

# GRIT: Faster and Better Image captioning Transformer Using Dual Visual Features

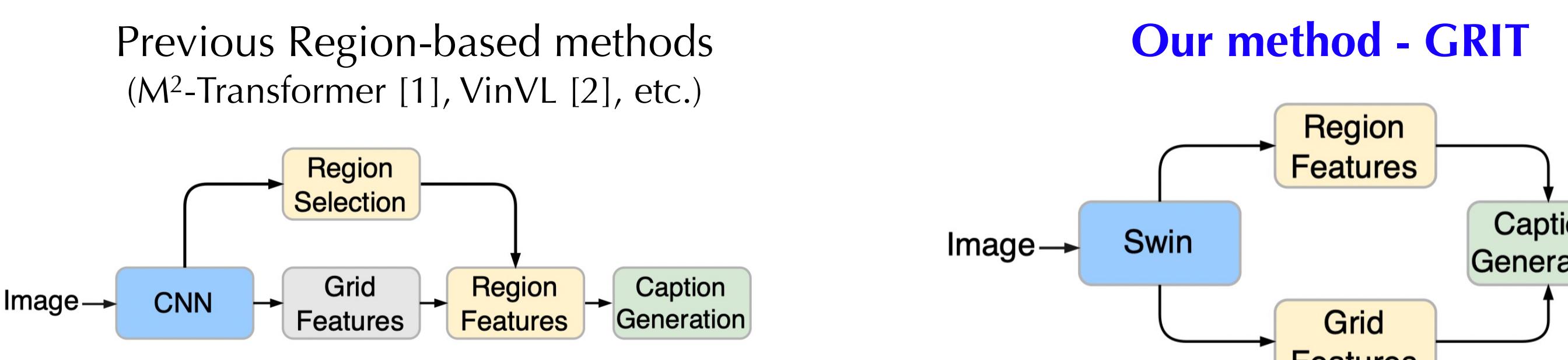
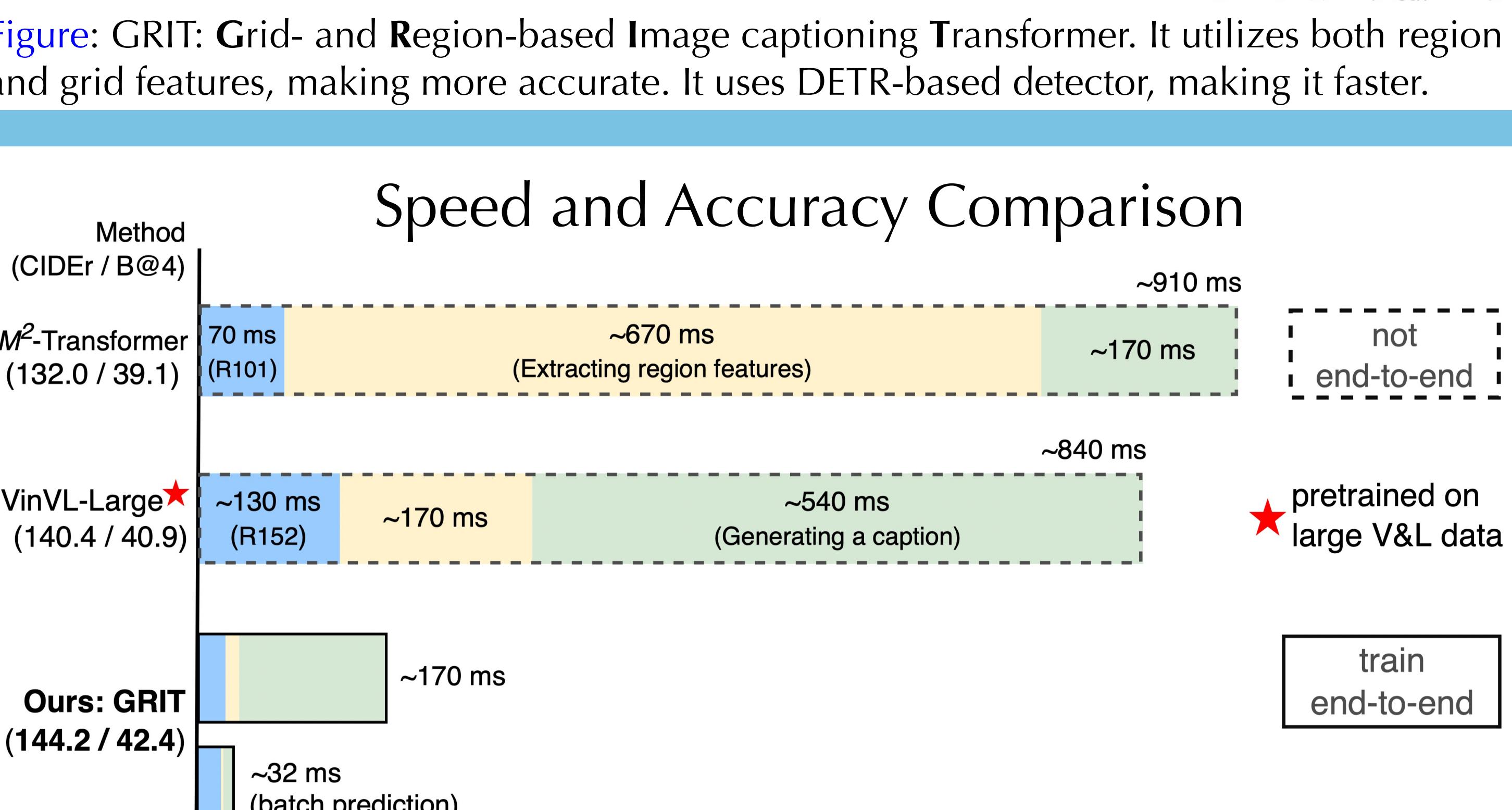
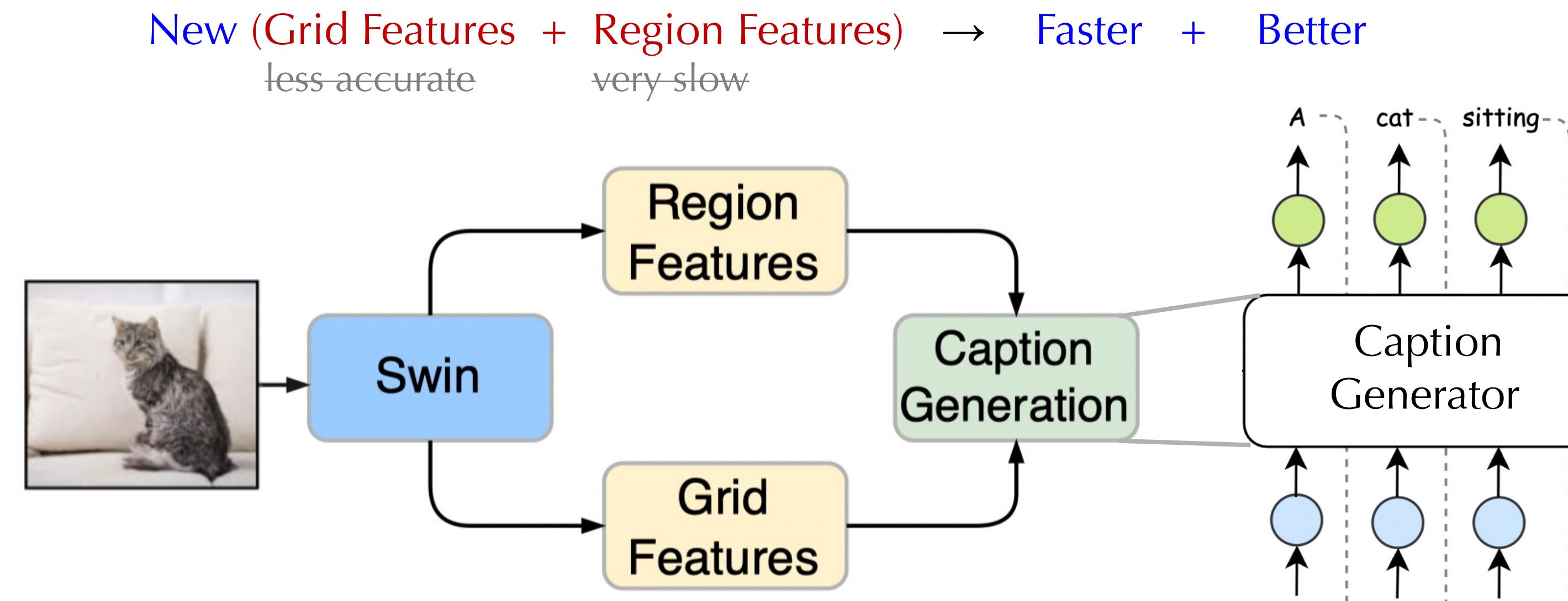
Van-Quang Nguyen<sup>1</sup>, Masanori Suganuma<sup>1,2</sup>, Takayuki Okatani<sup>1,2</sup>
<sup>1</sup>GSIS, Tohoku University, <sup>2</sup>RIKEN Center for AIP

## Abstract

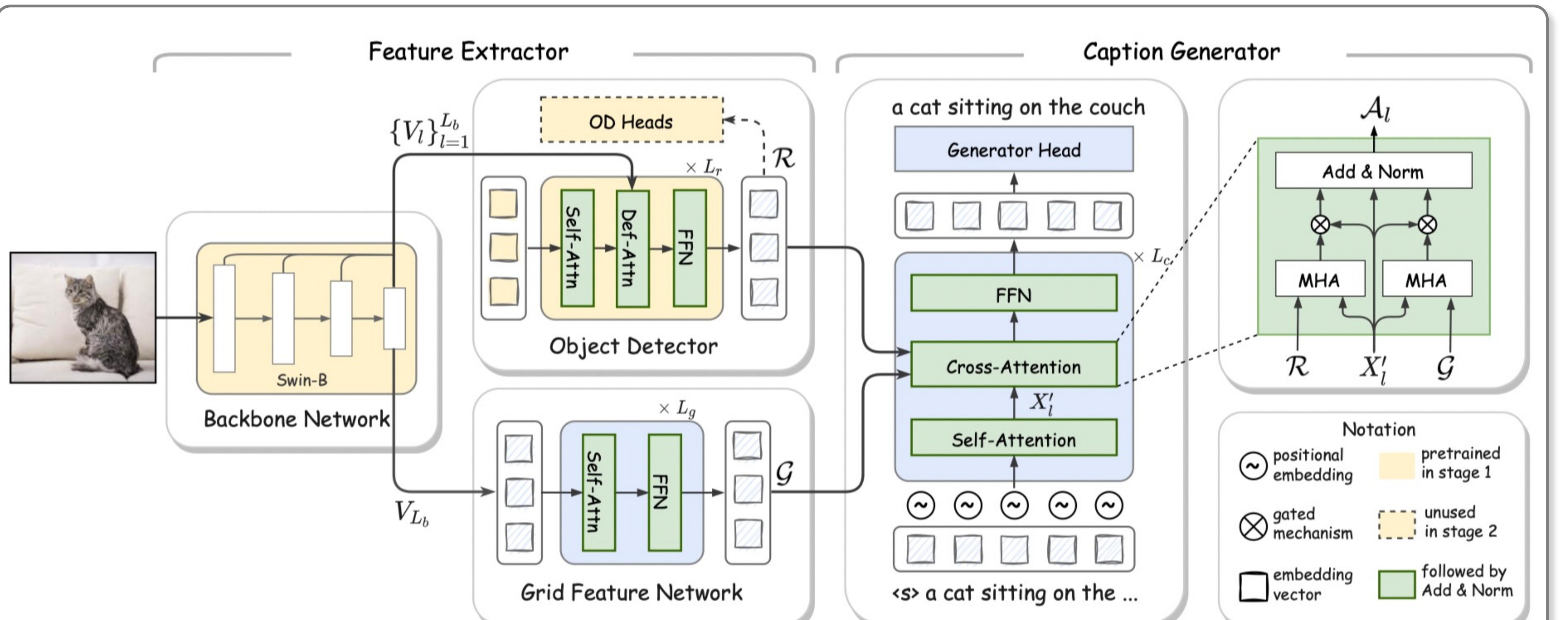
We propose Transformer-only neural architecture named **GRIT** that effectively utilizes dual visual features to generate better captions:

- GRIT replaces the CNN-based detector employed in previous methods with a DETR-based one, making it computationally faster.
- Its monolithic design consisting only of Transformers enables end-to-end training of the model.

**Question: How do we extract and fuse good visual representations?**

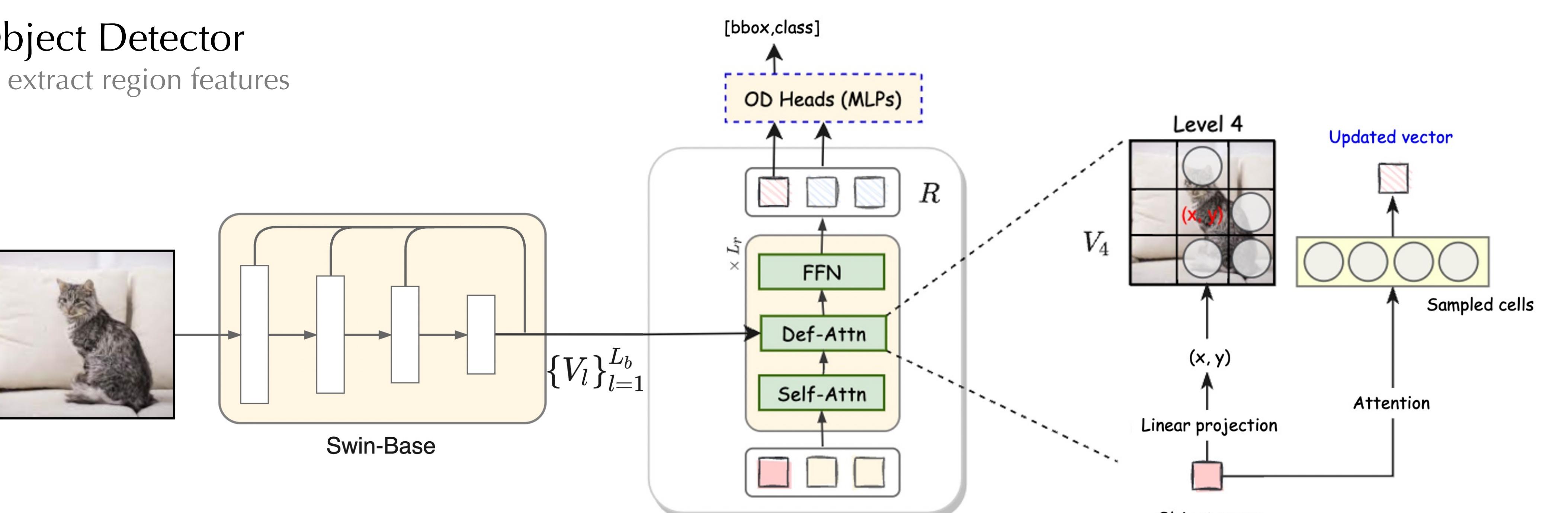


## Our Method



### Object Detector

To extract region features

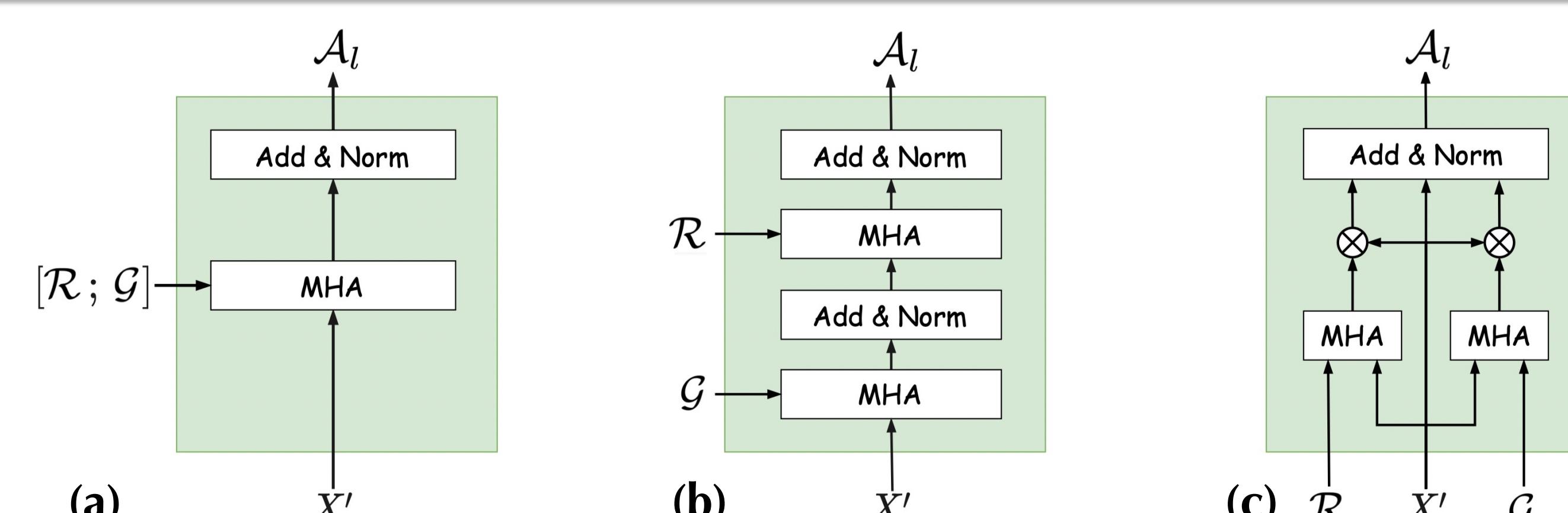


### Cross-Attention Mechanism

The mechanism to use the dual visual features

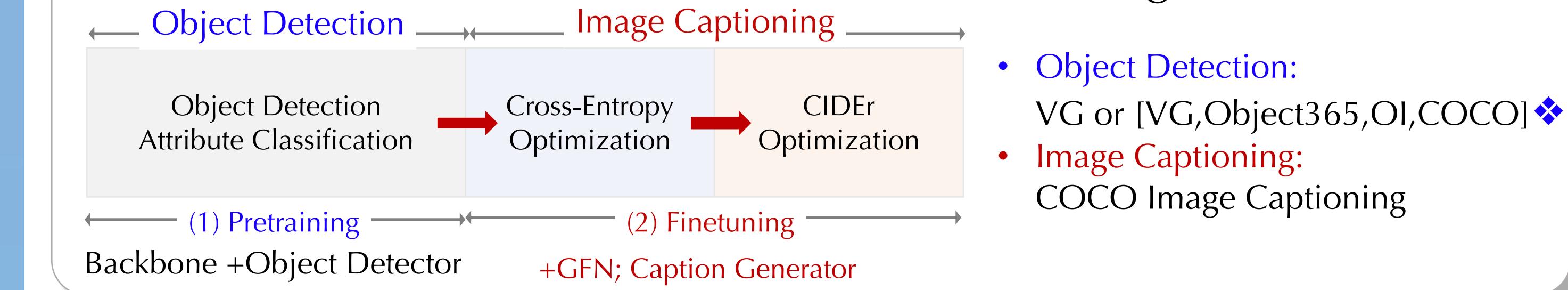
- Concatenated Cross-Attention
- Sequential Cross-Attention
- Parallel Cross-Attention

→ Parallel Cross-Attention yields the best!



## Experimental Results

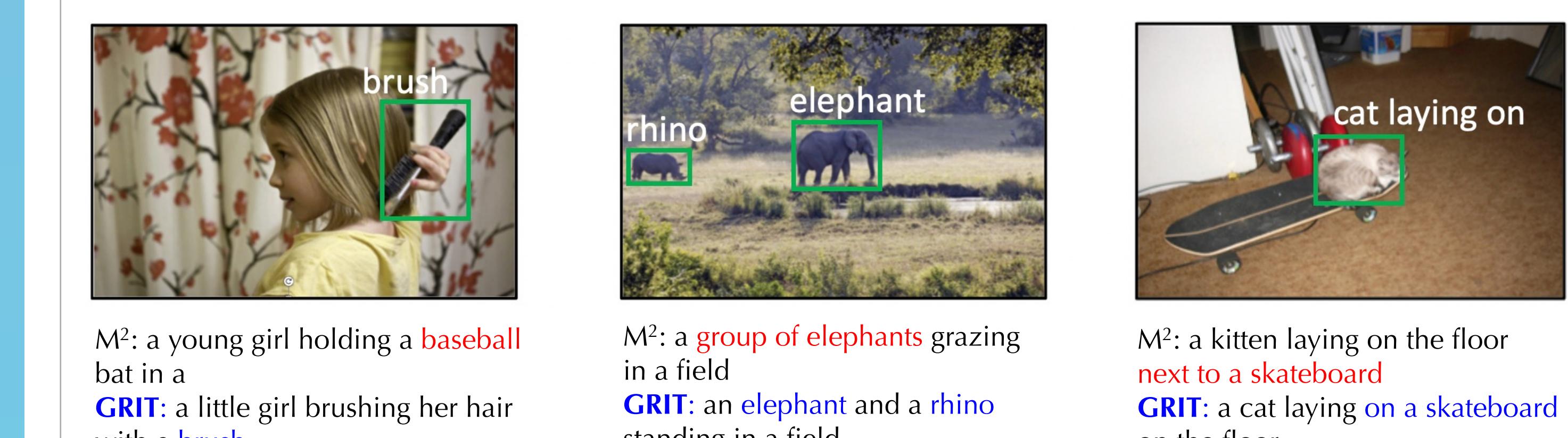
### End-to-End Training



### Results on the Karpathy test split of COCO

| Method  | V. E.<br>Type | # VL<br>Data | Performance Metrics |             |             |             |              |
|---|---------------|--------------|---------------------|-------------|-------------|-------------|--------------|
|   |               |              | B@1                 | B@4         | M           | R           | C            |
| w/ VL pretraining   |               |              |                     |             |             |             |              |
| UVLP  | R             | 3.0M         | -                   | 39.5        | 29.3        | -           | 129.3        |
| Oscar <sub>base</sub>   | R             | 6.5M         | -                   | 40.5        | 29.7        | -           | 137.6        |
| VinVL <sup>†</sup> <sub>large</sub>   | R             | 8.9M         | -                   | <b>41.0</b> | 31.1        | -           | 140.9        |
| SimVLM <sub>huge</sub>  | G             | 1.8B         | -                   | 40.6        | <b>33.7</b> | -           | <b>143.3</b> |
| w/o VL pretraining  |               |              |                     |             |             |             |              |
| SAT   | G             | -            | -                   | 31.9        | 25.5        | 54.3        | 106.3        |
| RSTNet  | G             | -            | 81.8                | 40.1        | 29.8        | 59.5        | 135.6        |
| Up-Down   | R             | -            | 79.8                | 36.3        | 27.7        | 56.9        | 120.1        |
| AoA   | R             | -            | 80.2                | 38.9        | 29.2        | 58.8        | 129.8        |
| $\mathcal{M}^2$ Transformer   | R             | -            | 80.8                | 39.1        | 29.2        | 58.6        | 131.2        |
| TCIC  | R             | -            | 81.8                | 40.8        | 29.5        | 59.2        | 135.4        |
| Dual Global   | R+G           | -            | 81.3                | 40.3        | 29.2        | 59.4        | 132.4        |
| DLCT  | R+G           | -            | 81.4                | 39.8        | 29.5        | 59.1        | 133.8        |
| <b>GRIT</b>   | R+G           | -            | 83.5                | 41.9        | 30.5        | 60.5        | 142.2        |
| <b>GRIT</b> (  4 datasets) | R+G           | -            | <b>84.2</b>         | <b>42.4</b> | <b>30.6</b> | <b>60.7</b> | <b>144.2</b> |
| <b>GRIT</b>   | R+G           | -            | <b>84.2</b>         | <b>42.4</b> | <b>30.6</b> | <b>60.7</b> | <b>144.2</b> |

### Qualitative Examples



## References & Our Code

- [1] Cornia, M., Stefanini, M., Baraldi, L., Cucchiara, R.: Meshed-memory transformer for image captioning. In CVPR, 2020.
- [2] Zhang, P., Li, X., Hu, X., Yang, J., Zhang, L., Wang, L., Choi, Y., Gao, J.: Vinvl: Revisiting visual representations in vision-language models. In CVPR, 2021.

Code and pretrained models given at: <https://www.github.com/davidnvq/grit>
