



## Motivation

### Cross-modal retrieval

- Searching for data when the query and database have different modalities (image and text).

### Ambiguity problem

- Even a single image often contains various contexts.
- Visual manifestations of a caption vary significantly.



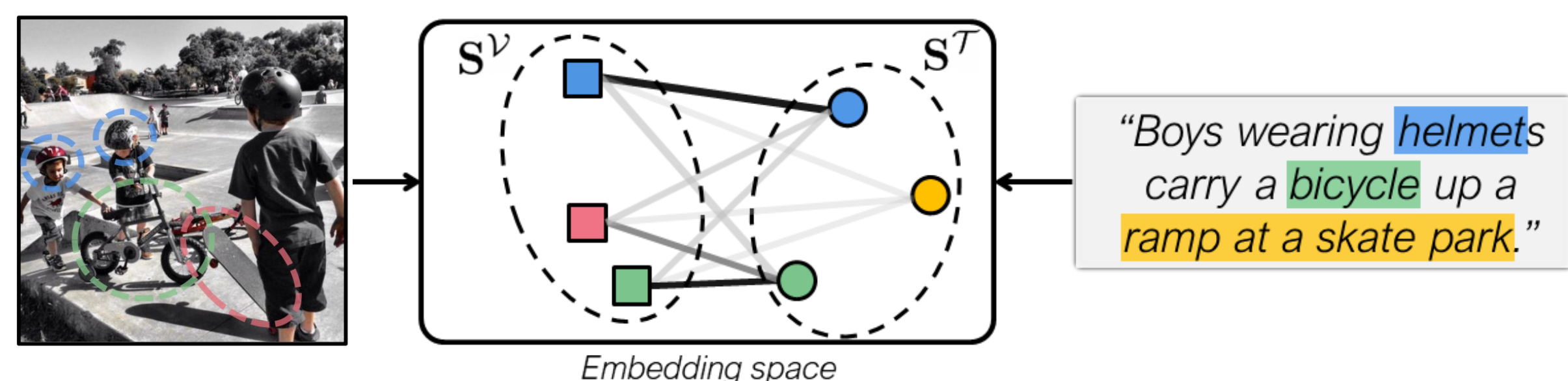
"Boys wearing **helmets** carry a **bicycle** up a ramp at a skate park."

"Small children stand near **bicycles** at a skate park."

"A group of young children riding **bikes** and **skateboards**."

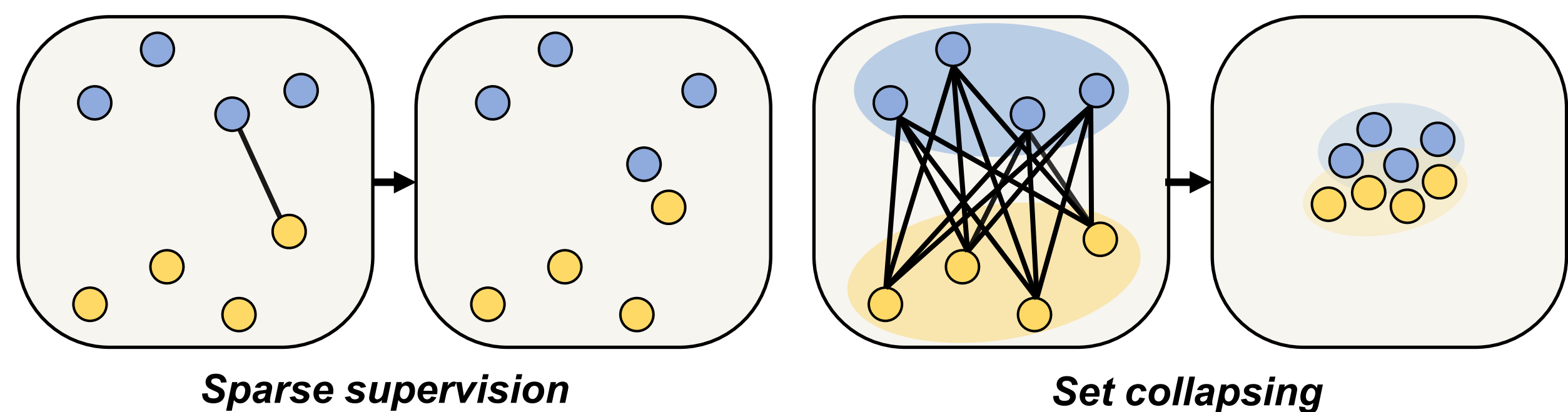
### Previous work: Set-based embedding

- Represent the data with the **set of embedding vectors** (embedding set) [1,2].
- Ambiguity of the data is addressed by elements of the embedding set, which represent diverse semantics.

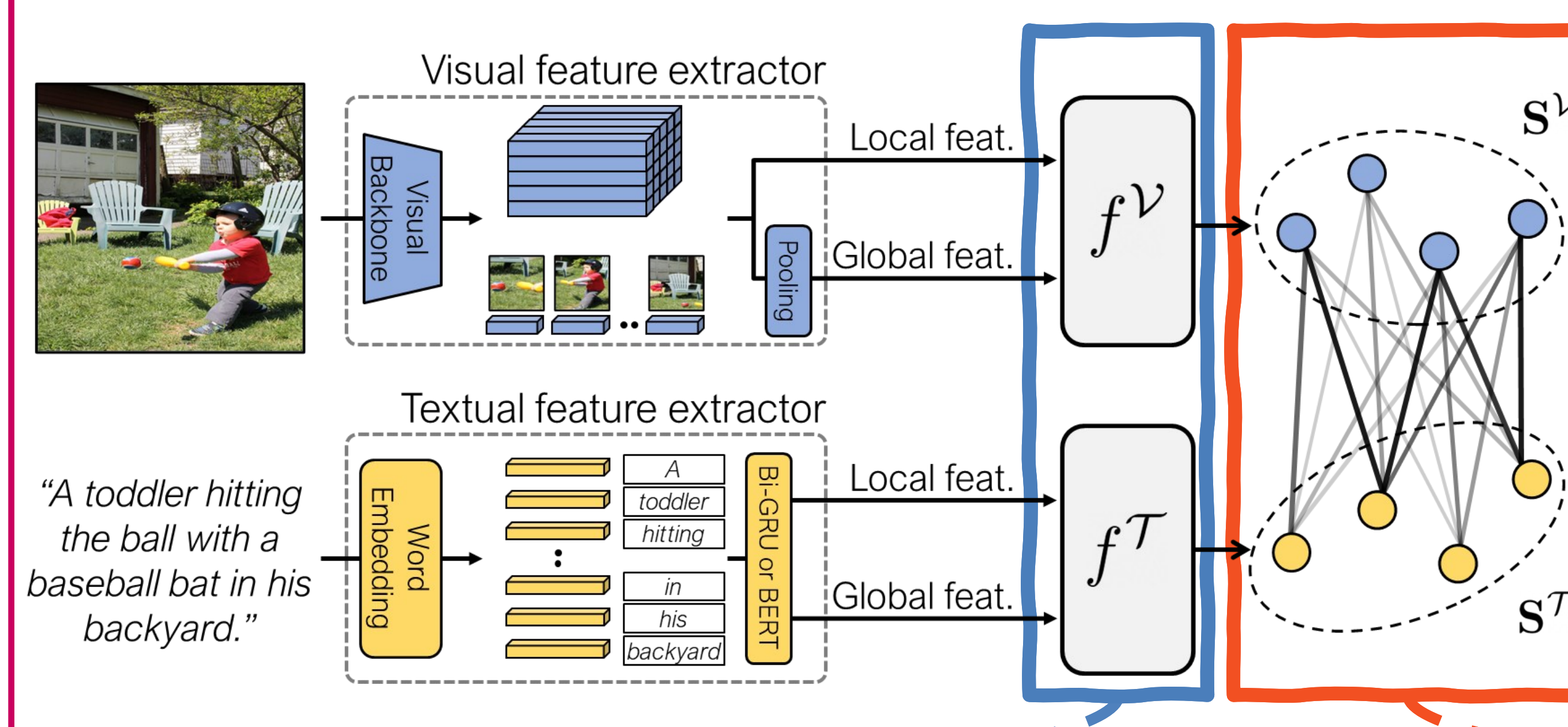


### Drawbacks of previous set-based embedding

- Sparse supervision** → An embedding set most of whose elements remain untrained.
- Set collapsing** → An embedding set with a small variance which does not encode sufficient ambiguity.



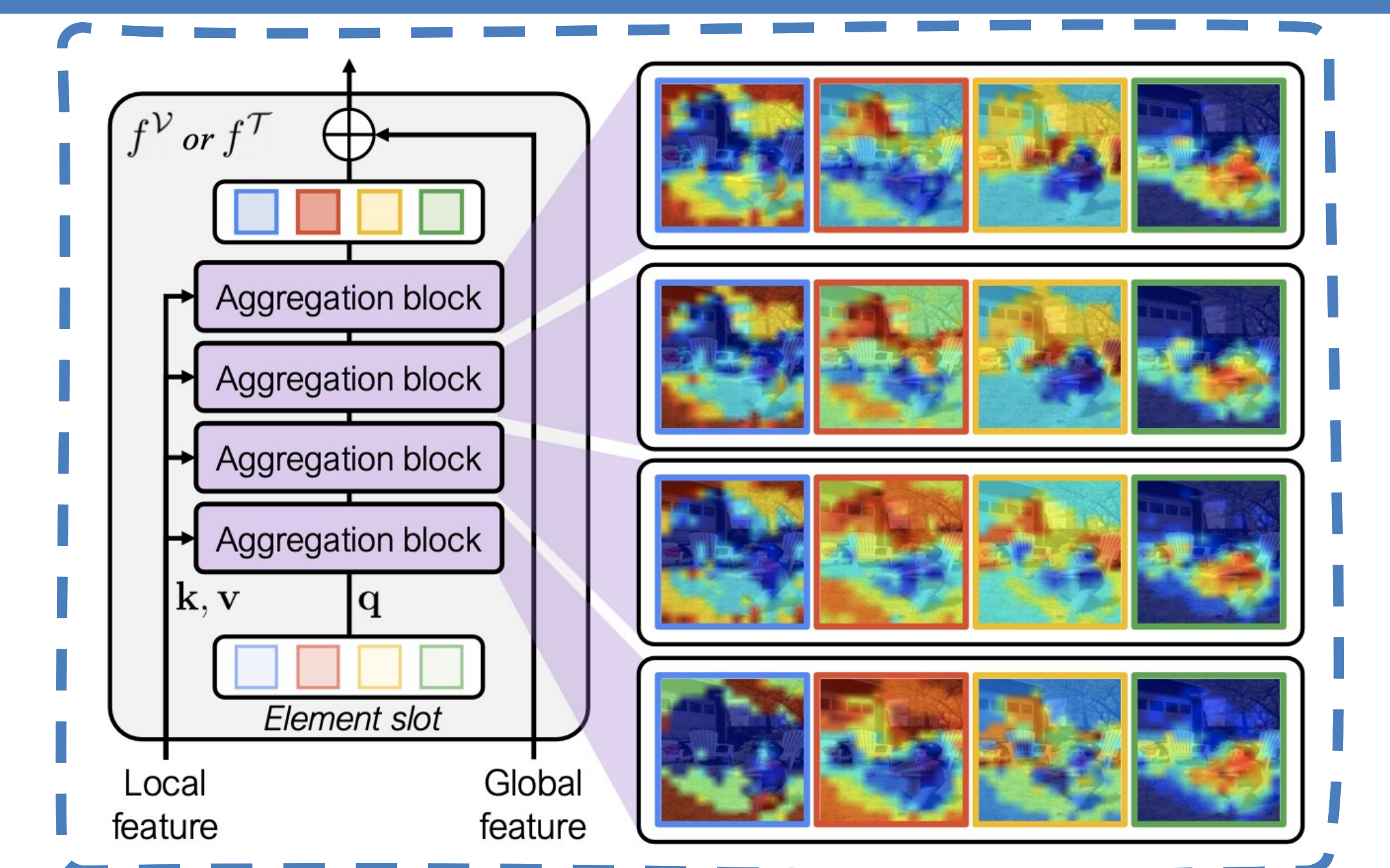
## Our solutions



- Smooth-Chamfer similarity:** Similarity function between sets that provides *dense supervision without collapsing*.
- Set-prediction module:** The module captures *diverse semantic ambiguity* of input, motivated by slot-attn [3].

### Set-prediction module

- Element slots **compete** for aggregating input, progressively transformed into embedding set.
- Competition between slots** makes element encode substantially different semantics.



$$\text{attn} = \text{softmax} \left( \frac{1}{\sqrt{D}} k(\text{inputs}) \cdot q(\text{slots})^T, \text{axis}='slots' \right)$$

Slot-attn [3] based attention scheme (Ours)

$$\text{attn} = \text{softmax} \left( \frac{1}{\sqrt{D}} k(\text{inputs}) \cdot q(\text{slots})^T, \text{axis}='inputs' \right)$$

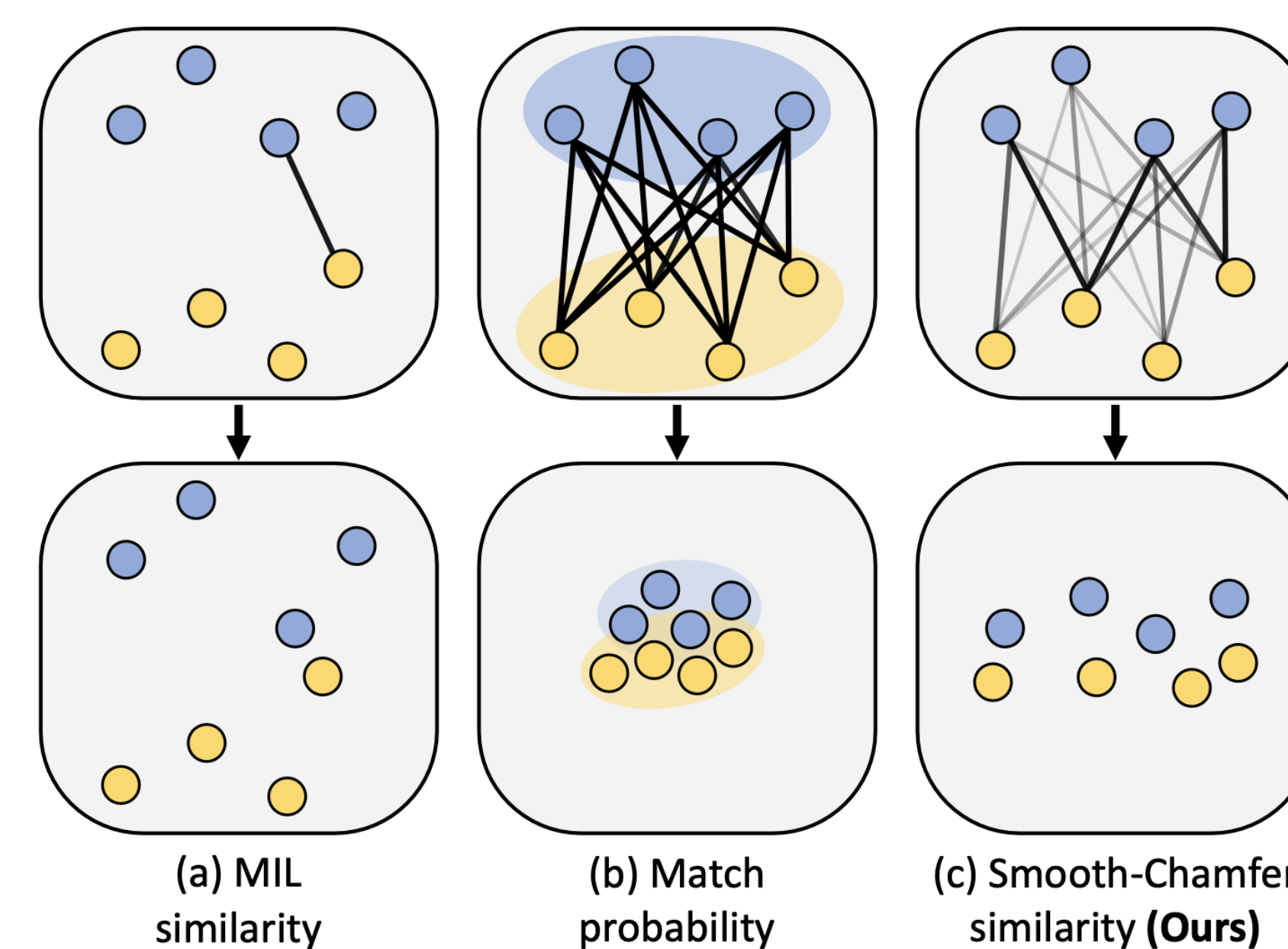
Conventional transformer attention scheme

### Smooth-Chamfer similarity

Proposed SC similarity associates

- every possible pair → *Resolves sparse supervision*
- with different degree of weights. → *Resolves set collapsing*

$$s_{SC}(\mathbf{S}_1, \mathbf{S}_2) = \frac{1}{2\alpha |\mathbf{S}_1|} \sum_{x \in \mathbf{S}_1} \text{LSE}(\alpha c(x, y))_{y \in \mathbf{S}_2} + \frac{1}{2\alpha |\mathbf{S}_2|} \sum_{y \in \mathbf{S}_2} \text{LSE}(\alpha c(x, y))_{x \in \mathbf{S}_1}$$



## Experiments

### Achieved SOTA on COCO, Flickr30K, CxC, and ECCV-caption

Method	CA	Image-to-Text			Text-to-Image			Image-to-Text			Text-to-Image			RSUM	
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10		
<b>Faster R-CNN + Bi-GRU</b>															
SCAN <sup>†</sup> [30]	✓	72.7	94.8	98.4	58.8	88.4	94.8	507.9	50.4	82.2	90.0	38.6	69.3	80.4	410.9
VSRN <sup>†</sup> [31]	✗	76.2	94.8	98.2	62.8	89.7	95.1	516.8	53.0	81.1	89.4	40.5	70.6	81.1	415.7
CAAN [53]	✓	75.5	95.4	98.5	61.3	89.7	95.2	515.6	52.5	83.3	90.9	41.2	70.3	82.9	421.1
IMRAM <sup>†</sup> [6]	✓	76.7	95.6	98.5	61.7	89.1	95.0	516.6	53.7	83.2	91.0	39.7	69.1	79.8	416.5
SGRAM <sup>†</sup> [14]	✓	79.6	96.2	98.5	63.2	90.7	96.1	524.3	57.8	-	91.6	41.9	-	81.3	-
VSE <sub>oc</sub> [27]	✗	78.5	96.0	98.7	61.7	90.3	95.6	520.8	56.6	83.6	91.4	39.3	69.9	81.1	421.9
NAAF <sup>†</sup> [52]	✓	80.5	96.5	98.8	64.1	90.7	96.5	527.2	58.9	85.2	92.0	42.5	70.9	81.4	430.9
Ours <sup>†</sup>	✗	79.8	96.2	98.6	63.6	90.7	95.7	524.6	58.8	84.9	91.5	41.1	72.0	82.4	430.7
Ours <sup>†</sup>	✗	80.6	96.3	98.8	64.7	91.4	96.2	<b>528.0</b>	60.4	86.2	92.4	42.6	73.1	83.1	<b>437.8</b>
<b>ResNeXt-101 + BERT</b>															
VSE <sub>oc</sub> [27]	✗	84.5	98.1	99.4	72.0	93.9	97.5	545.4	66.4	89.3	94.6	51.6	79.3	87.6	468.9
VSE <sub>oc</sub> <sup>†</sup> [27]	✗	85.6	98.0	99.4	73.1	94.3	97.7	548.1	68.1	90.2	95.2	52.7	80.2	88.3	474.8
Ours <sup>†</sup>	✗	86.3	97.8	99.4	72.4	94.0	97.6	547.5	69.1	90.7	95.6	52.1	79.6	87.8	474.9
Ours <sup>†</sup>	✗	86.6	98.2	99.4	73.4	94.5	97.8	<b>549.9</b>	71.0	91.8	96.3	53.4	80.9	88.6	<b>482.0</b>

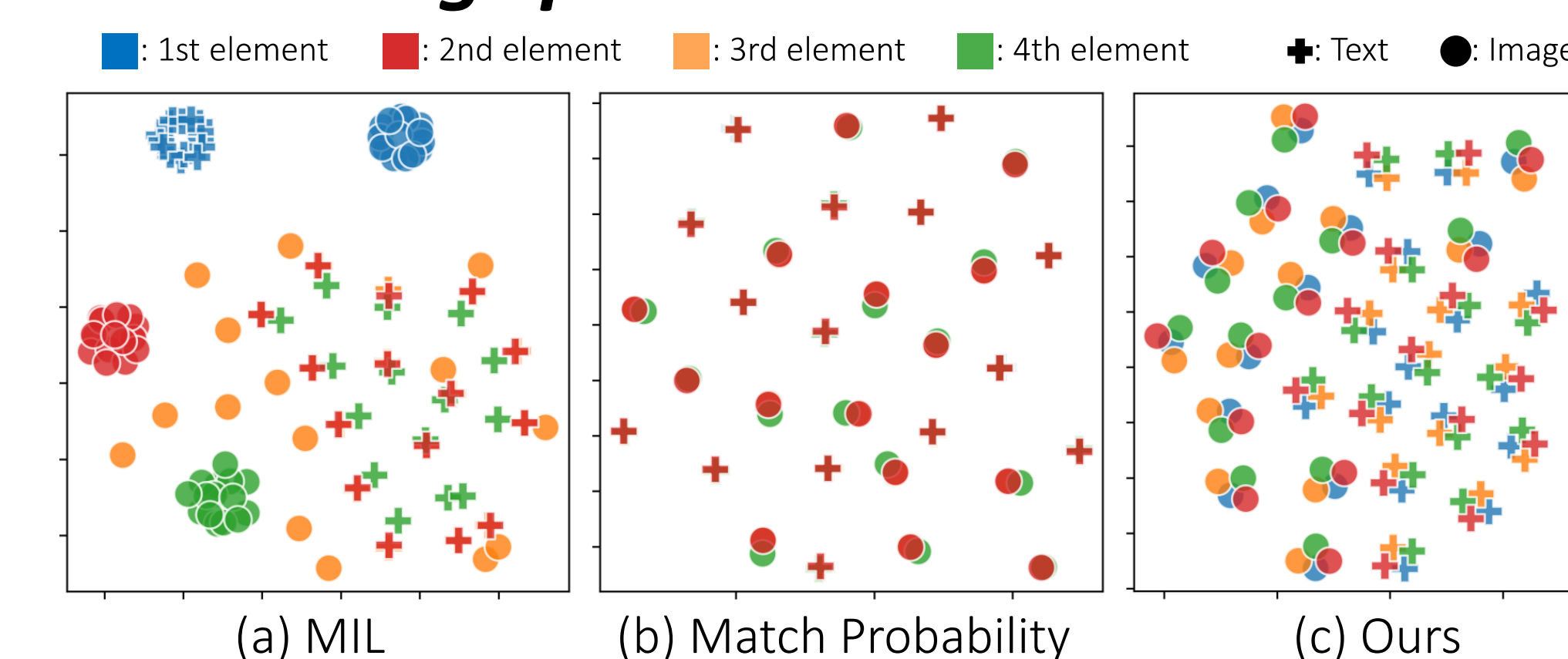
### Embedding set elements & their nearest caption



### Ablation studies: similarity and model

Similarity	Arch.	RSUM	Setting	log(Var.)	RSUM
MIL	Ours	491.7	PIE-Net	-7.35	483.3
MP	Ours	490.5	Ours \w MP	-5.27	490.5
Ours (Chamfer)	Ours	499.6	Transformer	-2.27	496.1
Ours (S-Chamfer)	PIE-Net	483.3	Ours	-2.13	<b>500.8</b>
Ours (S-Chamfer)	Ours	<b>500.8</b>			

### Embedding space visualization



Our method successfully resolves sparse supervision & set collapsing issues.