

Defining Diversity

We group 3DVG prompts by their *targets/anchors, attributes, relationships*, and other language structures.

Find the **food storage** which does not have a **green, rectangular** object **on top** of it.

Imperative Target Negation Color Shape Anchor Vert. Coref.

Target Reference

- Type (generic, coarse-grained, fine-grained)
- Coreferences
- Noun phrase position

Attribute Understanding

- Counts (total, target, anchor)
- Types

Relationship Understanding

- Counts (total, target, anchor)
- Types

Anchor Reference

- Type (single, multiple, non-object, agent-based)

Diversity Statistics

- Lexical bigrams

Attribute Types

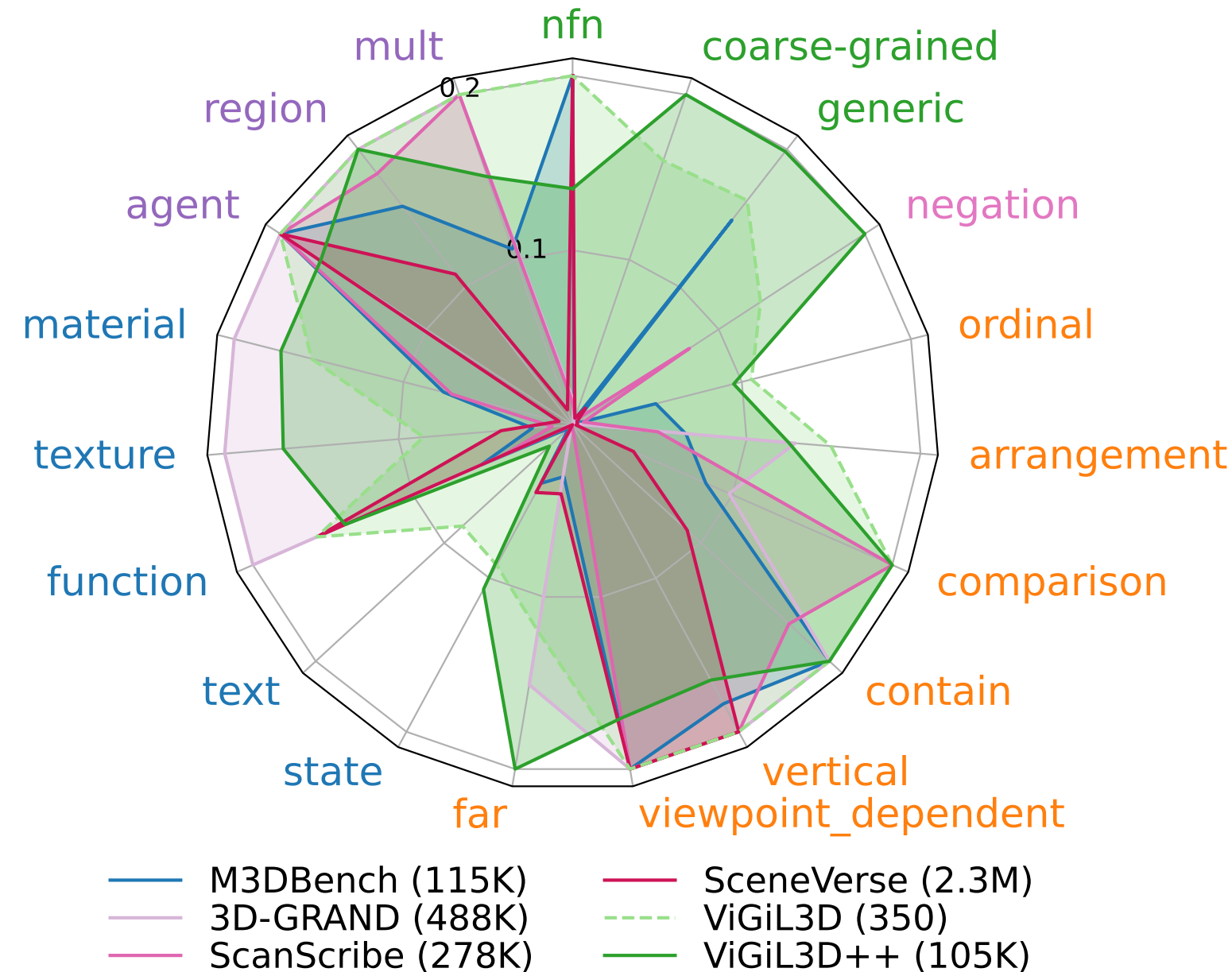
Attribute Types	Relationship Types
Color	Near
Function	Contain
Size	Far
Texture	Arrangement
Shape	Style
Number	Text Label
Material	State
	Directional
	Ordinal
	Vertical
	Comparison

Structural Patterns

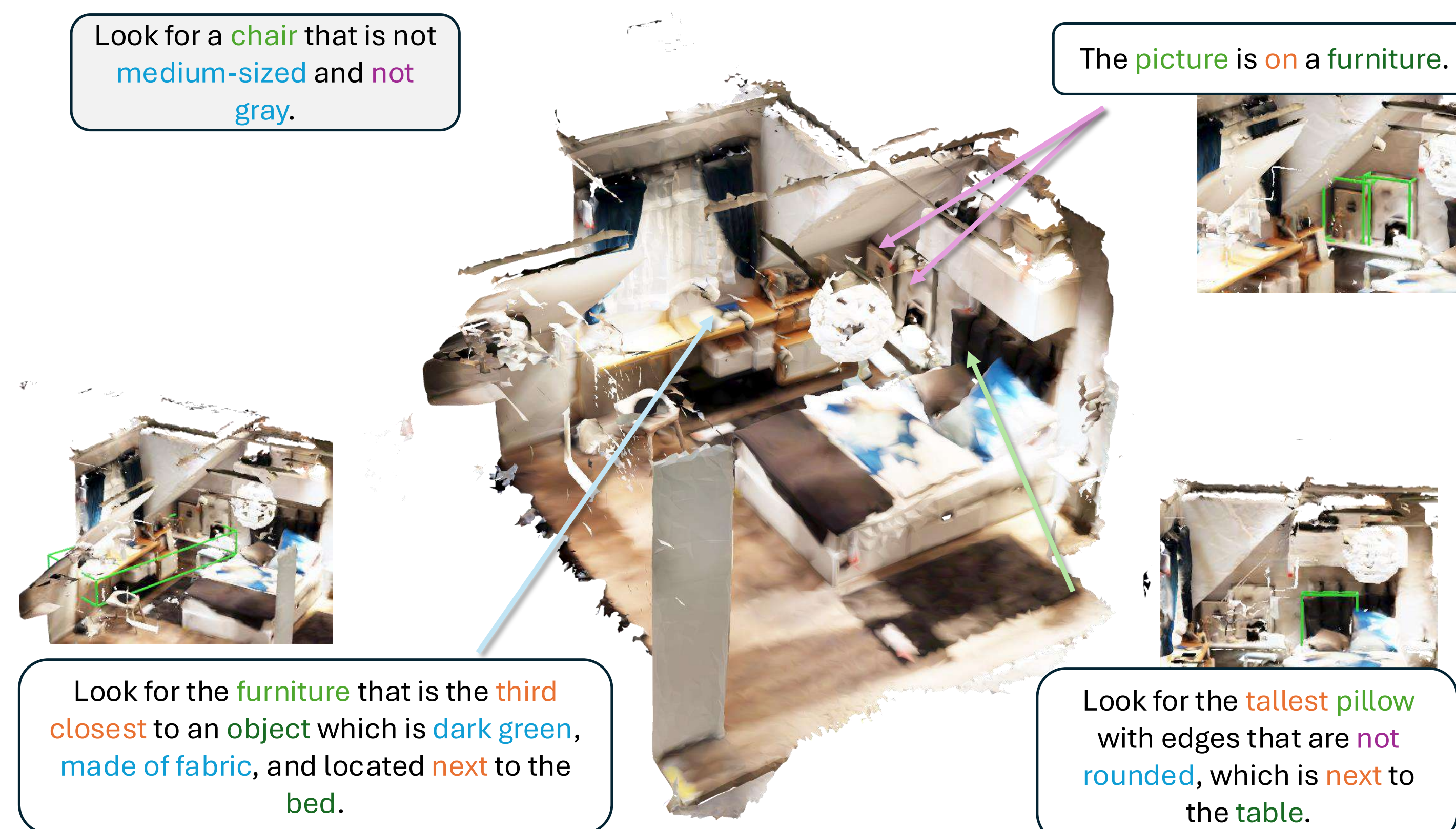
- Sentence structure (e.g. imperative, inverted)
- Logical operators (e.g. negation)

How diverse are existing datasets?

ViGiL3D++ has equal or greater linguistic diversity than existing scaled datasets.

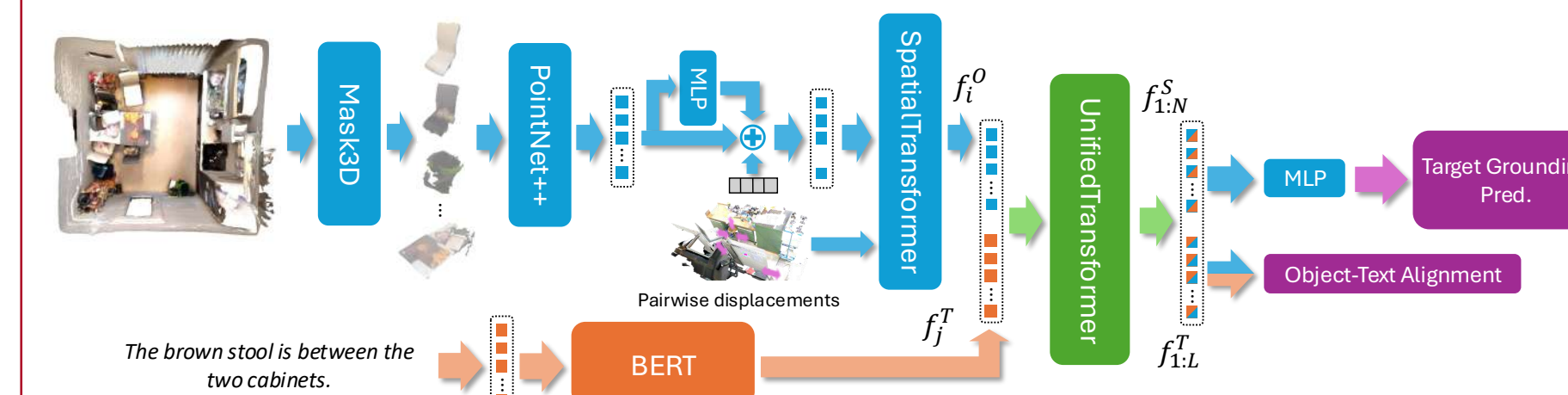


How important is the diversity of the 3DVG descriptions we use for training and evaluation?



V3DM Model

We train a model based on 3D-VisTA [9] to test the utility of ViGiL3D++.



Training on Diverse Data

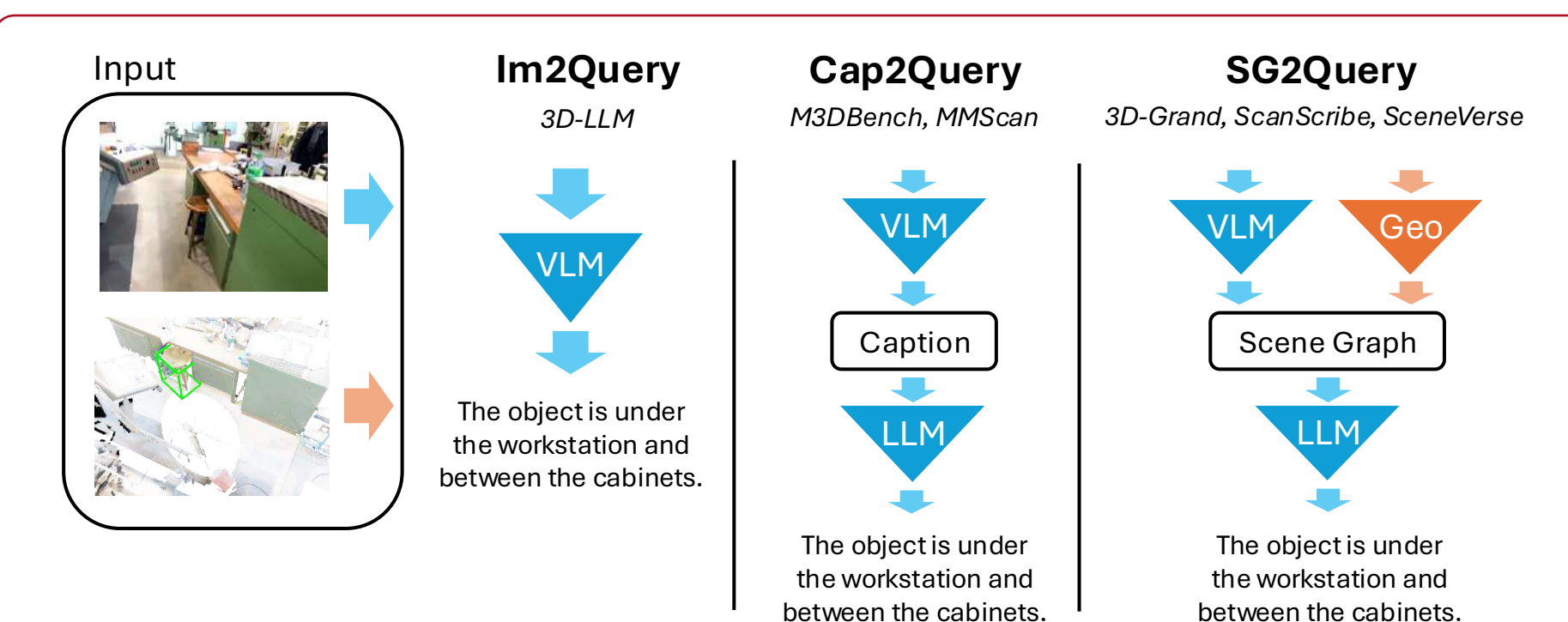
Training on ViGiL3D++ yields better performance on 3DVG benchmarks.

Model	Data		ViGiL3D [6]				Acc on ScanRefer [1]
	Pretrain	Finetune	Acc	F1	Uniq	Multi	Overall
ZSVG3D [8]	—	—	18.9	12.2	73.9	24.9	34.4
3D-VisTA [9]	ScanScribe	ScanRefer	15.1	15.1	84.6	48.2	55.2
3D-GRAND [7]	3D-GRAND	ScanRefer	16.0	9.2	73.3	33.0	40.8
GPS [3]	SceneVerse	ScanRefer	24.1	14.7	88.1	46.4	54.5
V3DM	ScanScribe	ViGiL3D++ _{SR}	<u>22.9</u>	15.7	88.7	<u>47.6</u>	55.6

Insights

- Training and evaluating on diverse language descriptions is important.
- VLMs do not reliably reason about spatial constraints (preserving consistency and generating with diversity).
- Encoding intuition about spatial relationships into concrete scene graphs is challenging due to the contextual information required.

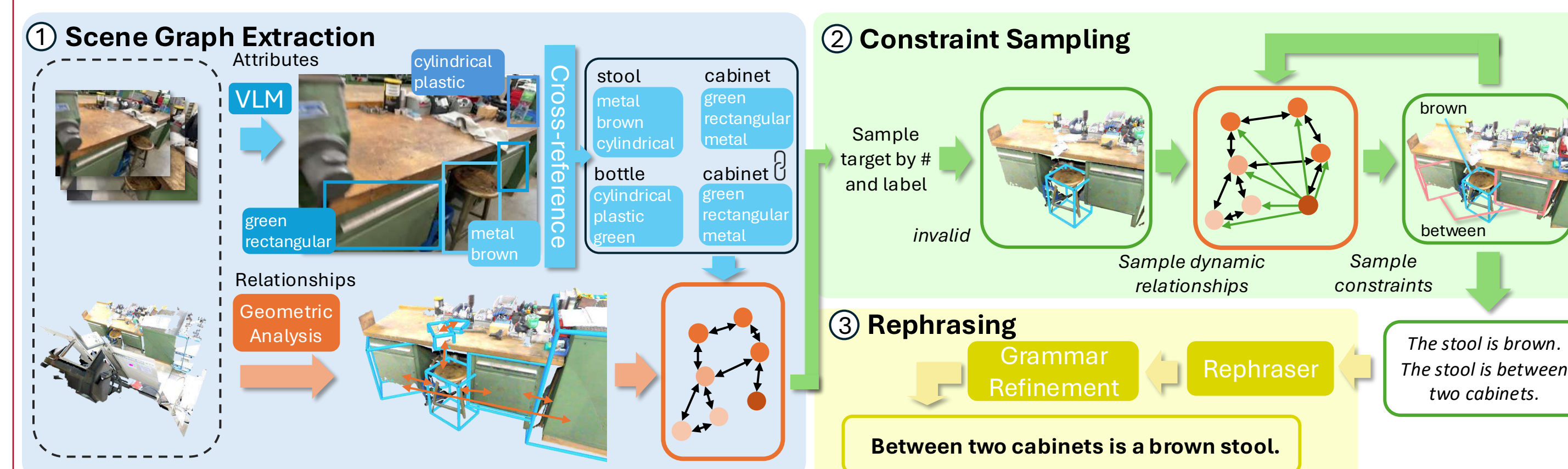
Baselines



Can VLMs be used for 3DVG query generation?

ViGiL3D++:

- Use *cross-referencing* to extract a more consistent scene graph
- Use *constraint sampling* for diversity and logical consistency in queries



Acknowledgements: This work was funded in part by a CIFAR AI Chair and the NSERC Discovery Grant, and enabled by support from the Digital Research Alliance of Canada and a CFI/BCKDF JELF.

References:

- D. Z. Chen et al. ScanRefer. In *ECCV*, 2020.
- Y. Hong et al. 3D-LLM. *NeurIPS*, 2023.
- B. Jia et al. SceneVerse. In *ECCV*, 2024.
- M. Li et al. M3DBench. In *ECCV*, 2024.
- R. Lyu et al. MMScan. *NeurIPS*, 2024.
- A. T. Wang et al. ViGiL3D. In *ACL*, 2025.
- J. Yang et al. 3D-GRAND. In *CVPR*, 2025.
- Z. Yuan et al. ZSVG3D. In *CVPR*, 2024.
- Z. Zhu et al. 3D-VisTA. In *ICCV*, 2023.