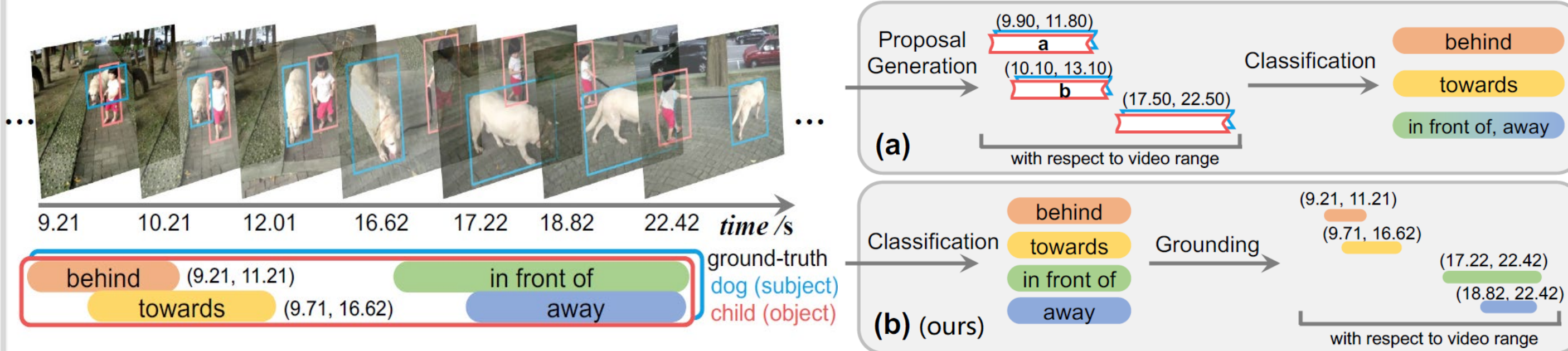
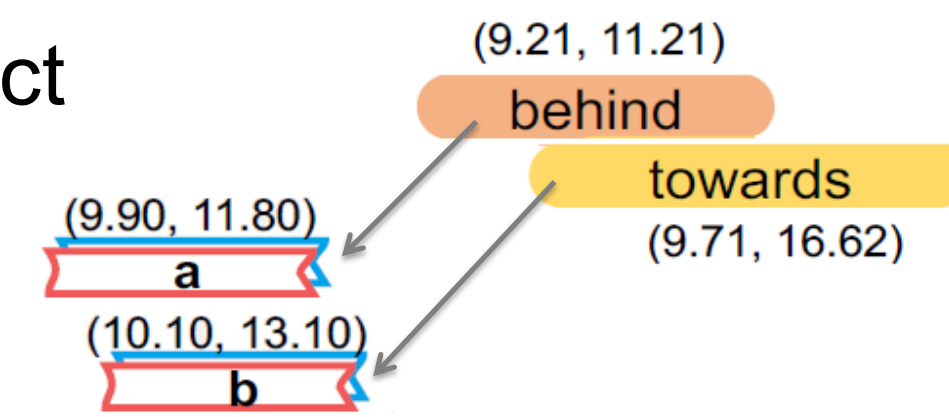


## Detecting Video Visual Relations



### Previous Work *subj-obj pair proposal* $\rightarrow$ *relation classification*

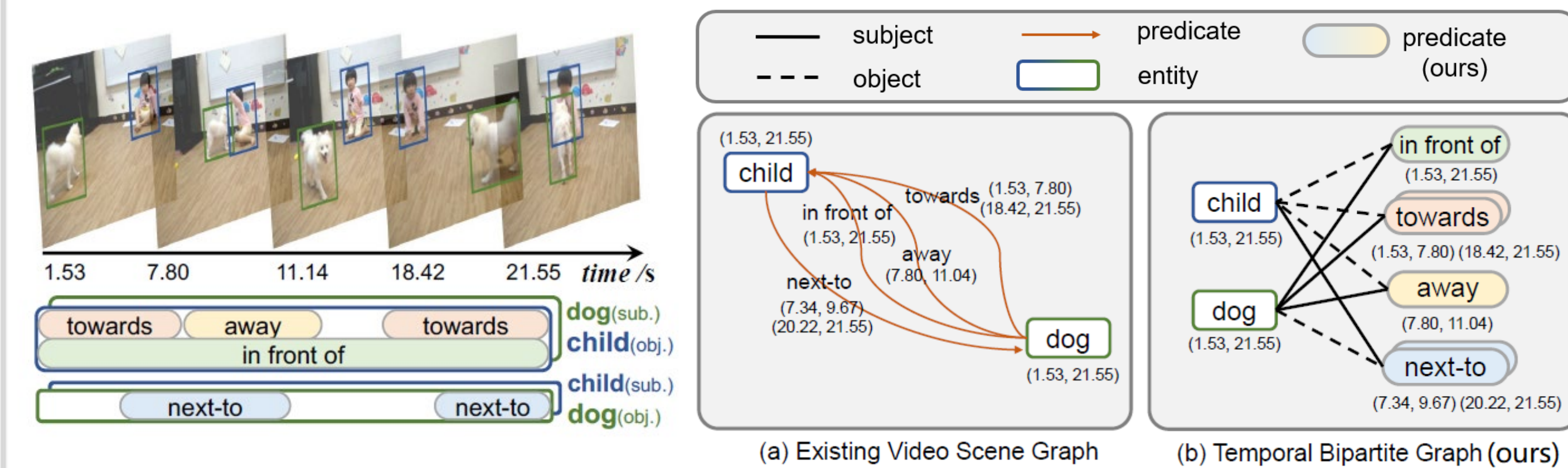
- Label assignment is partially correct
- Discard relation context
- Upper-bounded by proposals



### Ours *relation classification* $\rightarrow$ *temporal grounding*

- A new classification-then-grounding framework
- A novel **B**ipartite **G**raph based model BIG
- Reformulate video scene graphs as temporal bipartite graphs

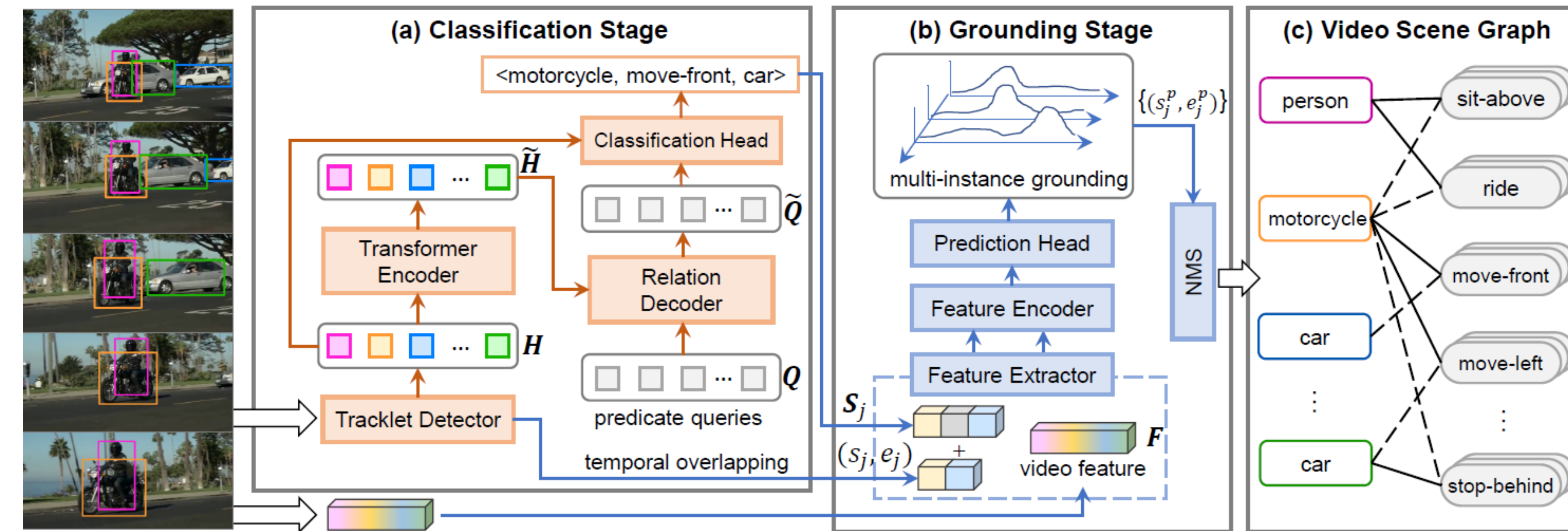
## Temporal Bipartite Graphs



### Advantages

- it avoids enumerating all entity pairs for predicate (relation) prediction
- it is easier to model entity pairs with multiple predicates
- it can be easily extended to more general relations with more semantic roles

## Classification-Then-Grounding



### Classification Stage

- A Transformer-based model
- Classify the categories of all entity & predicate nodes
- Learn the edges of the bipartite graph

R-norm	$F_*$	RelDet (%)		RelTag (%)	
		mAP	R@50	P@1	P@5
✓	✓	7.98	7.71	61.65	51.10
✓	✓	8.02	7.36	61.65	51.68
✓	✓	8.29	7.92	64.42	51.70

Table 5. Ablations of BIG-C for the R-norm and  $F_*$  of RaCA module on VidOR.

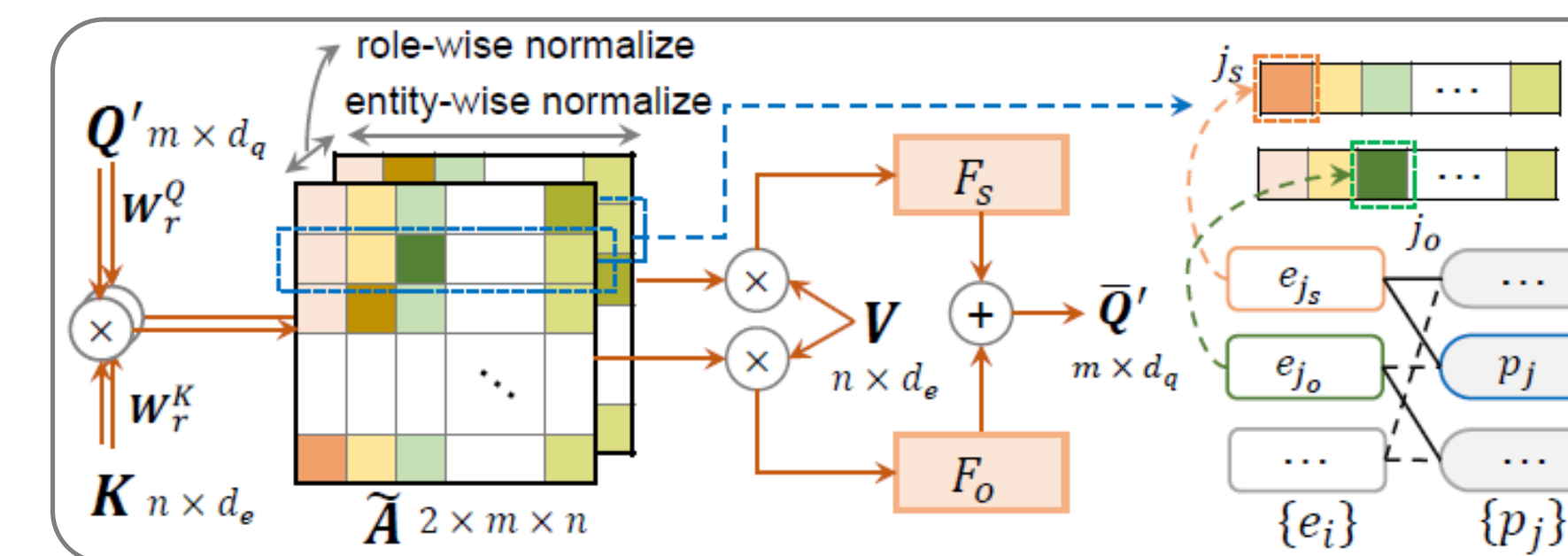
### Grounding Stage

- Localize the temporal location of each predicate node
- Take triplet category as language query

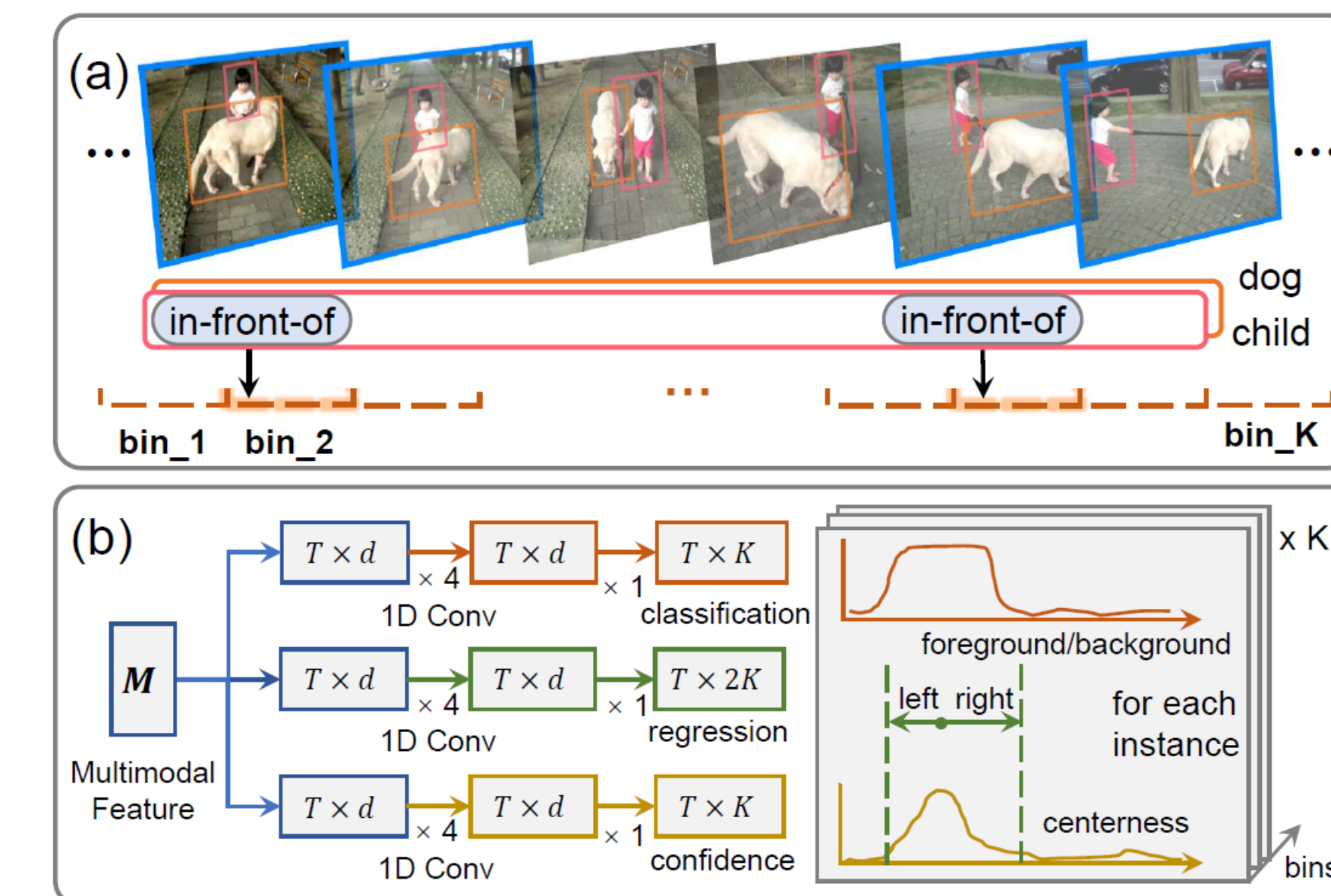
#Bins	$fR_S@K$ (%)			$fR_M@K$ (%)		
	50	100	150	50	100	150
1	12.96	15.59	16.76	5.53	6.86	7.46
5	13.07	15.83	17.26	5.75	7.20	8.05
10	13.04	15.89	17.61	5.75	7.30	8.25

Table 6. Ablations for multi-instance grounding with different number of bins on VidOR.

### Role-aware Cross-Attention



### Multi-instance Grounding



## Experiments Results

### Compare with SOTA on VidOR dataset

Models	Detector	Features				RelDet			RelTag		
		Visual	Lang	Motion	Mask	mAP	R@50	R@100	P@1	P@5	P@10
Liu et al. [24] <sub>CVPR'20</sub>	RefineDet	RoI+I3D <sub>r</sub>		✓		6.85	8.21	9.90	51.20	40.73	—
Chen et al. [13] <sub>ICCV'21</sub>	Faster R-CNN	RoI+I3D <sub>r</sub>	✓	✓		10.04	8.94	10.69	61.52	50.05	38.48
Chen et al. [13] <sub>ICCV'21</sub>	Faster R-CNN	RoI+I3D <sub>r</sub>	✓	✓	✓	11.21	9.99	11.94	68.86	55.16	43.40
IVRD [22] <sub>MM'21</sub>	Faster R-CNN	RoI		✓		7.42	7.36	9.41	53.40	42.70	—
Chen et al. [13] <sub>ICCV'21</sub>	Faster R-CNN	RoI		✓		<b>8.93</b>	7.38	9.22	56.89	44.76	34.07
VidVRD-II [34] <sub>MM'21</sub>	Faster R-CNN	RoI		✓		<b>8.65</b>	<b>8.59</b>	<b>10.69</b>	57.40	44.54	33.30
<b>BIG-C (Ours)</b>	MEGA	RoI		✓		8.03	7.60	9.39	<b>62.25</b>	<b>50.96</b>	<b>40.30</b>
<b>BIG (Ours)</b>	MEGA	RoI+I3D <sub>r</sub>		✓		8.28	<b>7.74</b>	<b>9.82</b>	<b>62.13</b>	<b>51.25</b>	<b>40.48</b>
VRU <sup>+</sup> 19-top1 [37] <sub>MM'19</sub>	FGFA		✓	✓		6.56	6.89	8.83	51.20	40.73	—
MHA [36] <sub>MM'20</sub>	FGFA		✓	✓		6.59	6.35	8.05	50.72	41.56	—
VRU <sup>+</sup> 20-top1 [46] <sub>MM'20</sub>	CascadeRCNN	RoI	✓	✓	✓	<b>9.93</b>	<b>9.12</b>	—	<b>67.43</b>	—	—
Chen et al. [13] <sub>ICCV'21</sub>	Faster R-CNN	RoI	✓	✓	✓	<b>9.54</b>	<b>8.49</b>	<b>10.17</b>	59.24	47.24	35.99
<b>BIG-C (Ours)</b>	MEGA	RoI	✓	✓		8.29	7.92	9.65	64.42	<b>51.70</b>	<b>41.05</b>
<b>BIG (Ours)</b>	MEGA	RoI+I3D <sub>r</sub>	✓	✓		8.54	8.03	<b>10.04</b>	<b>64.42</b>	<b>51.80</b>	<b>40.96</b>

Table 3. Performance (%) on VidOR of SOTA models. The **Best** and **second best** are marked in according formats. **Visual**: I3D<sub>r</sub> and I3D<sub>f</sub> denote region-level and frame-level I3D features, respectively. **Lang**: The word embeddings of entity categories. **Motion**: It refers to the relative motion feature of entity pairs [34]. **Mask**: It means the localization mask of entities [46].

### Compare with SOTA on VidVRD dataset

Models	Features		RelDet			RelTag		
	Visual	Motion	mAP	R@50	R@100	P@1	P@5	P@10
VidVRD [35] <sub>MM'17</sub>	iDT	✓	8.58	5.54	6.37	43.00	28.90	20.80
GSTEG [40] <sub>CVPR'19</sub>	iDT	✓	9.52	7.05	8.67	51.50	39.50	28.23
VRD-GCN [30] <sub>MM'19</sub>	iDT	✓	16.26	8.07	9.33	57.50	41.00	28.50
MHA [36] <sub>MM'20</sub>	iDT	✓	19.03	9.53	10.38	57.50	41.40	29.45
IVRD [22] <sub>MM'21</sub>	RoI	✓	22.97	12.40	14.46	68.83	49.87	35.57
VidVRD-II [34] <sub>MM'21</sub>	RoI	✓	29.37	19.63	22.92	70.40	53.88	40.16
Liu et al. [24] <sub>CVPR'20</sub>	RoI+I3D <sup>+</sup>	✓	18.38	11.21	13.69	60.00	43.10	32.24
Chen et al. [13] <sub>ICCV'21</sub>	RoI+I3D	✓	20.08	13.73	16.88	62.50	49.20	38.45
Liu et al. [24] <sub>CVPR'20</sub>	RoI <sup>+</sup>		14.01	8.47	<b>11.00</b>	56.50	36.70	26.60
TRACE [39] <sub>ICCV'21</sub>	RoI		15.06	7.67	10.32	—	—	—
<b>BIG-C (Ours)</b>	RoI <sup>+</sup>		<b>17.56</b>	<b>9.59</b>	10.92	<b>56.50</b>	<b>44.30</b>	<b>32.35</b>
Liu et al. [24] <sub>CVPR'20</sub>	RoI+I3D <sup>+</sup>		14.81	9.14	<b>11.39</b>	55.50	38.90	28.90
TRACE [39] <sub>ICCV'21</sub>	RoI+I3D		17.57	9.08	11.15	<b>61.00</b>	<b>45.30</b>	<b>33.50</b>
<b>BIG-C (Ours)</b>	RoI+I3D <sup>+</sup>		<b>17.67</b>	<b>9.63</b>	11.29	56.00	43.80	32.85
<b>BIG (Ours)</b>	RoI <sup>+</sup>		<b>26.08</b>	<b>14.10</b>	<b>16.25</b>	<b>73.00</b>	<b>55.10</b>	<b>40.00</b>

## Bipartite Graph Visualization

