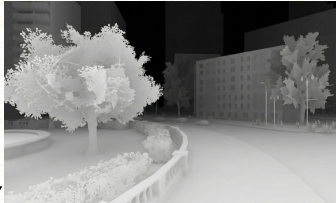


Semantic understanding



Instruction
Fly to the front
of the tree ahead
on the left.

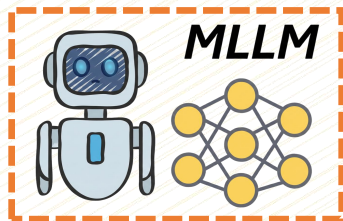
3D coordinate
determination



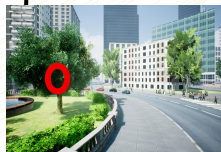
No collision
trajectory



Input



Output

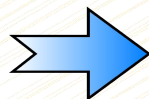
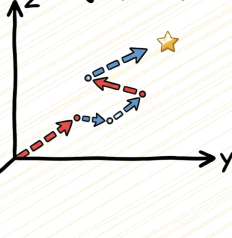


2D Target
Position

Un-projection

(x, y)

(u, v, z)



3D Target
Position

Ego-Planner

