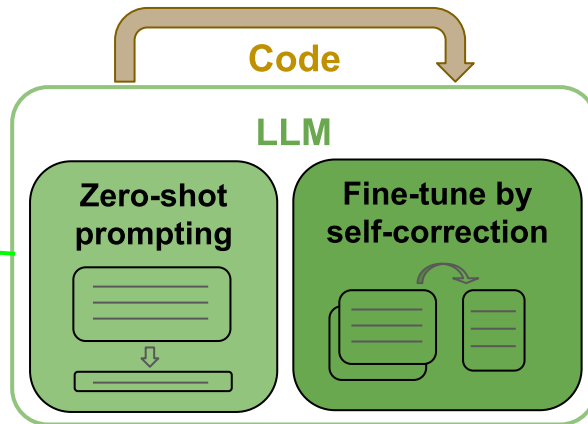


Answer
ID/bbox

3D Object Detector

Spatial-Semantic info
of *all* objects



Code

LLM

Zero-shot
prompting

Fine-tune by
self-correction

Reason

Spatial-Semantic info
of *relevant* objects

Object Filter

"Hey Baymax, can you
help grab me *the chair
in the corner of the
room, between the
white and yellow
desks?*"

