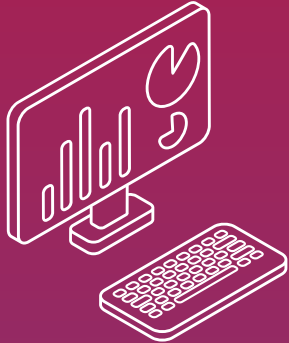# VL-BEIT

Parsa Sharifi

- vision-language foundation model
- bidirectional multimodal Transformer
- visual question answering, visual reasoning, and image-text retrieval

- Backbone model
- Pre-training

- Mixture Of-Modality-Experts (MOME) Transformer
- Multi-head self-attention layer
- A feed-forward expert layer(pool)
- Hard Routing Mechanism

- Pretraining Tasks
  - Masked Language Modeling
  - Masked Image Modeling
  - Masked Vision-Language Modeling
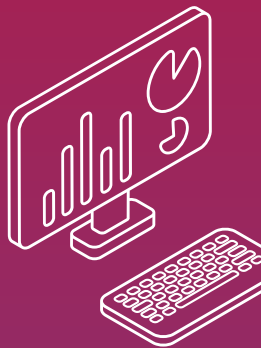
Attention is [M] we [M] ?
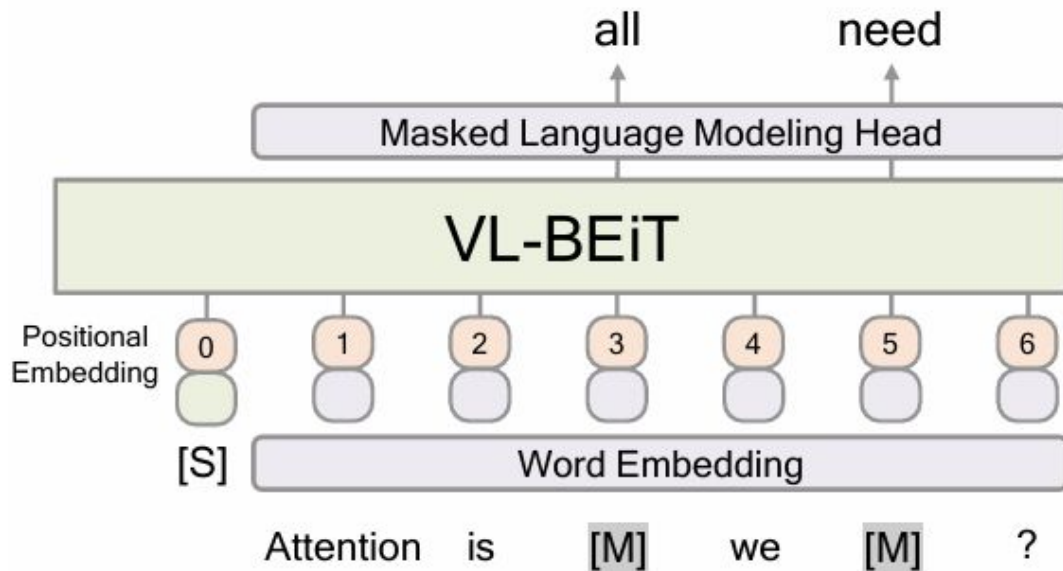
- Input Representations

$$H^v = \left[ v_{\texttt{[I\_CLS]}}, v_1, \ldots, v_N \right] + V_{pos}.$$

$$H^w = \left[ w_{\texttt{[T\_CLS]}}, w_1, \ldots, w_M, w_{\texttt{[T\_SEP]}} \right] + T_{pos}.$$
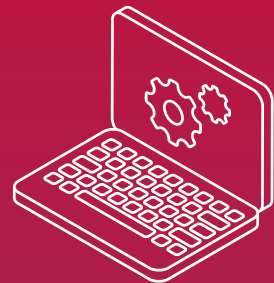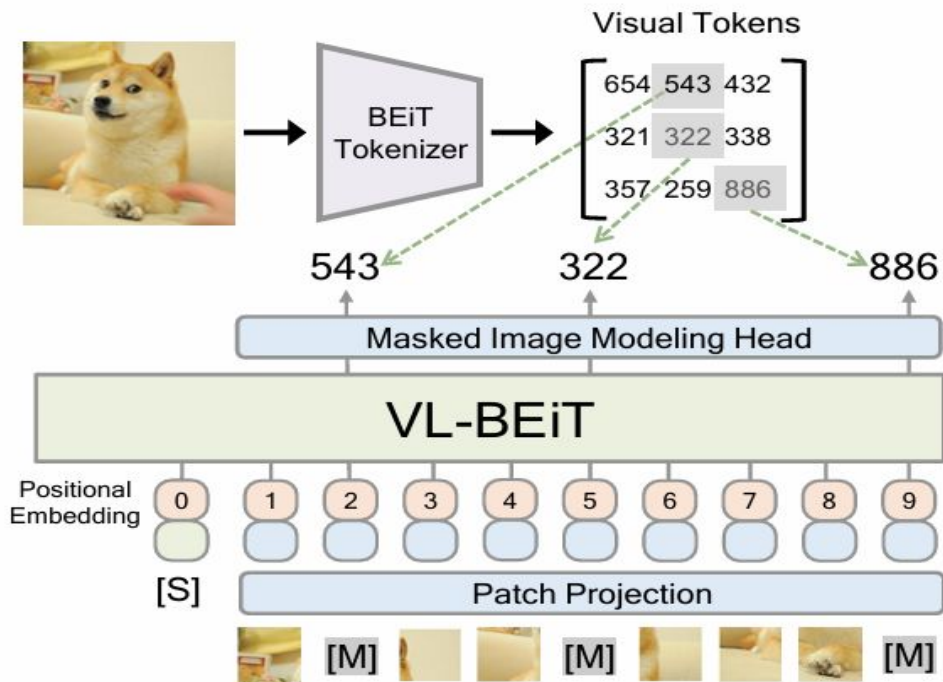
$$H^{vl} = \left[ H^w; H^v \right]$$

(a) Masked Language Modeling

(b) Masked Image Modeling

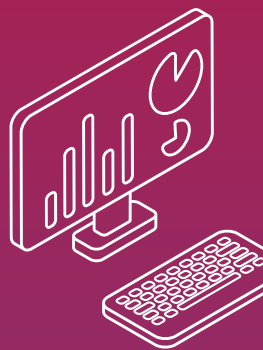(c) Masked Vison-Language Modeling

| Model | VQA | | NLVR2 | |
| --- | --- | --- | --- | --- |
| | test-dev | test-std | dev | test-P |
| *Base-size models pretrained on the same data* | | | | |
| UNITER | 72.70 | 72.91 | 77.18 | 77.85 |
| VILLA | 73.59 | 73.67 | 78.39 | 79.30 |
| UNIMO | 73.79 | 74.02 | - | - |
| ViLT | 71.26 | - | 75.70 | 76.13 |
| ALBEF | 74.54 | 74.70 | 80.24 | 80.50 |
| VLMo | 76.64 | 76.89 | **82.77** | **83.34** |
| VL-BEiT | **77.53** | **77.75** | 81.93 | 82.66 |

| Model | COCO | | Flickr30K | |
|---|---|---|---|---|
| | TR | IR | TR | IR |
| *Fusion encoder* | | | | |
| UNITER | 64.4 | 50.3 | 85.9 | 72.5 |
| VILLA | - | - | 86.6 | 74.7 |
| ViLT | 61.5 | 42.7 | 83.5 | 64.4 |
| *Dual encoder* | | | | |
| VLMo | 74.8 | 57.2 | 92.3 | 79.3 |
| *Dual encoder + Fusion encoder reranking* | | | | |
| ALBEF | 73.1 | 56.8 | 94.3 | 82.8 |
| VL-BEiT | **79.5** | **61.5** | **95.8** | **83.9** |