

### TL;DR

Previous:

1. Segment whole 3D scene
2. Retrieve the most suitable fixed-shape mask

Mask incorrect/  
different granularity  
=>  
Can't find / Bad masks

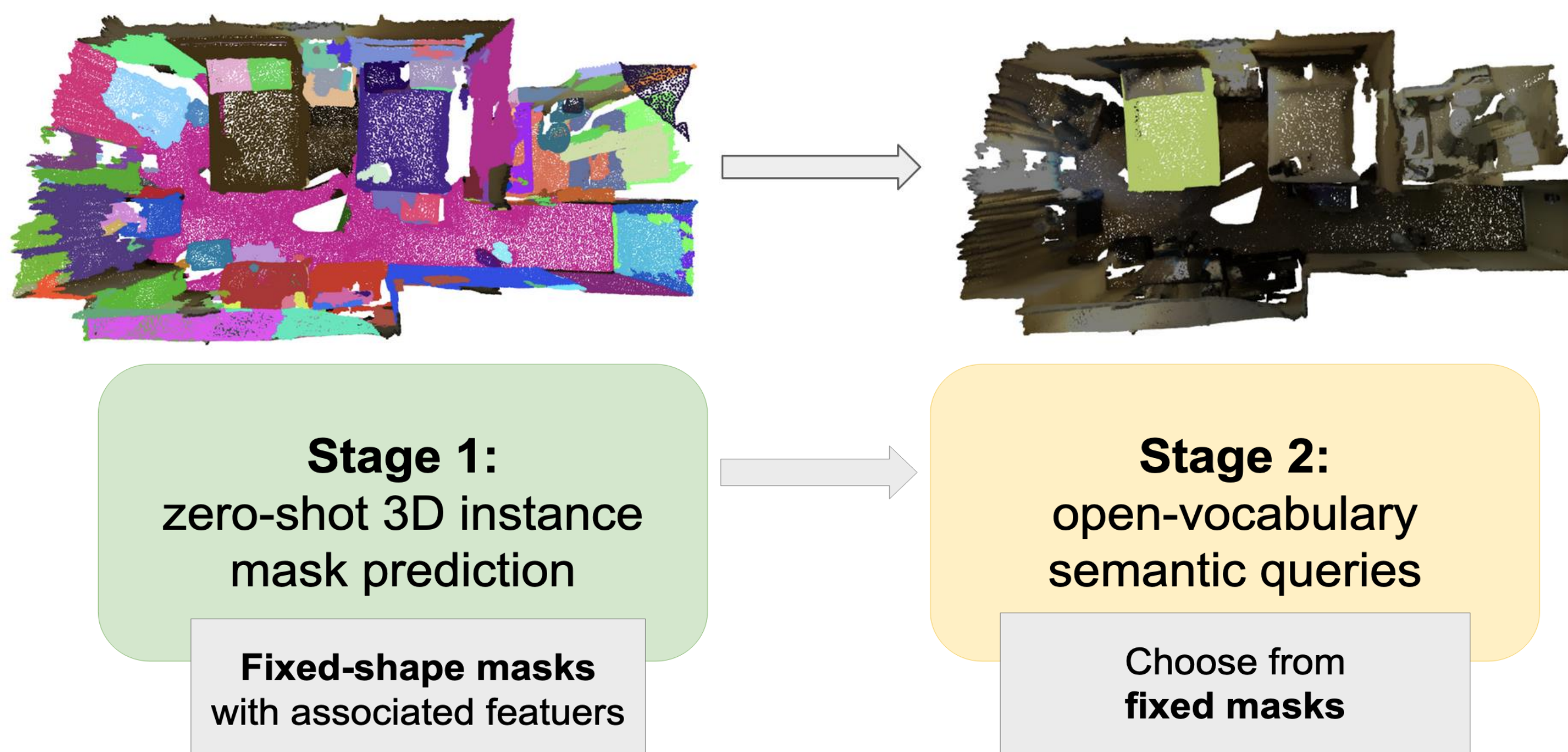
Ours:

1. Segment whole 3D scene
2. Query-Aware mask from 2D
3. Refine fixed masks from step 1

Good masks!  
This is what I want.

### Background:

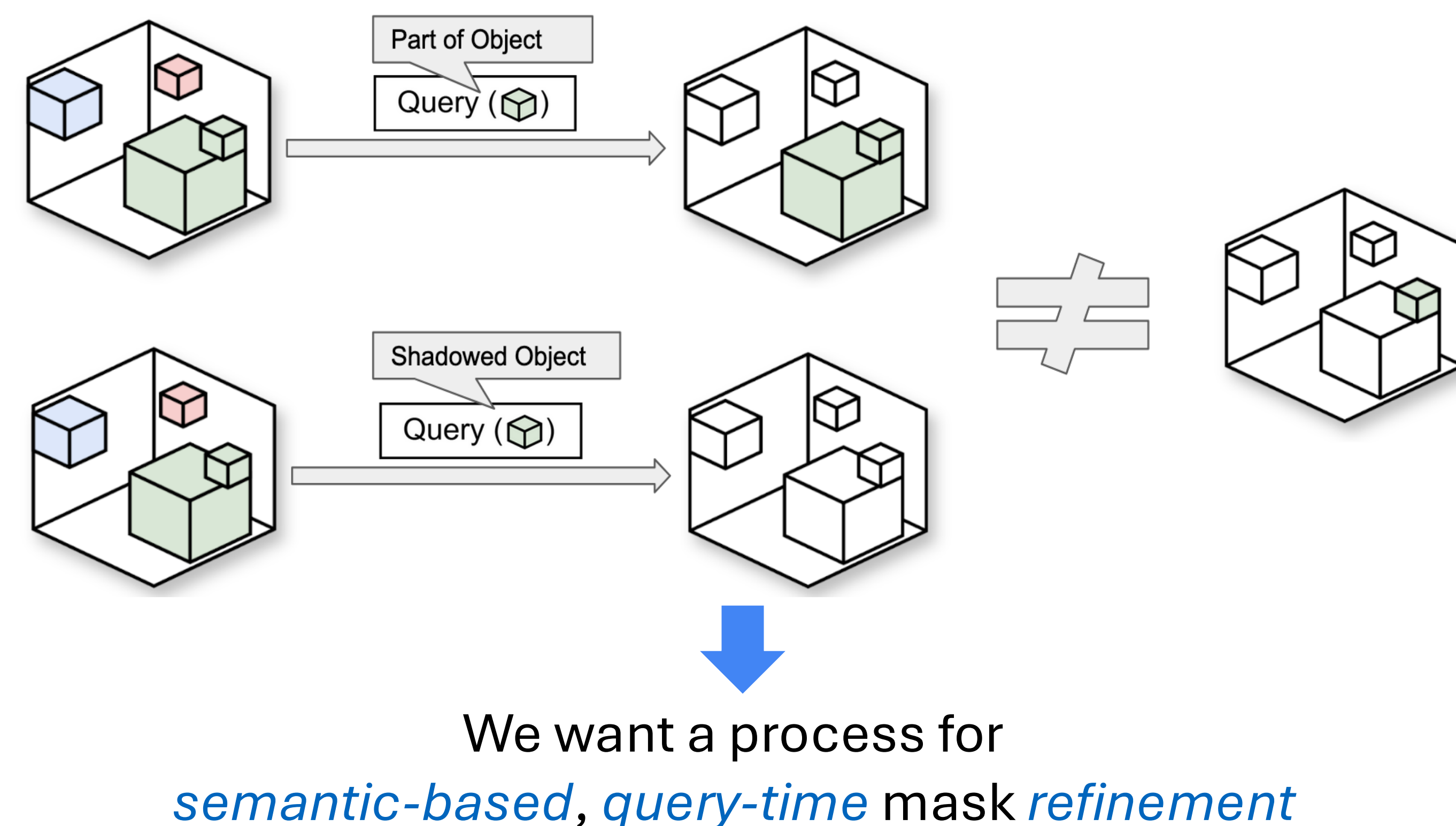
#### Previous Open-Voc 3D instance segmentation



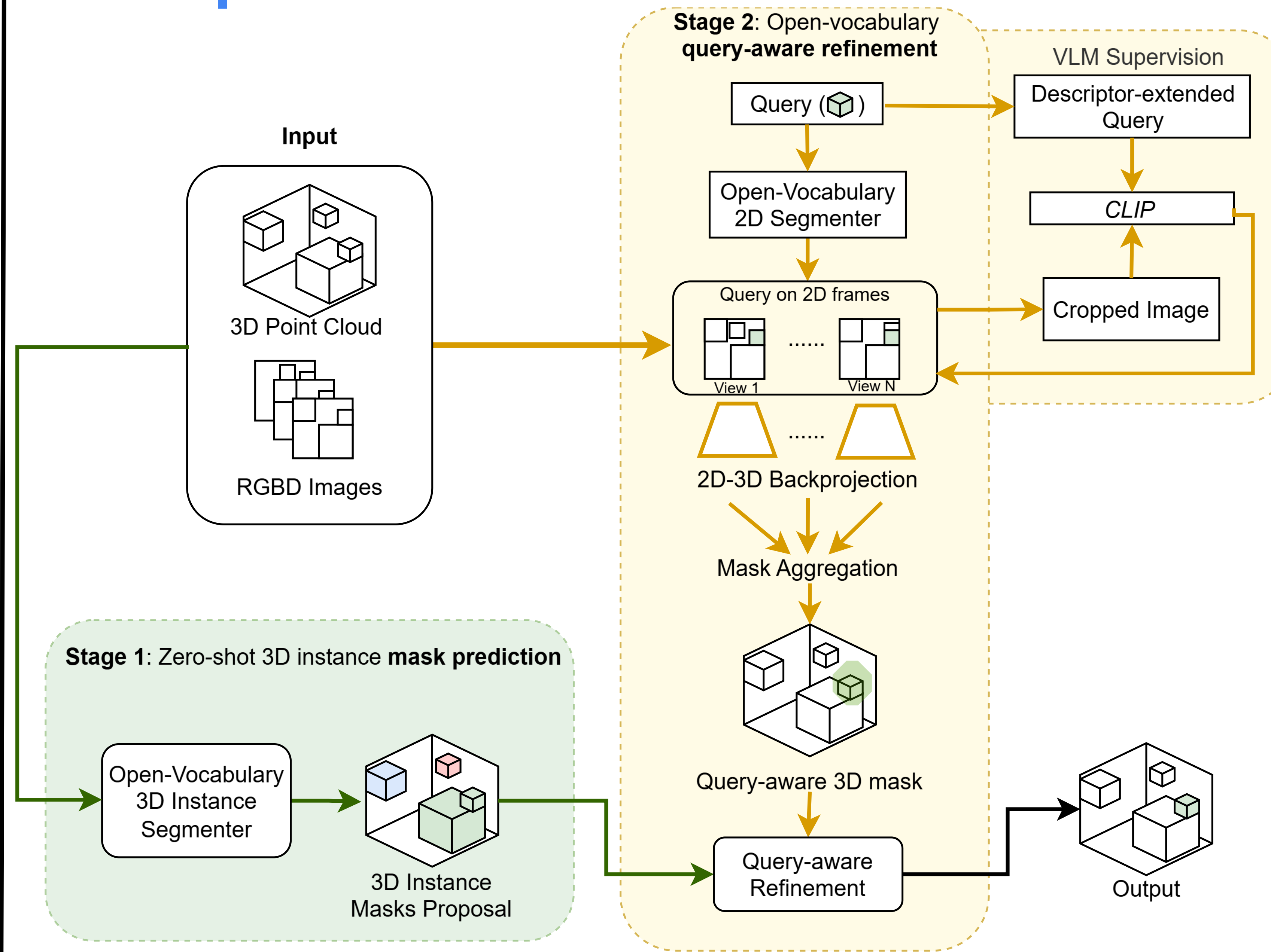
Text queries are not used to their **full extent**: Current methods treat the query phase as **purely a retrieval process**

### Motivation:

Want to find exactly what you want?  
**NO, YOU CAN'T**



### Our Pipeline:



Even if the model shadows small objects, there's still a **second chance** for it to be rediscovered utilizing the semantic information entailed by the text prompts.

### Quantitative Results:

Scannet200 classes*	Methods	AP	AP50	AP25
Overall	Open3DIS	23.7	28.2	31.2
	Open3DIS + BeyondFF	<b>27.4</b>	<b>33.3</b>	<b>39.6</b>
Head	Open3DIS	27.0	32.5	35.5
	Open3DIS + BeyondFF	<b>29.4</b>	<b>36.5</b>	<b>42.0</b>
Common	Open3DIS	21.3	24.9	26.5
	Open3DIS + BeyondFF	<b>26.9</b>	<b>32.0</b>	<b>38.6</b>
Tail	Open3DIS	22.4	26.6	31.1
	Open3DIS + BeyondFF	<b>25.4</b>	<b>30.8</b>	<b>37.8</b>
Base	Open3DIS	23.6	28.7	32.4
	Open3DIS + BeyondFF	<b>27.8</b>	<b>34.6</b>	<b>43.2</b>
Novel	Open3DIS	23.7	28.0	30.8
	Open3DIS + BeyondFF	<b>27.2</b>	<b>32.8</b>	<b>38.3</b>

\* Results tested on 120 classes, distributed as follows: Head:40, Com-mon:40, Tail:40, Base:90, Novel:30

### Qualitative Results:

