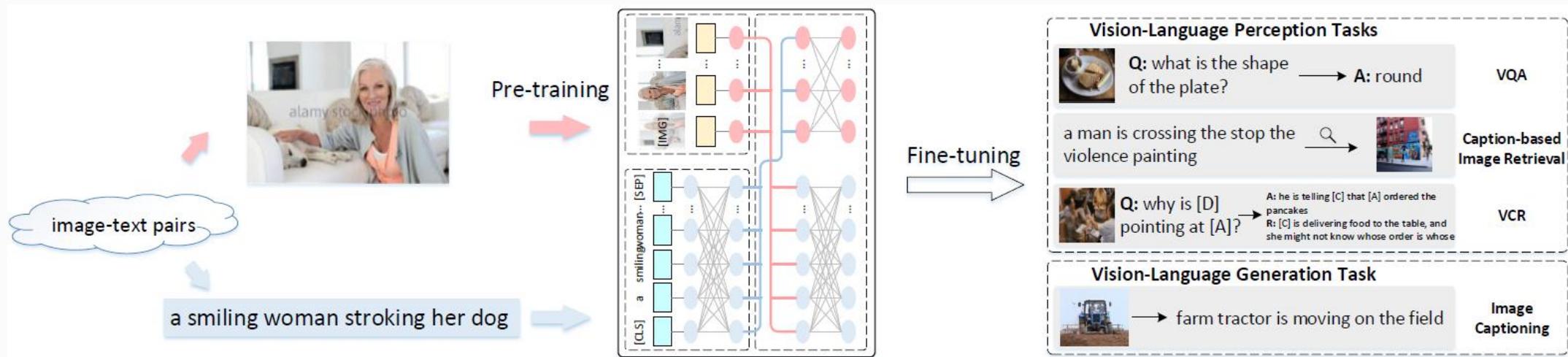
Model	Group	B@4		METEOR		ROUGE-L		CIDEr-D	
		c5	c40	c5	c40	c5	c40	c5	c40
X-LAN	Pan, et al., CVPR'20	40.3	72.4	29.6	39.2	59.5	75.0	131.1	133.5
HIP	Yao, et al., ICCV'19	39.3	71.0	28.8	38.1	59.0	74.1	127.9	130.2
AoANet	Huang, et al., ICCV'19	39.4	71.2	29.1	38.5	58.9	74.5	126.9	129.6
GCN-LSTM	Yao, et al., ECCV'18	38.7	69.7	28.5	37.6	58.5	73.4	125.3	126.5
RFNet	Jiang, et al., ECCV'18	38.0	69.2	28.2	37.2	58.2	73.1	122.9	125.1
Up-Down	Anderson, et al., CVPR'18	36.9	68.5	27.6	36.7	57.1	72.4	117.9	120.5
LSTM-A	Yao, et al., , ICCV'17	35.6	65.2	27	35.4	56.4	70.5	116	118
Watson Multimodal	Rennie, et al., CVPR'17	34.4	63.6	26.8	35.3	55.9	70.4	112.3	114.6
G-RMI	Liu, et al., ICCV'17	33.1	62.4	25.5	33.9	55.1	69.4	104.2	107.1
MetaMind/VT_GT	Lu, et al., CVPR'17	33.6	63.7	26.4	35.9	55	70.5	104.2	105.9
DLTC@MSR	Gan, et al., CVPR'17	33.1	63.1	25.7	34.8	54.3	69.6	100.3	101.3
reviewnet	Yang, et al., NIPS'16	31.3	59.7	25.6	34.7	53.3	68.6	96.5	96.9

Code:



Direction III: vision-language pre-training

- Vision-language Pre-training
 - Pre-train multi-modal encoder representation on large-scale vision-language benchmarks
 - Fine-tune multi-modal encoder on vision-language downstream tasks (e.g., VQA, VCR...)

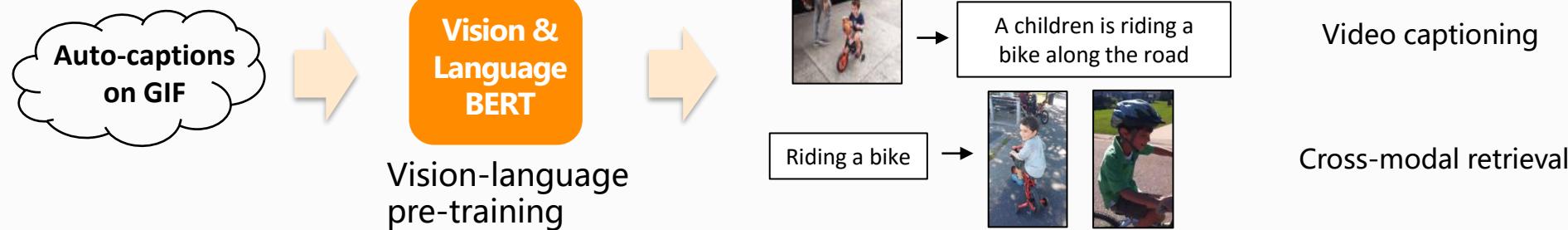


- Dataset
 - Image-Language: Conceptual Captions [Sharma, ACL18, Google]
 - Video-Language: **Auto-captions on GIF** [Pan, Arxiv20, JD AI Research]
- Technology / Open questions
 - How to design a good encoder-decoder structure and proxy tasks?
 - How to pre-train a structure for both VL understanding/generation downstream tasks?

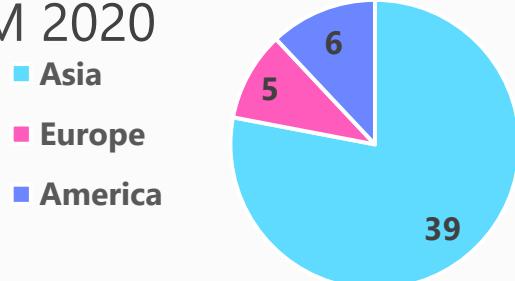
Pre-training for Video Captioning Challenge

<http://www.auto-video-captions.top/2020/>

- Task Description
 - Leverage **Auto-CapTIONs (ACTION 1.0) on GIF** benchmark to boost research on an emerging task of visual-language pre-training for downstream tasks (e.g., cross-modal retrieval, and video captioning)



- Challenge
 - 50 teams registered challenge
 - 10 teams submitted results
 - Awards will be announced at ACM MM 2020



Rank	Team	Organization	B@4	M	C	S
1	Old Boys	Tsinghua University Beijing University of Posts and Telecommunications Shanghai Ocean University	21.14	17.38	24.42	5.65
2	SYSU-CS	Sun Yat-Sen University	20.41	17.02	23.80	5.39
3	IVIPC-King	University of Electronic Science and Technology of China	18.24	16.46	21.36	5.25



双排扣大衣略带一点中性的帅气，加上legging，曲线美尽显。



黑色风衣是经典必备的基础款，显瘦显气质，绝对不会穿出错，搭配白色衬衫和深色小脚裤，简单的搭配很显气质，知性而利落。

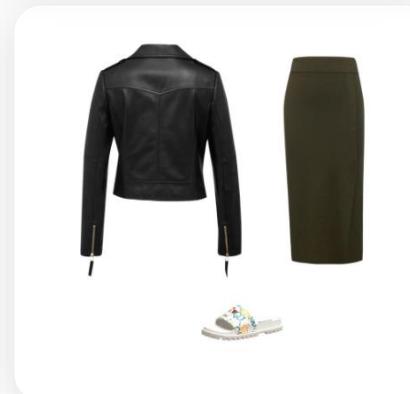
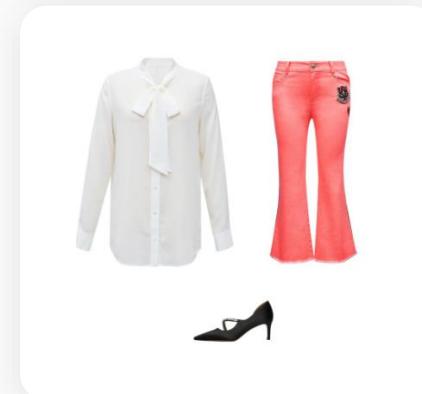
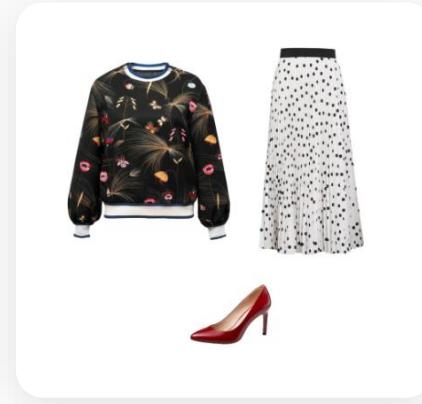
Scenario I: Fashion Collocation in JD.com 京东搭配购



User
Recommend fashion collocation based on scene and personalized needs

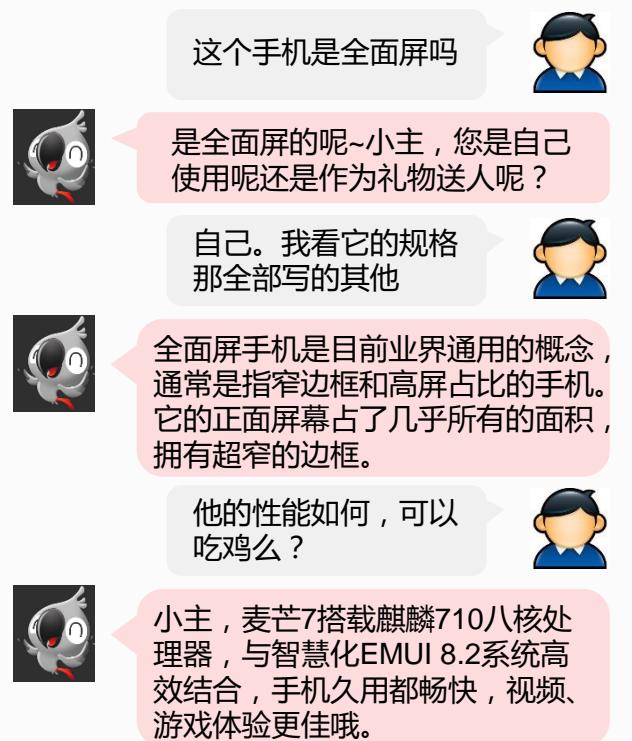
POP Shop
Automatically generate fashion product marketing plan and improve efficiency

Platform
Improve user experience, enhance user browsing time and click



Scenario II: Task-oriented multi-modal dialogue

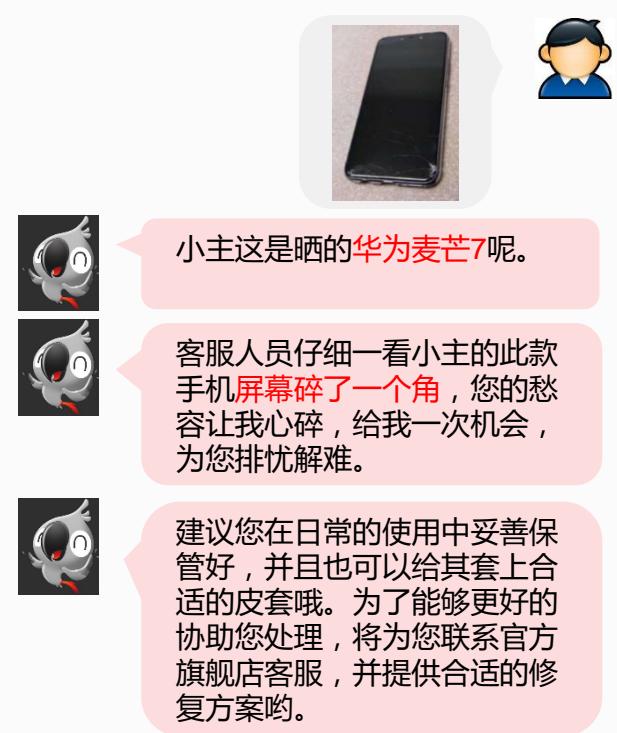
- Task-oriented dialogue for E-commerce customer service
 - Traditional dialogue: only focus on textual (or voice) modality
 - Multi-modal dialogue: capture semantics across textual and visual (image, video) modalities



Traditional task-oriented dialogue



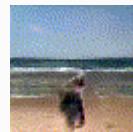
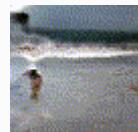
Task-oriented multi-modal dialogue



Language to Vision



Direction I: Caption to Video



beach



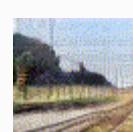
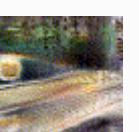
a long hair man who wears a black vest is playing Guitar



golf

digit 8 is moving left
and right digit 2 is up and
down and digit 0 is
left and right

[Yu, arxiv'20]

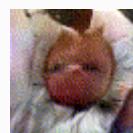
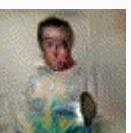
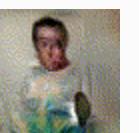


train station



a cook puts noodles
into some boiling
water

a person is
cutting beef



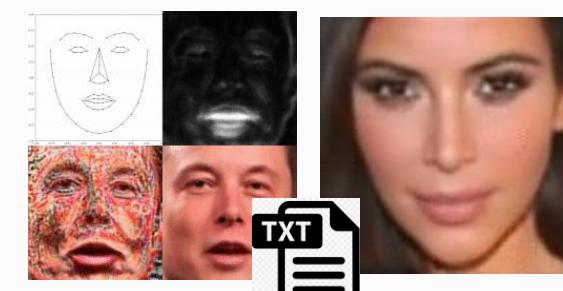
baby

[Vondrick, NIPS'16]

[Pan, Yao, Mei, ACM MM'17]

Talking head/full-body
Generation

[Chen, CVPR'19, Wang & Mei, MM'20]



Direction II: Multi-objects Scene

an old clock next to a light post in front of a steeple



a herd of sheep grazing on a lush green field



a fruit stand display with bananas and kiwi



a young girl eating a slice of pizza



A woman, man and a dog standing in the snow

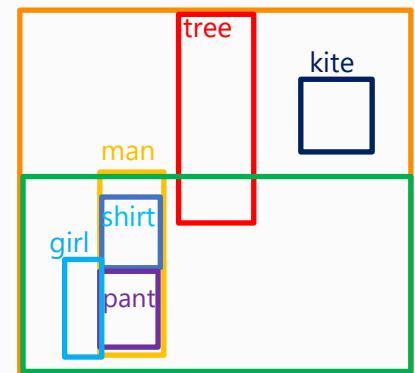
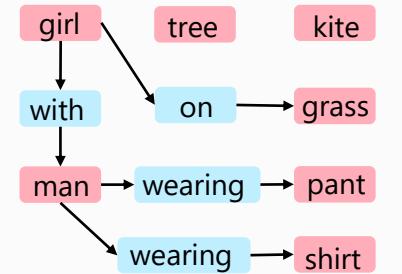


Two rectangular slices of pizza with multiple toppings



[Xu, Zhang, Huang, Zhang, Gan, Huang, He, CVPR'18]

[Hinz, Heinrich, Wermter, ICLR'19]



[Hua, Bai, Zhang, Mei, arxiv'20]

Trends & Application

- Trend 1 : Generative Pre-Training
 - Background : Lang (GPT-3, promising) + Vision (ImageGPT, emerging)
 - Lang-to-vision is naturally *Generative Pre-Training* for high-level V&L tasks
- Trend 2 : Knowledge
 - Background : Knowledge has been successful in language domain (dialog, QA)
 - Incorporating knowledge (besides data statistics) in V&L is important for boosting model generalization, explainability
- Applications
 - 文本到图像的检索
 - 小说/文章的自动配图/插画
 - 数字虚拟人 – 从“相似性”到“互补性”转变的V&L
 - 文本驱动的视觉形象 (文本 → 表情、唇形、动作)
 - 用于播报、对话场景，如客服/虚拟主播/带货