

Image-Text Pairs: 6.5M

- ① Masked Token Loss ② Contrastive Loss

Image-Text Representation



Word-Tag-Region Triplet

Pre-training



Fine-tuning

Understanding

- VQA
- GQA
- NLVR2
- Image-Text Retrieval
- Text-Image Retrieval

Generation

- Image Captioning
- Novel Object Captioning