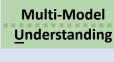
"What is the boy holding in his hand?"

Internet Vision-language Datasets
"Pick up the corn."





"A blue wrist





UP-VLA

Generation

Future Prediction



Robotic Action Datasets



Action

