

# Vision and Language Pretrained Models

Luke Ye

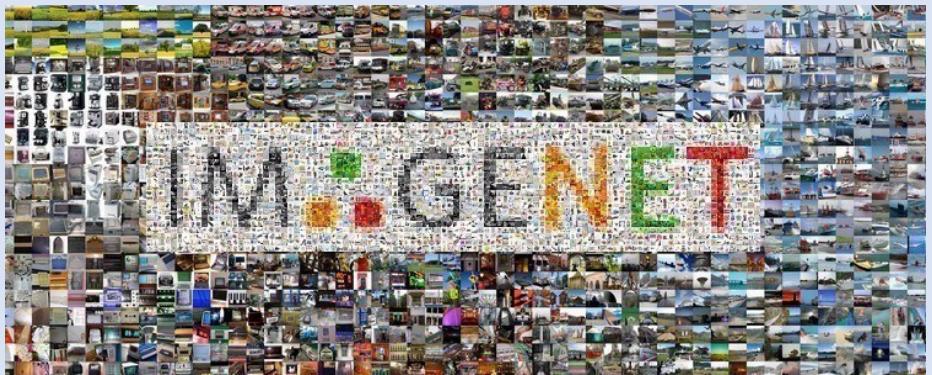
52194506006

# Outline

- Background
  - Pretrained Model
  - Vision and Language Pretrain
- Recent Advances in Vision and Language Pretrained Models
- Summary and Thinking

## Background

# Pretrained Model



Object  
Detection

Semantic  
Segmentation

Instance  
Segmentation

...

BERT (Devlin et al, 2018)



Wikipedia



BookCorpus



Question  
Answering

Sentence  
Classification

Natural  
Language  
Inference

...

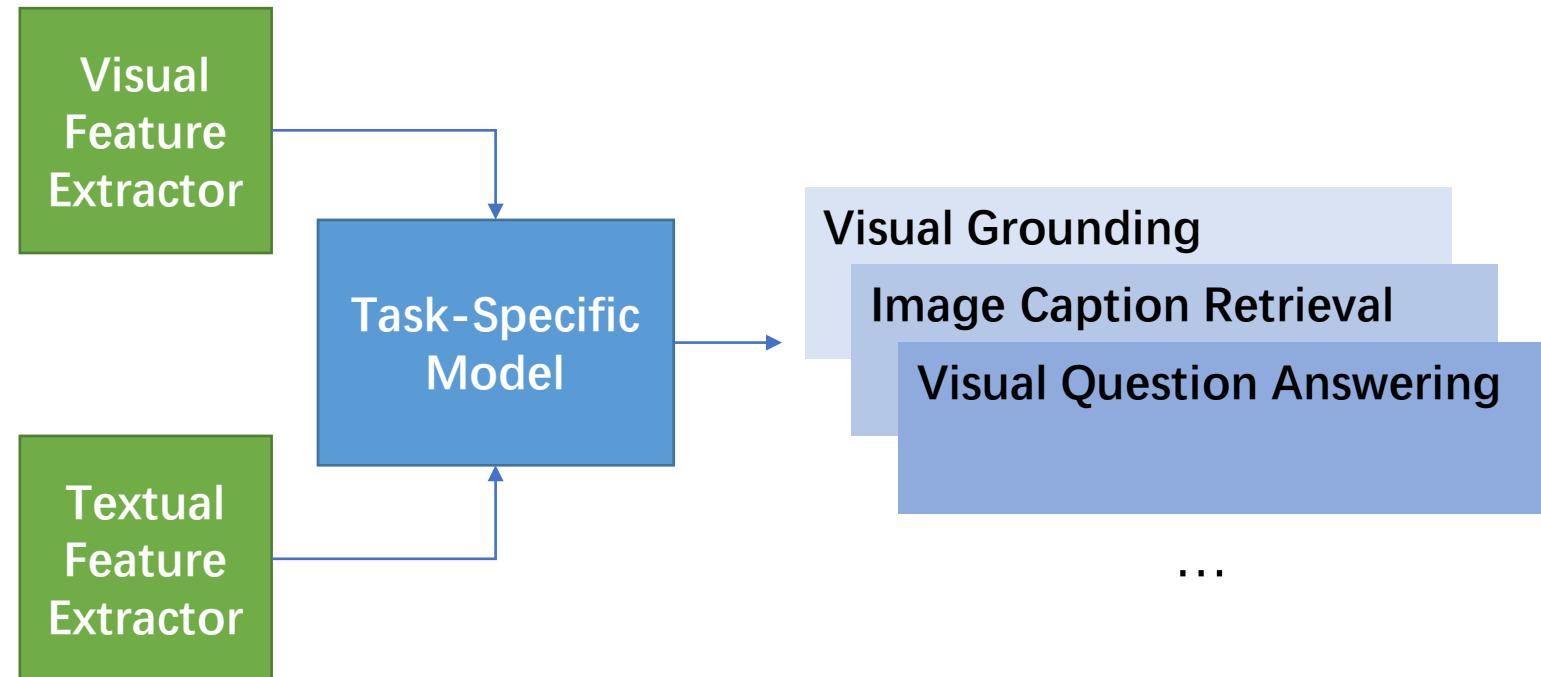
## Background

# Vision and Language Pretrain



How many slices of pizza are there?

A [puppy](#) with a [tie](#) is sitting at [table](#) with a [cake](#).



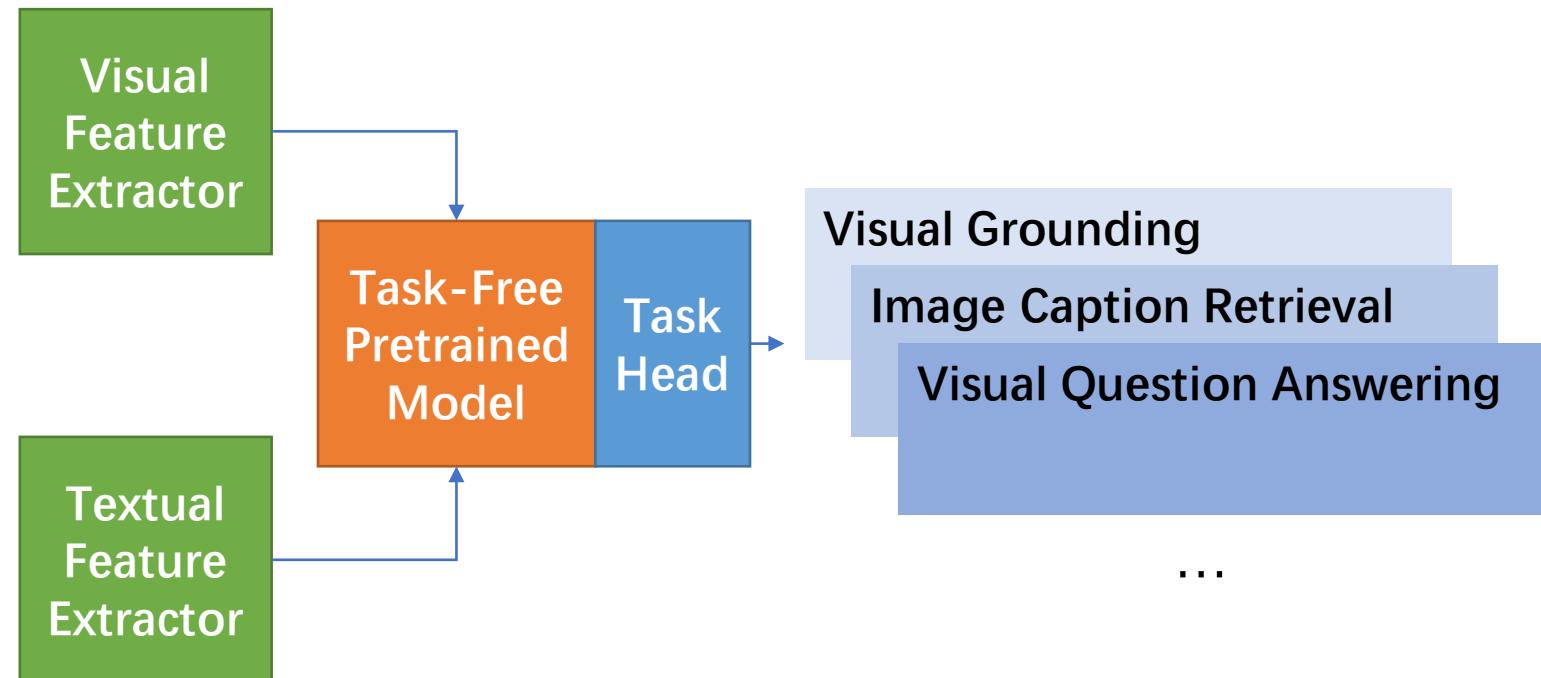
## Background

# Vision and Language Pretrain



How many slices of pizza are there?

A [puppy](#) with a [tie](#) is sitting at [table](#) with a [cake](#).



# Recent Advances in Vision and Language Pretrained Models

## ViLBERT

### Data

- Conceptual Captions

*3.3 million image-caption pairs automatically scraped from alt-text enabled web images*



**Alt-text:** A Pakistani worker helps to clear the debris from the Taj Mahal Hotel November 7, 2005 in Balakot, Pakistan.



**Conceptual Captions:** a worker helps to clear the debris.

**Alt-text:** Musician Justin Timberlake performs at the 2017 Pilgrimage Music & Cultural Festival on September 23, 2017 in Franklin, Tennessee.

**Conceptual Captions:** pop artist performs at the festival in a city.

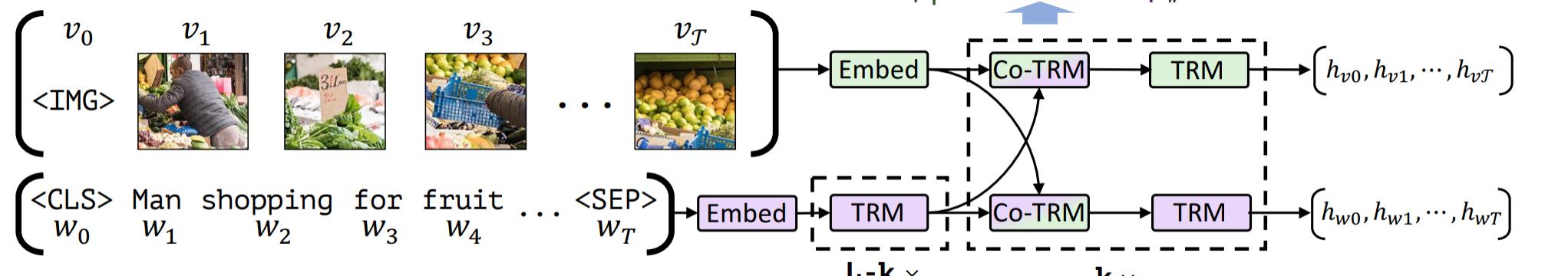
### Task

- Masked multi-modal modelling
- Multi-modal alignment prediction

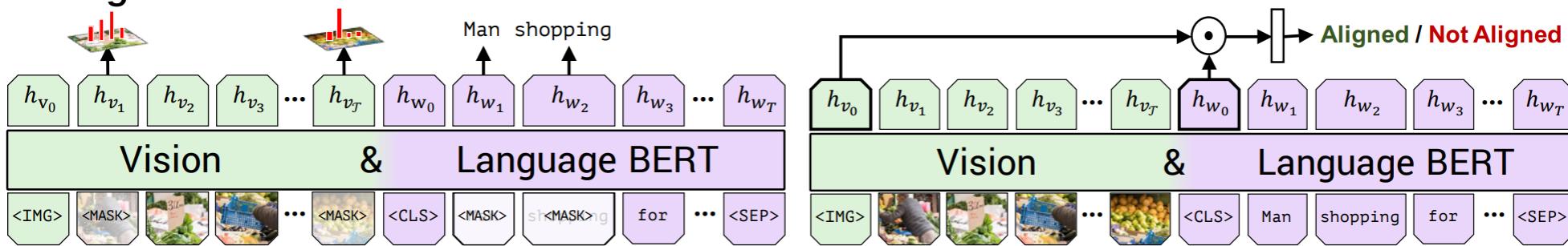
# Recent Advances in Vision and Language Pretrained Models

## ViLBERT

### Model Structure



### Training Task



(a) Masked multi-modal learning

(b) Multi-modal alignment prediction

# Recent Advances in Vision and Language Pretrained Models

## ViLBERT

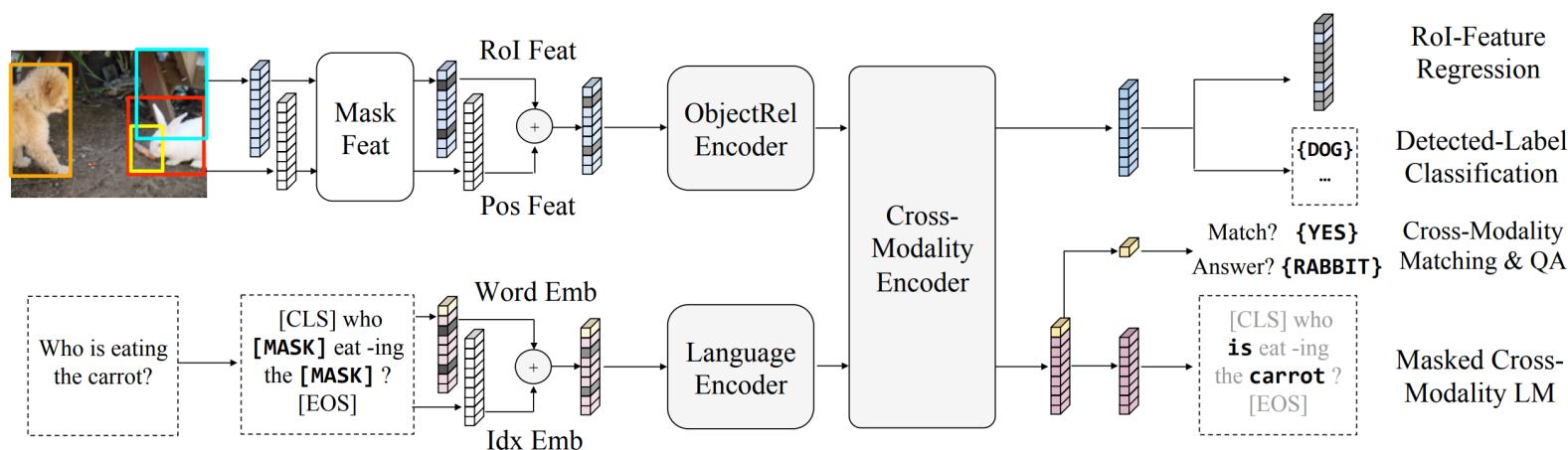
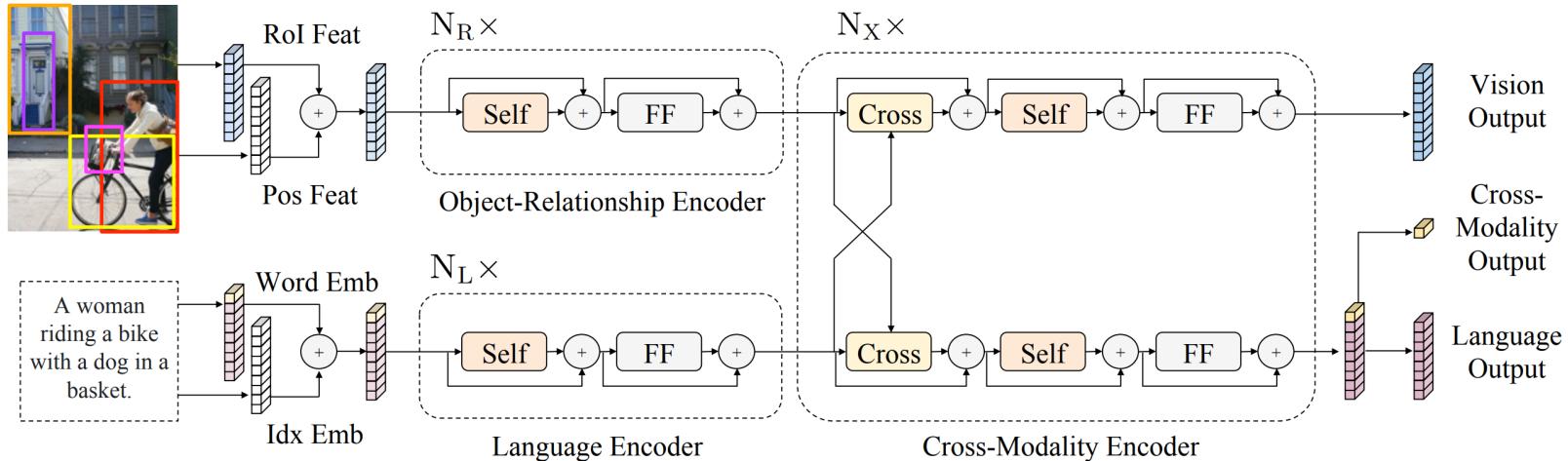
	VQA [3]			VCR [25]			RefCOCO+ [32]			Image Retrieval [26]			ZS Image Retrieval		
Method	test-dev (test-std)	Q→A	QA→R	Q→AR	val	testA	testB	R1	R5	R10	R1	R5	R10		
DFAF [36]	70.22 (70.34)	-	-	-	-	-	-	-	-	-	-	-	-		
SOTA	R2C [25]	-	63.8 (65.1)	67.2 (67.3)	43.1 (44.0)	-	-	-	-	-	-	-	-		
	MAttNet [33]	-	-	-	-	65.33	71.62	56.02	-	-	-	-	-		
	SCAN [35]	-	-	-	-	-	-	48.60	77.70	85.20	-	-	-		
	Single-Stream <sup>†</sup>	65.90	68.15	68.89	47.27	65.64	72.02	56.04	-	-	-	-	-		
Ours	Single-Stream	68.85	71.09	73.93	52.73	69.21	75.32	61.02	-	-	-	-	-		
	ViLBERT <sup>†</sup>	68.93	69.26	71.01	49.48	68.61	75.97	58.44	45.50	76.78	85.02	0.00	0.00		
	ViLBERT	<b>70.55 (70.92)</b>	<b>72.42 (73.3)</b>	<b>74.47 (74.6)</b>	<b>54.04 (54.8)</b>	<b>72.34</b>	<b>78.52</b>	<b>62.61</b>	<b>58.20</b>	<b>84.90</b>	<b>91.52</b>	<b>31.86</b>	<b>61.12</b>	<b>72.80</b>	

† : 不使用Conceptual Captions预训练

Single-Stream: 使用单个BERT模型处理视觉和文本输入

# Recent Advances in Vision and Language Pretrained Models

## LXMERT



### Pre-training Datasets

- MS COCO
  - Visual Genome
  - VQA v2.0
  - GQA
  - VG-QA
- For QA Task

# Recent Advances in Vision and Language Pretrained Models

## LXMERT

### Comparison

Method	VQA				GQA			NLVR <sup>2</sup>		(Test-set result)
	Binary	Number	Other	Accu	Binary	Open	Accu	Cons	Accu	
Human	-	-	-	-	91.2	87.4	89.3	-	96.3	
Image Only	-	-	-	-	36.1	1.74	17.8	7.40	51.9	
Language Only	66.8	31.8	27.6	44.3	61.9	22.7	41.1	4.20	51.1	
State-of-the-Art	85.8	53.7	60.7	70.4	76.0	40.4	57.1	12.0	53.5	
<b>LXMERT</b>	<b>88.2</b>	<b>54.2</b>	<b>63.1</b>	<b>72.5</b>	<b>77.8</b>	<b>45.0</b>	<b>60.3</b>	<b>42.1</b>	<b>76.2</b>	

70.92 (ViLBERT)

### Ablation studies

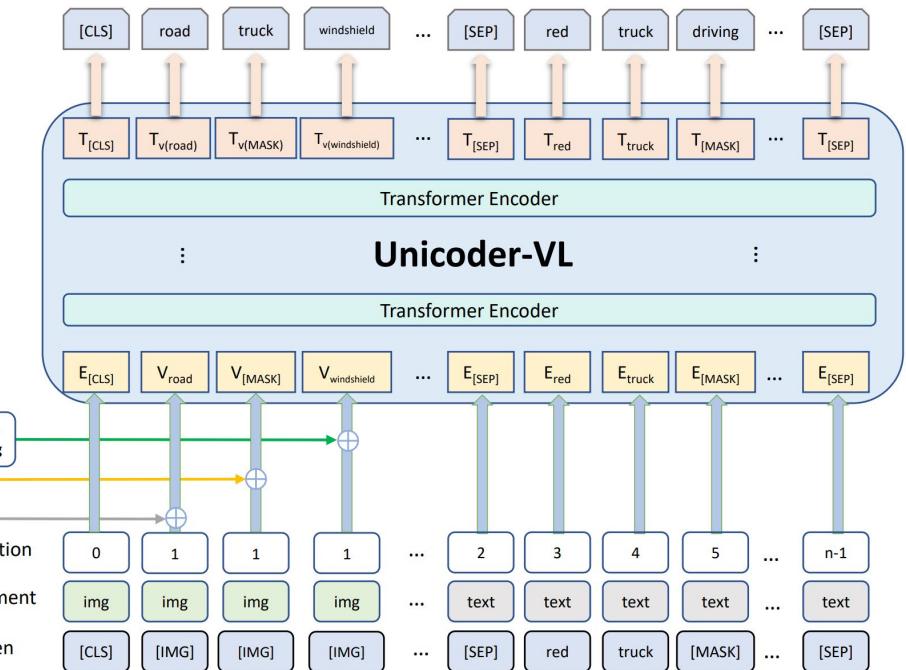
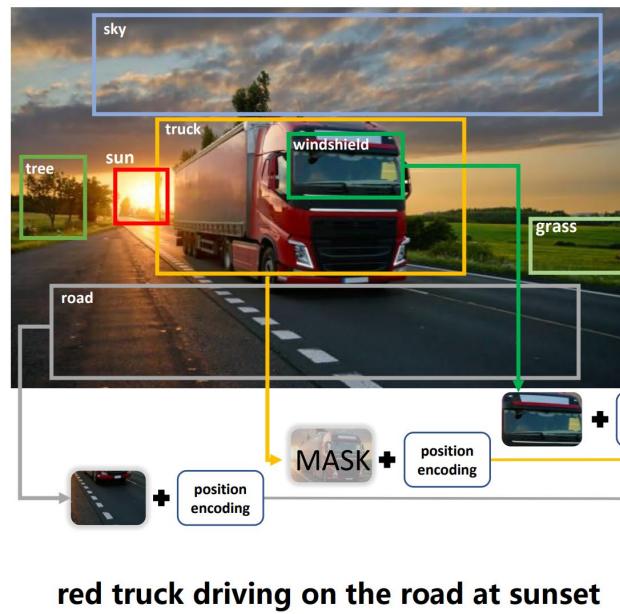
Method	VQA	GQA	NLVR <sup>2</sup>	Method	VQA	GQA	NLVR <sup>2</sup>	(Dev-set result)
1. P20 + DA	68.0	58.1	-	1. No Vision Tasks	66.3	57.1	50.9	
2. P20 + FT	68.9	58.2	72.4	2. Feat	69.2	59.5	72.9	
3. P10+QA10 + DA	69.1	59.2	-	3. Label	69.5	59.3	73.5	
<b>4. P10+QA10 + FT</b>	<b>69.9</b>	<b>60.0</b>	<b>74.9</b>	<b>4. Feat + Label</b>	<b>69.9</b>	<b>60.0</b>	<b>74.9</b>	

# Recent Advances in Vision and Language Pretrained Models

## Unicoder-VL

### Improvement

- use a single-stream architecture
- use object label classification instead minimize KL divergence

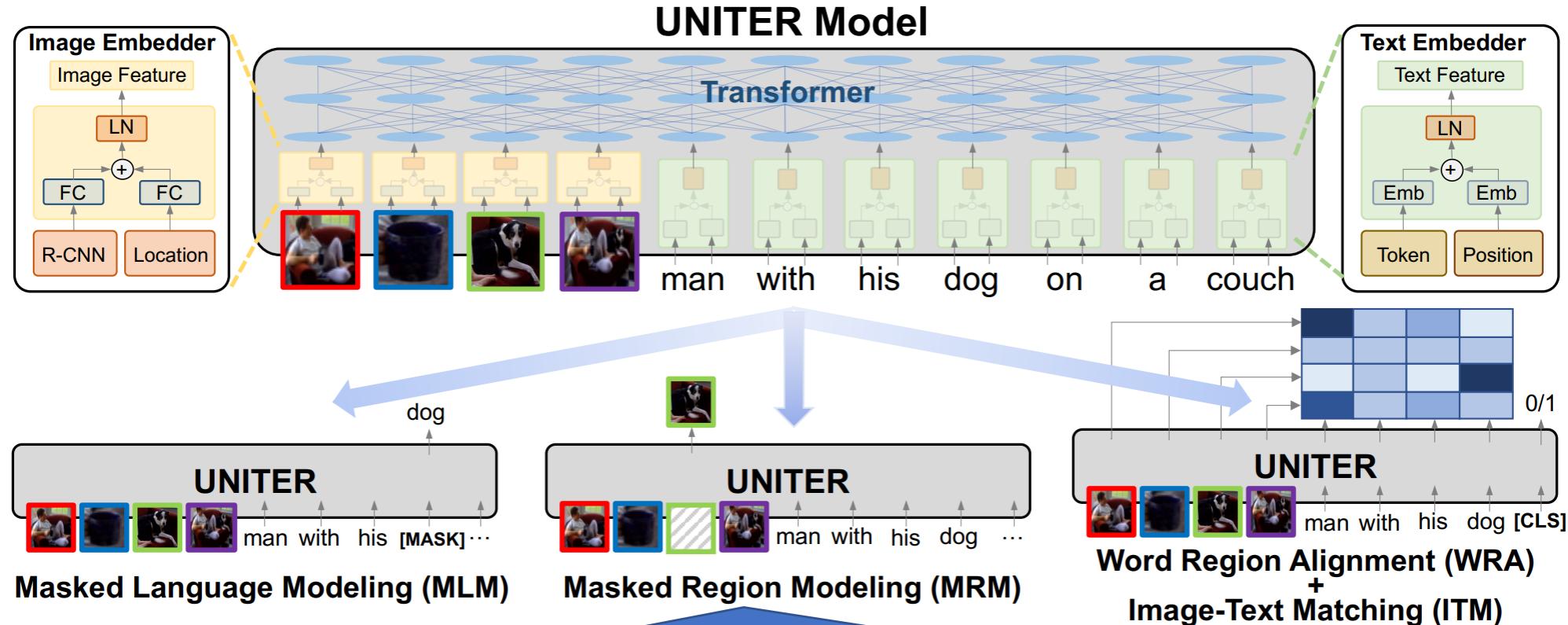


# Recent Advances in Vision and Language Pretrained Models

## Unicoder-VL

Methods	MSCOCO									Flickr30k								
	Sentence Retrieval			Image Retrieval			Sentence Retrieval			Image Retrieval								
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10						
1K Test set																		
DVSA (Karpathy and Fei-Fei 2015)	38.4	69.9	80.5	27.4	60.2	74.8	22.2	48.2	61.4	15.2	37.7	50.5						
m-CNN (Ma et al. 2015)	42.8	73.1	84.1	32.6	68.6	82.8	33.6	64.1	74.9	26.2	56.3	69.6						
DSPE (Wang, Li, and Lazebnik 2016)	50.1	79.7	89.2	39.6	75.2	86.9	40.3	68.9	79.9	29.7	60.1	72.1						
VSE++ (Faghri et al. 2017)	64.7	-	95.9	52.0	-	92.0	52.9	79.1	87.2	39.6	69.6	79.5						
SCAN (Lee et al. 2018)	72.7	94.8	98.4	58.8	88.4	94.8	67.4	90.3	95.8	48.6	77.7	85.2						
SCG (Shi et al. 2019)	76.6	96.3	99.2	61.4	88.9	95.1	71.8	90.8	94.8	49.3	76.4	85.6						
PFAN (Wang et al. 2019)	76.5	96.3	99.0	61.6	89.6	95.2	70.0	91.8	95.0	50.4	78.7	86.1						
ViLBERT (Lu et al. 2019) <sup>†</sup>	-	-	-	-	-	-	-	-	-	58.2	84.9	91.5						
UNITER (Chen et al. 2019) <sup>†</sup>	-	-	-	-	-	-	84.7	<b>97.1</b>	99.0	71.5	<b>91.2</b>	<b>95.2</b>						
Unicoder-VL (zero-shot)	54.4	82.8	90.6	43.4	76.0	87.0	64.3	85.8	92.3	48.4	76.0	85.2						
Unicoder-VL (w/o pre-training)	75.1	94.3	97.8	63.9	91.6	96.5	73.0	89.0	94.1	57.8	82.2	88.9						
Unicoder-VL	<b>84.3</b>	<b>97.3</b>	<b>99.3</b>	<b>69.7</b>	<b>93.5</b>	<b>97.2</b>	<b>86.2</b>	96.3	<b>99.0</b>	<b>71.5</b>	90.9	94.9						
5K Test set																		
SCAN (Lee et al. 2018)	50.4	82.2	90.0	38.6	69.3	80.4	-	-	-	-	-	-						
SCG (Shi et al. 2019)	56.6	84.5	92.0	39.2	68.0	81.3	-	-	-	-	-	-						
UNITER (Chen et al. 2019) <sup>†</sup>	<b>63.3</b>	87.0	<b>93.1</b>	<b>48.4</b>	<b>76.7</b>	<b>85.9</b>	-	-	-	-	-	-						
Unicoder-VL	62.3	<b>87.1</b>	92.8	46.7	76.0	85.3	-	-	-	-	-	-						

# UNITER

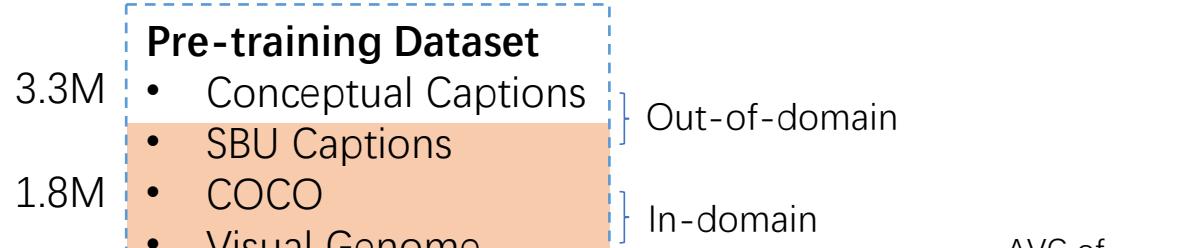


- Masked Region Classification (MRC)
- Masked Region Feature Regression (MRFR)
- Masked Region Classification with KL-divergence (MRC-kl)

# Recent Advances in Vision and Language Pretrained Models

## UNITER

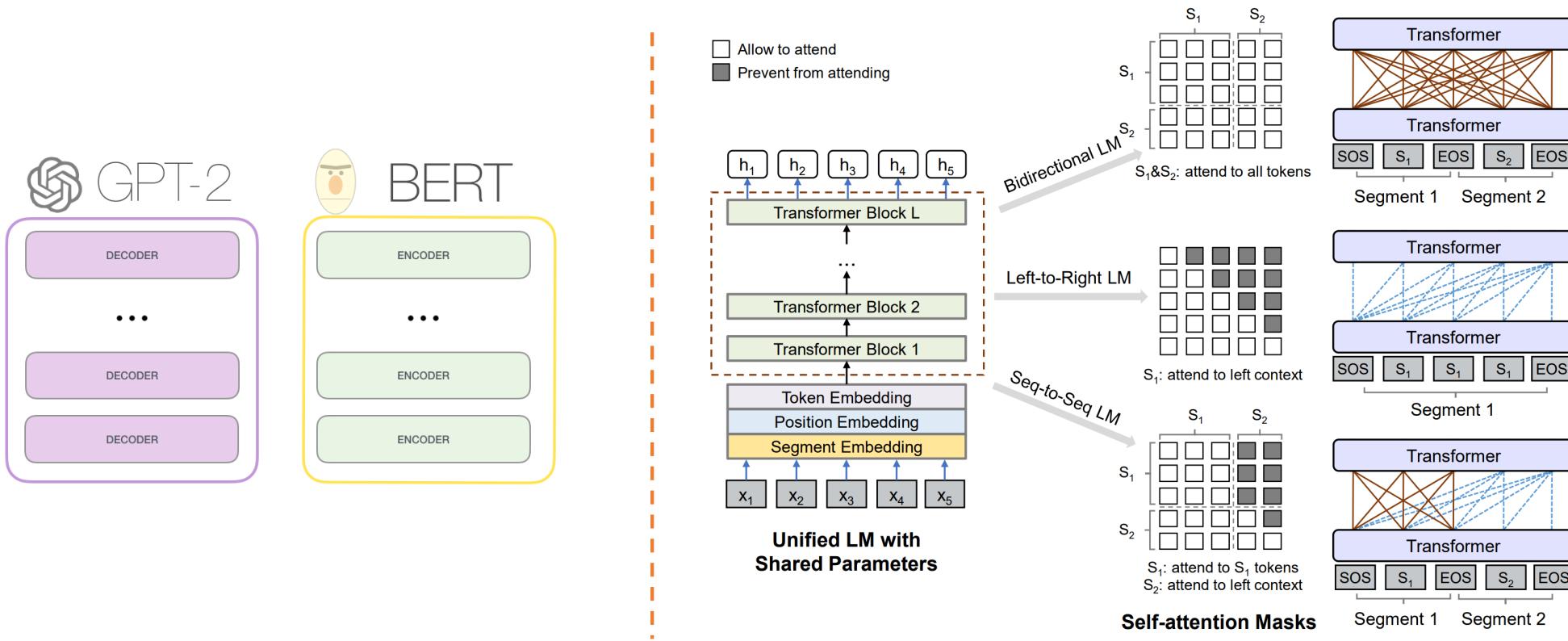
Tasks	12 layer						24 layer	
	SOTA	ViLBERT	VLBERT (Large)	Unicoder -VL	VisualBERT	LXMERT	UNITER Base	Large
VQA	test-dev	70.63	70.55	71.79	-	70.80	72.42	72.70 <b>73.82</b>
	test-std	70.90	70.92	72.22	-	71.00	72.54	72.91 <b>74.02</b>
VCR	Q→A	72.60	73.30	<b>75.80</b>	-	71.60	-	75.00 <b>77.30</b>
	QA→R	75.70	74.60	78.40	-	73.20	-	77.20 <b>80.80</b>
	Q→AR	55.00	54.80	<b>59.70</b>	-	52.40	-	58.20 <b>62.80</b>
NLVR <sup>2</sup>	dev	54.80	-	-	-	67.40	74.90	77.18 <b>79.12</b>
	test-P	53.50	-	-	-	67.00	74.50	77.85 <b>79.98</b>
SNLI-VE	val	71.56	-	-	-	-	-	78.59 <b>79.39</b>
	test	71.16	-	-	-	-	-	78.28 <b>79.38</b>
ZS IR (Flickr)	R@1	-	31.86	-	48.40	-	-	66.16 <b>68.74</b>
	R@5	-	61.12	-	76.00	-	-	88.40 <b>89.20</b>
	R@10	-	72.80	-	85.20	-	-	92.94 <b>93.86</b>
IR (Flickr)	R@1	48.60	58.20	-	71.50	-	-	72.52 <b>75.56</b>
	R@5	77.70	84.90	-	91.20	-	-	92.36 <b>94.08</b>
	R@10	85.20	91.52	-	95.20	-	-	96.08 <b>96.76</b>
IR (COCO)	R@1	38.60	-	-	48.40	-	-	50.33 <b>52.93</b>
	R@5	69.30	-	-	76.70	-	-	78.52 <b>79.93</b>
	R@10	80.40	-	-	85.90	-	-	87.16 <b>87.95</b>
ZS TR (Flickr)	R@1	-	-	-	64.30	-	-	80.70 <b>83.60</b>
	R@5	-	-	-	85.80	-	-	<b>95.70</b> <b>95.70</b>
	R@10	-	-	-	92.30	-	-	<b>98.00</b> 97.70
TR (Flickr)	R@1	67.90	-	-	<b>86.20</b>	-	-	85.90 <b>87.30</b>
	R@5	90.30	-	-	96.30	-	-	97.10 <b>98.00</b>
	R@10	95.80	-	-	<b>99.00</b>	-	-	98.80 <b>99.20</b>
TR (COCO)	R@1	50.40	-	-	62.30	-	-	64.40 <b>65.68</b>
	R@5	82.20	-	-	87.10	-	-	87.40 <b>88.56</b>
	R@10	90.00	-	-	92.80	-	-	93.08 <b>93.76</b>
Ref-COCO	val	87.51	-	-	-	-	91.64	<b>91.84</b>
	testA	89.02	-	-	-	-	92.26	<b>92.65</b>
	testB	87.05	-	-	-	-	90.46	<b>91.19</b>
	val <sup>d</sup>	77.48	-	-	-	-	81.24	<b>81.41</b>
	testA <sup>d</sup>	83.37	-	-	-	-	86.48	<b>87.04</b>
Ref-COCO+VG	testB <sup>d</sup>	70.32	-	-	-	-	73.94	<b>74.17</b>
	val	75.38	-	80.31	-	-	83.66	<b>84.25</b>
	testA	80.04	-	83.62	-	-	86.19	<b>86.34</b>
	testB	69.30	-	75.45	-	-	78.89	<b>79.75</b>
	val <sup>d</sup>	68.19	72.34	72.59	-	-	75.31	<b>75.90</b>
Ref-COCO	testA <sup>d</sup>	75.97	78.52	78.57	-	-	81.30	<b>81.45</b>
	testB <sup>d</sup>	57.52	62.61	62.30	-	-	65.58	<b>66.70</b>
	val	81.76	-	-	-	-	86.52	<b>87.85</b>
Ref-COCOg	test	81.75	-	-	-	-	86.52	<b>87.73</b>
	val <sup>d</sup>	68.22	-	-	-	-	74.31	<b>74.86</b>
	test <sup>d</sup>	69.46	-	-	-	-	74.51	<b>75.77</b>



Pre-training Data	Pre-training Tasks	Meta-Sum	VQA		IR		TR		NLVR <sup>2</sup>		Ref-COCO+	
			test-dev	val	(Flickr)	val	(Flickr)	dev	val <sup>d</sup>			
None	1 None	314.34	67.03	61.74	65.55	51.02	68.73					
Wiki + BookCorpus	2 MLM (text only)	346.24	69.39	73.92	83.27	50.86	68.80					
In-domain (COCO+VG)	3 MRFR	344.66	69.02	72.10	82.91	52.16	68.47					
	4 ITM	385.29	70.04	78.93	89.91	74.08	72.33					
	5 MLM	386.10	71.29	77.88	89.25	74.79	72.89					
	6 MLM + ITM	393.04	71.55	81.64	91.12	75.98	72.75					
	7 MLM + ITM + MRC	393.97	71.46	81.39	91.45	76.18	73.49					
	8 MLM + ITM + MRFR	396.24	71.73	81.76	92.31	76.21	74.23					
	9 MLM + ITM + MRC-kl	397.09	71.63	82.10	92.57	76.28	74.51					
	10 MLM + ITM + MRC-kl + MRFR	399.97	71.92	83.73	92.87	76.93	74.52					
	11 MLM + ITM + MRC-kl + MRFR + WRA	400.93	72.47	83.72	93.03	76.91	74.80					
	MLM + ITM + MRC-kl + MRFR	396.51	71.68	82.31	92.08	76.15	74.29					
	12 (w/o cond. mask)											
Out-of-domain (SBU+CC)	13 MLM + ITM + MRC-kl + MRFR + WRA	396.91	71.56	84.34	92.57	75.66	72.78					
	14 MLM + ITM + MRC-kl + MRFR + WRA	<b>405.24</b>	<b>72.70</b>	<b>85.77</b>	<b>94.28</b>	<b>77.18</b>	<b>75.31</b>					
In-domain + Out-of-domain												

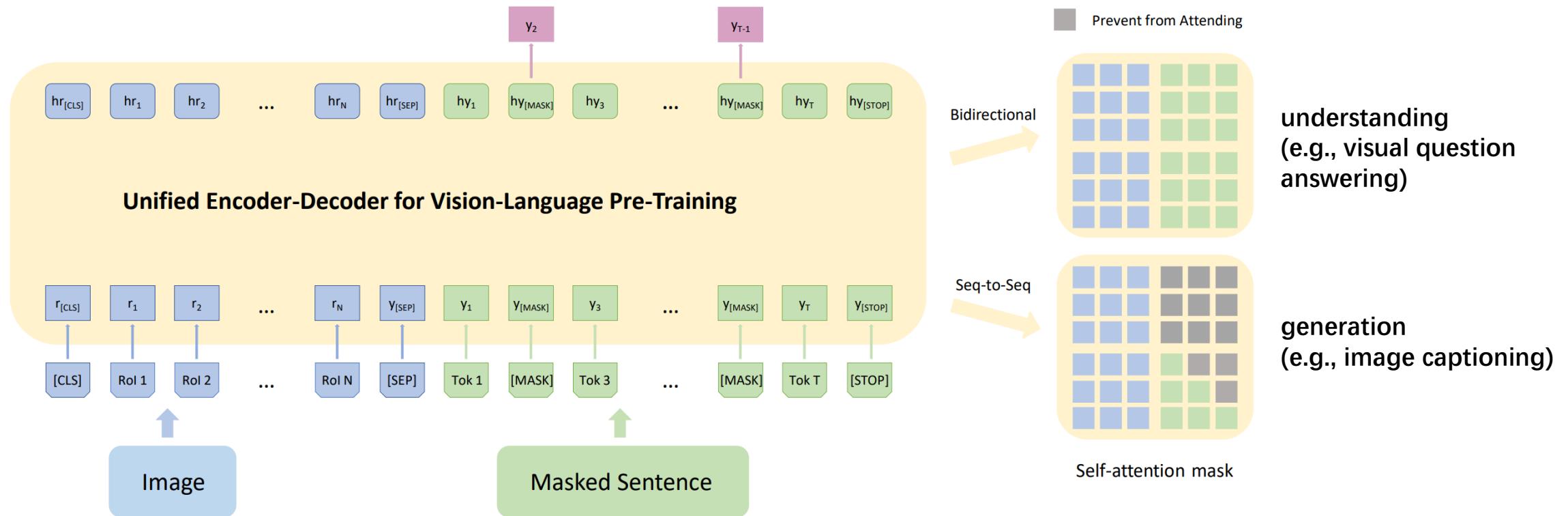
# Recent Advances in Vision and Language Pretrained Models

## Unified VLP - Background



# Recent Advances in Vision and Language Pretrained Models

## Unified VLP



# Recent Advances in Vision and Language Pretrained Models

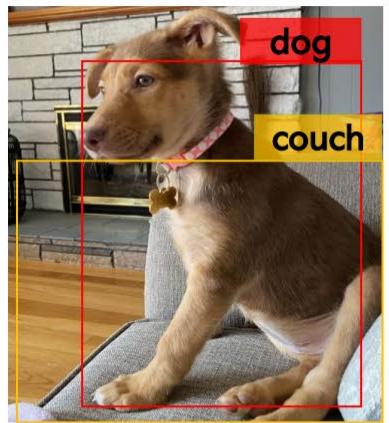
## Unified VLP

Method	COCO				VQA 2.0 (Test-Standard)				Flickr30k			
	B@4	M	C	S	Overall	Yes/No	Number	Other	B@4	M	C	S
BUTD (Anderson et al. 2018)	36.2	27.0	113.5	20.3	65.7	-	-	-	27.3	21.7	56.6	16.0
NBT (with BBox) (Lu et al. 2018)	34.7	27.1	107.2	20.1	-	-	-	-	27.1	21.7	57.5	15.6
GCN-LSTM (spa) (Yao et al. 2018)	<b>36.5</b>	27.8	115.6	20.8	-	-	-	-	-	-	-	-
GCN-LSTM (sem)	<b>36.8</b>	27.9	116.3	20.9	-	-	-	-	-	-	-	-
GVD (Zhou et al. 2019)	-	-	-	-	-	-	-	-	26.9	22.1	60.1	16.1
GVD (with BBox)	-	-	-	-	-	-	-	-	27.3	22.5	62.3	16.5
BAN (Kim, Jun, and Zhang 2018)	-	-	-	-	70.4	85.8	<b>53.7</b>	<b>60.7</b>	-	-	-	-
DFAF (Gao et al. 2019)	-	-	-	-	70.3	-	-	-	-	-	-	-
AoANet* (Huang et al. 2019)	37.2	28.4	119.8	21.3	-	-	-	-	-	-	-	-
ViLBERT* (Lu et al. 2019)	-	-	-	-	70.9	-	-	-	-	-	-	-
LXMERT* (Tan and Bansal 2019)	-	-	-	-	72.5	88.2	54.2	63.1	-	-	-	-
<b>Ours</b>												
w/o VLP pre-training (baseline)	35.5	28.2	114.3	21.0	70.0	86.3	52.2	59.9	27.6	20.9	56.8	15.3
seq2seq pre-training only	<b>36.5</b>	<b>28.4</b>	<b>117.7</b>	<b>21.3</b>	70.2	86.7	52.7	59.9	<b>31.1</b>	<b>23.0</b>	<b>68.5</b>	<b>17.2</b>
bidirectional pre-training only	36.1	28.3	116.5	21.2	<b>71.3</b>	<b>87.6</b>	<b>53.5</b>	<b>61.2</b>	<b>30.5</b>	22.6	63.3	16.9
Unified VLP	<b>36.5</b>	<b>28.4</b>	<b>116.9</b>	<b>21.2</b>	<b>70.7</b>	<b>87.4</b>	52.1	60.5	30.1	<b>23.0</b>	<b>67.4</b>	<b>17.0</b>

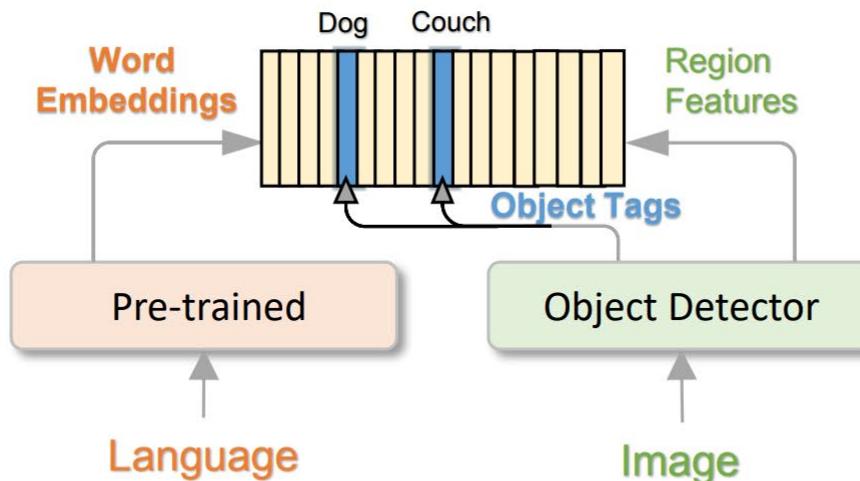
\* indicates unpublished works

# Recent Advances in Vision and Language Pretrained Models

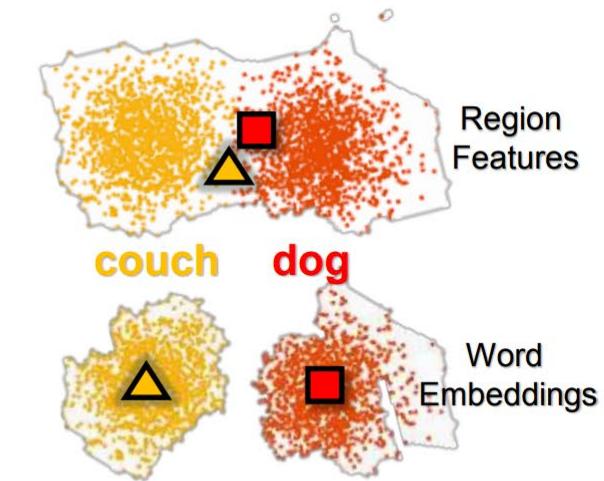
## Oscar



a) Image-text pair



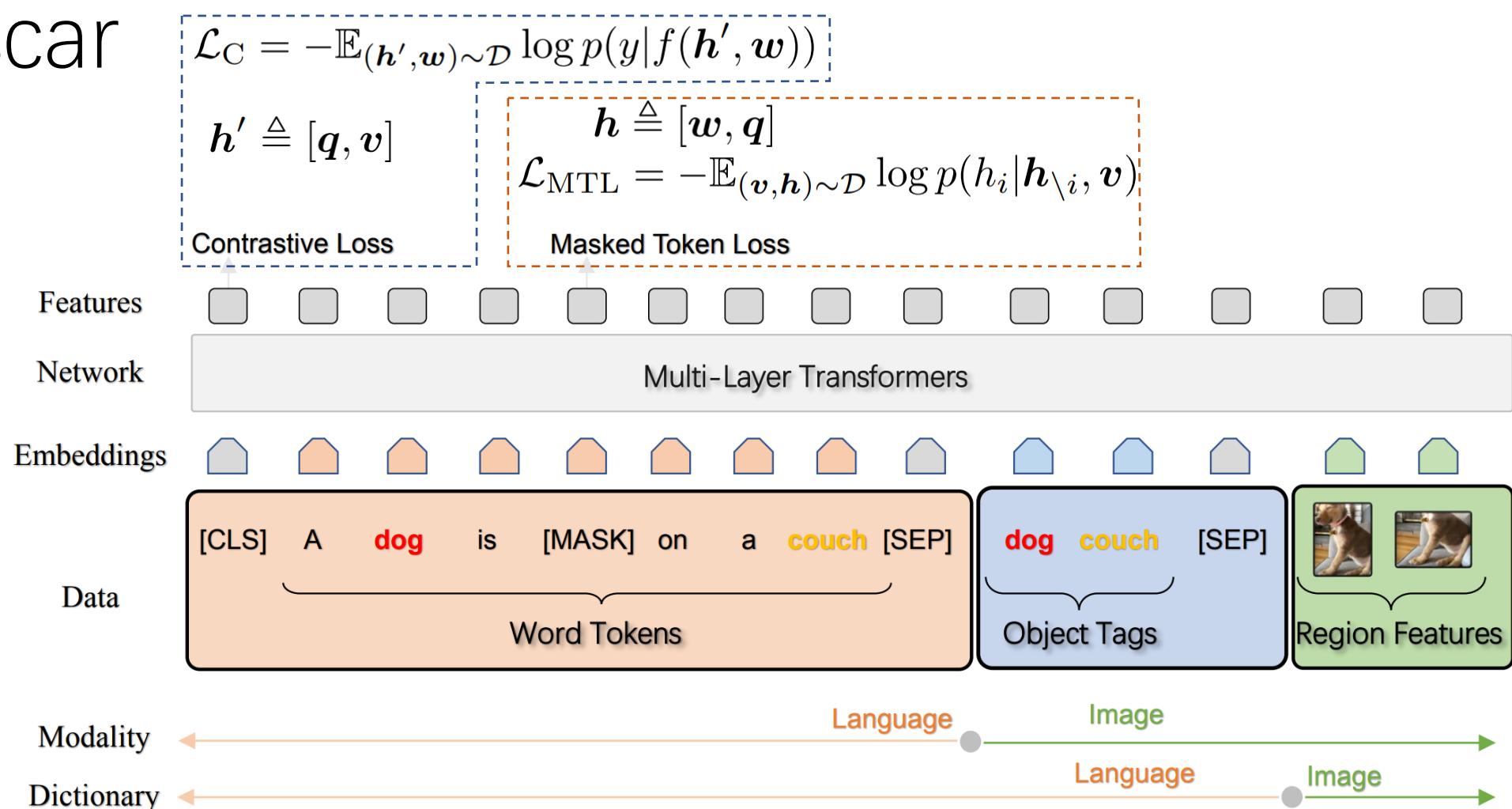
(b) Objects as anchor points



(c) Semantics spaces

# Recent Advances in Vision and Language Pretrained Models

Oscar



The image-text pair was represent as a triple [ **word tokens** , **object tags** , **region features** ]

# Recent Advances in Vision and Language Pretrained Models

# Oscar

Method	Size	Text Retrieval			Image Retrieval			Text Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
1K Test Set													
DVSA [14]	-	38.4	69.9	80.5	27.4	60.2	74.8	-	-	-	-	-	-
VSE++ [7]	-	64.7	-	95.9	52.0	-	92.0	41.3	-	81.2	30.3	-	72.4
DPC [46]	-	65.6	89.8	95.5	47.1	79.9	90.0	41.2	70.5	81.1	25.3	53.4	66.4
CAMP [42]	-	72.3	94.8	98.3	58.5	87.9	95.0	50.1	82.1	89.7	39.0	68.9	80.2
SCAN [18]	-	72.7	94.8	98.4	58.8	88.4	94.8	50.4	82.2	90.0	38.6	69.3	80.4
SCG [33]	-	76.6	96.3	99.2	61.4	88.9	95.1	56.6	84.5	92.0	39.2	68.0	81.3
PFAN [41]	-	76.5	96.3	99.0	61.6	89.6	95.2	-	-	-	-	-	-
Unicoder-VL [19]	B	84.3	97.3	99.3	69.7	93.5	97.2	62.3	87.1	92.8	46.7	76.0	85.3
12-in-1 [24]	B	-	-	-	65.2	91.0	96.2	-	-	-	-	-	-
UNITER [5]	B	-	-	-	-	-	-	63.3	87.0	93.1	48.4	76.7	85.9
UNITER [5]	L	-	-	-	-	-	-	66.6	89.4	94.3	51.7	78.4	86.9
OSCAR	B	88.4	<b>99.1</b>	<b>99.8</b>	75.7	95.2	98.3	70.0	91.1	95.5	54.0	80.8	88.5
OSCAR	L	<b>89.8</b>	98.8	99.7	<b>78.2</b>	<b>95.8</b>	<b>98.3</b>	<b>73.5</b>	<b>92.2</b>	<b>96.0</b>	<b>57.5</b>	<b>82.8</b>	<b>89.8</b>

(a) Image-text retrieval

Method	ViLBERT	VL-BERT	VisualBERT	LXMERT	12-in-1	UNITER <sub>B</sub>	UNITER <sub>L</sub>	OSCAR <sub>B</sub>	OSCAR <sub>L</sub>
Test-dev	70.63	70.50	70.80	72.42	73.15	72.27	<b>73.24</b>	73.16	<b>73.61</b>
Test-std	70.92	70.83	71.00	72.54	-	72.46	73.40	<b>73.44</b>	<b>73.82</b>

(b) VQA

Method	MAC	VisualBERT	LXMERT	12-in-1	UNITER <sub>B</sub>	UNITER <sub>L</sub>	OSCAR <sub>B</sub>	OSCAR <sub>L</sub>
Dev	50.8	67.40	74.90	-	77.14	<b>78.40</b>	78.07	<b>79.12</b>
Test-P	51.4	67.00	74.50	78.87	77.87	<b>79.50</b>	78.36	<b>80.37</b>

(c) NLVR2

Method	MAC	VisualBERT	LXMERT	12-in-1	UNITER <sub>B</sub>	UNITER <sub>L</sub>	OSCAR <sub>B</sub>	OSCAR <sub>L</sub>
Dev	50.8	67.40	74.90	-	77.14	<b>78.40</b>	78.07	<b>79.12</b>
Test-P	51.4	67.00	74.50	78.87	77.87	<b>79.50</b>	78.36	<b>80.37</b>

(c) NLVR2

Method	Test-dev		Test-std		Method	cross-entropy		optimization		Method	CIDEr optimization		
	B@4	M	C	S		B@4	M	C	S		B@4	M	
LXMERT [39]	60.00	-	60.33	-	BUTD [2]	36.2	27.0	113.5	20.3	36.3	27.7	120.1	21.4
MMN [4]	-	-	60.83	-	VLP [47]	36.5	28.4	117.7	21.3	39.5	29.3	129.3	23.2
12-in-1 [24]	-	-	60.65	-	AoANet [11]	37.2	28.4	119.8	21.3	38.9	29.2	129.8	22.4
NSM [12]	-	-	<b>63.17</b>	-	OSCAR <sub>B</sub>	36.5	<b>30.3</b>	<b>123.7</b>	<b>23.1</b>	<b>40.5</b>	<b>29.7</b>	<b>137.6</b>	<b>22.8</b>
OSCAR <sub>B</sub>	61.19	-	61.23	-	OSCAR <sub>L</sub>	<b>37.4</b>	<b>30.7</b>	<b>127.8</b>	<b>23.5</b>	<b>41.7</b>	<b>30.6</b>	<b>140.0</b>	<b>24.5</b>
OSCAR <sub>B</sub> *	<b>61.58</b>	-	<b>61.62</b>	-									

(d) GQA

Method	in-domain		near-domain		out-of-domain		overall	
	CIDEr	SPICE	CIDEr	SPICE	CIDEr	SPICE	CIDEr	SPICE
UpDown [1]	78.1	11.6	57.7	10.3	31.3	8.3	55.3	10.1
UpDown + CBS [1]	80.0	12.0	73.6	11.3	66.4	9.7	73.1	11.1
UpDown + ELMo + CBS [1]	79.3	<b>12.4</b>	73.8	11.4	71.7	9.9	74.3	11.2
OSCAR <sub>B</sub>	79.6	12.3	66.1	11.5	45.3	9.7	63.8	11.2
OSCAR <sub>B</sub> + CBS	80.0	12.1	80.4	<b>12.2</b>	75.3	<b>10.6</b>	79.3	<b>11.9</b>
OSCAR <sub>B</sub> + SCST + CBS	<b>83.4</b>	12.0	<b>81.6</b>	12.0	<b>77.6</b>	<b>10.6</b>	<b>81.1</b>	11.7
OSCAR <sub>L</sub>	79.9	<b>12.4</b>	68.2	11.8	45.1	9.4	65.2	11.4
OSCAR <sub>L</sub> + CBS	78.8	12.2	78.9	<b>12.1</b>	77.4	10.5	78.6	<b>11.8</b>
OSCAR <sub>L</sub> + SCST + CBS	<b>85.4</b>	11.9	<b>84.0</b>	11.7	<b>80.3</b>	10.0	<b>83.4</b>	11.4

(f) Evaluation on NoCaps Val. Models are trained on COCO only without pre-training.

## Training and inference tricks in image caption task

- Constrained beam search (CBS)
- self-critical sequence training (SCST)

# Recent Advances in Vision and Language Pretrained Models

Oscar

Oscar

Baseline(no tags)

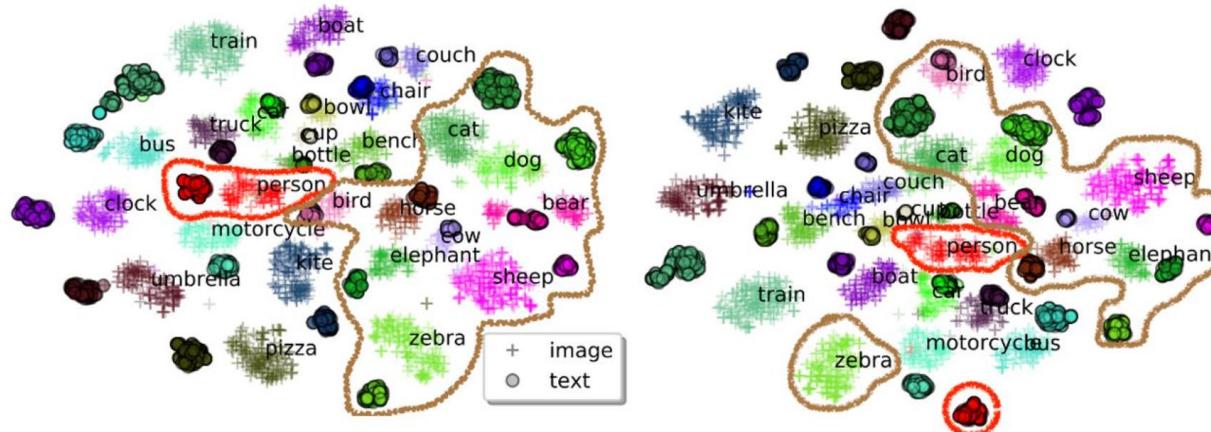


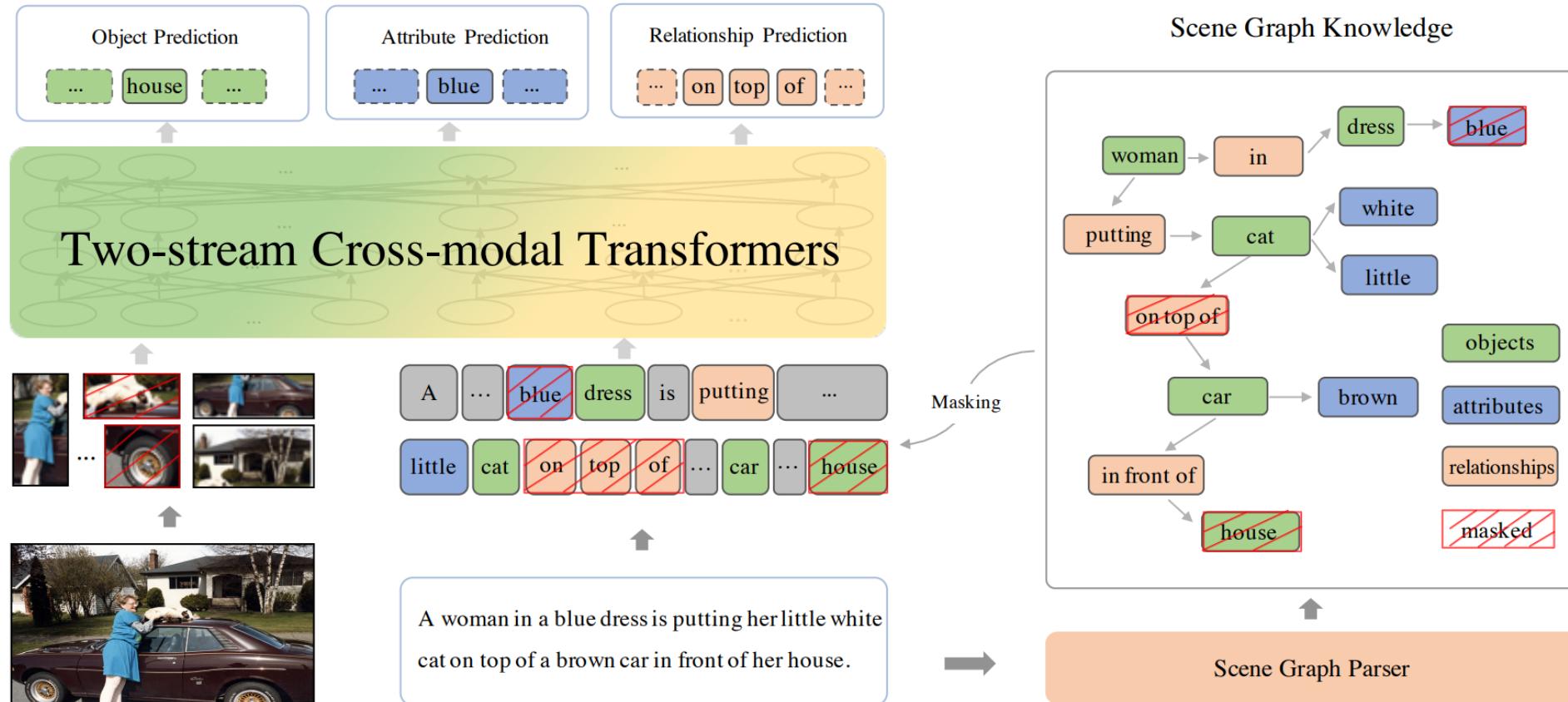
Figure 4: 2D visualization using t-SNE. The points from the same object class share the same color. Oscar (left) improves the cross-domain alignment over the baseline without object tags (right). Red and grey curves cover the objects of the same and related semantics, respectively.

Table 4: Results with various pre-training schemes.

Pre-train	VQA dev	Text Retrieval			Image Retrieval			Image Captioning				
		R@1	R@5	R@10	R@1	R@5	R@10	B@4	M	C	S	
visual genome open images	BASELINE (NO TAGS)	70.93	84.4	98.1	99.5	73.1	94.5	97.9	34.5	29.1	115.6	21.9
	OSCAR <sup>VG</sup>	71.70	88.4	99.1	99.8	75.7	95.2	98.3	36.4	30.3	123.4	23.0
	OSCAR <sup>OI</sup>	71.15	85.9	97.9	99.5	72.9	94.3	97.6	35.3	29.6	119.5	22.6

# Recent Advances in Vision and Language Pretrained Models

## ERNIE-ViL



# Recent Advances in Vision and Language Pretrained Models

## ERNIE-VIL

Models	VCR			RefCOCO+					
	Q→A	QA→R	Q→AR	val	testA	testB			
Out-of-domain	ViLBERT-base	72.42 (73.3)	74.47 (74.6)	54.04 (54.8)	72.34	78.52	62.61		
	Unicoder-VL-base	72.6 (73.4)	74.5 (74.4)	54.4 (54.9)	-	-	-		
	VLBERT-base	73.8 (-)	74.4 (-)	55.2 (-)	71.60	77.72	60.99		
	UNITER-base	-	-	-	72.78	-	-		
	VLBERT-large	75.5 (75.8)	77.9 (78.4)	58.9 (59.7)	72.59	78.57	62.30		
	ERNIE-ViL-base	74.37 (77.0)	79.65 (80.3)	61.24 (62.1)	74.02	80.33	<b>64.74</b>		
	ERNIE-ViL-large	<b>78.52(79.2)</b>	<b>83.37(83.5)</b>	<b>65.81(66.3)</b>	<b>74.24</b>	<b>80.97</b>	64.70		
Out-of-domain + in-domain	UNITER-base	74.56 (75.0)	77.03 (77.2)	57.76 (58.2)	75.31	81.30	65.58		
	VILLA-base	75.54 (76.4)	78.78 (79.1)	59.75 (60.6)	76.05	81.65	65.70		
	UNITER-large	77.22 (77.3)	80.49 (80.8)	62.59 (62.8)	75.90	81.45	66.70		
	VILLA-large	78.45 (78.9)	82.57 (82.8)	65.18 (65.7)	<b>76.17</b>	81.54	66.84		
	ERNIE-ViL-large	<b>78.62</b> (-)	<b>83.42</b> (-)	<b>65.95</b> (-)	75.95	<b>82.07</b>	<b>66.88</b>		
Models	VQA		IR-Flickr30K			TR-Flickr30K			
	test-dev	test-std	R@1	R@5	R@10	R@1	R@5		
Out-of-domain	ViLBERT-base	70.55	70.92	58.20	84.90	91.52	-	-	-
	Unicoder-VL-base	-	-	71.50	90.90	94.90	86.20	96.30	99.00
	VLBERT-base	71.16	-	-	-	-	-	-	-
	UNITER-base	71.56	-	-	-	-	-	-	-
	VLBERT-large	71.79	72.22	-	-	-	-	-	-
	ERNIE-ViL-base	72.62	72.85	74.44	92.72	95.94	86.70	97.80	99.00
	ERNIE-ViL-large	<b>73.78</b>	<b>73.96</b>	<b>75.10</b>	<b>93.42</b>	<b>96.26</b>	<b>88.70</b>	<b>97.30</b>	<b>99.10</b>
Out-of-domain + in-domain	UNITER-base	72.70	72.91	72.52	92.36	96.08	85.90	97.10	98.80
	OSCAR-base	73.16	73.61	-	-	-	-	-	-
	VILLA-base	73.59	73.67	74.74	92.86	95.82	86.60	97.90	<b>99.20</b>
	12-in-1-base	73.15	-	67.90	-	-	-	-	-
	UNITER-large	73.82	74.02	75.56	94.08	96.76	87.30	<b>98.00</b>	<b>99.20</b>
	OSCAR-large	73.44	73.82	-	-	-	-	-	-
	VILLA-large	74.69	74.87	76.26	<b>94.24</b>	<b>96.84</b>	87.90	97.50	98.80
	ERNIE-ViL-large	<b>74.75</b>	<b>74.93</b>	<b>76.70</b>	93.58	96.44	<b>88.10</b>	<b>98.00</b>	<b>99.20</b>

# Visual Dialog

Task	Num of sentences
Dialogue	up to 10 rounds sequences
Caption	
Image-Text Retrieval	1 sentence
VQA	2 sentences (question-answer pair)

## Captioning

Two people are in a wheelchair and one is holding a racket.



## Visual Dialog

- Q: How many people are on wheelchairs ?  
A: Two  
Q: What are their genders ?  
A: One male and one female  
Q: Which one is holding a racket ?  
A: The woman



## Visual Dialog

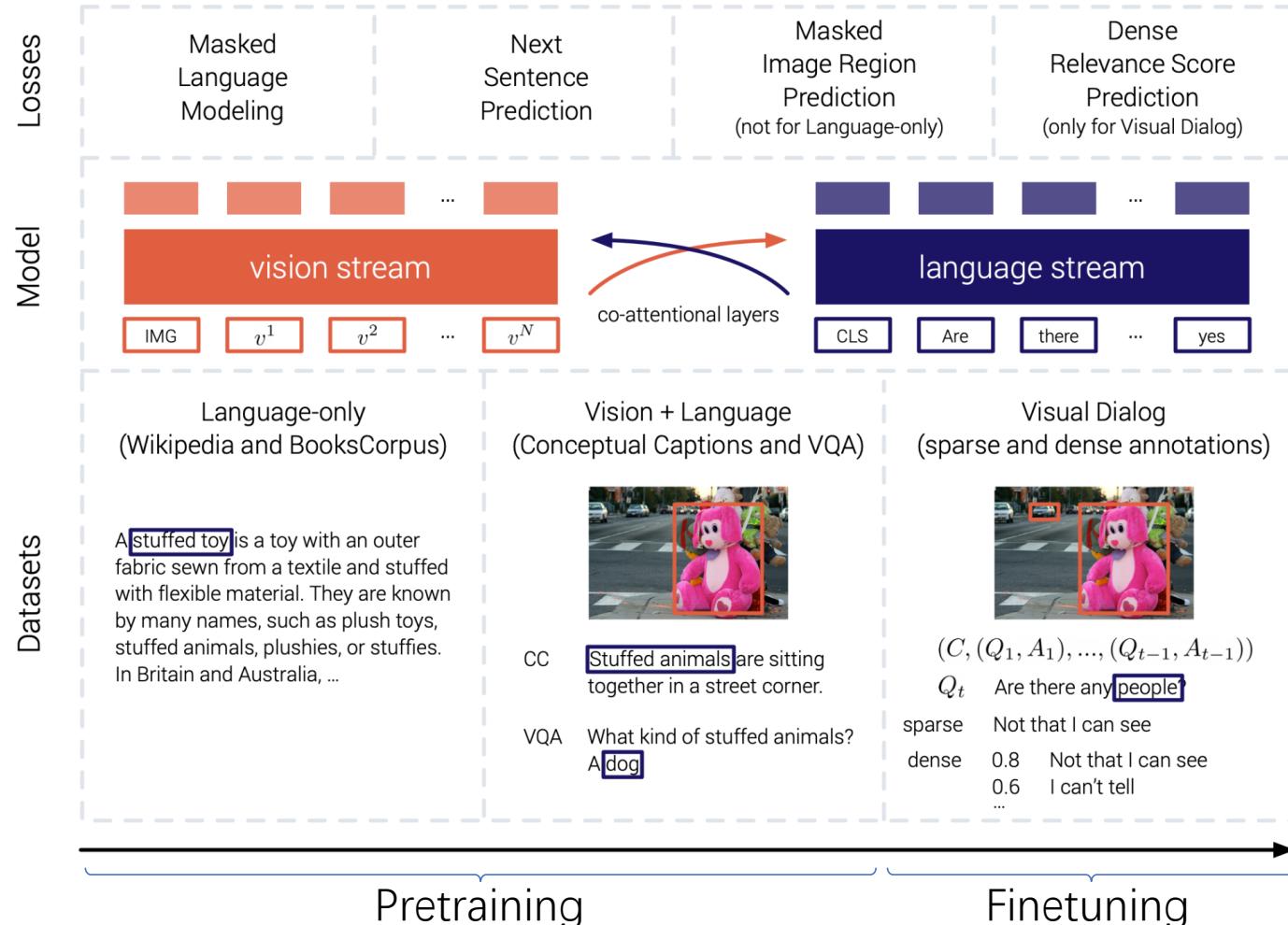
- Q: What is the gender of the one in the white shirt ?  
A: She is a woman  
Q: What is she doing ?  
A: Playing a Wii game  
Q: Is that a man to her right  
A: No, it's a woman

## Candidate answers:

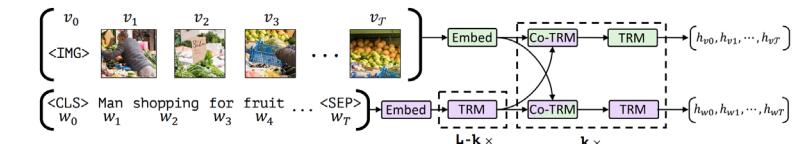
$$A_t = \{A_t^{(1)}, \dots, A_t^{(100)}\}$$

# Recent Advances in Vision and Language Pretrained Models

## VisDial-BERT



## ViLBERT



# Recent Advances in Vision and Language Pretrained Models

## VisDial-BERT

	Model	NDCG ↑	MRR ↑	R@1 ↑	R@5 ↑	R@10 ↑	MR ↓
BERT	NSP	56.17	63.37	49.17	80.62	89.42	4.21
	NSP + MLM	57.26	64.40	50.30	81.60	90.43	4.01
ViLBERT	w/ L-only	62.64	67.86	54.54	84.34	92.36	3.44
	w/ CC [3]	60.80	67.13	53.59	84.39	92.49	3.44
+dense	w/ CC [3]+VQA [4]	64.94	<b>69.10</b>	<b>55.88</b>	<b>85.50</b>	<b>93.29</b>	<b>3.25</b>
	CE	<b>75.24</b>	52.22	39.92	65.05	80.63	6.17
	CE + NSP	69.24	65.88	53.41	80.92	90.18	4.24

Table 3: Results on VisDial v1.0 val. ↑ indicates higher is better

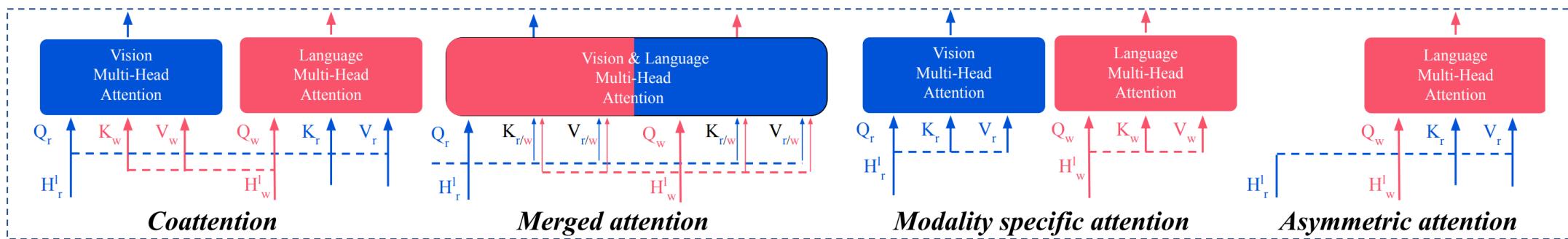
Model	NDCG ↑	MRR ↑	R@1 ↑	R@5 ↑	R@10 ↑	MR ↓
Published Results	GNN [15]	52.82	61.37	47.33	77.98	87.83
	CorefNMN [12]	54.70	61.50	47.55	78.10	88.80
	RvA [14]	55.59	63.03	49.03	80.40	89.83
	HACAN [22]	57.17	64.22	50.88	80.63	89.45
	NMN [12]	58.10	58.80	44.15	76.88	86.88
	DAN [17]	57.59	63.20	49.63	79.75	89.35
	DAN <sup>†</sup> [17]	59.36	64.92	51.28	81.60	90.88
	ReDAN [18]	61.86	53.13	41.38	66.07	74.50
	ReDAN <sup>†</sup> [18]	64.47	53.74	42.45	64.68	75.68
	DualVD [25]	56.32	63.23	49.25	80.23	89.70
	FGA [16]	56.93	66.22	52.75	82.92	91.08
	DL-61 [23]	57.32	62.20	47.90	80.43	89.95
	DL-61 <sup>†</sup> [23]	57.88	63.42	49.30	80.77	90.68
	MReal - BDAI* [24]	74.02	52.62	40.03	68.85	79.15
Leaderboard Entries	LF	45.31	55.42	40.95	72.45	82.83
	HRE	45.46	54.16	39.93	70.45	81.50
	MN	47.50	55.49	40.98	72.30	83.30
	MN-Att	49.58	56.90	42.43	74.00	84.35
	LF-Att	51.63	60.41	46.18	77.80	87.30
	MS ConvAI	55.35	63.27	49.53	80.40	89.60
	USTC-YTH	56.47	61.44	47.65	78.13	87.88
	UET-VNU	57.40	59.50	45.50	76.33	85.82
	square	60.16	61.26	47.15	78.73	88.48
	MS D365 AI	64.47	53.73	42.45	64.68	75.68
Ours	w/ CC [3]+VQA [4]	63.87	<b>67.50</b>	<b>53.85</b>	<b>84.68</b>	<b>93.25</b>
	CE	<b>74.47</b>	50.74	37.95	64.13	80.00
	CE + NSP	68.08	63.92	50.78	79.53	89.60

Table 4: Summary of results on VisDial v1.0 test-std. ↑ indicates higher is better. ↓ indicates lower is better. † denotes ensembles  
\* denotes the winning team of the 2019 Visual Dialog Challenge.

## Summary and Thinking

# Summary

- Dataset and data processing
  - Automatically collecting datasets
  - Processing images into a sequence involves defining “visual tokens”
- Multimodal Attention



- Loss Functions
    - masked region modeling loss
    - masked sub-word modeling loss
    - image-text matching loss
- minimizes the KL-divergence  
minimizes the cross-entropy loss
- $$\begin{aligned} & -y \log(\sigma(s_\theta(r^m, w^m))) \\ & -(1-y) \log(1 - \sigma(s_\theta(r^m, w^m))) \end{aligned}$$
 binary classification loss
- $$-\log \left( \frac{e^{s_\theta(r^m, w^m)}}{e^{s_\theta(r^m, w^m)} + \sum_{(\hat{r}, \hat{w}) \sim \mathcal{N}} e^{s_\theta(\hat{r}^m, \hat{w}^m)}} \right)$$
 contrastive loss

## Summary and Thinking

# Thinking and further work

- Exploring cognitive understanding level task
- Knowledge enhanced multimodal pretrain



**Title: Why am I so tired?**  
**Generated caption: A woman sitting on top of a bed in a room.**

intent	without event	get to her destination quickly
	with event	get to work on time
before	ground truth	save time by eating lunch while walking
	without event	walk out of the building
	with event	pick up the sandwich
	ground truth	purchase the sandwich at a deli
		take the subway to her current location
		be late for work
		be eating a sandwich while walking to work
after	without event	walk into a building
	with event	eat the sandwich
	ground truth	cross the street at a crosswalk
		throw away the trash from the sandwich
		look down at the sandwich thinking it's terrible
		finish eating the sandwich anyway