# GLOMA: **G**rounded **L**ocation for **O**bject **Ma**nipulation

Yifan (Brandon) Yang, Mohammad Yasar, Tariq Iqbal

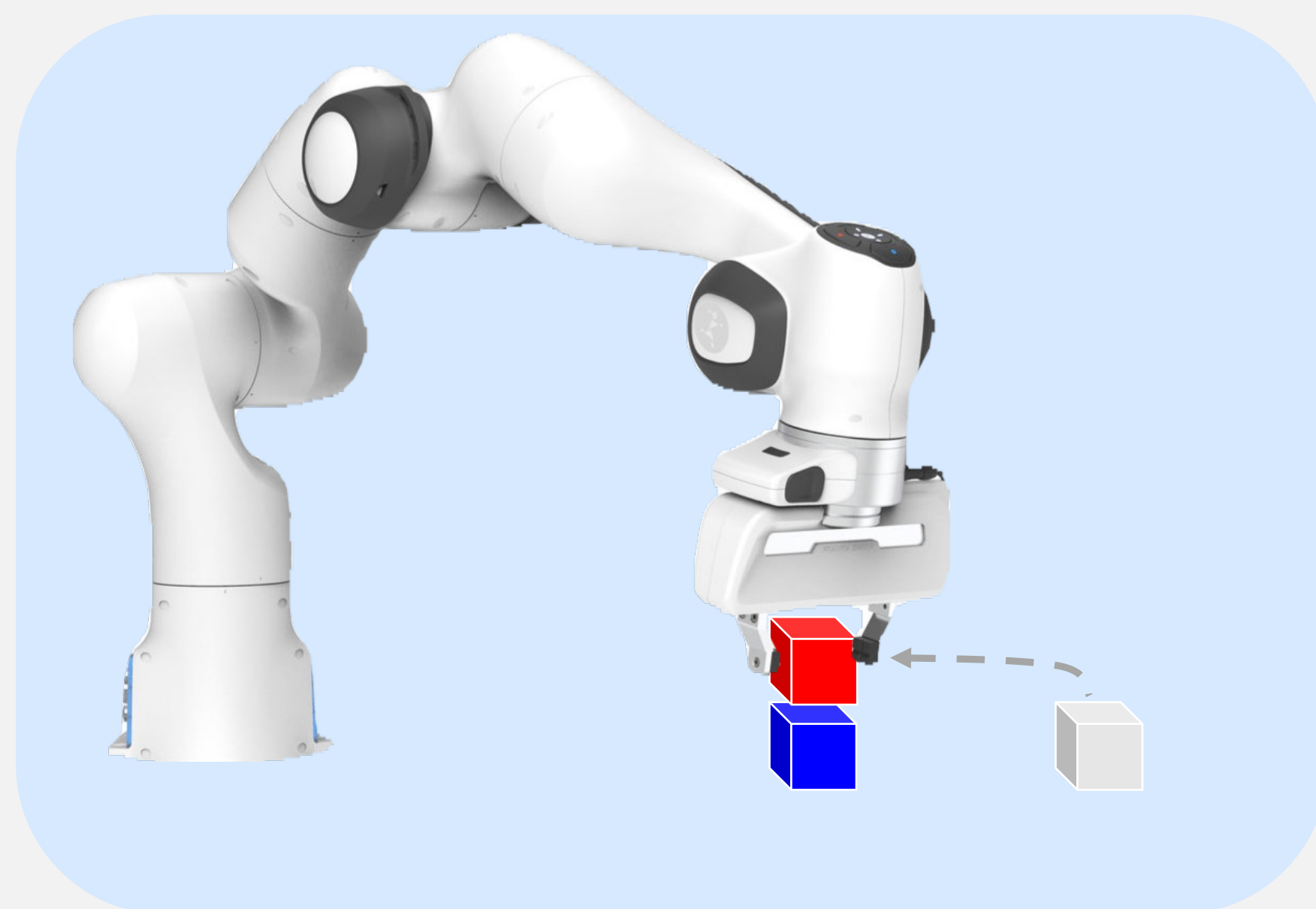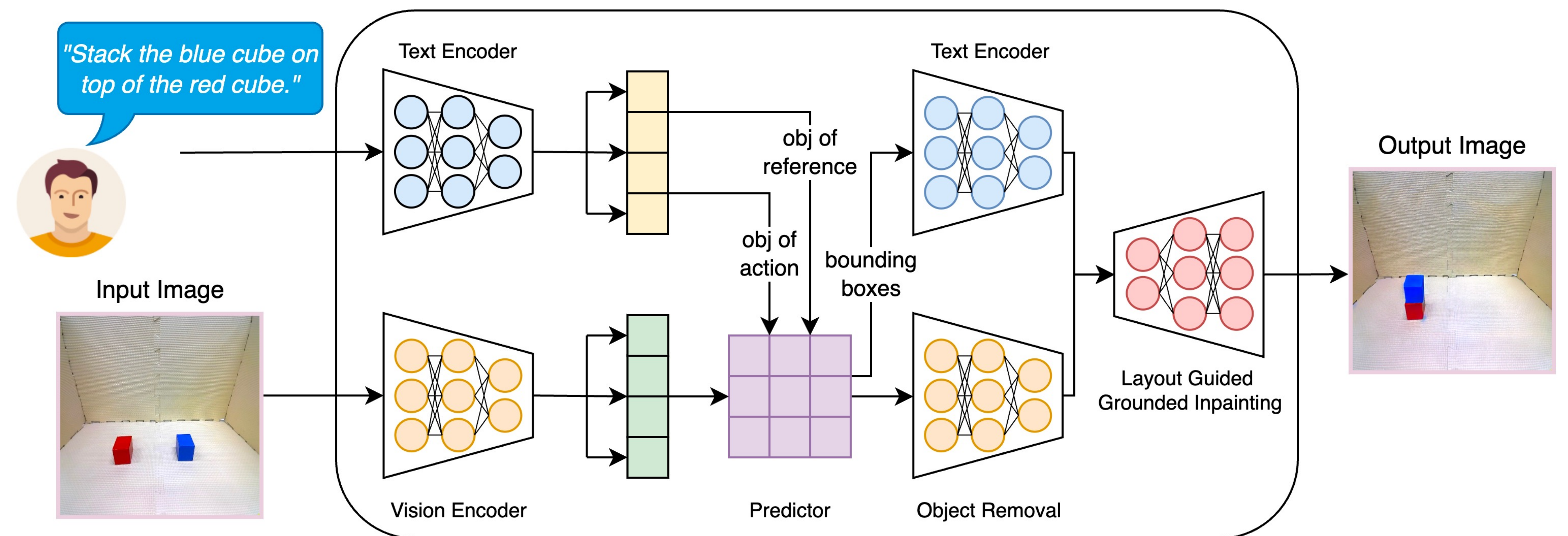Collaborative Robotics Lab

UNIVERSITY of VIRGINIA

## Motivation



- Solving Robotics Tasks with goal-conditioned RL.
- Require goal state represented as an image.

- Existing diffusion models proficient at style transfers, text-to-image prompts.
- Current models unable to achieve real-time object location manipulation.

## Architecture



## Quantitative Results

| Model | Task | SSIM ↑ | VGG16 feature similarity ↑ | FID Score ↓ |
|-------|------|--------|----------------------------|-------------|
| LEDITS | Stacking | 0.74 | 0.48 | |
| | Moving | 0.76 | 0.59 | 299.085 |
| | Other Objects | 0.73 | 0.26 | |
| ControlNet | Stacking | 0.29 | 0.27 | |
| | Moving | 0.31 | 0.29 | 426.548 |
| | Other Objects | 0.37 | 0.22 | |
| InstructPix2Pix | Stacking | 0.84 | 0.60 | |
| | Moving | 0.80 | 0.49 | 308.686 |
| | Other Objects | 0.76 | 0.13 | |
| **GLOMA (Ours)** | Stacking | **0.86** | **0.78** | |
| | Moving | **0.86** | **0.75** | **190.119** |
| | Other Objects | **0.83** | **0.48** | |

## Qualitative Results