



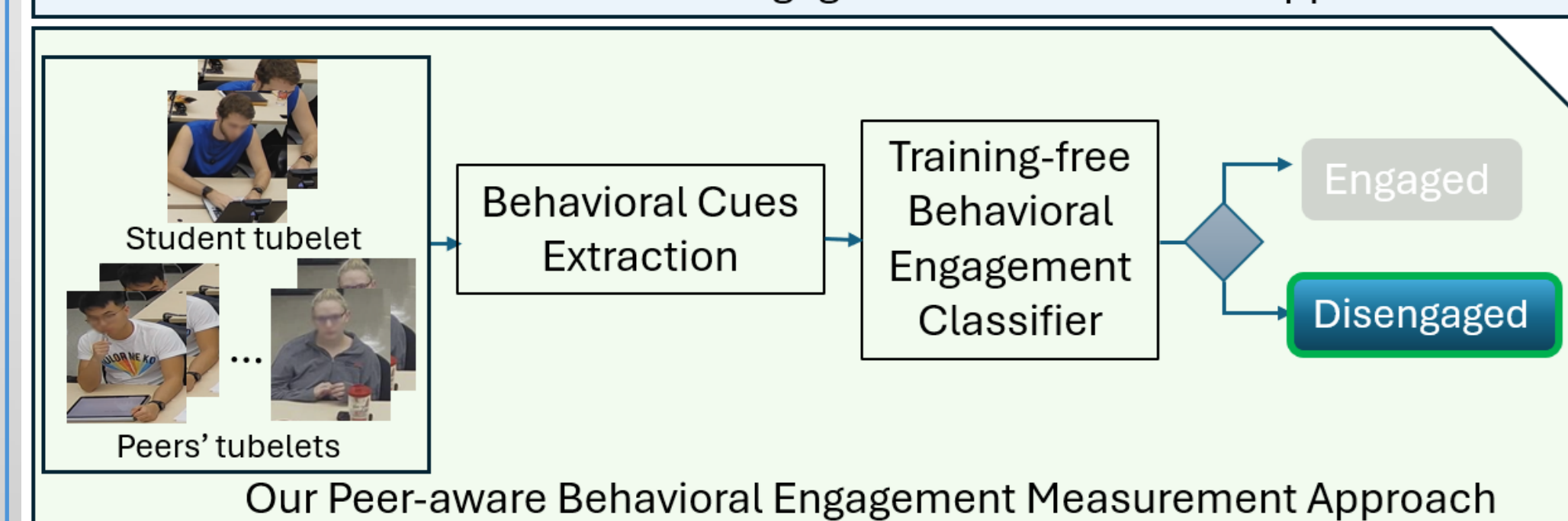
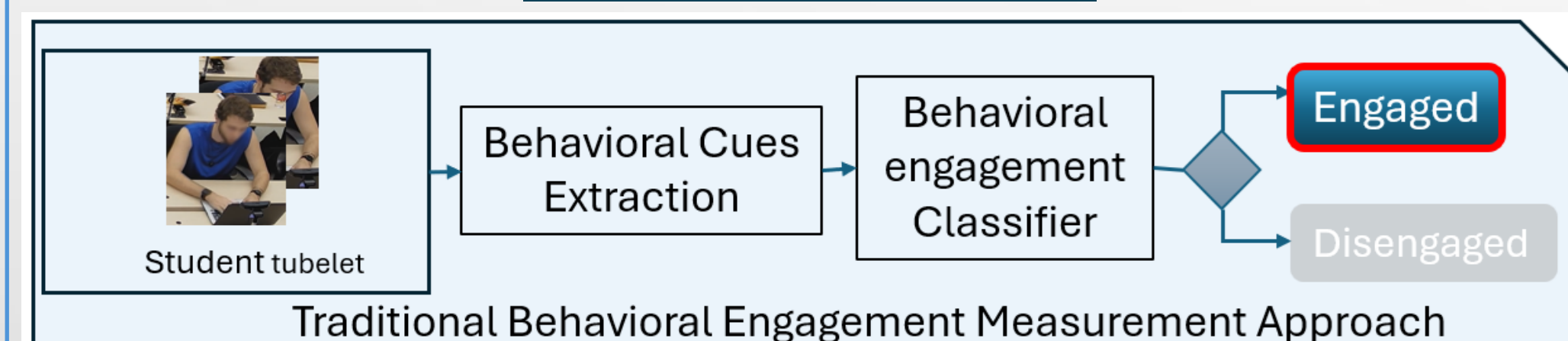
Motivation

Understanding student behavior in the classroom is essential to improving both pedagogical quality and student engagement.

Existing engagement measurement approaches:

- Require extensive annotated data
- Ignore classroom context:
- Ignore the temporal order of student actions and their duration.

Context Matters!

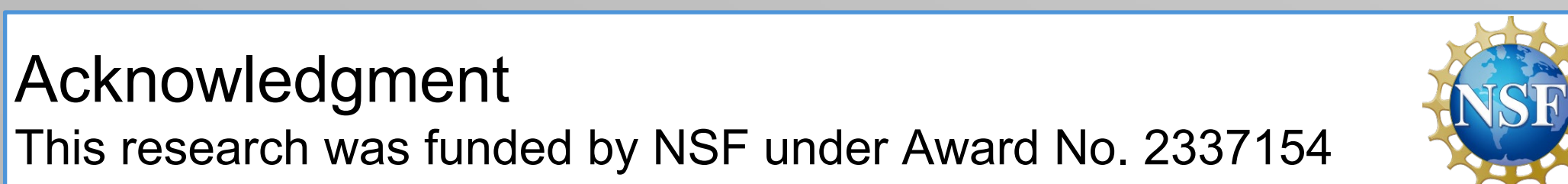


Action Dictionary

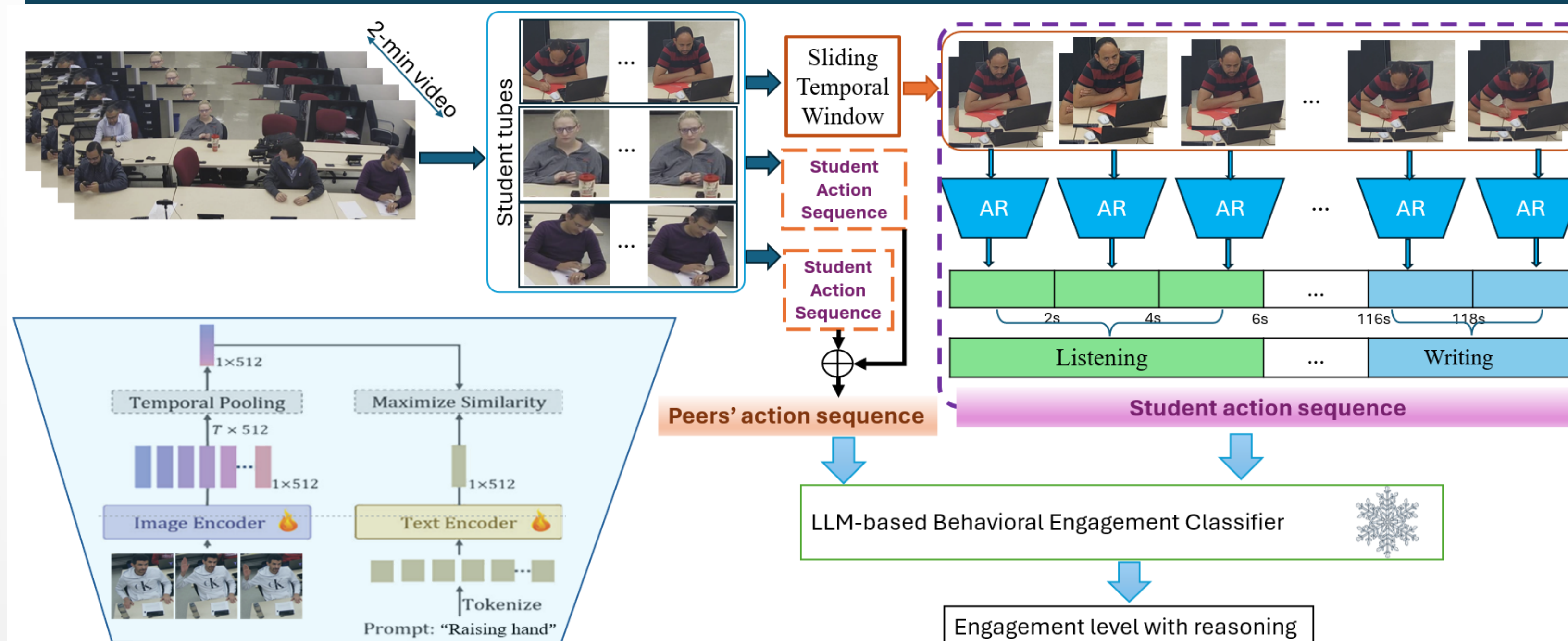


Acknowledgment

This research was funded by NSF under Award No. 2337154 (Farag, PI; Thomas Tretter (Co-PI), Measuring Student Engagement in Introductory Engineering STEM Classes, 10/15/2024 -9/30/2027)



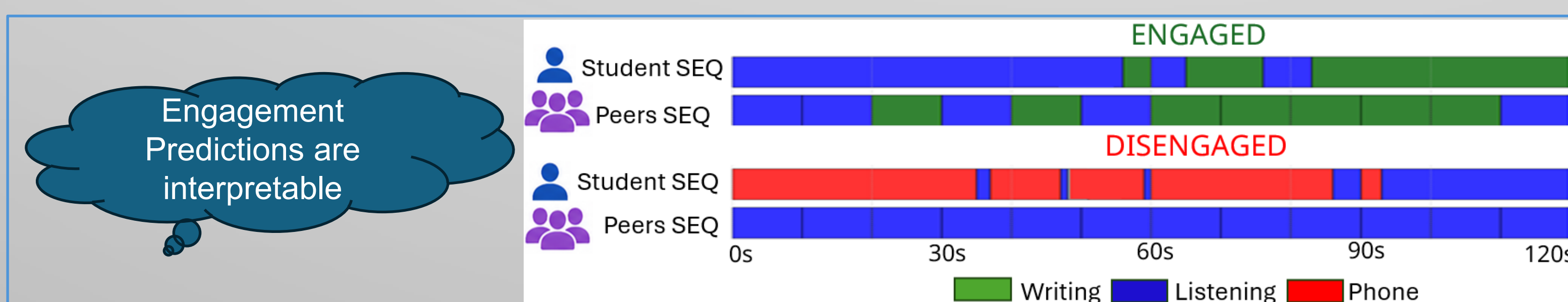
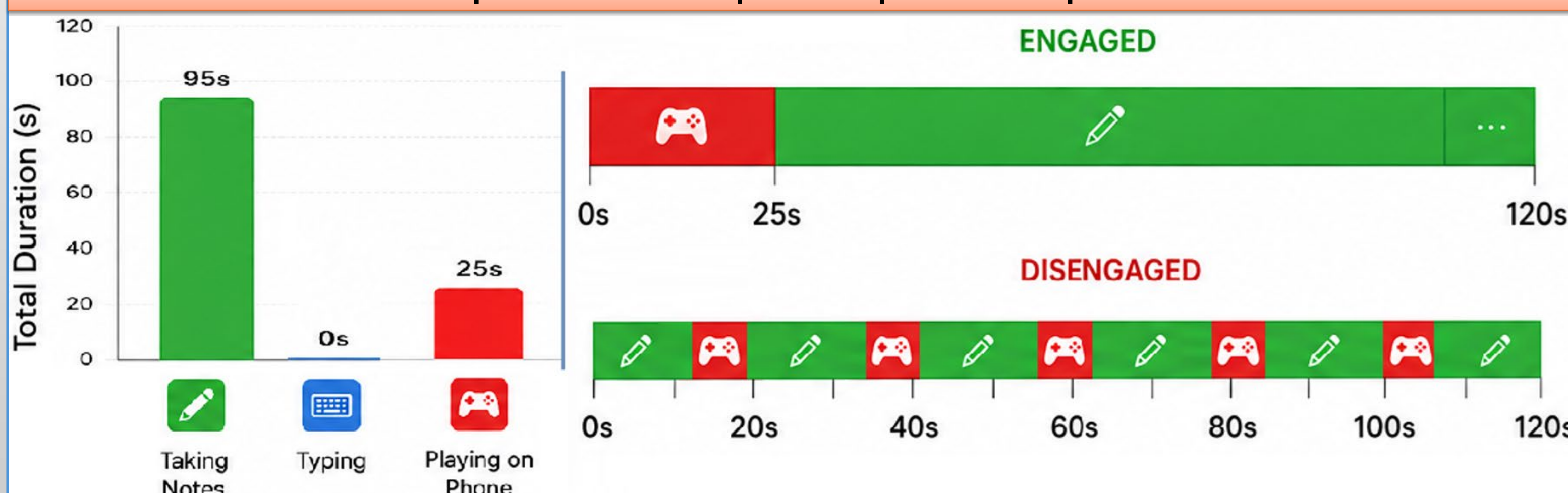
Proposed Framework



We propose a novel framework integrating Vision-Language Models & Large Language Models for student engagement classification. It consists of:

- Few-Shot VLM Action Recognition
- VLM-based Action Parsing
- Peer-Aware LLM -based Engagement Classification

Sequences better capture frequent interruptions



Engagement Predictions are interpretable

Experimental Results

We collected a behavioral engagement dataset with two subsets of annotated videos:

- Student actions -13 action categories: 208 training, 46 testing trimmed clips.
- Student engagement- 455 annotated two-minute student video clips with engagement and temporal action labels.

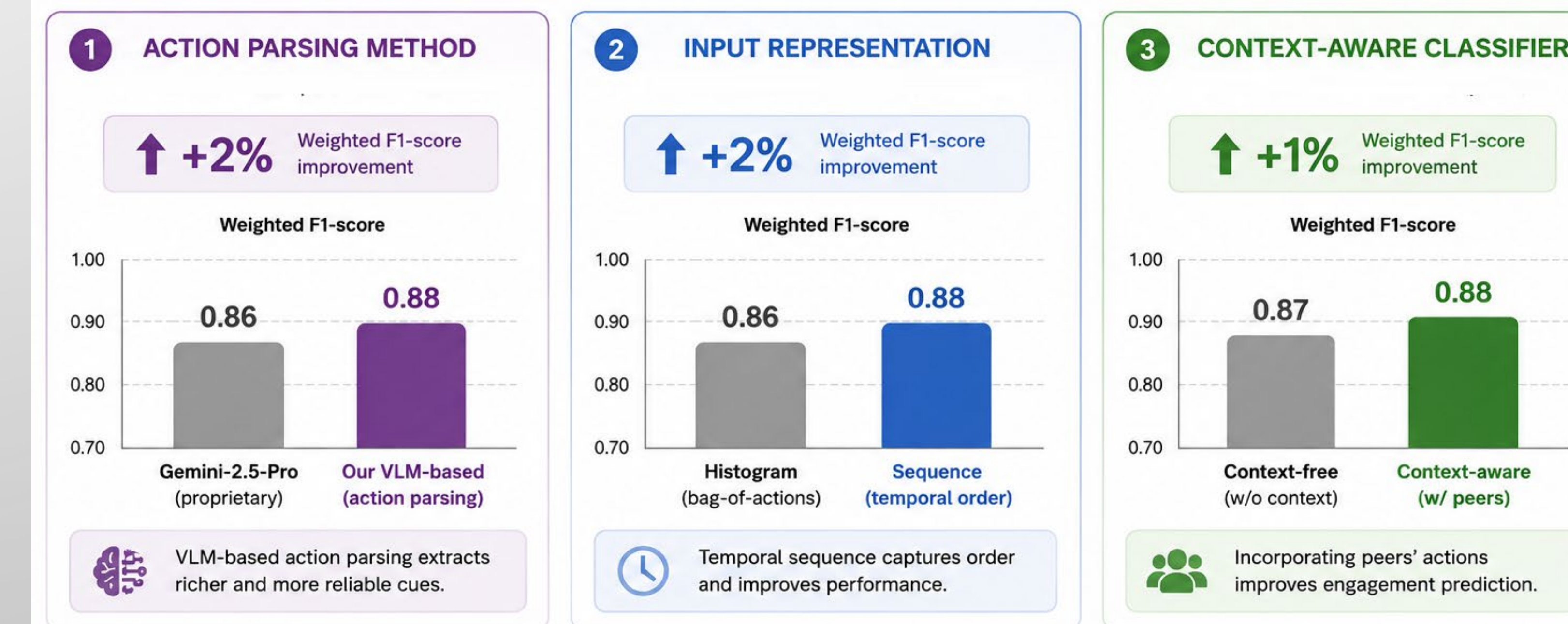
Table 1. Comparison of action recognition models (%).

Model	Student Action Dataset			
	K = 2	K = 4	K = 8	K = 16
XCLIP	43.2	79.6	88.6	97.9
VIFI-CLIP+	47.3	81.8	93.2	97.9

Table 2. Comparison of temporal action segmentation methods (%).

TAS	Accuracy	Edit Score	F1@[10,25,50]
Gemini-2.5-flash	57.2	37.4	[39.0, 34.5, 25.3]
Gemini-2.5-Pro	69.8	51.8	[58.2, 55.2, 45.2]
Our VLM-based	<u>67.0</u>	<u>45.7</u>	<u>[48.3, 43.9, 31.4]</u>

Engagement Classification: Key Comparisons (Gemma-2-9B)



KEY TAKEAWAY

Leveraging VLM-based action parsing, temporal action sequences, and classroom context consistently improves engagement classification performance.