

Large-scale Pre-training for Grounded Video Caption Generation

Evangelos Kazakos¹, Cordelia Schmid², Josef Sivic¹

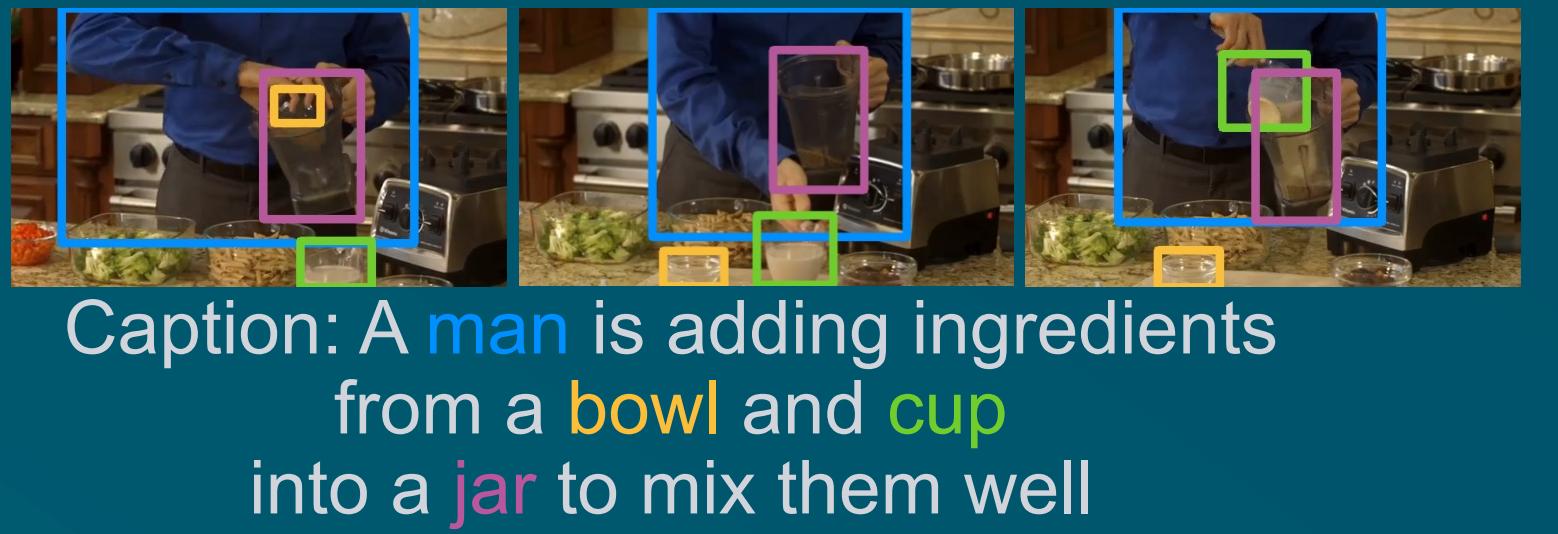
¹Czech Institute of Informatics, Robotics and Cybernetics at the Czech Technical University in Prague

²Inria, École normale supérieure, CNRS, PSL Research University



OVERVIEW

Goal: Simultaneous captioning & grounding of nouns to bounding boxes



Caption: A man is adding ingredients from a **bowl** and **cup** into a **jar** to mix them well

Motivation



Contributions

- Automatic annotation method
- HowToGround1M dataset
- Manually annotated iGround dataset
- The GROVE grounded video caption generation model
- SOTA on 5 grounding datasets

AUTOMATIC ANNOTATION METHOD

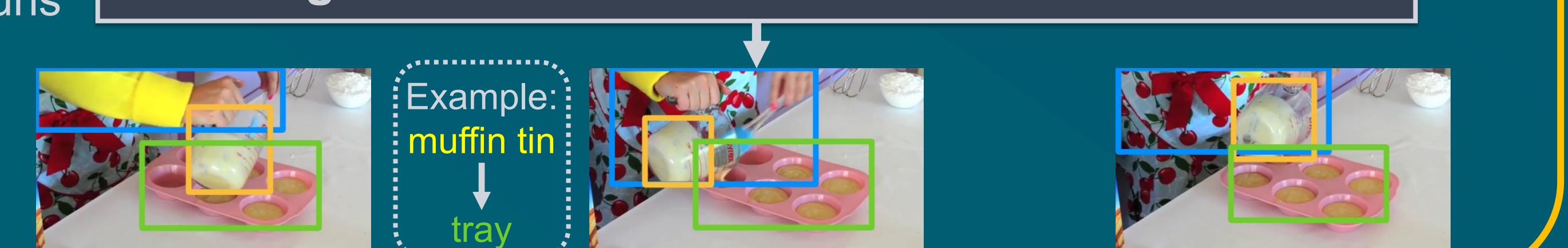
Stage 1: Frame-level grounding with image VLM



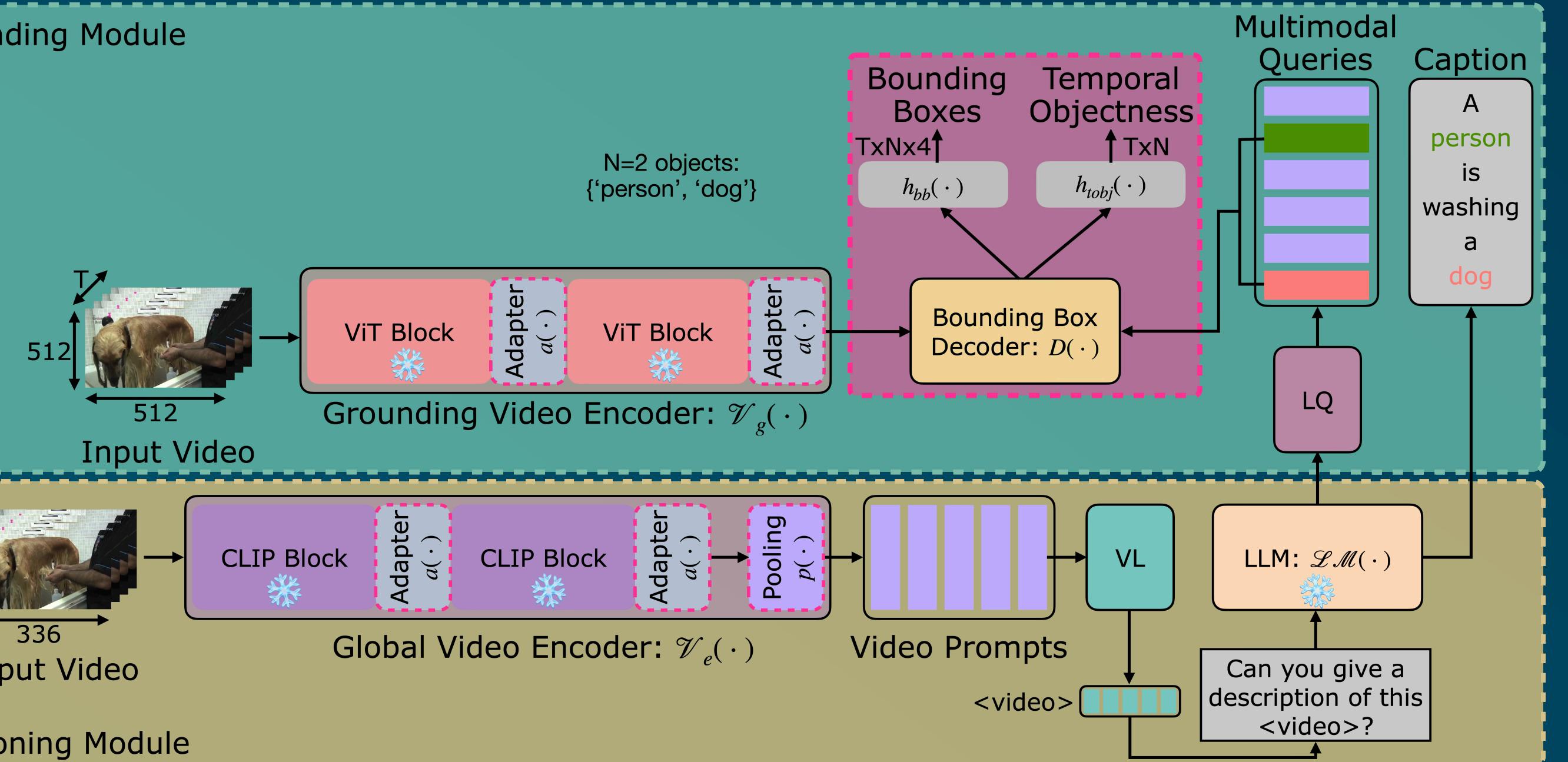
Stage 2: Caption aggregation with LLM



Stage 3: Connect boxes across frames with LLM

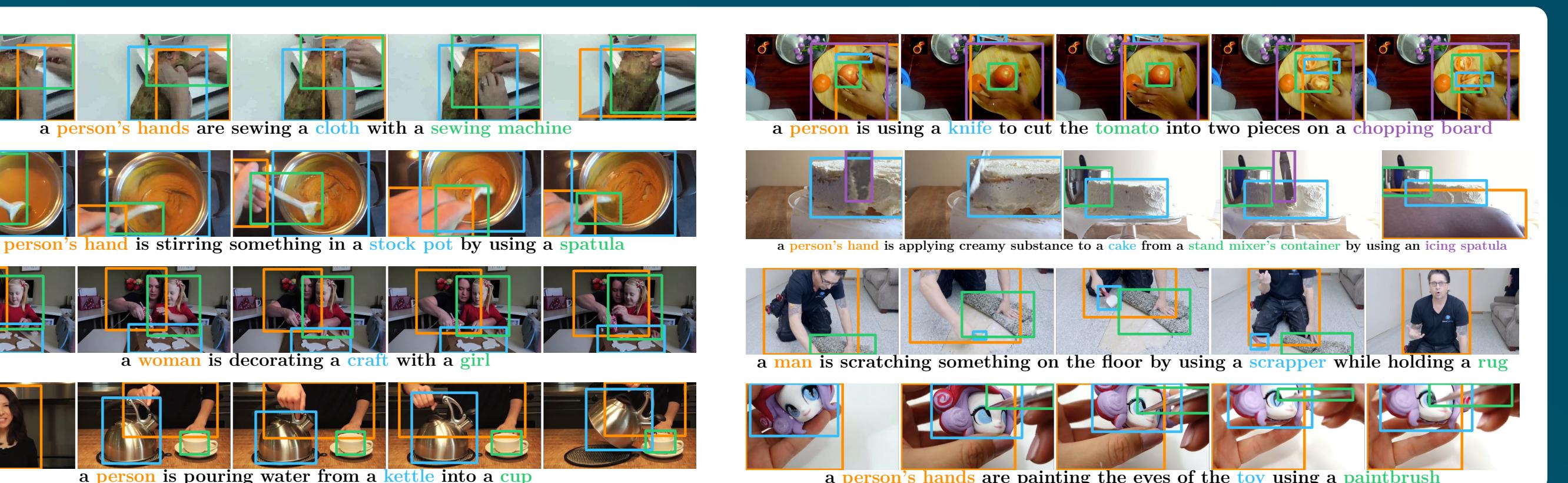


THE GROVE MODEL



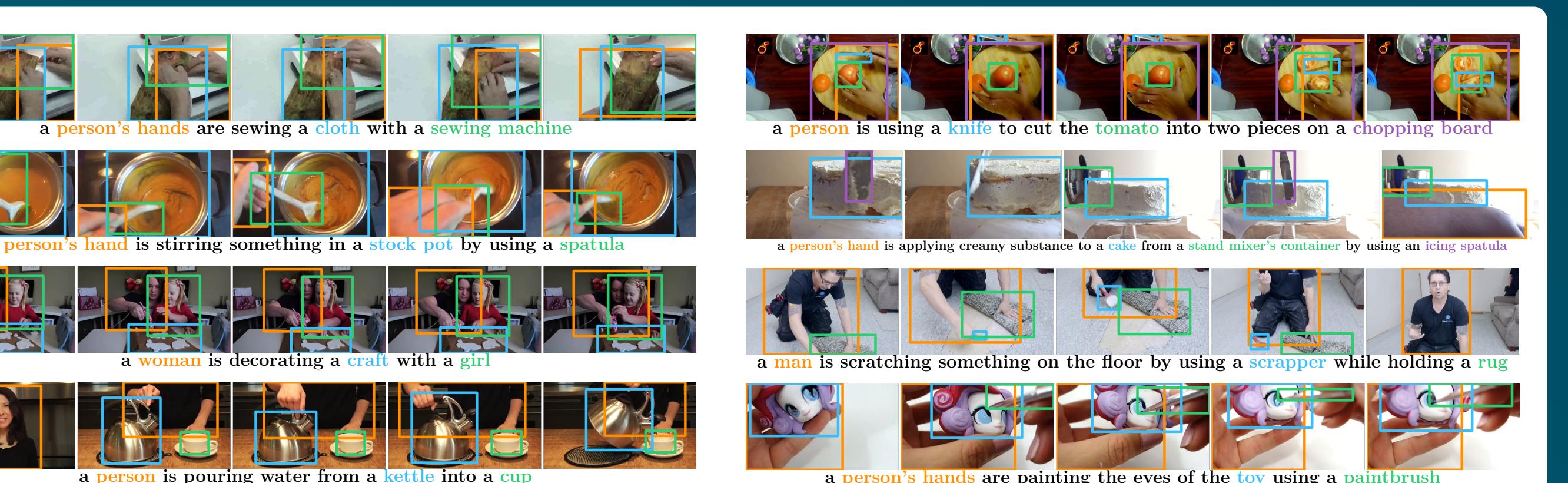
- 3D convolutional adapters for spatio-temporal modeling ($a(\cdot)$)
- Captioning module generates caption and tags noun phrases
- Cross-attention of embedded noun phrases with frame-level features ($D(\cdot)$)
- Temporal objectness head for predicting the presence of an object in a frame ($h_{tobj}(\cdot)$)

HowToGround1M DATASET



- Automatically annotated
- 1M videos
- 80.1M bounding boxes
- Ideal for pre-training

iGround DATASET



- Manually annotated
- 3500 videos
- Train/val/test: 2000/500/1000
- Ideal for fine-tuning and evaluation

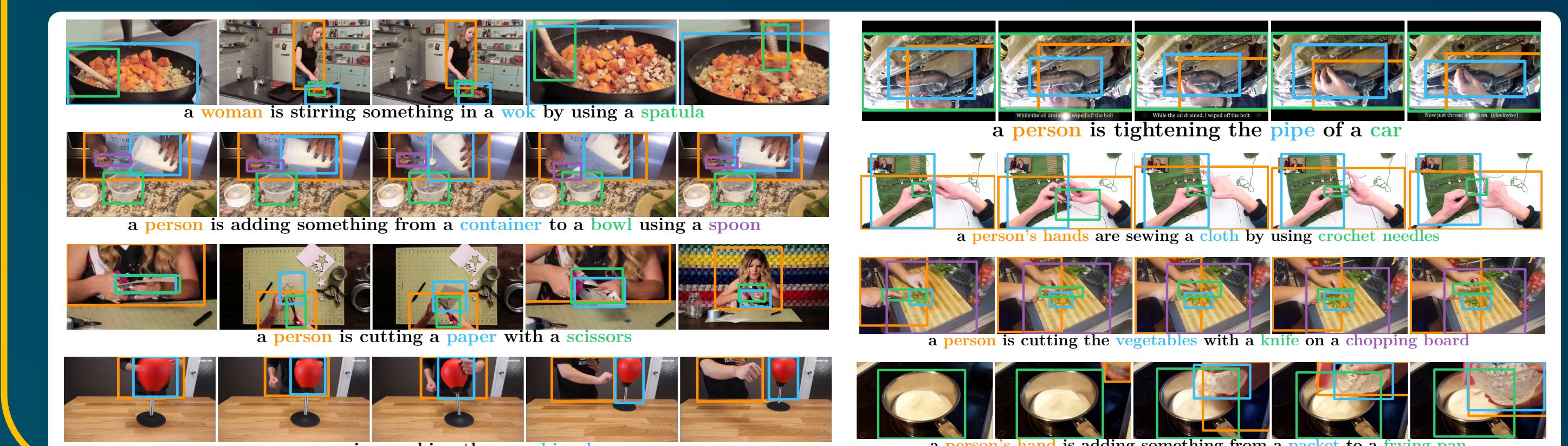
RESULTS

TAKEAWAYS

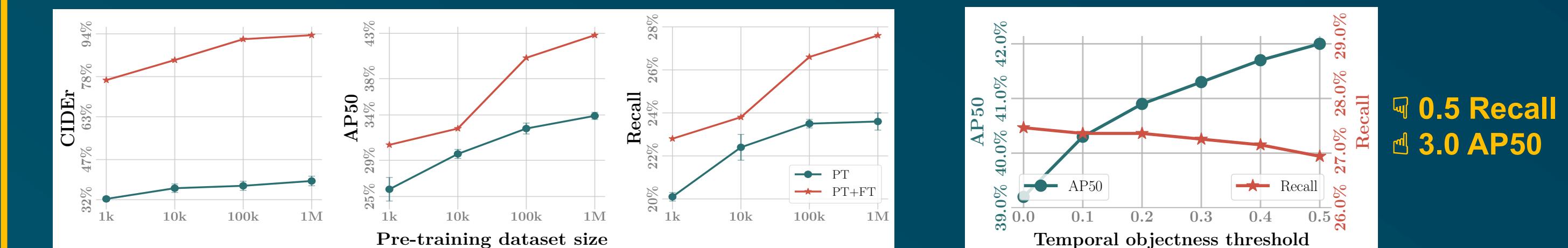
- GROVE >> GLaMM in single-frame setup
- GROVE >> automatic annotation
- Pre-training is crucial for grounding
- Pre-training + fine-tuning 🔥

- SOTA in VidSTG
- Best performance without fine-tuning (interrogative)
- ++ SOTA in ANet-Entities, GroundingYT, YouCook-Inter. (paper)

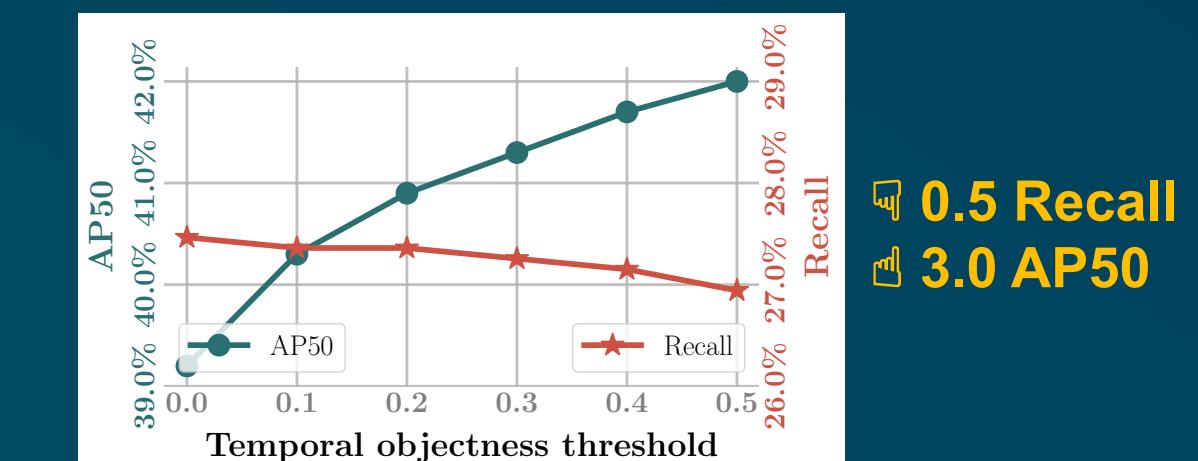
Qualitative results (iGround)



ABLATIONS



- **Scaling:** Performance scales with pre-training dataset size
- Fine-tuning benefits from scaling



- **Temporal objectness:** Little sacrifice in recall for in AP50