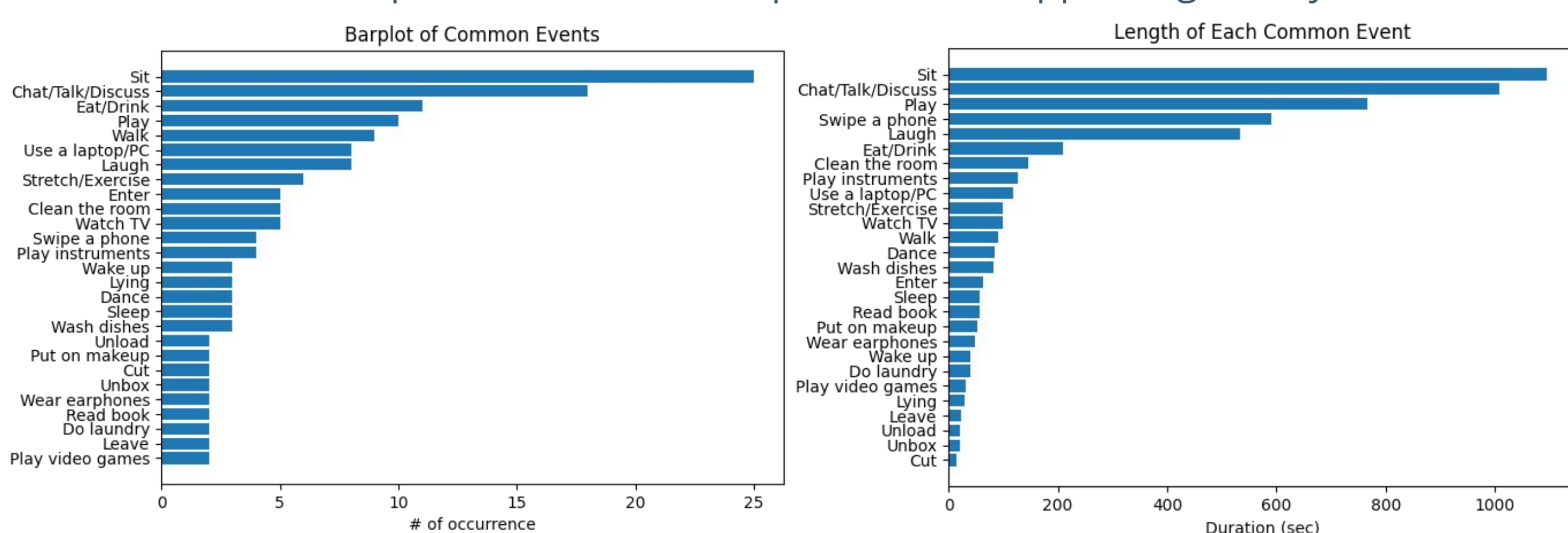


Introduction

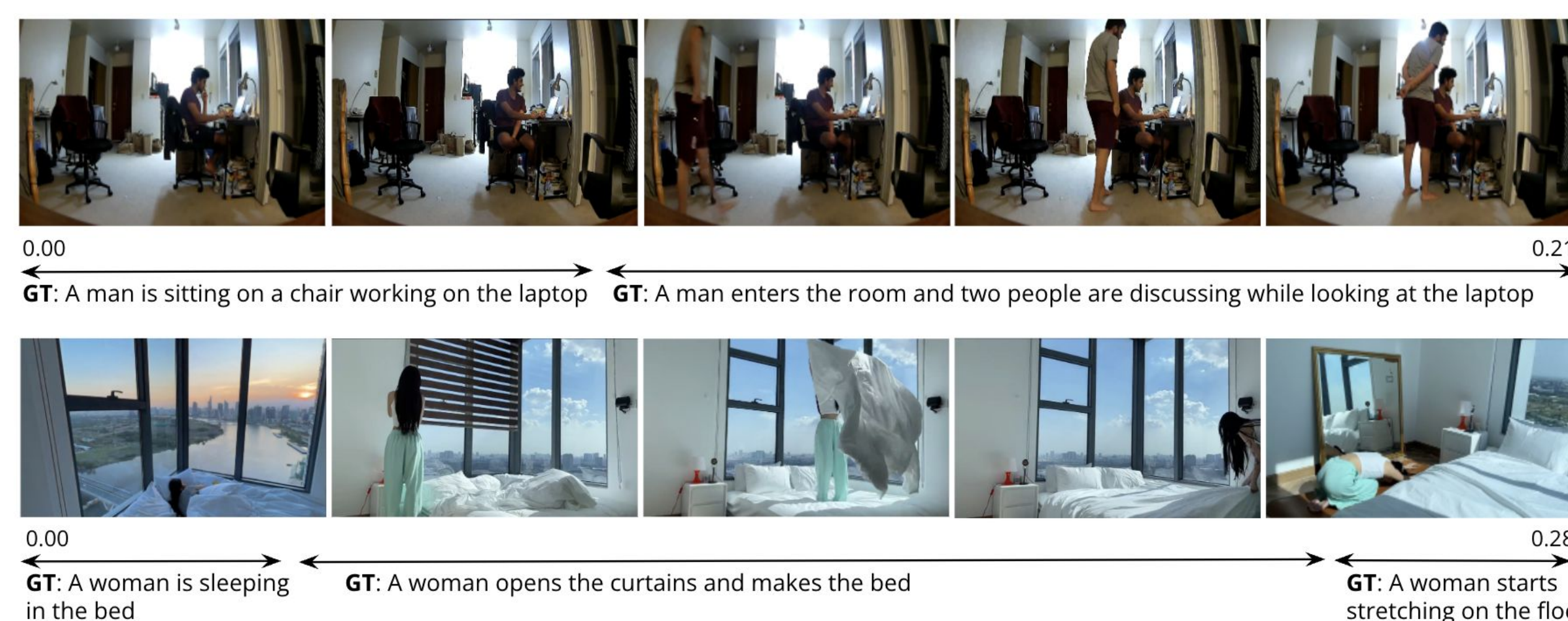
- The team is focusing on exploring technical solutions of video captioning for indoor/outdoor scenarios and super long video (> 1 hour), video content understanding by bridging visual and language information and leveraging video data to generate captions to describe the semantic content in video recordings.
- To generate the captions for the super long video, the team worked to do in-depth research on academia-leading solutions on dense video captioning and improve state-of-the-art method, PDVC model, to generate descriptions on Wyze device captured videos.
- Different from normal video captioning, the task of dense video captioning involves both detecting and describing events in a video; we can describe the super long video well.
- Our model and software developed will be used to create meaningful captions that help users monitor their daily events and identify situations that need their attention.

Dataset

- Our dataset consists of indoor home scenarios:
 - 100 videos within 20~30 sec (25 from Wyze cam + 75 from Internet)
 - 10-min video*1 + 1-hr videos*5 (all from Wyze cam)
 - 27 Most common events
- The total duration of the dataset is 349.25 minutes and the average length is 197.69 seconds. The 27 most frequently occurring events are shown in the barplot. The most common event is "Sit", which appears in 25 videos. Next comes the event "Chat/Talk/Discuss", with 18 occurrences. On the rightmost side of the barplot are the least frequent events, appearing in only 2 videos.

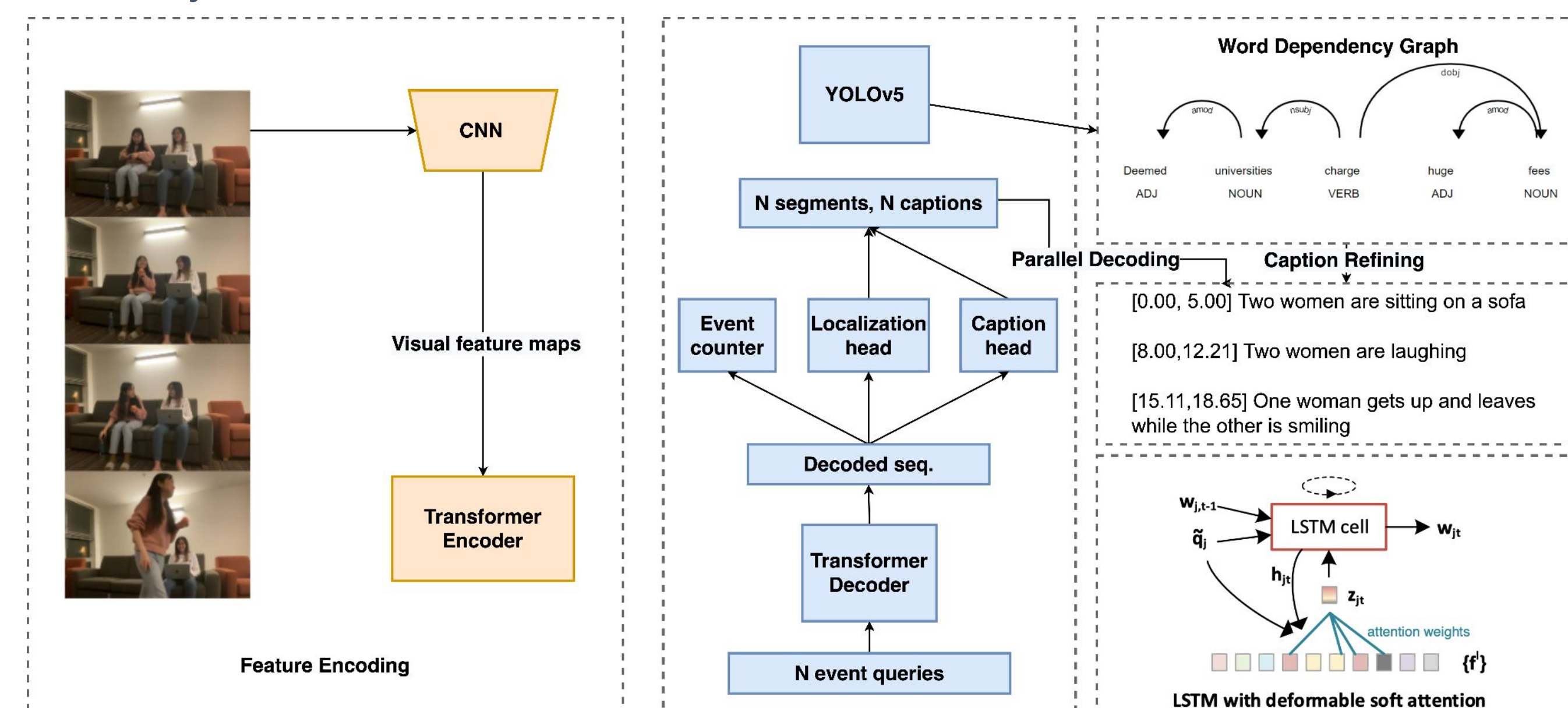


DATASET EXAMPLES:



Enhanced PDVC Model

- Our proposed solution is an end-to-end machine learning model with an Encoder-Decoder type architecture.
- The encoder is a 3D CNN which extracts multi-scale frame features from the pre-trained feature extractor.
- The decoder has three parallel heads to generate captions, predict temporal boundaries for events in a video and count the number of events in the input.
- The generated caption is refined using the YOLO Object Detection Model for obtaining the correct object class in the video

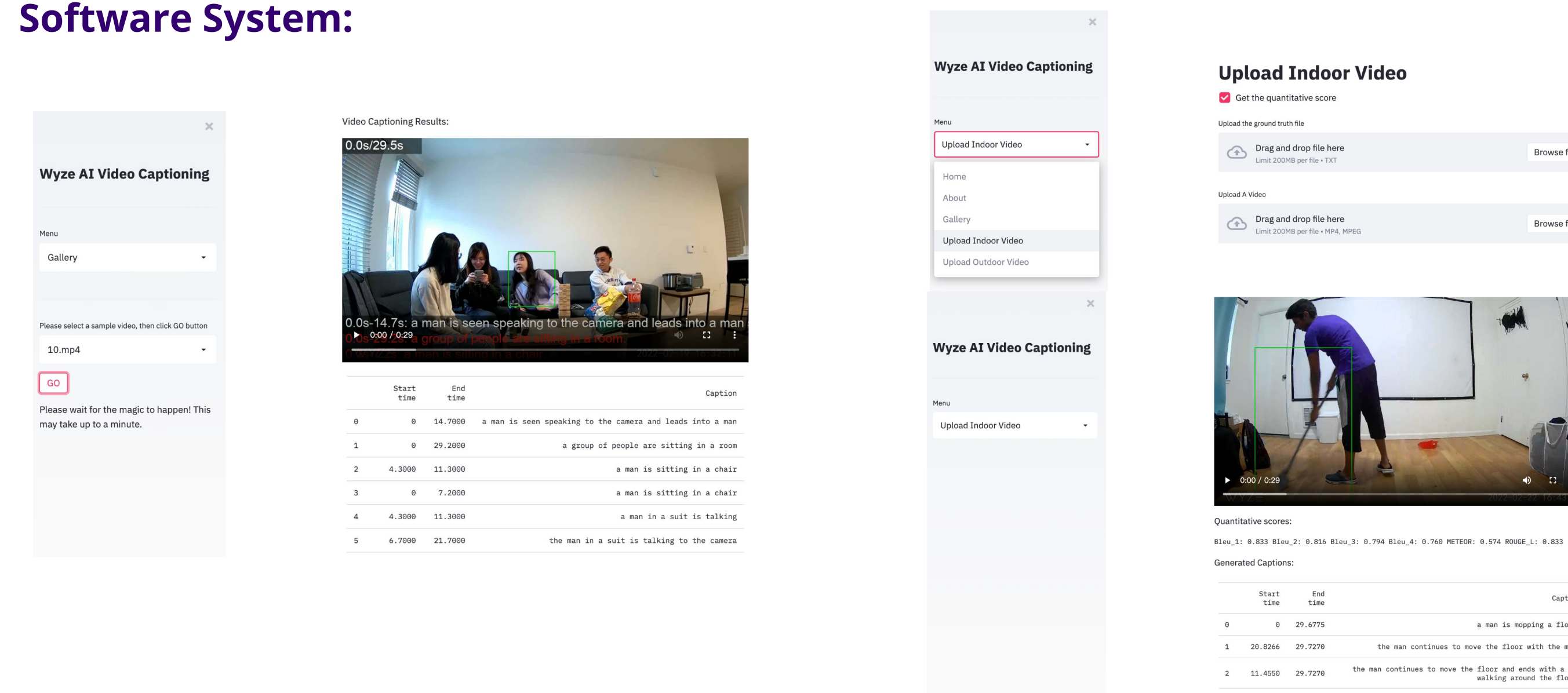


- Dense captioning employs learning of two tasks simultaneously
 - Localization of events
 - Captioning of events
- PDVC uses the same intermediate features so that information is shared between the two tasks.

Software System

- Using Streamlit with Ngrok and Colab to build the website.
- Features:
 - Gallery: User can test some sample videos.
 - Upload indoor/outdoor video: User can upload indoor/outdoor videos and get the quantitative results through their mobile and computer.

Software System:



Results

EVALUATION DATASET:

- A dataset designed to evaluate the performance of the video caption generated on about 50 indoor videos and 10 outdoor videos using the metrics mentioned below.

METRIC NAME	METHODOLOGY
BLEU	N-GRAM PRECISION
ROUGE	N-GRAM RECALL
METEOR	N-GRAM WITH SYNONYM MATCHING
CIDEr	TF-IDF WEIGHTED N-GRAM SIMILARITY

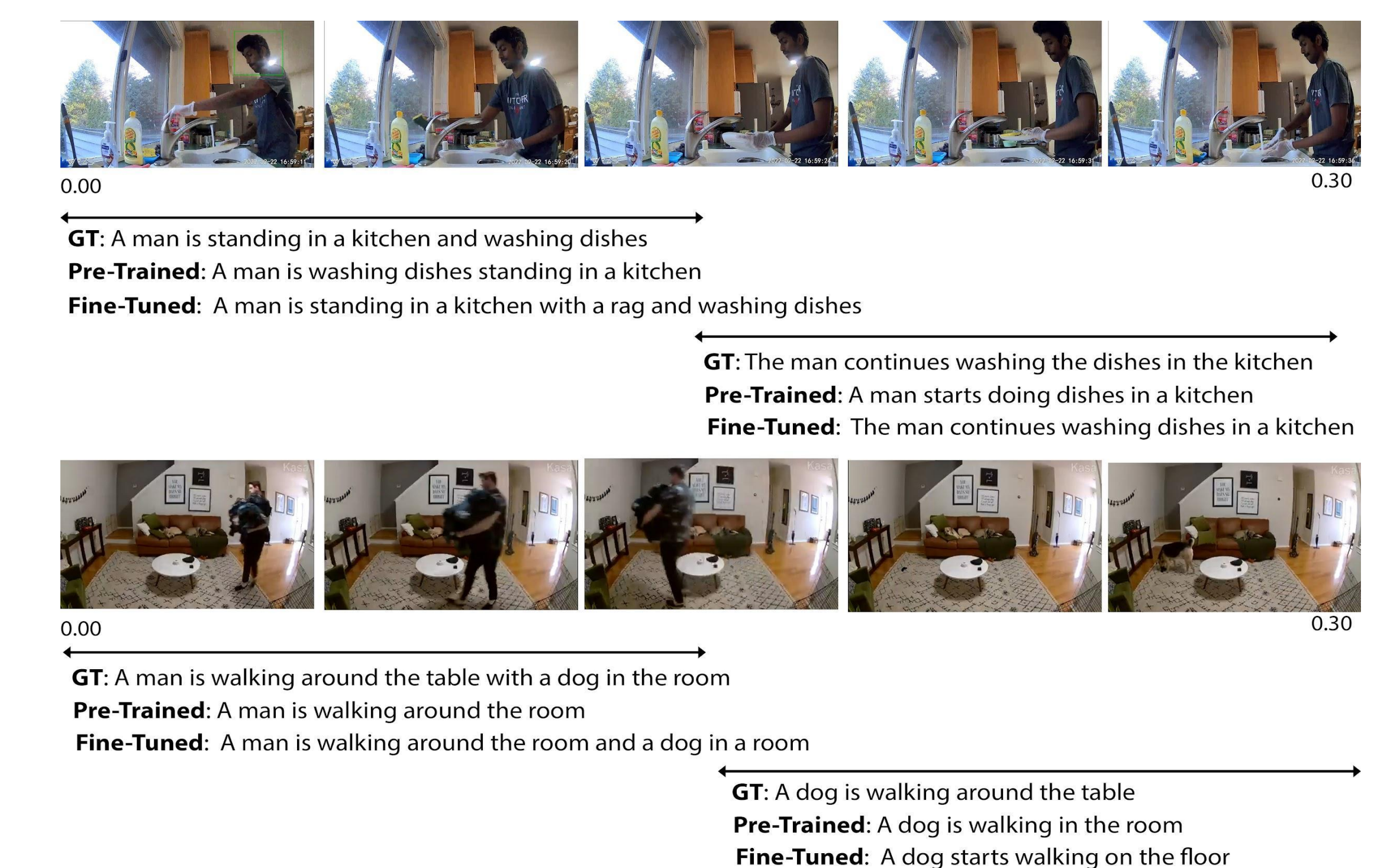
QUANTITATIVE RESULTS: (INDOOR)

MODEL	BLEU_1	BLEU_2	BLEU_3	BLEU_4	METEOR	ROUGE_L	CIDEr
Pre-Trained	0.683	0.606	0.543	0.488	0.375	0.683	2.009
Fine-Tuned	0.785	0.725	0.669	0.620	0.424	0.749	2.854

QUANTITATIVE RESULTS: (OUTDOOR)

MODEL	BLEU_1	BLEU_2	BLEU_3	BLEU_4	METEOR	ROUGE_L	CIDEr
Pre-Trained	0.538	0.357	0.222	0.132	0.244	0.484	0.445
Fine-Tuned	0.776	0.698	0.626	0.570	0.446	0.756	2.521
Fine_Tuned with YOLO	0.897	0.875	0.850	0.825	0.607	0.920	2.956

QUALITATIVE RESULTS:



Future Work and References

- Future Work
 - Enhancement of the PDVC Model Architecture.
 - Creating an enhanced UI capable of user authentication and authorization.
 - Integration of the system with the Wyze Hardware for seamless integration as a single system.
 - Generating better captions for anomaly, vehicle based and outdoor scenarios.
- References
 - Alwassel, Humam, Silvio Giancola, and Bernard Ghanem. "Tsp: Temporally-sensitive pretraining of video encoders for localization tasks." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021.
 - Wang, Teng, et al. "End-to-End Dense Video Captioning with Parallel Decoding." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021.