

Input:
-Text
-Image
-Video



(e.g., **LAMP**, **RL-VLM-F**,
iVideoGPT)

LAMP



(e.g., **Eureka**,
Text2Reward)

LLMs



Text

VLMs



Text



**RL
Agent**

State /Action

Diffusion Reward



Diffusion

(e.g., **Diffusion
Reward**, **Diffusion-
QL**)

Input:
-Text



Eureka

Future State
/Image
/Video

WMS/VPMs



(e.g.,
UniSim,
GenRL)



Input:
-Text
-Image
-Sensory
-Trajectory