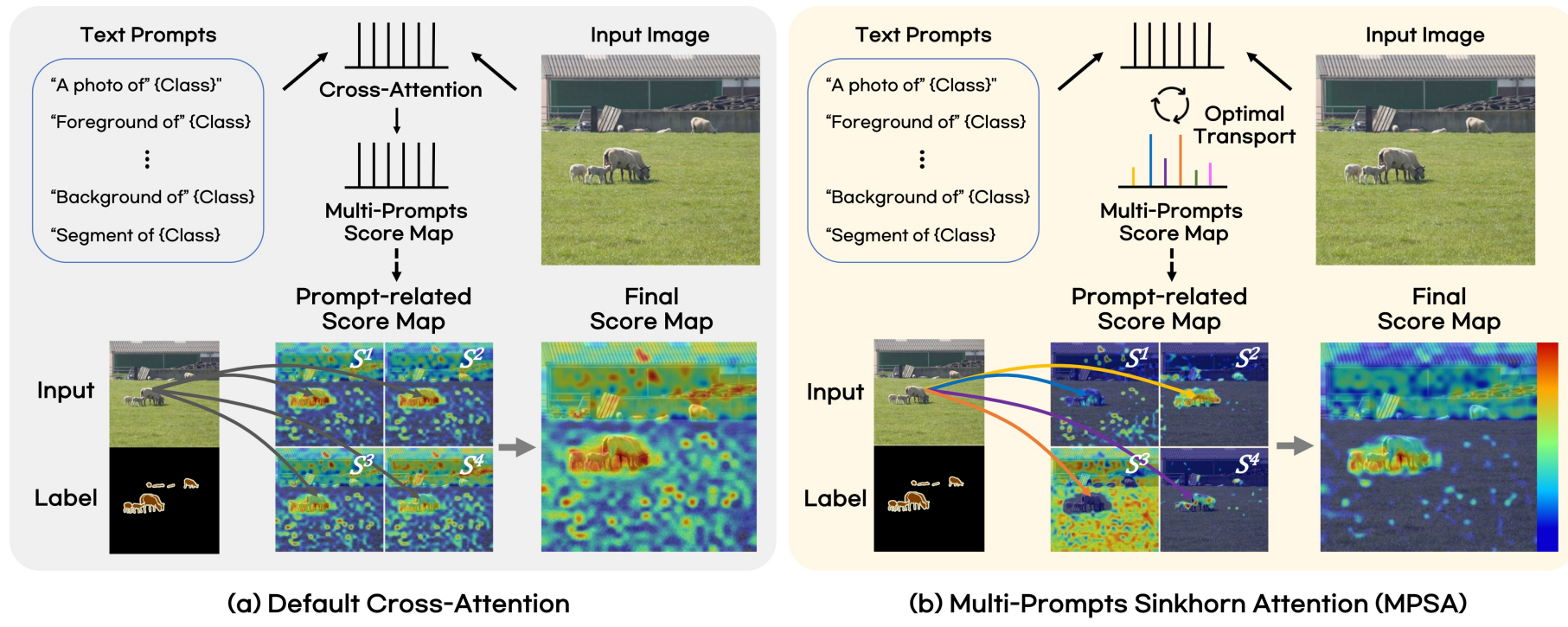


Introduction

Despite the success of CLIP transfer learning for zero-shot semantic segmentation (ZS3), **challenges remain in closely aligning text embeddings with pixel embeddings.**

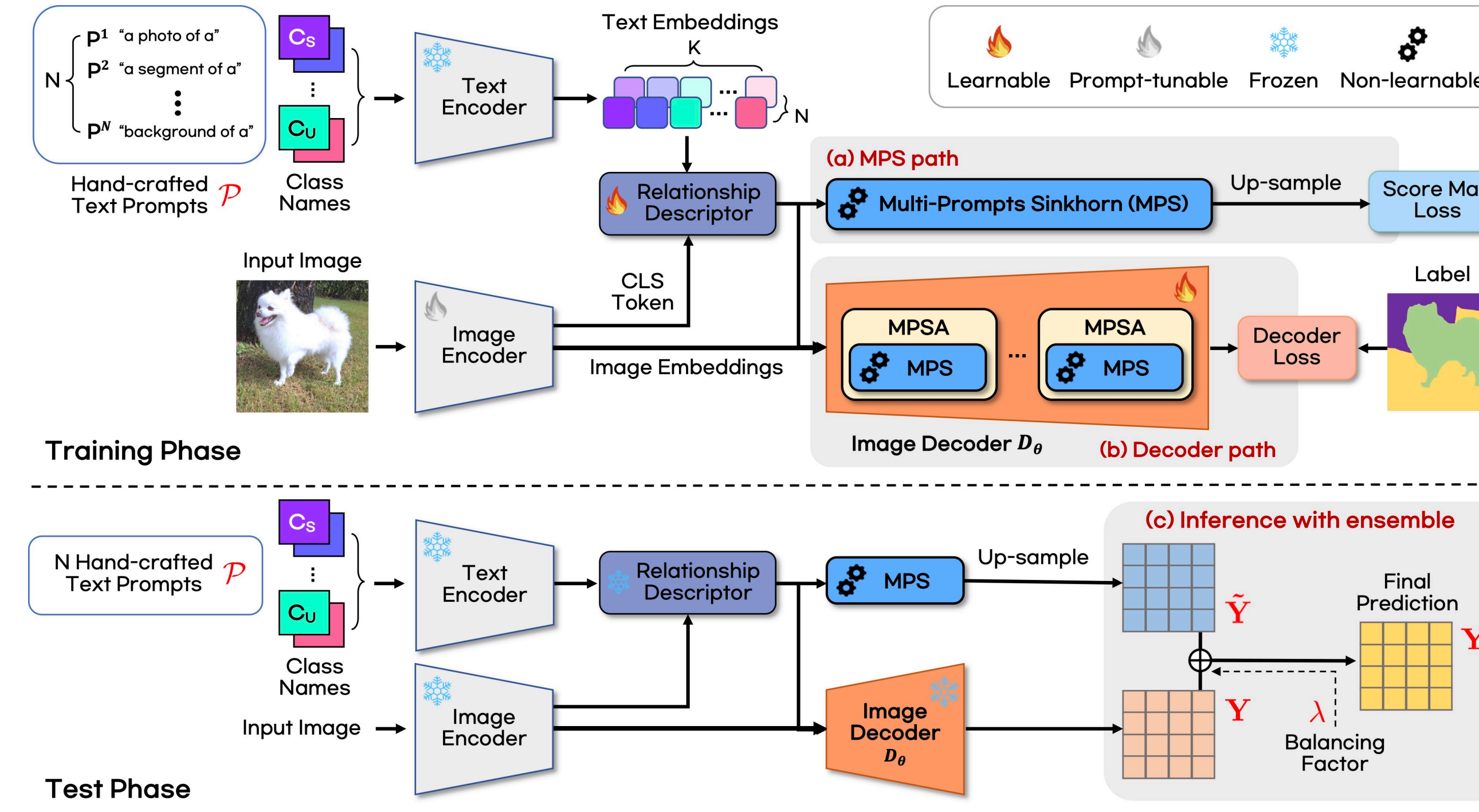
To address this issue, we propose **OTSeg**, a multimodal attention mechanism aimed at enhancing the effectiveness of matching multiple text prompts with corresponding pixel embeddings.

Specifically, we propose **Multi-Prompt Sinkhorn (MPS)** based on the **Optimal Transport**, and along with its extension, **Multi-Prompts Sinkhorn Attention (MPSA)** which effectively replaces cross-attention mechanisms within Transformer.



Overall Pipeline of OTSeg+

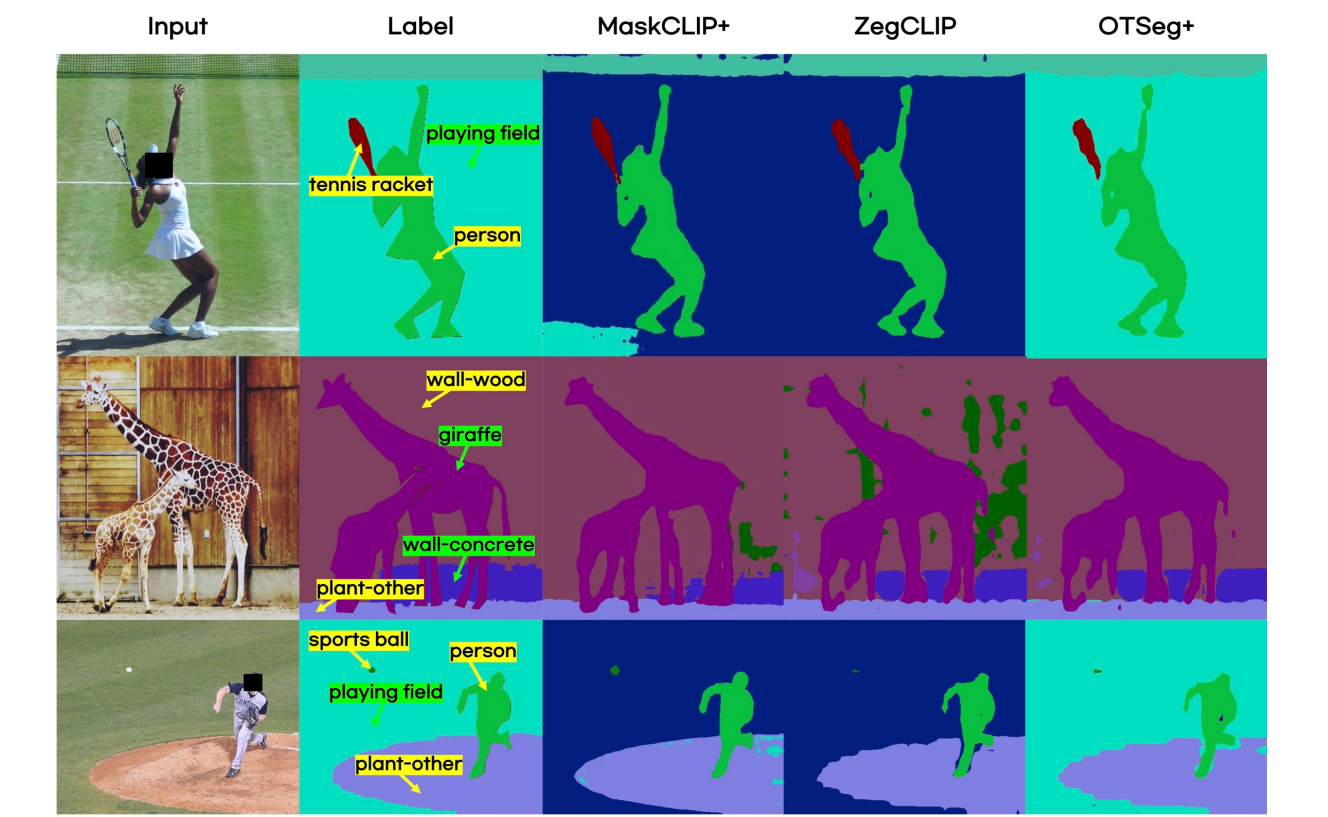
OTSeg+ ensembles (a) MPS path output and (b) Decoder path output to yield final prediction.



Comparison with Other Methods

OTSeg+ achieves **state-of-the-art performance across 3 benchmarks.**

Methods	VOC 2012		PASCAL Context			COCO-Stuff164K		
	mIoU(U)	mIoU(S)	hIoU	mIoU(U)	mIoU(S)	hIoU	mIoU(U)	mIoU(S)
Inductive setting								
ZegFormer [10]	63.6	86.4	73.3	-	-	-	33.2	36.6
Zsseg [35]	72.5	83.5	77.6	-	-	-	36.3	37.8
ZegCLIP [38]	77.8	91.9	84.3	54.6	46.0	49.9	41.4	40.2
OTSeg	78.1	92.1	84.5	56.7	53.0	54.8	41.4	41.4
OTSeg+	81.6	93.3	87.1	60.4	55.2	57.7	41.8	41.5
Transductive setting								
Zsseg [35]	78.1	79.2	79.3	-	-	-	43.6	39.6
MaskCLIP+ [37]	88.1	86.1	87.4	66.7	48.1	53.3	54.7	45.0
FrostSeg [25]	82.6	91.8	86.9	-	-	-	49.1	42.2
MVP-SEG+ [14]	87.4	89.0	88.0	67.5	48.7	54.0	55.8	39.9
ZegCLIP [38]	89.9	92.3	91.1	68.5	46.8	55.6	59.9	40.7
OTSeg	94.3	94.2	94.2	66.7	53.4	59.3	60.7	41.8
OTSeg+	94.3	94.3	94.4	67.0	54.0	59.8	62.6	41.4
Fully-supervised								
ZegCLIP [38]	90.9	92.4	91.6	78.7	46.5	56.9	63.2	40.7
OTSeg	94.4	94.0	94.2	78.1	55.2	64.7	64.0	41.8
OTSeg+	95.0	94.1	94.6	78.4	54.5	65.5	63.2	41.5

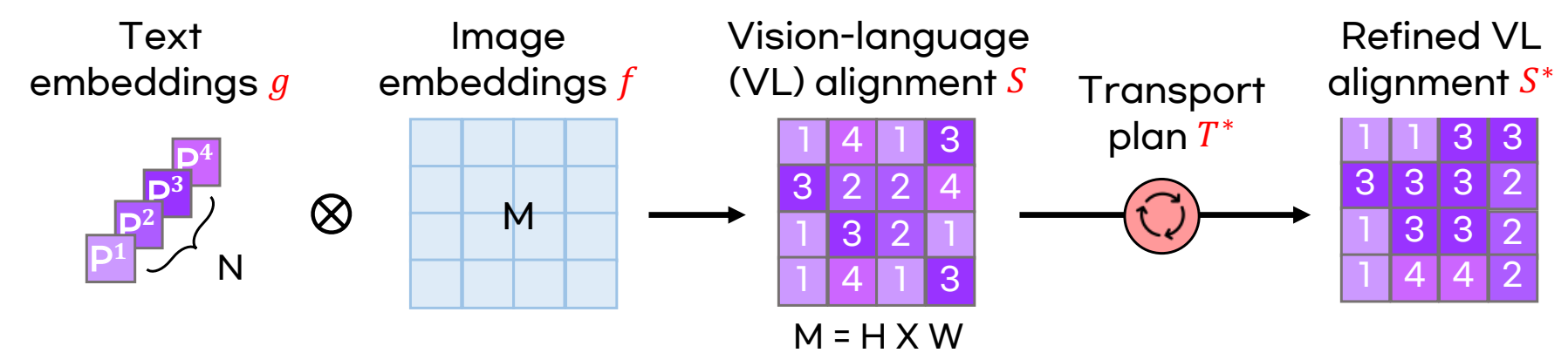


Multi-Prompt Sinkhorn (MPS)

- Inspired by Sinkhorn algorithm, we propose **MPS** to allocate **M-image pixel embeddings** to **N-text prompt embeddings**.
- The discrete optimal transport plan T^* tries to maximize the Vision-language (VL) alignment S by minimizing the total cost C .

$$S^* = \text{MPS}(S) = \mathcal{M}(T^* \odot S),$$

$$\text{where } T^* = \text{Sinkhorn}\left(\frac{C}{\epsilon}\right), C := 1 - S$$



- The total cost C is the matrix multiplication of the text embeddings g and the image embeddings f .
- Minimizing the total cost C corresponds to maximizing the VL alignment S to yield the refined VL alignment S^* .

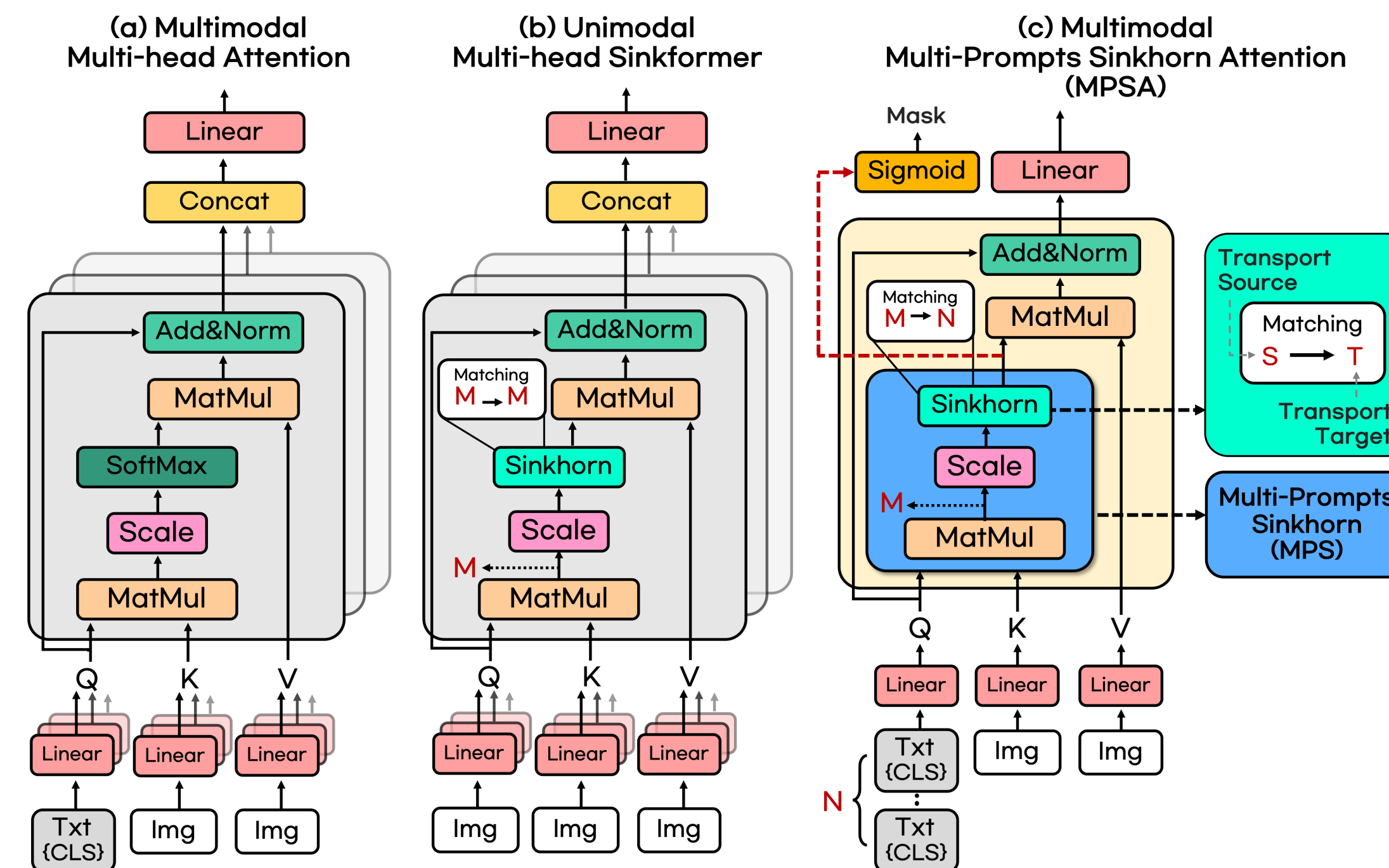
$$T^* = \underset{T \in \mathbb{R}^{M \times N}}{\text{argmin}} \sum_{i=1}^M \sum_{j=1}^N T_{ij} C_{ij} - \epsilon H(T),$$

$$\text{where } C_{ij} = 1 - \frac{f_i g_j^T}{\|f_i\|_2 \|g_j\|_2} = 1 - S_{ij}$$

Multi-Prompt Sinkhorn Attention (MPSA)

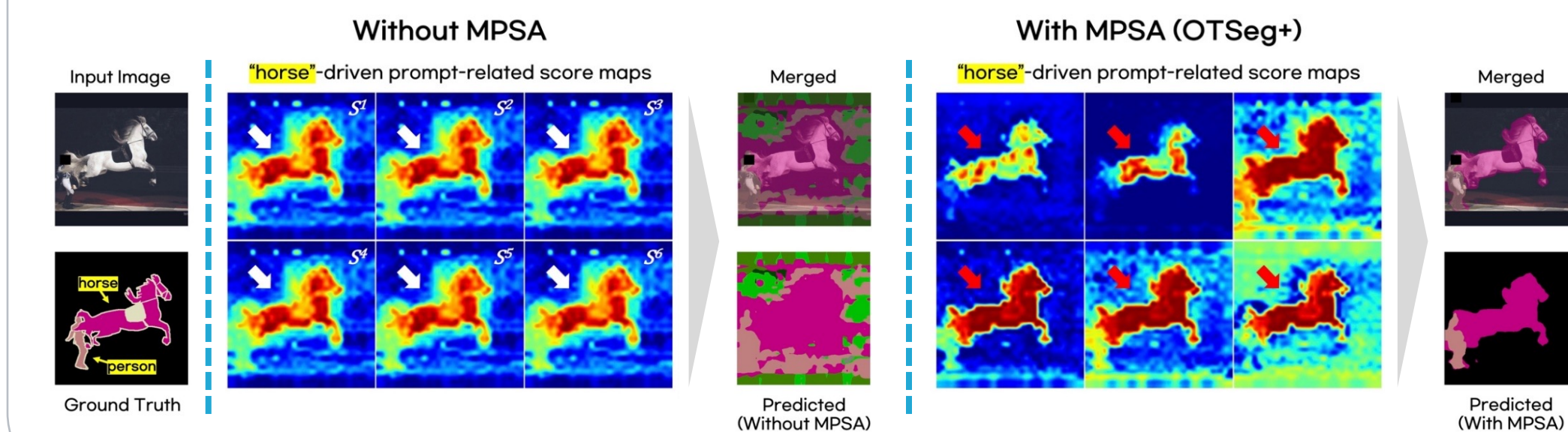
- MPS** can be extended into the cross-attention mechanism, referred as **MPSA** and integrated as a plugin module in each decoder layer.

$$\text{Cross-Attention}(Q, K, V) = \text{SoftMax}\left(\frac{QK^T}{\sqrt{d}}\right)V \xrightarrow{\text{Replace}} \text{MPSA}(Q, K, V) = \text{MPS}(QK^T)V$$



- OTSeg+ demonstrates **efficiency in both computational cost and attention mechanisms** by providing **diversely dispersed score maps**.

Method	# Parameter (M)	GFLOPS ↓	FPS ↑
Zsseg	61.1	1916.7	4.2
ZegFormer	60.3	1829.3	6.8
ZegCLIP	13.8	61.1	25.6
OTSeg	13.8	61.9 ^{-0.8}	23.6 ^{-2.0}
OTSeg+	13.8	61.9 ^{-0.8}	22.5 ^{-3.1}



Conclusion

- We introduce **OTSeg**, a novel multimodal matching framework for ZS3.
- Our proposed component **MPS** and **MPSA** can propose the way for new directions for future researches in **multimodal alignment and zero-shot learning requiring multi-conceptual semantic understandings of vision.**