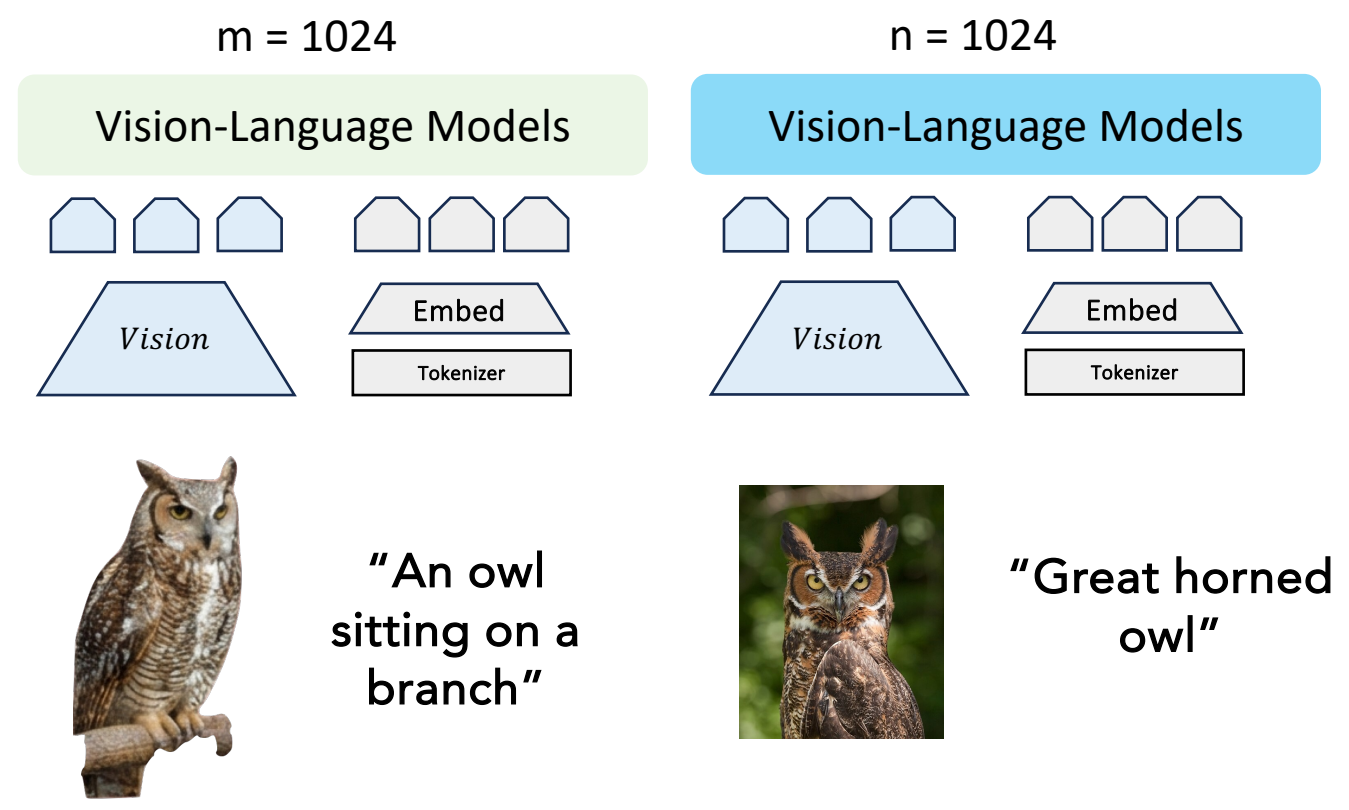
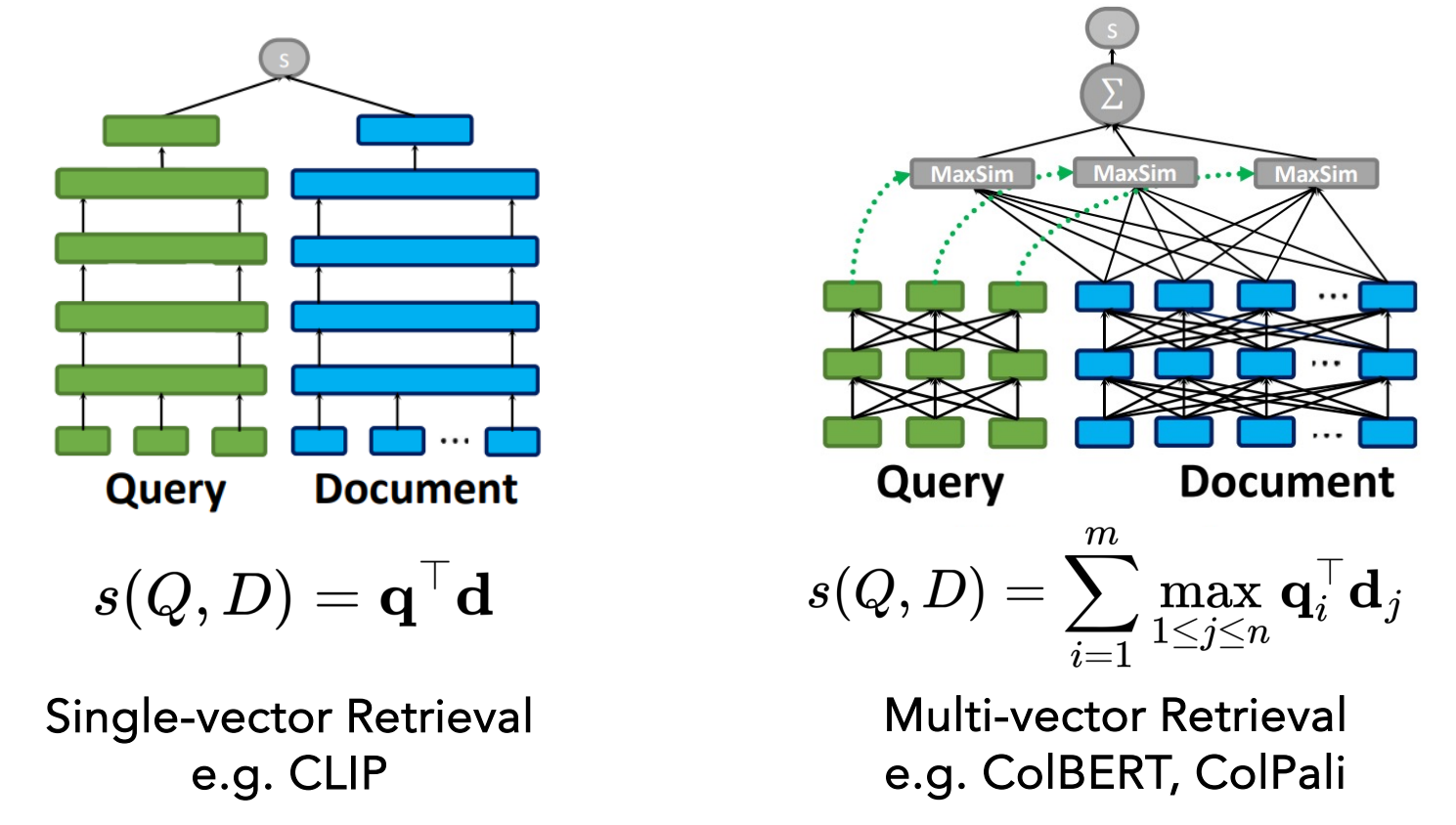


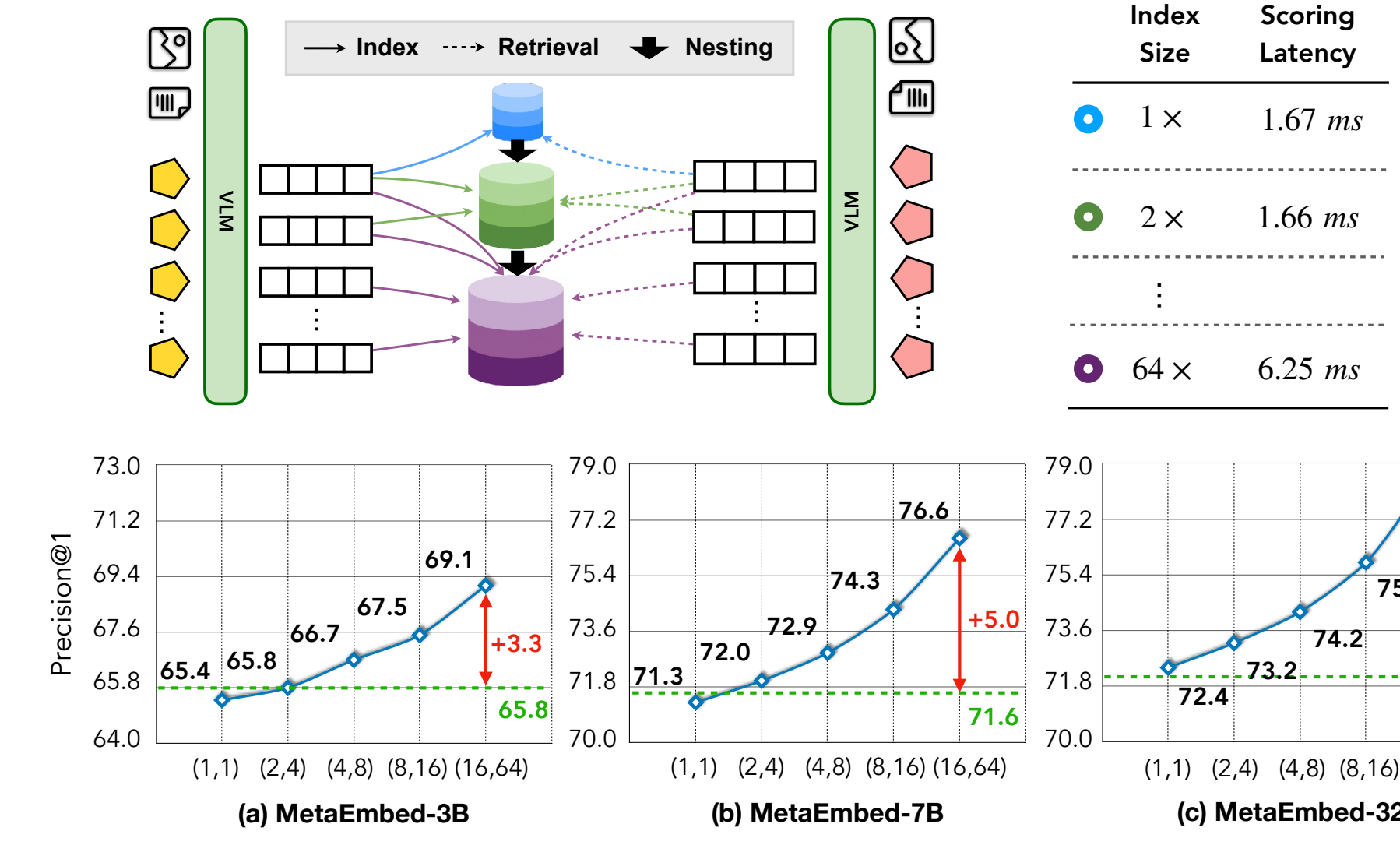
## The Problem: Multimodal Retrieval



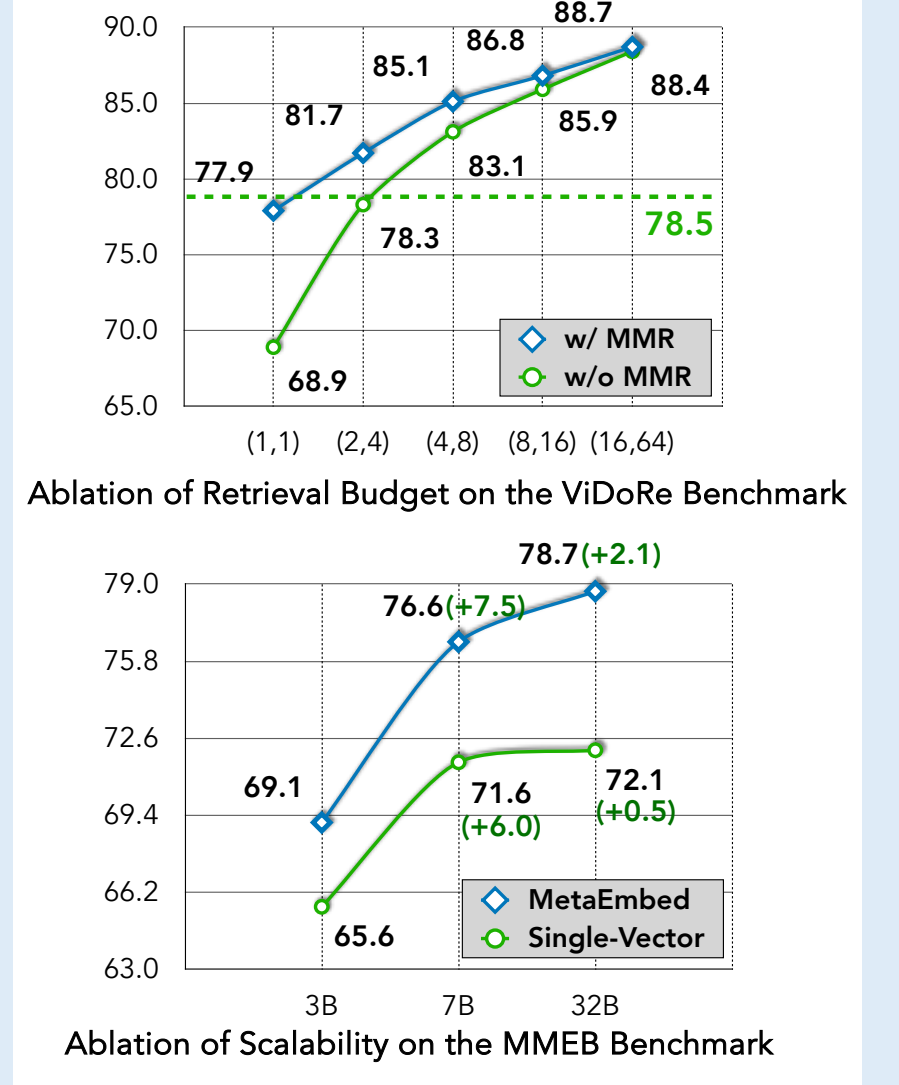
## Baseline Methods



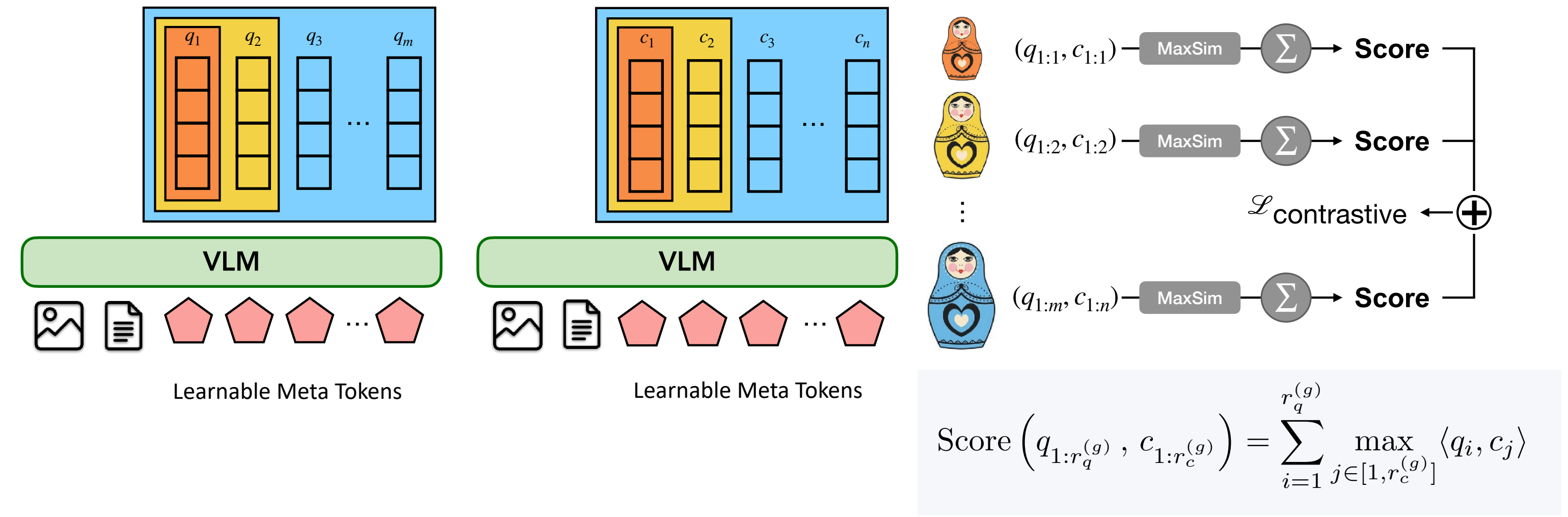
## Inference: Flexible Trade-offs: Latency vs Accuracy



## Ablations



## Key Idea: Matryoshka Multi-Vector Retrieval (MMR)



$$\mathcal{L}_{\text{contrastive}} = \mathcal{L}_{\text{NCE}}(\text{Score}(q_{11}, c_{11})) + \mathcal{L}_{\text{NCE}}(\text{Score}(q_{12}, c_{12})) + \dots + \mathcal{L}_{\text{NCE}}(\text{Score}(q_{1m}, c_{1n}))$$

- Contrastive losses on nested prefixes so the first few vectors form a coarse summary and additional vectors progressively refine it.

## Multimodal Retrieval Evaluation

- Models trained on MMEB and ViDoRe-Train.
- Best accuracy obtained when using 16 query vectors and 64 retrieval vectors.
- Competitive or superior accuracy even with a single vector.

Models	Size	Per Meta-Task Score				Average Score		
		Classification	VQA	Retrieval	Grounding	IND	OOD	Overall
Baseline Models								
CLIP	428M	55.2	19.7	53.2	62.2	47.6	42.8	45.4
MagicLens	613M	38.8	8.3	35.4	26.0	-	-	27.8
UniIR	428M	42.1	15.0	60.1	62.2	-	-	42.8
MM-EMBED	7B	48.1	32.2	63.8	57.8	-	-	50.0
GME	7B	56.9	41.2	67.8	53.4	-	-	55.8
mmE5	11B	67.6	62.7	71.0	89.7	72.4	66.6	69.8
MoCa-3B	3B	59.8	62.9	70.6	88.6	72.3	61.5	67.5
MoCa-7B	7B	65.8	64.7	75.0	<b>92.4</b>	74.7	67.6	71.5
B3-7B	7B	70.0	66.5	74.1	84.6	75.9	67.1	72.0
METAEMBED - PaliGemma Initialized								
METAEMBED-3B <sup>Gemma</sup>	3B	64.9	53.5	70.9	79.5	68.6	61.3	65.4
METAEMBED - Llama-3.2-Vision Initialized								
METAEMBED-11B	11B	66.4	42.1	74.3	<u>91.6</u>	65.7	64.3	65.1
METAEMBED - Qwen2.5-VL Initialized								
METAEMBED-3B	3B	62.7	68.1	71.9	79.6	73.5	63.8	69.1
METAEMBED-7B	7B	<b>71.3</b>	<b>74.2</b>	<b>78.7</b>	85.4	<b>81.8</b>	<b>70.0</b>	<b>76.6</b>
METAEMBED-32B	32B	<b>73.7</b>	<b>78.6</b>	<b>78.9</b>	88.1	<b>82.8</b>	<b>73.7</b>	<b>78.7</b>