# LEGO: Learning EGOcentric Action Frame Generation via Visual Instruction Tuning

Meta  Georgia Tech.  Bolin Lai[1,2]  Xiaoliang Dai[1]  Lawrence Chen[1]  Guan Pang[1]  James M. Rehg[3]  Miao Liu[1]  UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN  ECCV

[1]GenAI, Meta    [2]Georgia Tech    [3]UIUC

## Motivation

**User Prompt**

*"Can you provide instructions on how to **wash the trouser with the brush** in my current situation?"*

**Actions in Existing Egocentric Datasets**

start kettle, take out pip part, wipe off knife, open the tap, puncture the sandal, sew the glue can, place mug, open the glue can, **rub oil into dough**, **wash trouser with brush**, squeeze ball of dough, shake the spray adhesive, clean a sieve, **trim the flower**, turn the motherboard, clean rice, adjust gas, attach the wire, place spice, turn on microwave, tap excess water, open fridge, **pour the topsoil**, close plastic bag, pull out the thread

*"1. Check the Care Label: Check the care label on your trousers for any specific washing instructions.*

*2. Prep the Area: Find a clean and flat surface to work on, such as a table or countertop.*

*3. Spot Check Stains:"*
⋮

*"1. You should submerge the trouser in the water.*

*2. Use the brush to scrub the trouser, focusing on any stains or areas that may require extra attention.*

*3. Once the trouser is clean, you should rinse it."*
⋮

**(a) LLM Response**    **(b) Visual LLM Response**    **(c) Our model (LEGO) Response**

When a user asks for instructions on a task:

- *LLM* -- the answer is too generic and verbose, which is hard to follow.
- *Visual LLM* -- she still faces the challenge of parsing a written description.
- *LEGO (our method)* -- generates an image that provides visual guidance exactly in her situation from the egocentric viewpoint.
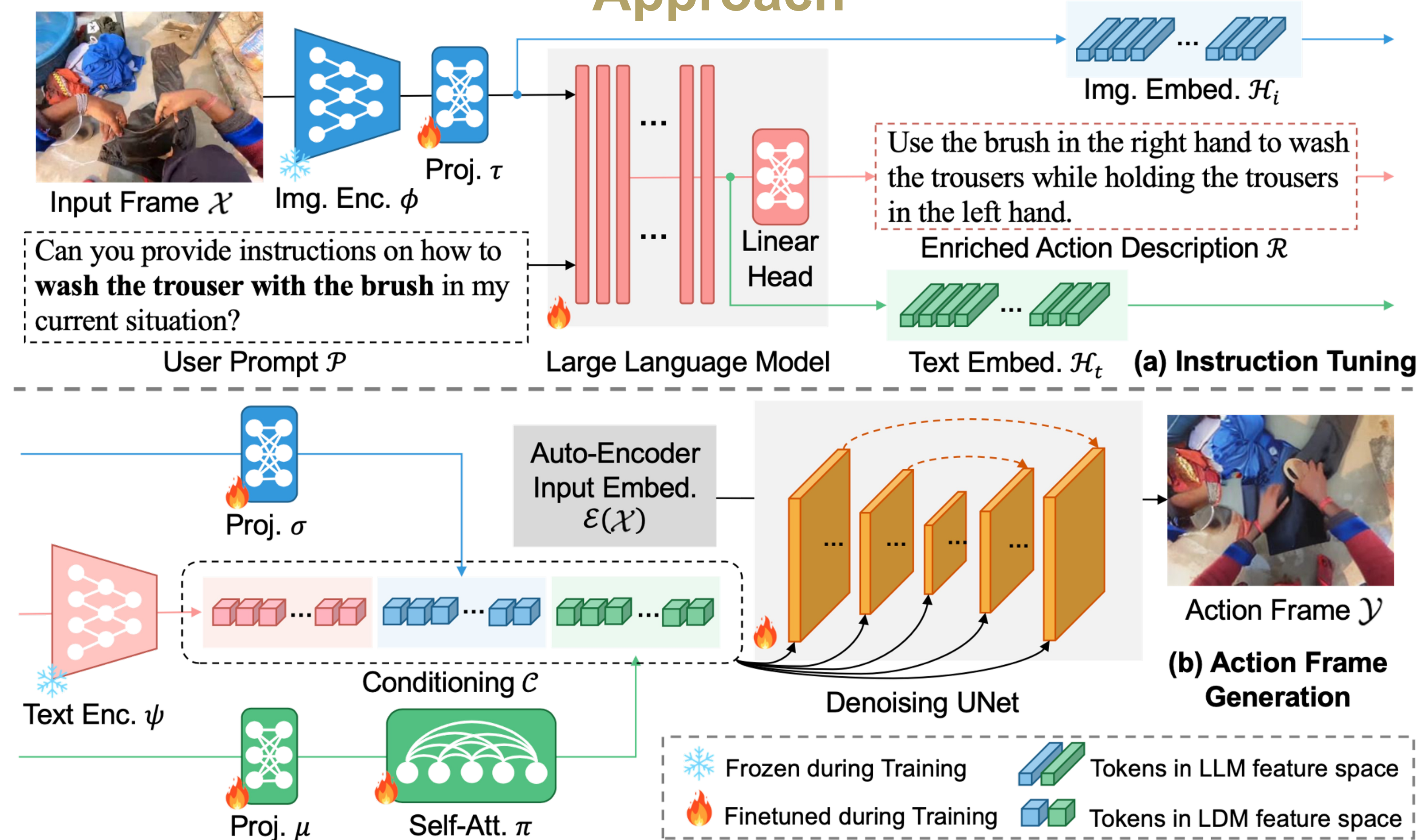
We thus propose a new task -- Egocentric Action Frame Generation,

**Input:** (1) User query of how to perform an action, (2) An image of current situation before an action happens. **Output:** An image in which the action is being performed.

## Challenges

- Action labels are short of necessary details for action frame generation.
- The off-the-shelf diffusion models are limited in action understanding due to domain gap.

⟶

- Enriching the action labels with LLM via visual instruction tuning.
- Leveraging finetuned LLM embeddings to improve egocentric action frame generation.

## Approach



Input Frame $\mathcal{X}$ — Img. Enc. $\phi$ — Proj. $\tau$ — Large Language Model — Linear Head — Img. Embed. $\mathcal{H}_i$

Use the brush in the right hand to wash the trousers while holding the trousers in the left hand.

Enriched Action Description $\mathcal{R}$

Can you provide instructions on how to **wash the trouser with the brush** in my current situation?

User Prompt $\mathcal{P}$    Text Embed. $\mathcal{H}_t$    **(a) Instruction Tuning**

Proj. $\sigma$ — Auto-Encoder Input Embed. $\mathcal{E}(\mathcal{X})$ — Denoising UNet — Action Frame $\mathcal{Y}$

Text Enc. $\psi$ — Proj. $\mu$ — Self-Att. $\pi$ — Conditioning $\mathcal{C}$    **(b) Action Frame Generation**

❄ Frozen during Training    Tokens in LLM feature space
🔥 Finetuned during Training    Tokens in LDM feature space

LEGO consists of two key components:

- *Visual Instruction Tuning* -- We finetune an LLM to generate detailed action descriptions which include information such as hands and spatial locations.
- *Action Frame Generation* -- We project image and text features from LLM to LDM space, and input them to a diffusion model as additional conditions to mitigate the domain gap.
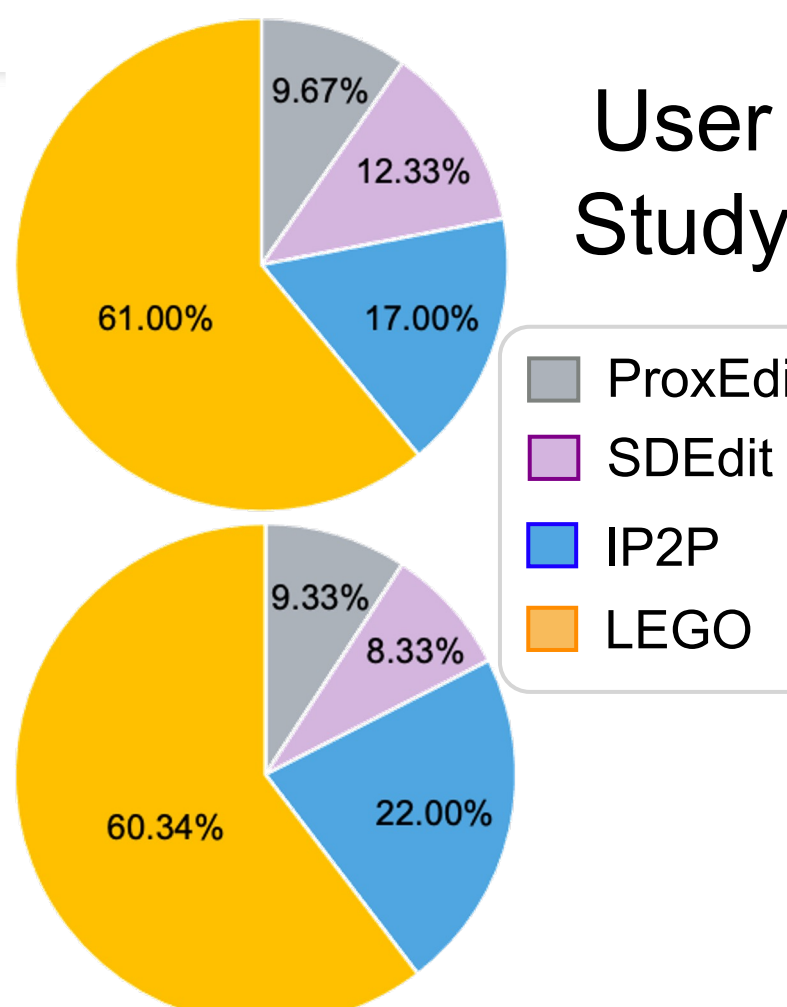
Input Frame
LEGO

*"cut a portion of clay mix with both hands"*  *"brush a wood with a brush"*  *"put tray in oven"*  *"close container"*  *"take soy milk"*

## Experiments and Results

| | Methods | EgoVLP | EgoVLP⁺ | CLIP | FID ↓ | PSNR | LPIPS ↓ |
|---|---|---|---|---|---|---|---|
| Ego4D | ProxEdit [26] | 44.51 | 72.68 | 68.17 | 33.01 | 11.88 | 40.90 |
| | SDEdit [59] | 50.07 | 72.90 | 73.35 | 33.35 | 11.81 | 41.60 |
| | IP2P [6] | 62.19 | 78.84 | 78.75 | 24.73 | 12.16 | 37.16 |
| | LEGO | 65.65 | 80.44 | 80.61 | 23.83 | 12.29 | 36.43 |
| E-Kitchens | ProxEdit [26] | 32.27 | 52.77 | 69.18 | 51.35 | 11.06 | 46.35 |
| | SDEdit [59] | 33.84 | 56.80 | 74.76 | 27.41 | 11.30 | 43.33 |
| | IP2P [6] | 42.97 | 61.06 | 77.03 | 20.64 | 11.23 | 40.82 |
| | LEGO | 45.89 | 62.66 | 78.63 | 21.57 | 11.33 | 40.36 |

**User Study**

9.67%, 12.33%, 17.00%, 61.00%
9.33%, 8.33%, 22.00%, 60.34%

ProxEdit  SDEdit  IP2P  LEGO

Input Frame    ProxEdit    SDEdit    InstructPix2Pix    LEGO

*"How to rinse the jacket inside the plastic bath?"*

*"How to take glass?"*

**Generating various actions in the same contexts:**

*"Can you provide instructions on how to {action} in my current situation?"*

*"open drawer"*  *"dry hands"*  *"cut cucumber"*

*"open microwave"*  *"pick up bowl"*  *"take knife"*

LEGO

## Contact