

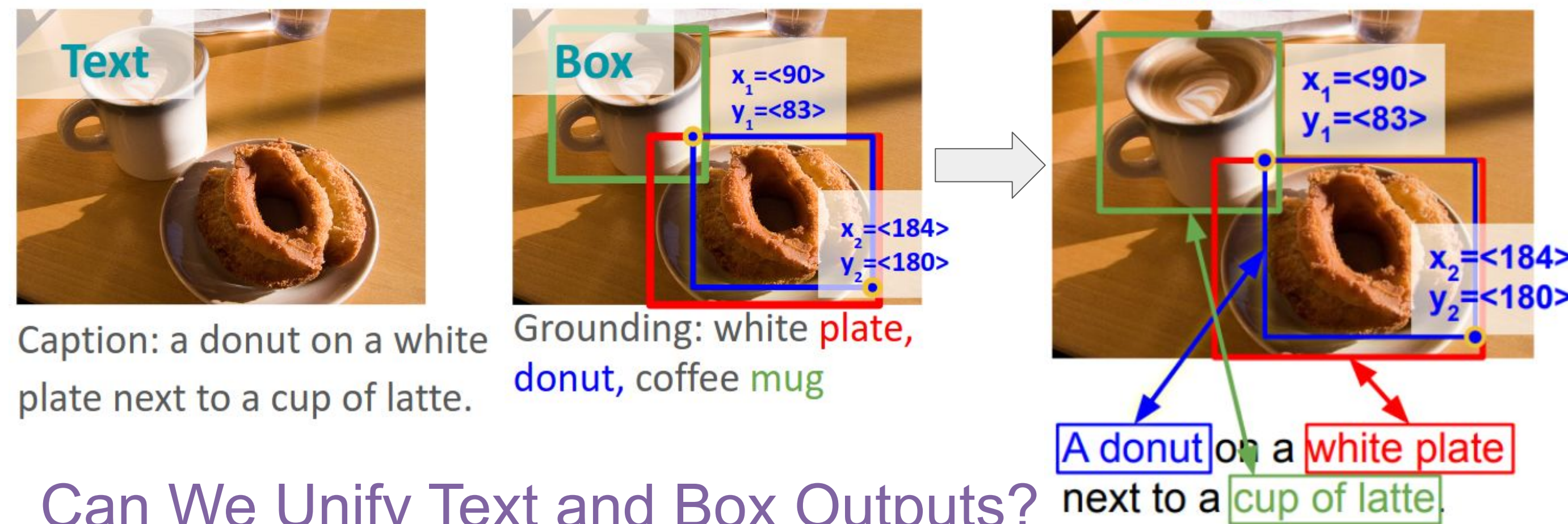
# UniTAB: Unifying Text and Box Outputs for Grounded Vision-Language Modeling

Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, Lijuan Wang



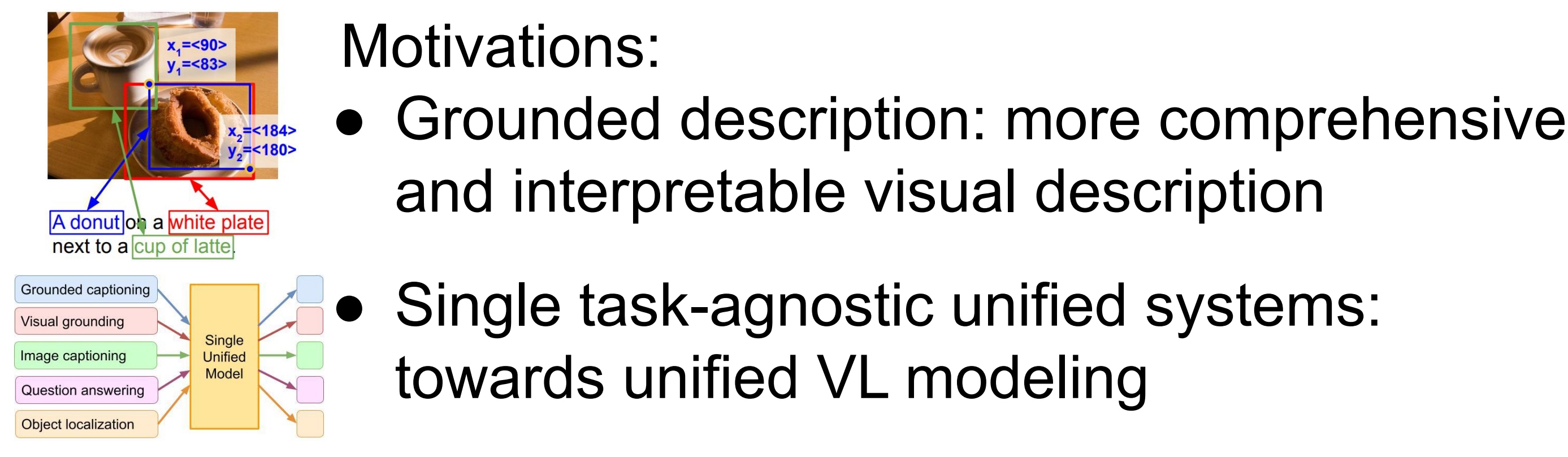
[github.com/microsoft/UniTAB](https://github.com/microsoft/UniTAB)

## Unifying Text and Box Outputs



### Can We Unify Text and Box Outputs?

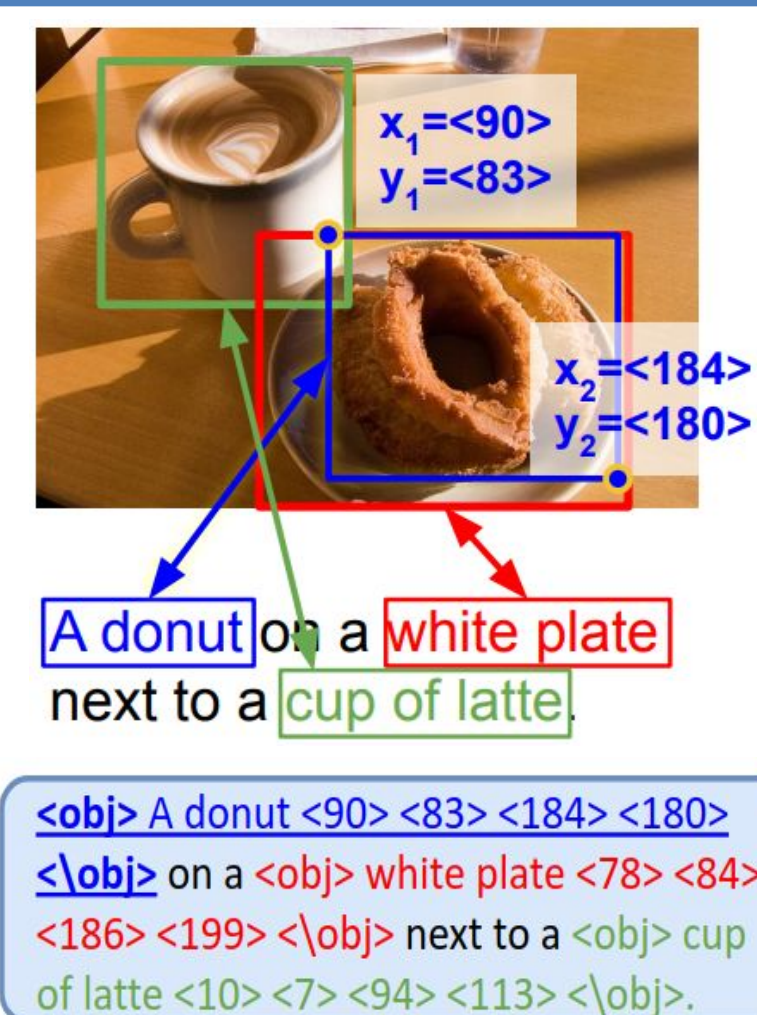
- Supporting both text and box outputs
- Representing word-box alignments



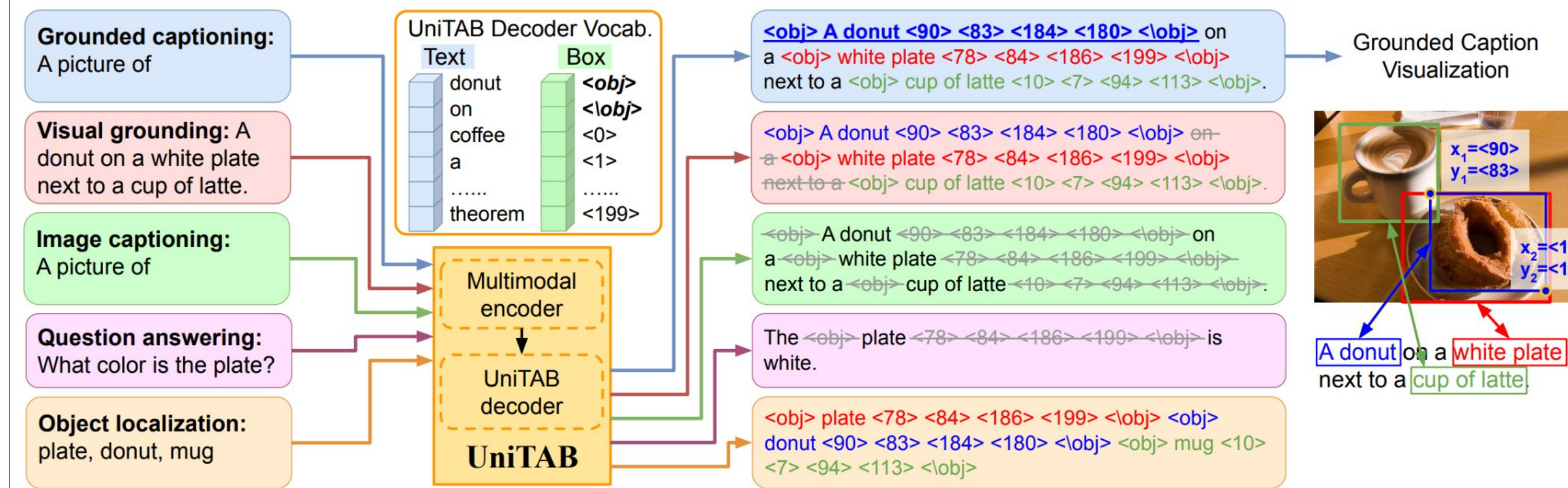
## Key Takeaways

### UniTAB: unifying text and box outputs for grounded VL modeling

- Grounded description ability
- Unified modeling for VL tasks
- Parameter efficient and generalizable

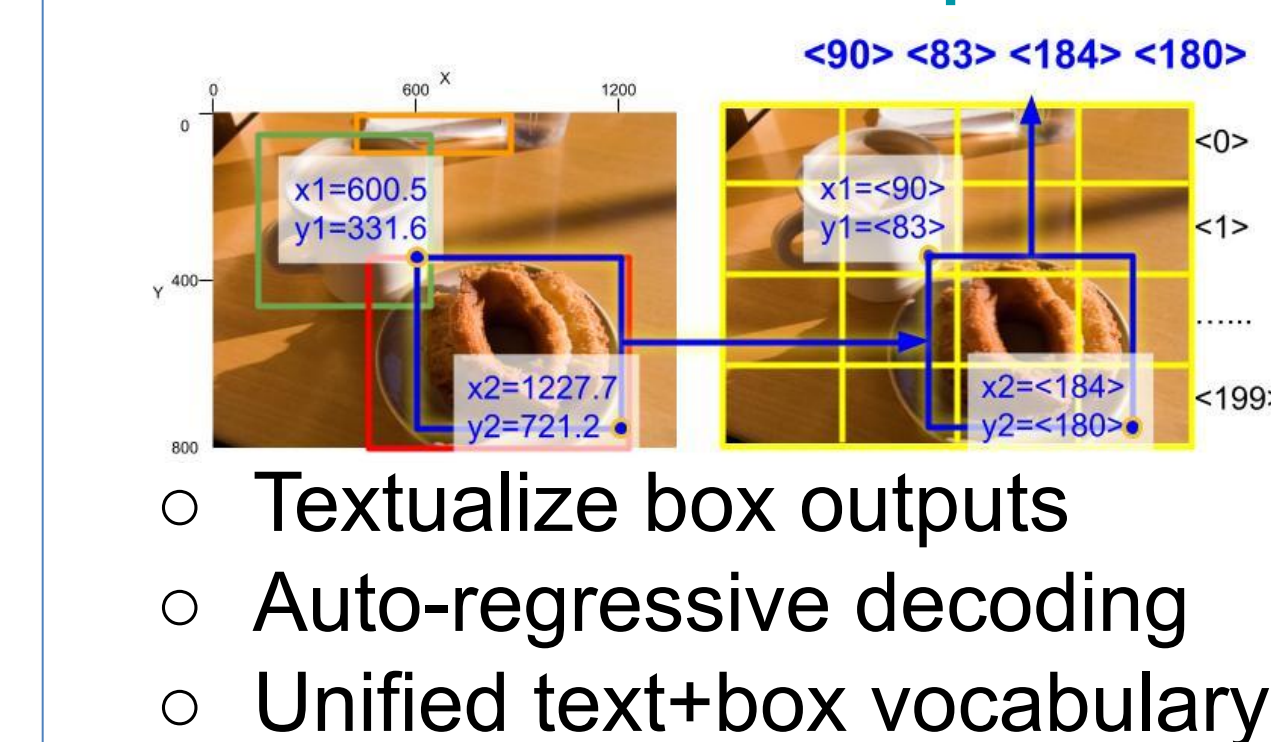


## How to Build a Shared Model for Text and Box?

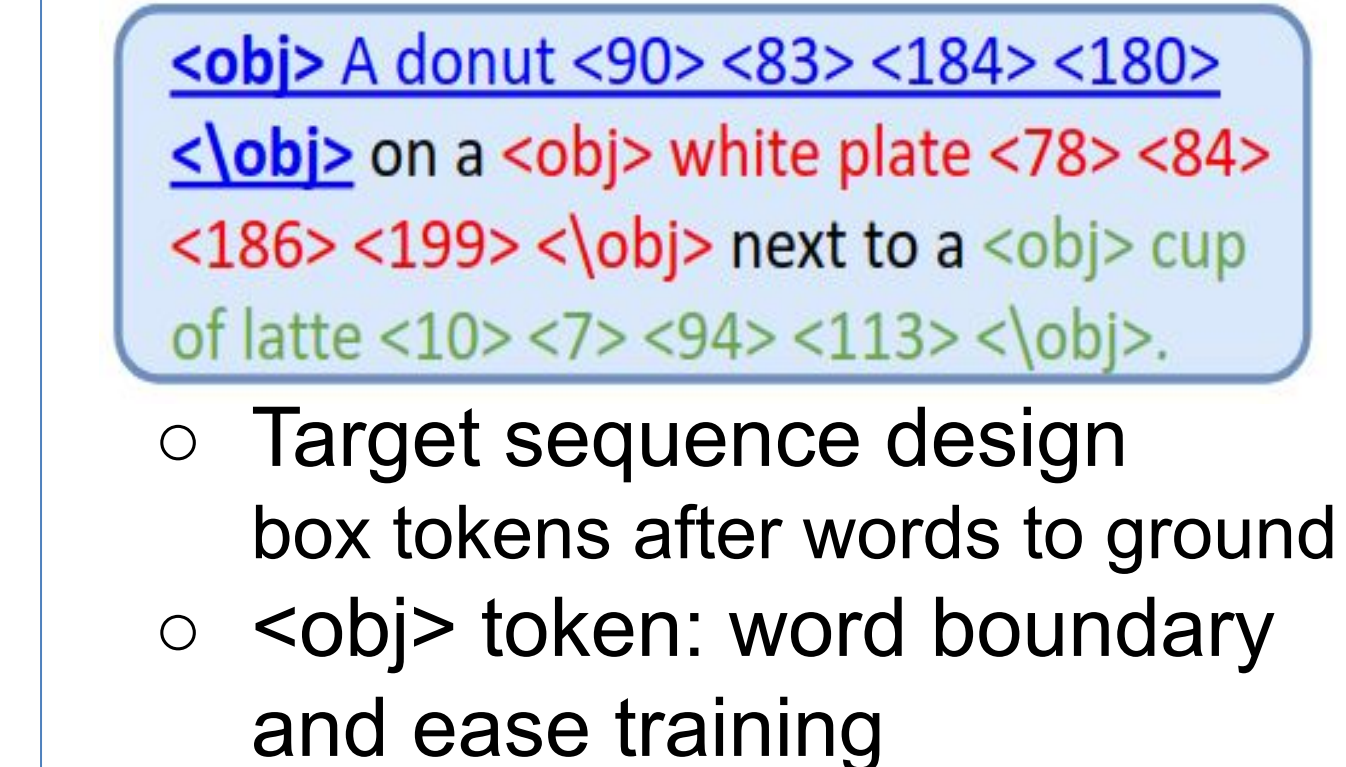


- Text and box outputs: unified decoding vocabulary
- Word-box alignments: <obj> token

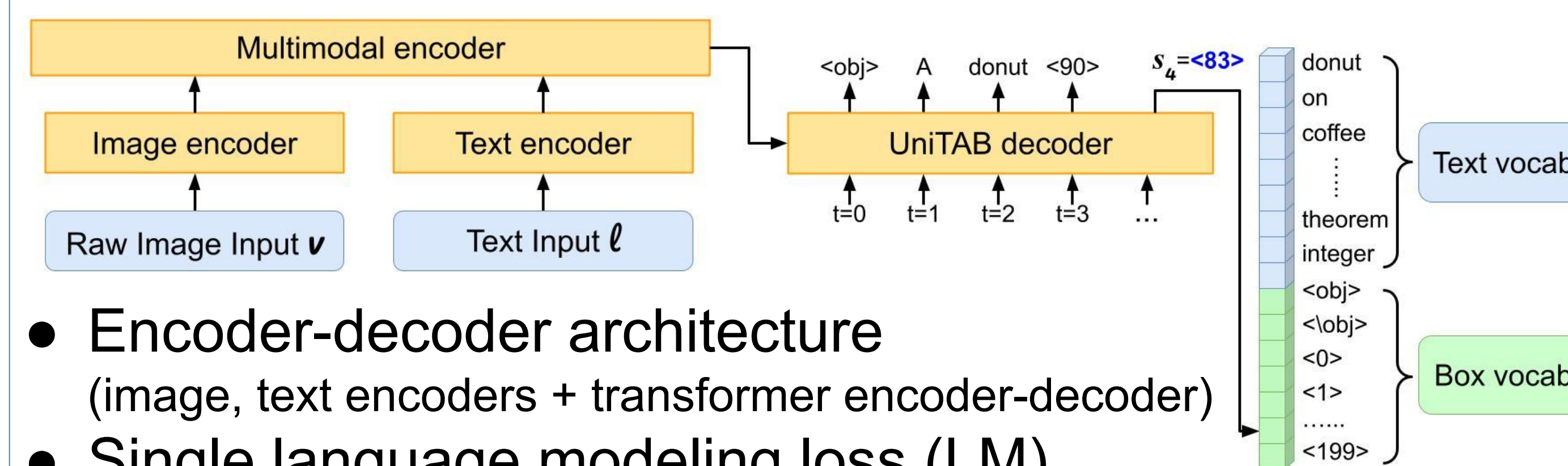
### Text and box outputs



### Word-box alignments



## UniTAB Framework and Training

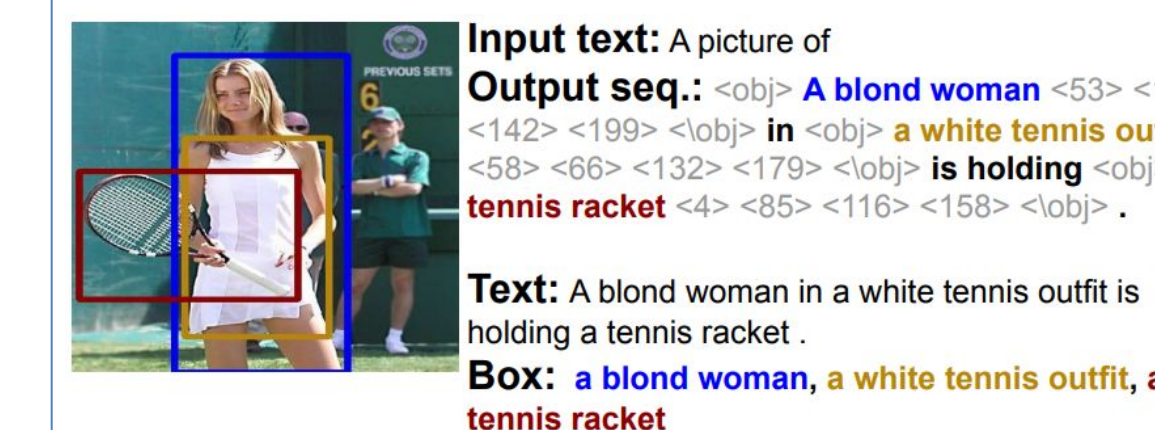


- Encoder-decoder architecture (image, text encoders + transformer encoder-decoder)
- Single language modeling loss (LM)

$$\mathcal{L}_{LM}(\theta) = -\sum_{t=1}^T \log P_{\theta}(s_t | s_{<t, v, l})$$

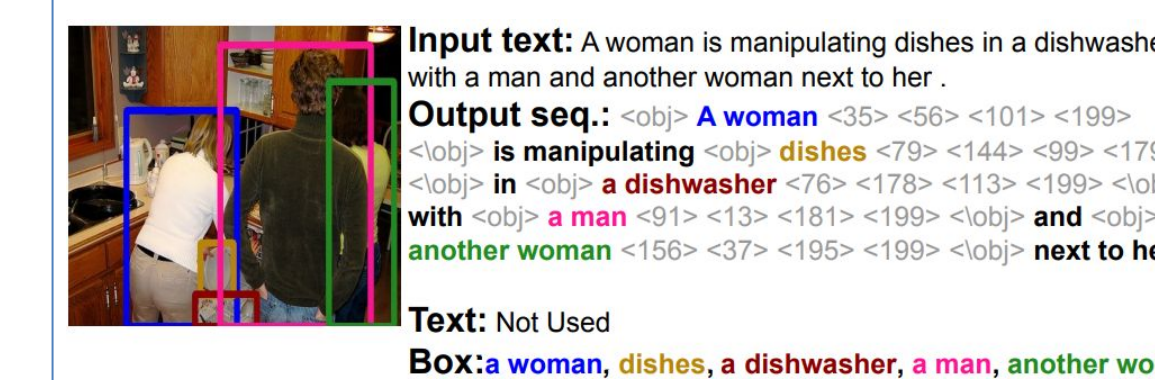
## Quantitative Results

- Text, box, alignment
- Grounded captioning
- Flickr30k Entities



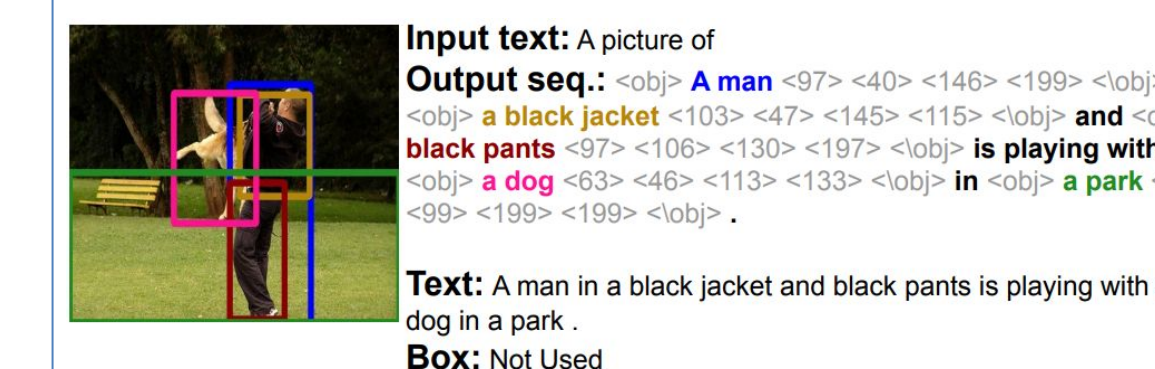
Method	Caption Eval.				Grounding Eval.	
	B@4	M	C	S	F1 <sub>all</sub>	F1 <sub>loc</sub>
NBT [49]	27.1	21.7	57.5	15.6	-	-
GVD [86]	27.3	22.5	62.3	16.5	7.55	22.2
Cyclical [50]	26.8	22.4	61.1	16.8	8.44	22.78
POS-SCAN [88]	30.1 <sup>†</sup>	22.6 <sup>†</sup>	69.3 <sup>†</sup>	16.8 <sup>†</sup>	7.17	17.49
Chen et al. [9]	27.2	22.5	62.5	16.5	7.91	21.54
UniTAB	<b>30.1</b>	<b>23.7</b>	<b>69.7</b>	<b>17.4</b>	<b>12.95</b>	<b>34.79</b>

- Box, alignment
- Visual grounding
- Refcoco/+g, Flickr30k



Method	Refcoco			Refcoco+			Refcocog	
	val	testA	testB	val	testA	testB	val-u	test-u
MAttNet [79]	76.40	80.43	69.28	64.93	70.26	56.00	66.67	67.01
FAOA [77]	72.05	74.81	67.59	55.72	60.37	48.54	59.03	58.70
TransVG [18]	81.02	82.72	78.35	64.82	70.70	56.94	68.67	67.73
ViLBERT [47]	-	-	-	72.34	78.53	62.61	-	-
UNITER [12]	81.41	87.04	74.17	75.90	81.45	66.70	74.02	68.67
VILLA [22]	82.39	87.48	74.84	76.17	81.54	66.84	76.18	76.71
MDETR [34]	86.75	89.58	81.41	79.52	84.09	70.62	81.64	80.89
UniTAB	<b>88.59</b>	<b>91.06</b>	<b>83.75</b>	<b>80.97</b>	<b>85.36</b>	<b>71.55</b>	<b>84.58</b>	<b>84.70</b>

- Text
- Image captioning, VQA
- MSCOCO, VQAv2



Method	#Pre-train	B@4			
		M	C	S	
Unified VLP [87]	3M	36.5	28.4	117.7	21.3
OSCAR [43]	4M	36.5	30.3	123.7	23.1
E2E-VLP [75]	180K	36.2	-	117.3	-
VL-T5 [13]	180K	34.5	28.7	116.5	21.9
VL-BART [13]	180K	35.1	28.7	116.6	21.5
UniTAB	200K	36.1	28.6	119.8	21.7

## Multi-task finetuning: UniTAB<sub>Shared</sub>

- A single set of parameters for all experimented VL tasks
- (1) Parameter efficient, (2) generalizing learned abilities



- Grounded description
- MSCOCO



- Object localization
- ImageNet