



# Vote2Cap-DETR: A Set-to-Set Perspective Towards 3D Dense Captioning

Winner Presentation of the Scan2Cap Challenge

**Sijin Chen**<sup>1</sup> Hongyuan Zhu<sup>2</sup> Xin Chen<sup>3</sup> Yinjie Lei<sup>4</sup> Gang YU<sup>3</sup> Tao Chen<sup>1†</sup> Taihao Li<sup>5</sup>

<sup>1</sup> Fudan University

<sup>2</sup> Institute for Infocomm Research (I2R) & Centre for Frontier AI Research (CFAR), A\*STAR, Singapore

<sup>3</sup> Tencent PCG

<sup>4</sup> Sichuan University

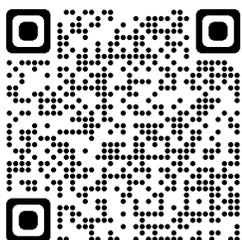
<sup>5</sup> Zhejiang Lab

† Corresponding author.

Paper



Project



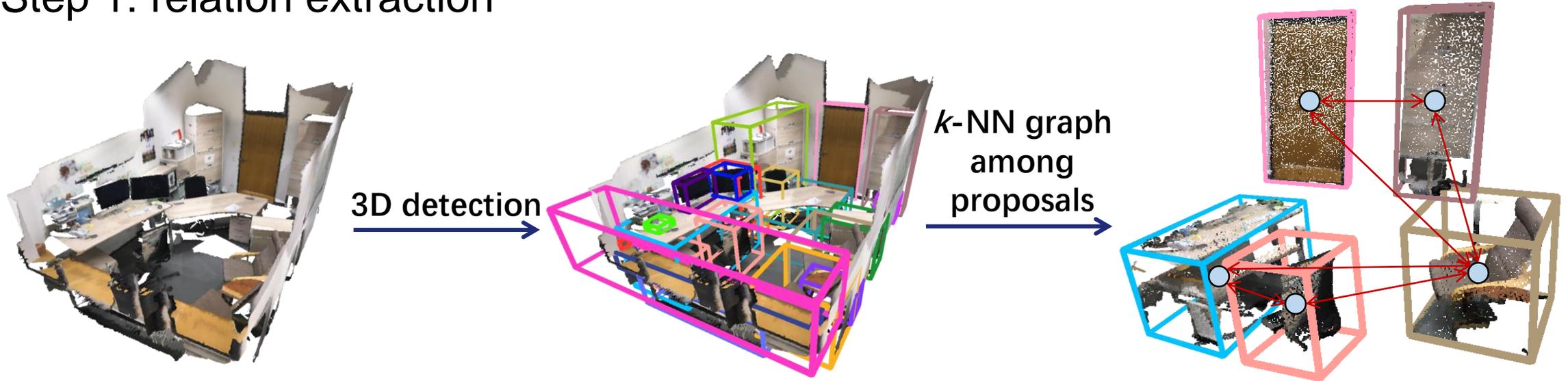
# Task Definition: 3D Dense Captioning



1. **accurate** localization of objects from a cluttered 3D scene;
2. **informative** and **object-centric** descriptions for each instance.

# Previous Explicit Approaches

Step 1: relation extraction



Step 2: IoU based proposal selection

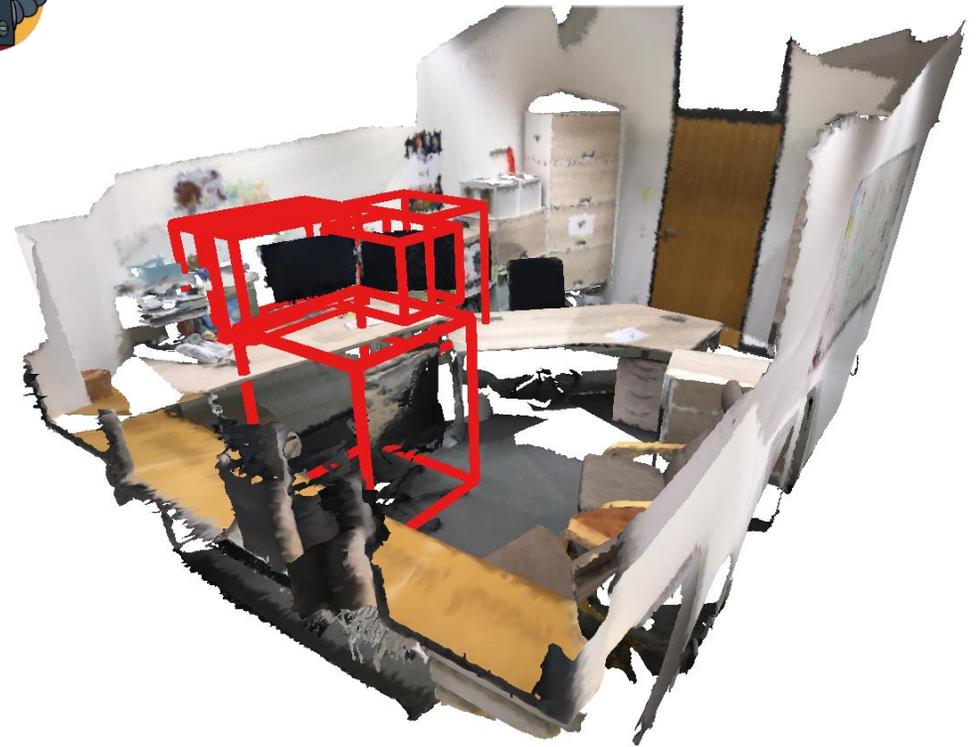


# Limitations

1. cumulative error caused by duplicated and inaccurate box proposals



How many monitors indeed?



# Limitations

2. hyper parameters hard to generalize to diverse 3D scenes

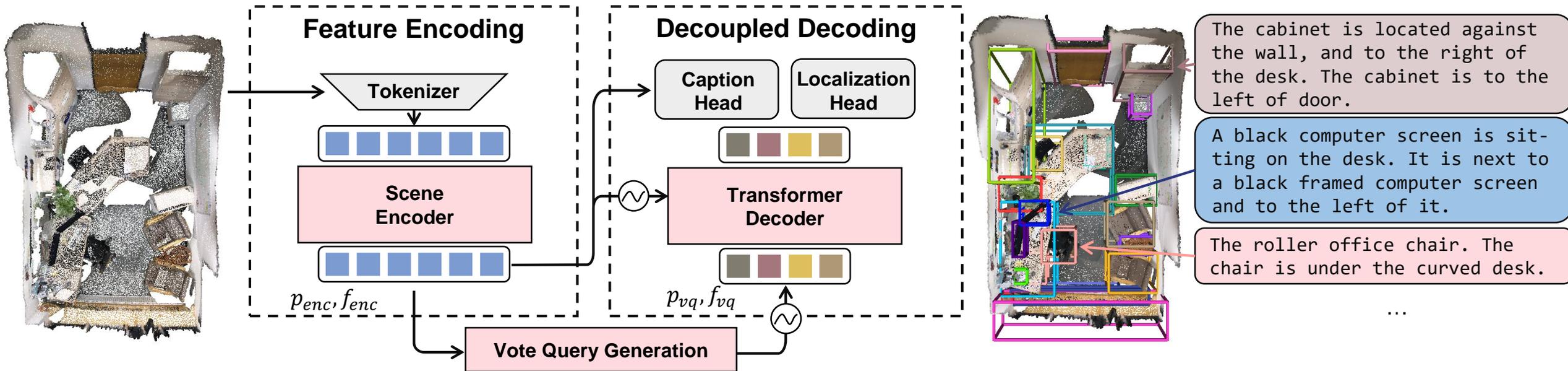


How to construct generalizable  $k$ -NN Graph?



# Vote2Cap-DETR: A Set-to-Set Approach

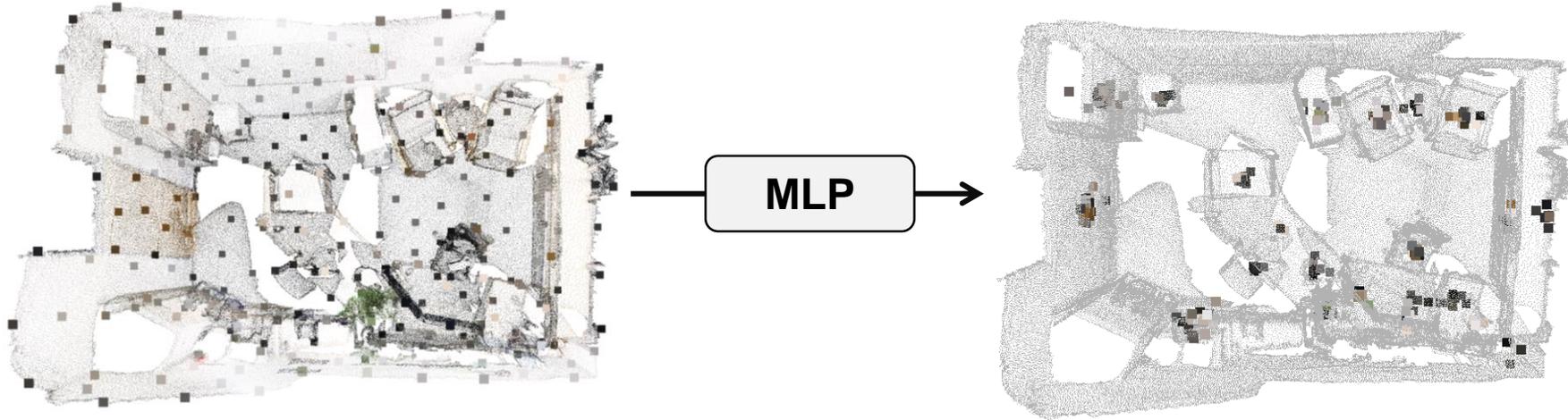
1. translate from a set of points to a set of “box-caption” proposals
2. learn query-to-query, query-to-scene interaction with decoder attention
3. set-to-set training, learning discriminative feature representations



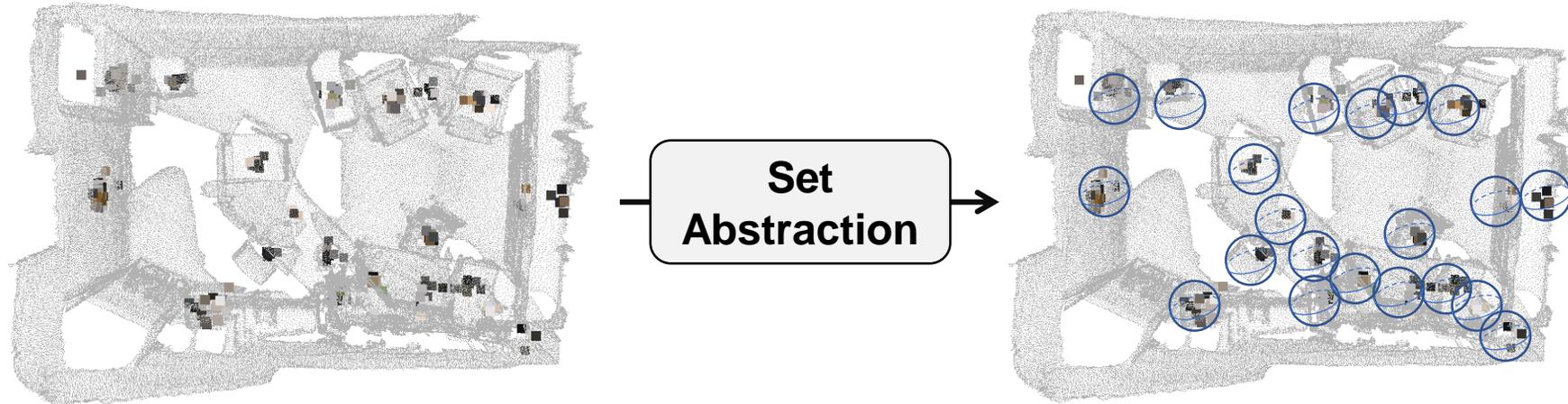
# Method: Component Designs

> design of the vote query

Step 1: vote shifting

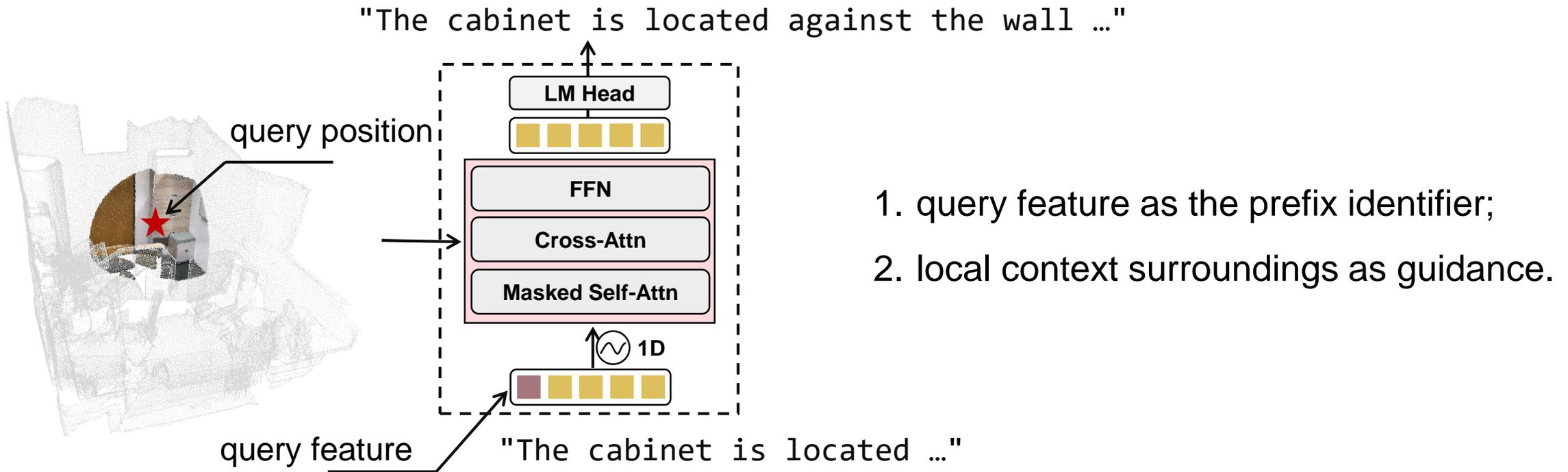


Step 2: local feature aggregation



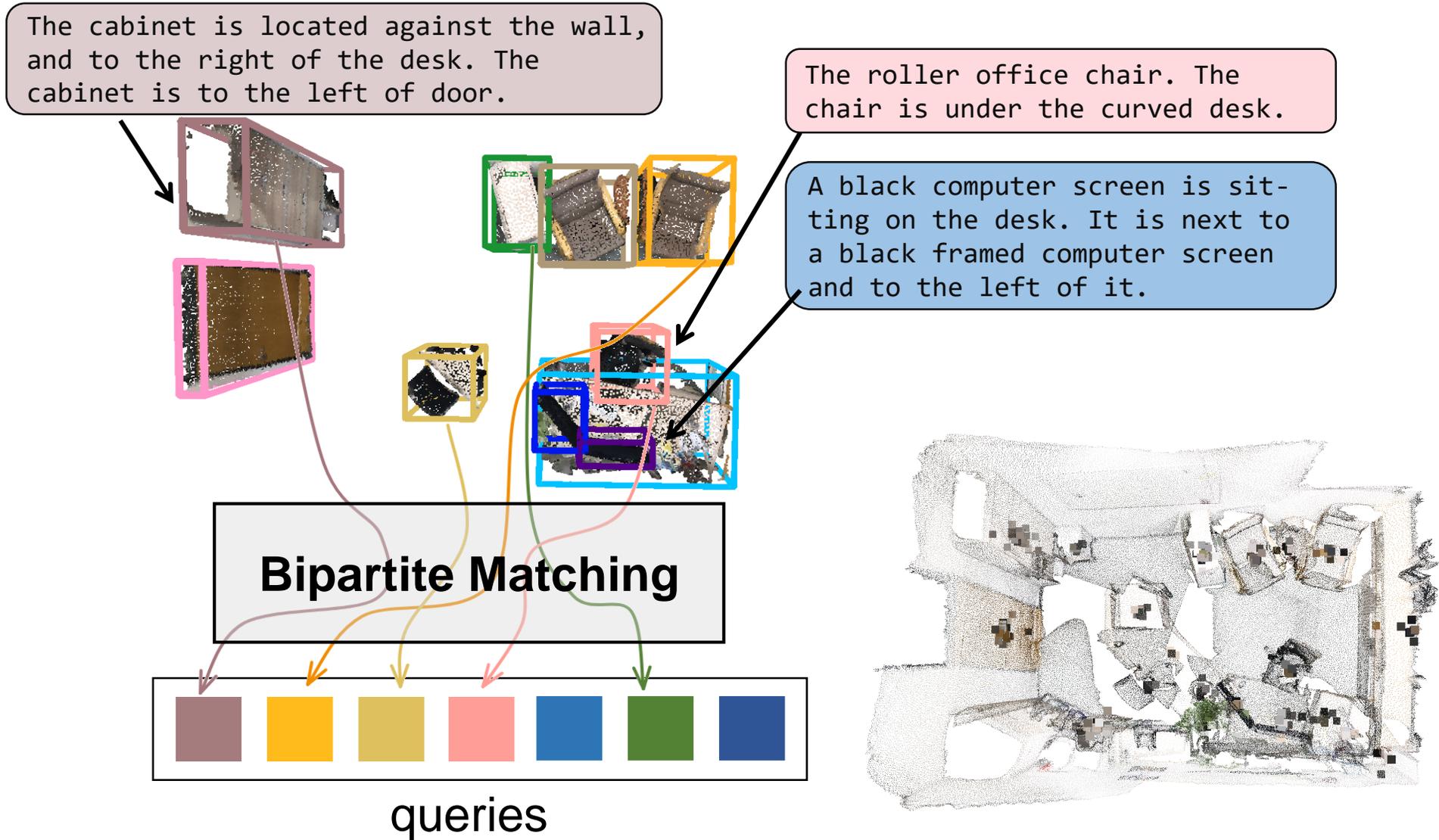
# Method: Component Designs

> design of the caption head



# Method: Set-to-Set Training

> matching queries to a set of “object-caption” pairs



# Quantitative Results

## > ScanRefer validation set

Method	$\mathcal{L}_{des}$	w/o additional 2D input								w/ additional 2D input							
		IoU = 0.25				IoU = 0.50				IoU = 0.25				IoU = 0.50			
		C $\uparrow$	B-4 $\uparrow$	M $\uparrow$	R $\uparrow$	C $\uparrow$	B-4 $\uparrow$	M $\uparrow$	R $\uparrow$	C $\uparrow$	B-4 $\uparrow$	M $\uparrow$	R $\uparrow$	C $\uparrow$	B-4 $\uparrow$	M $\uparrow$	R $\uparrow$
Scan2Cap [13]		53.73	34.25	26.14	54.95	35.20	22.36	21.44	43.57	56.82	34.18	26.29	55.27	39.08	23.32	21.97	44.78
MORE [20]		58.89	35.41	26.36	55.41	38.98	23.01	21.65	44.33	62.91	36.25	26.75	56.33	40.94	22.93	21.66	44.42
SpaCap3d [39]		58.06	35.30	26.16	55.03	42.76	25.38	22.84	45.66	63.30	36.46	26.71	55.71	44.02	25.26	22.33	45.36
3DJCG [4]	MLE	60.86	<b>39.67</b>	27.45	59.02	47.68	31.53	24.28	51.80	64.70	<b>40.17</b>	27.66	<b>59.23</b>	49.48	31.03	24.22	50.80
D3Net [7]		-	-	-	-	-	-	-	-	-	-	-	-	46.07	30.29	24.35	51.67
Ours		<b>71.45</b>	39.34	<b>28.25</b>	<b>59.33</b>	<b>61.81</b>	<b>34.46</b>	<b>26.22</b>	<b>54.40</b>	<b>72.79</b>	39.17	<b>28.06</b>	<b>59.23</b>	<b>59.32</b>	<b>32.42</b>	<b>25.28</b>	<b>52.53</b>
$\chi$ -Trans2Cap [43]		58.81	34.17	25.81	54.10	41.52	23.83	21.90	44.97	61.83	35.65	26.61	54.70	43.87	25.05	22.46	45.28
Scan2Cap [13]		-	-	-	-	-	-	-	-	-	-	-	-	48.38	26.09	22.15	44.74
D3Net [7]	SCST	-	-	-	-	-	-	-	-	-	-	-	-	62.64	35.68	<b>25.72</b>	<b>53.90</b>
Ours		<b>84.15</b>	<b>42.51</b>	<b>28.47</b>	<b>59.26</b>	<b>73.77</b>	<b>38.21</b>	<b>26.64</b>	<b>54.71</b>	<b>86.28</b>	<b>42.64</b>	<b>28.27</b>	<b>59.07</b>	<b>70.63</b>	<b>35.69</b>	25.51	52.28

## > Nr3D validation set

Method	$\mathcal{L}_{des}$	C@0.5 $\uparrow$	B-4@0.5 $\uparrow$	M@0.5 $\uparrow$	R@0.5 $\uparrow$
Scan2Cap [13]		27.47	17.24	21.80	49.06
SpaCap3d [39]		33.71	19.92	22.61	50.50
D3Net [7]	MLE	33.85	20.70	23.13	53.38
3DJCG [4]		38.06	22.82	23.77	52.99
Ours		<b>43.84</b>	<b>26.68</b>	<b>25.41</b>	<b>54.43</b>
$\chi$ -Tran2Cap [43]		33.62	19.29	22.27	50.00
D3Net [7]	SCST	38.42	22.22	24.74	54.37
Ours		<b>45.53</b>	<b>26.88</b>	<b>25.43</b>	<b>54.76</b>

# Qualitative Results



scene0011\_00

**3DJCG:** This is a rectangular **whiteboard**. It is on the wall.

**SpaCap3D:** The **whiteboard** is affixed to the wall. It is to the right of the window.

**Ours:** The **tv** is on the wall. It is to the right of the table.

**GT:** This is a big black tv. It is above a thin table.



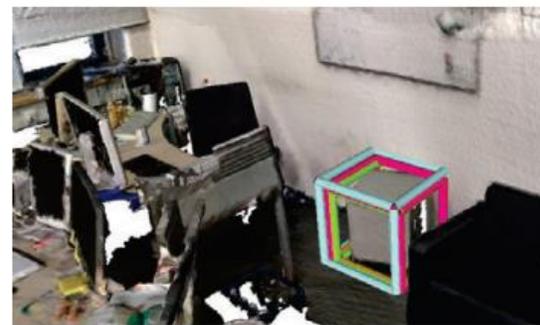
scene0015\_00

**3DJCG:** This is a **brown** table. It is in the **middle** of the room.

**SpaCap3D:** This is a **wooden** table. It is in the **center** of the room.

**Ours:** This is a **wooden** table. It is in the corner of the room.

**GT:** This is a small table with a wood look. It is the table closest to the front of the room in the upper left corner.



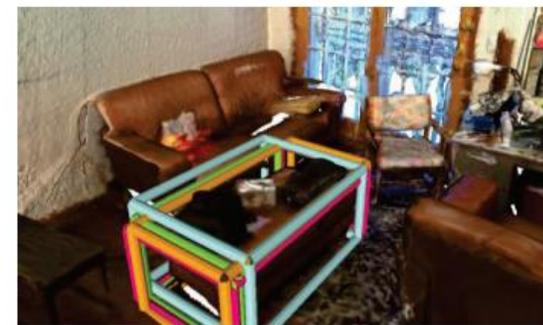
scene0025\_00

**3DJCG:** The is a small **brown** cabinet. It is to the right of the desk.

**SpaCap3D:** The cabinet is **below the desk**. It is **to the left of the chair**.

**Ours:** This is a **white** cabinet. It is to the right of the table.

**GT:** A white cabinet is sitting on the floor next to the wall. It is to the left of the couch.



scene0050\_00

**3DJCG:** This is a **brown** table. It is in front of the couch.

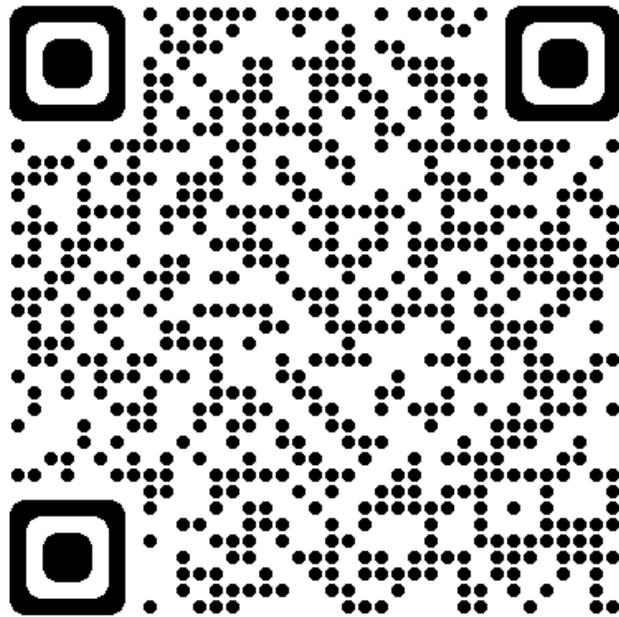
**SpaCap3D:** This is a **wooden** coffee table. It is in front of the couch.

**Ours:** This is a **brown ottoman**. It is to the right of the chair.

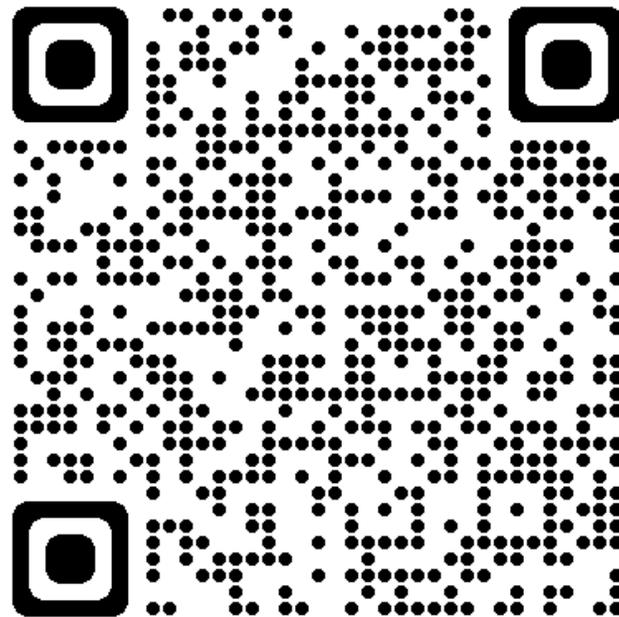
**GT:** This is a brown ottoman. It is in front of a couch.

# One More Thing...

Vote2Cap-DETR++



 Github



 Weights

