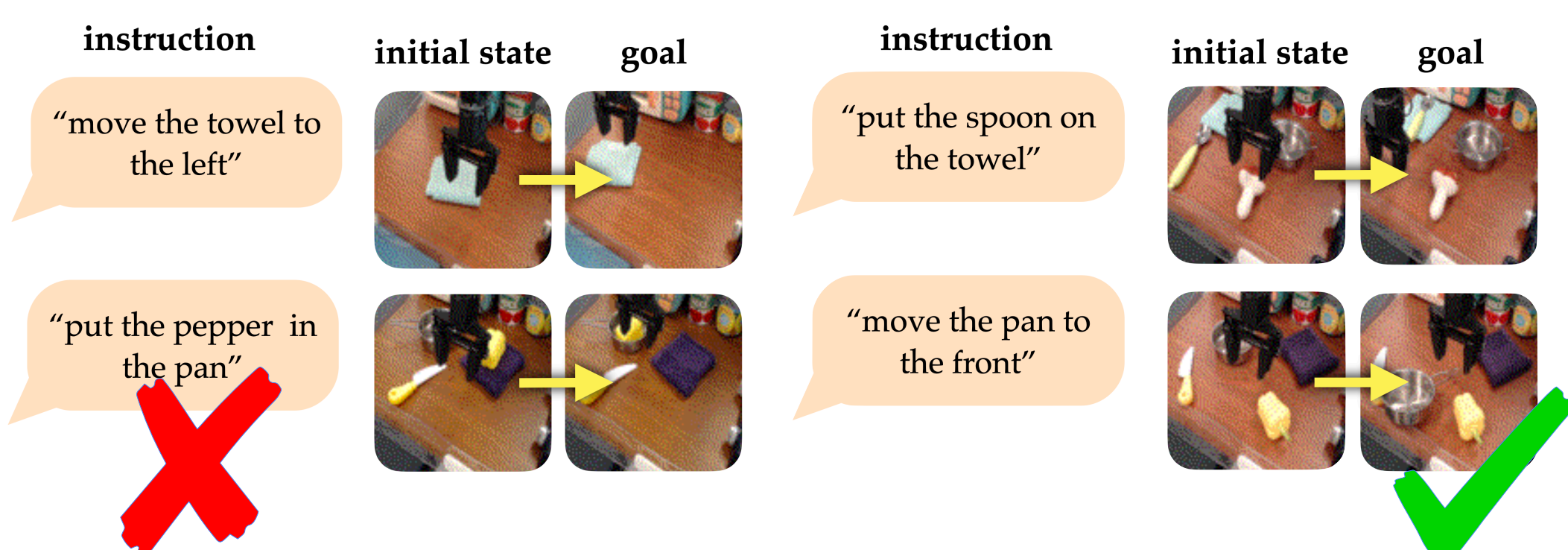


### Motivation

Goal-conditioned policies robustly generalize across tabletop manipulation skills and don't require task annotations. Language is a more natural form of task specification for humans, though it is harder to learn. How can we get the benefits of goal-conditioned learning when following language instructions?

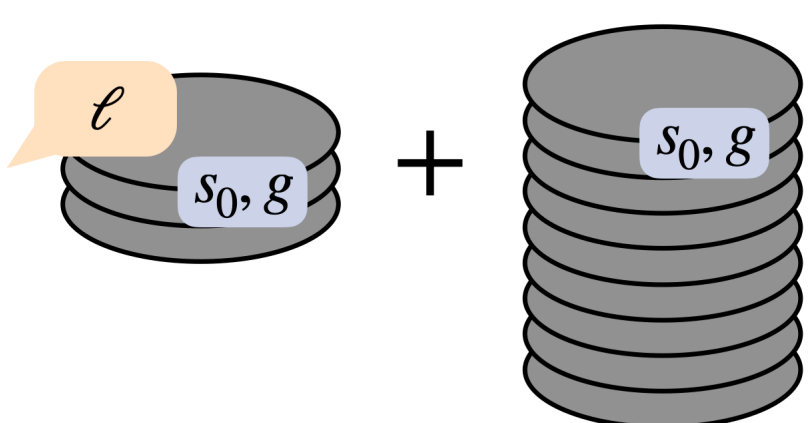
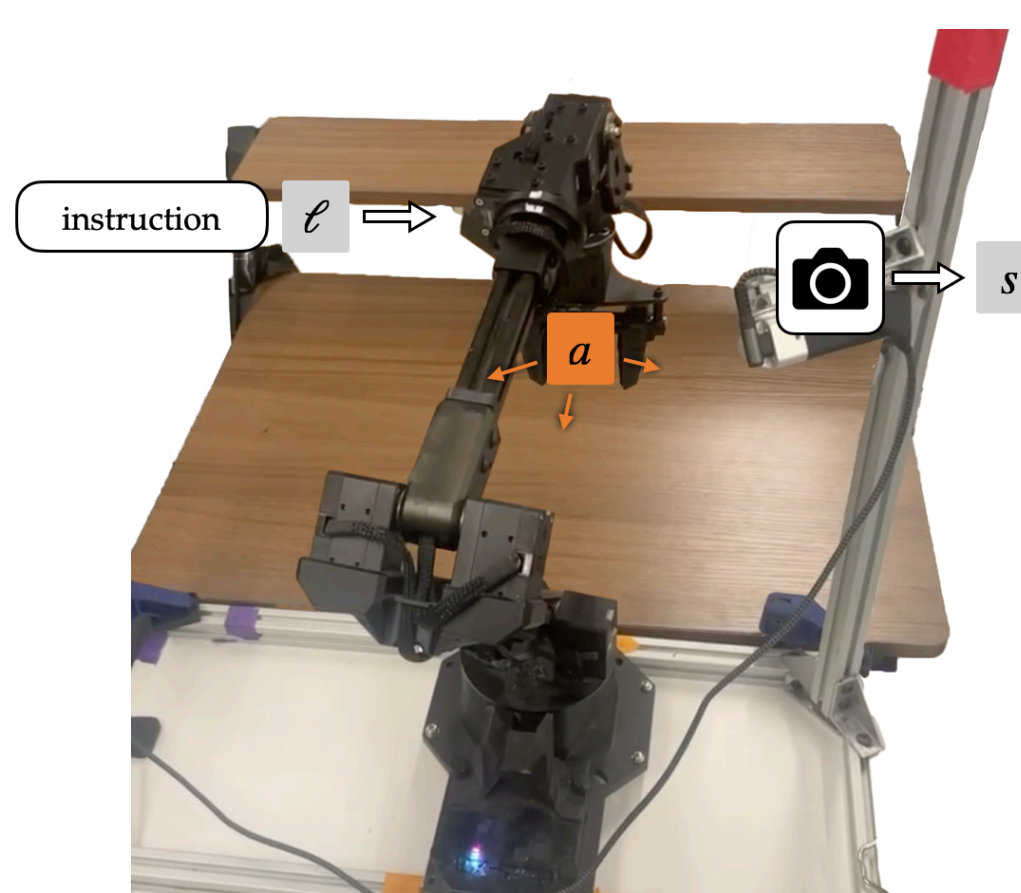


**Language:** less supervision, hard to learn

**Goals:** hindsight relabeling, better generalization

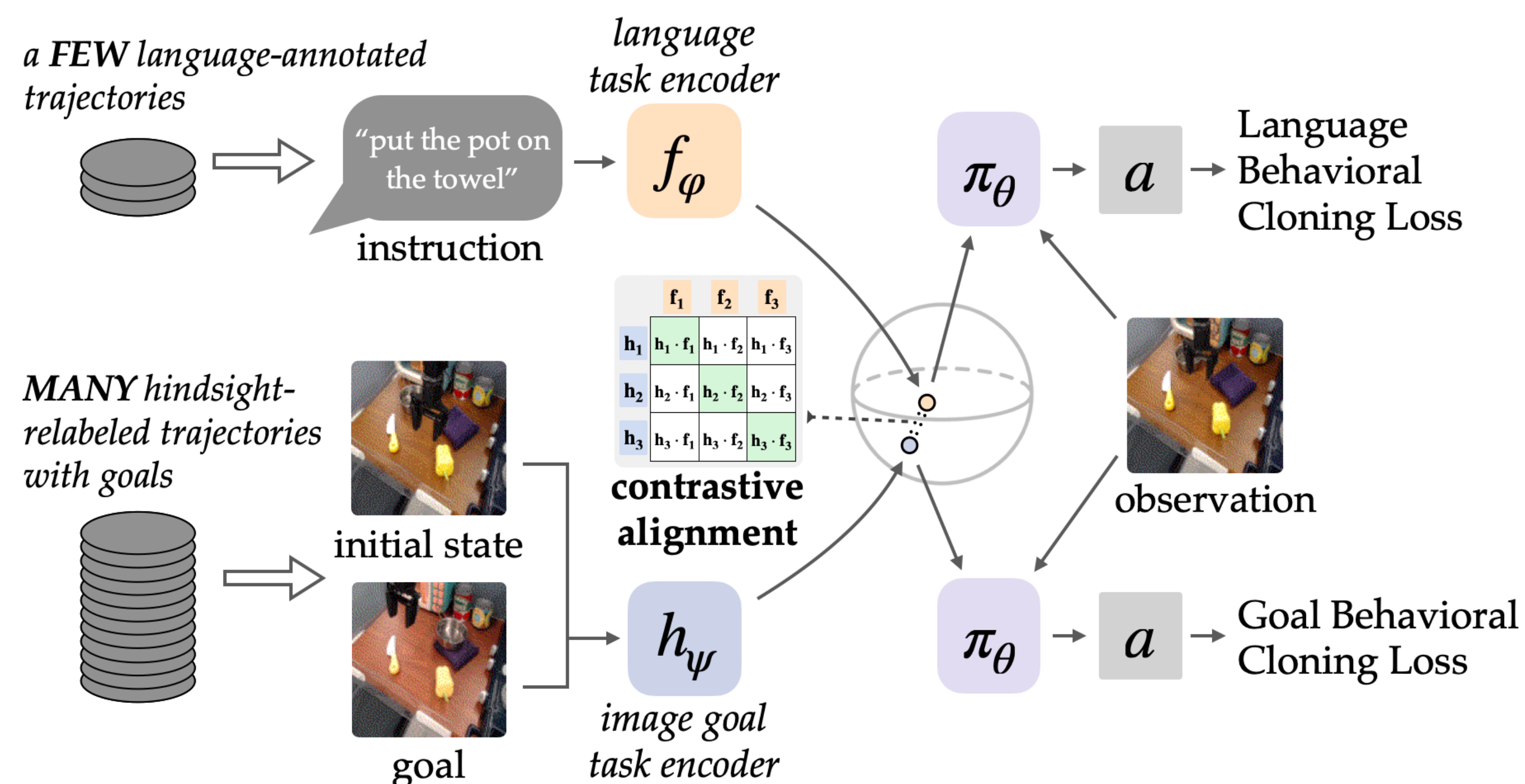
### Setup

We perform 6 DoF control of a WidowX arm for language-conditioned tabletop manipulation.



We have a few language-annotated trajectories along with many **unlabeled** trajectories. We need a **semi-supervised** approach to use both sources of data.

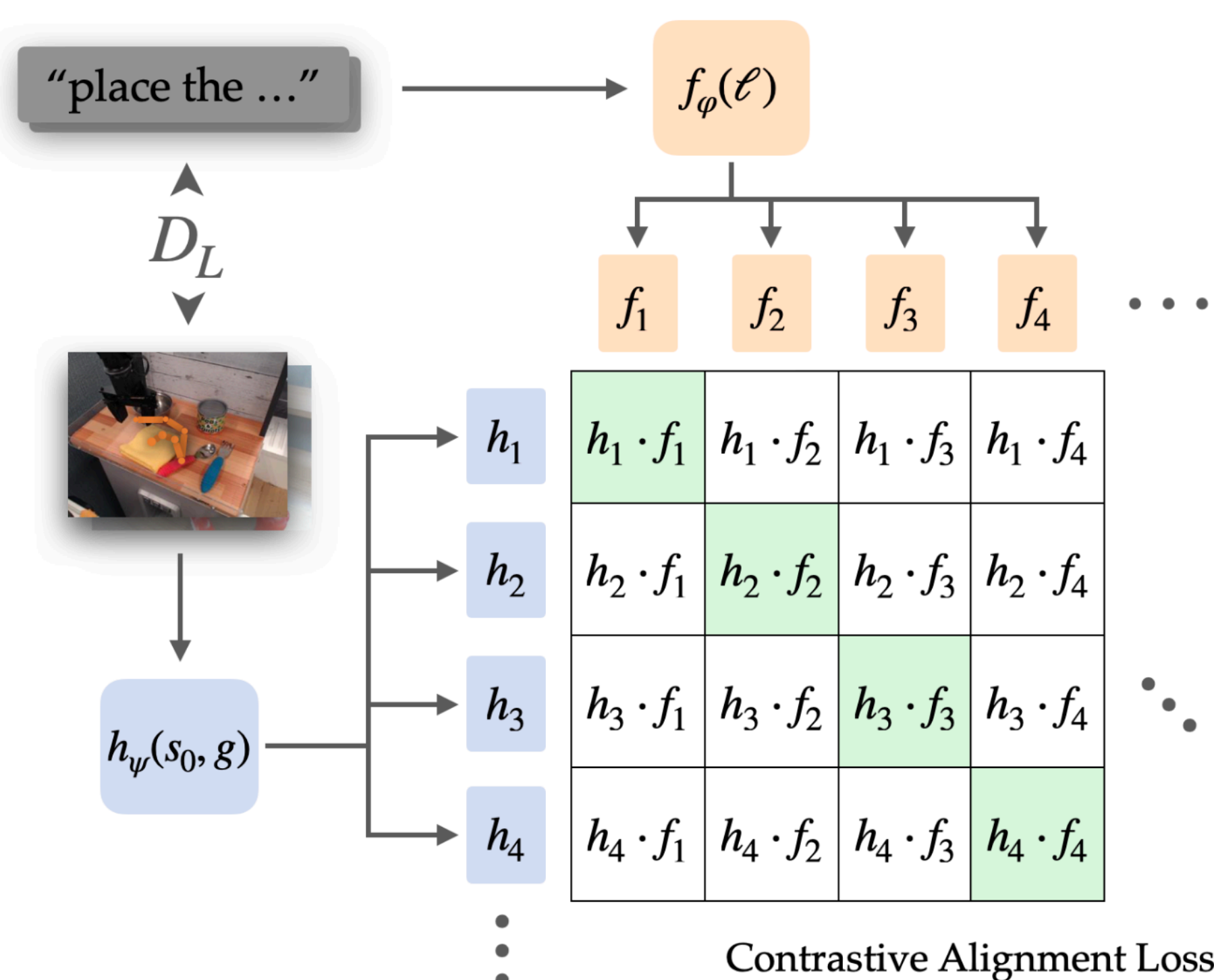
### Approach



### Goal Representations for Instruction Following (GRIF)

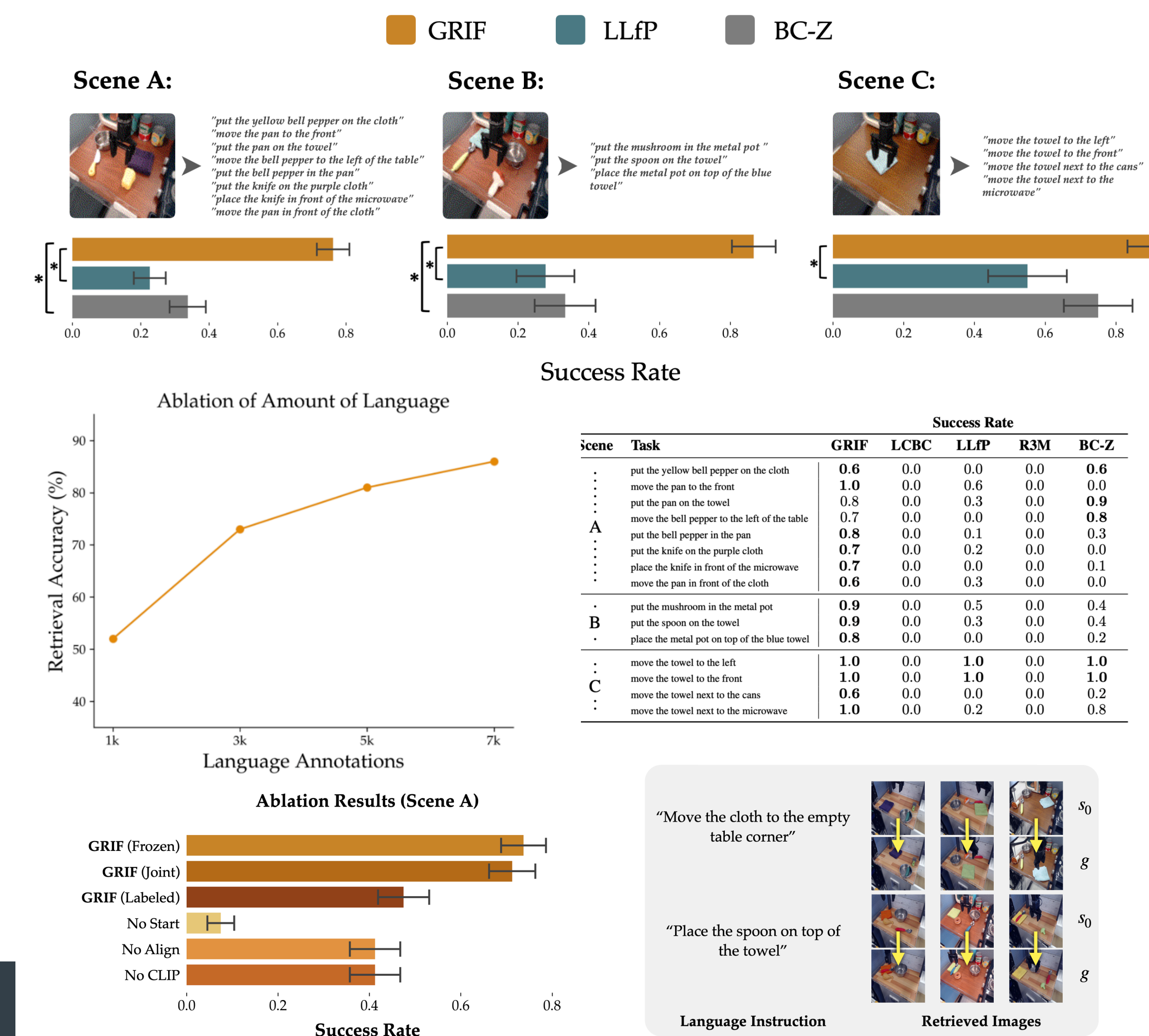
### Aligned Task Representations

Our approach aligns (start state, goal) pairs in the scene with language tasks using a **contrastive loss**. We learn a policy on top of these aligned representations using behavioral cloning.



With **hindsight relabeling**, we can use both unlabeled trajectories and language-annotated trajectories to train our policy. We also modify the **CLIP** architecture to leverage pre-trained language knowledge.

### Results



We outperform baseline and ablation approaches on **diverse unseen language instructions** in different scenes. In particular GRIF is better able to **ground** language instructions compared to baselines while performing well at manipulation.

### Summary

We train policies on **language-annotated** and **goal**-related trajectories

**Contrastive** language-goal task alignment enables robust grounding

Our approach outperforms baselines on **unseen** instructions