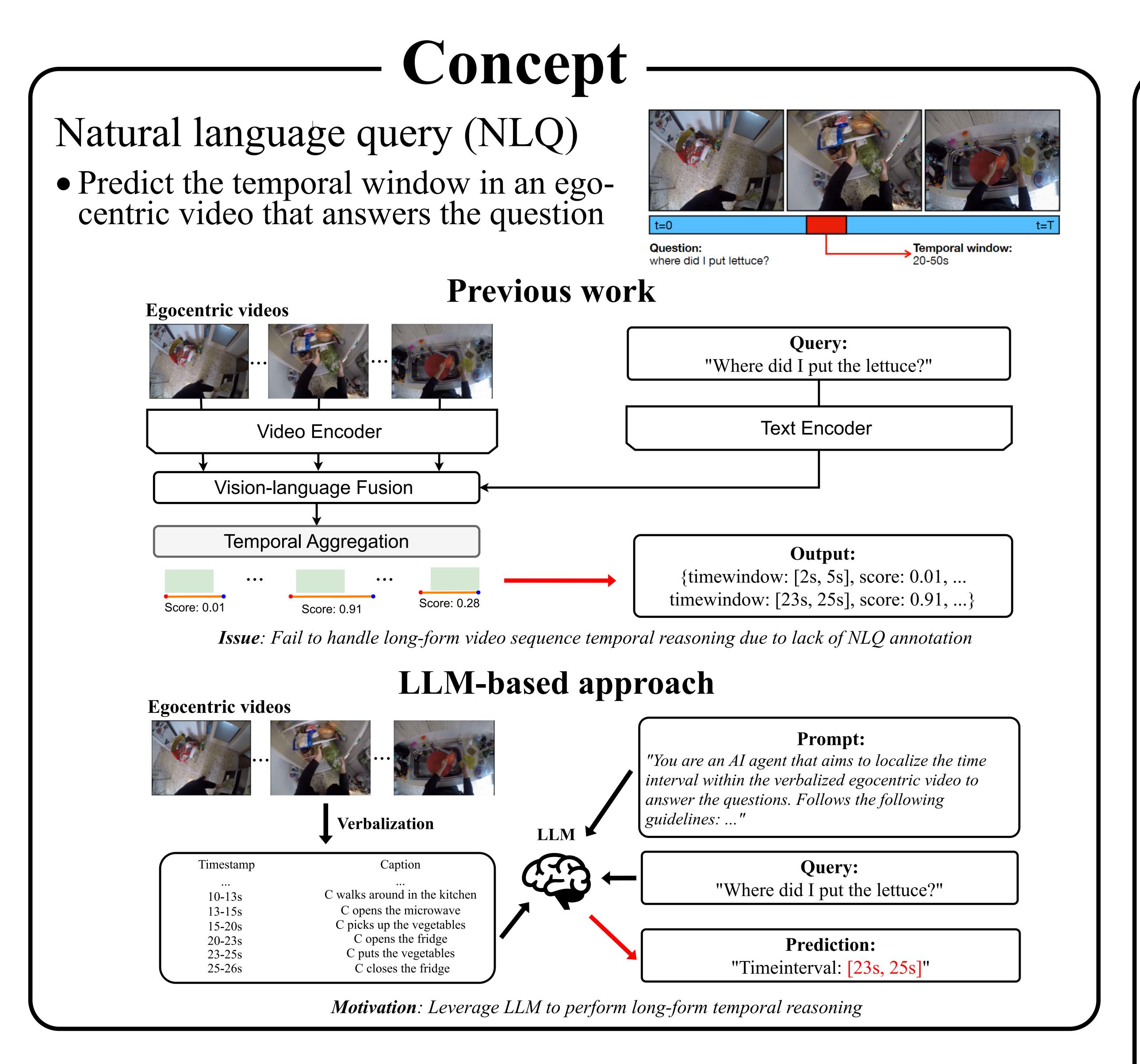


Retrieval Augmented Generation for Natural Language Query in Egocentric Video

Jiahao Nick Li, Li Gu, Omid Reza Heidari

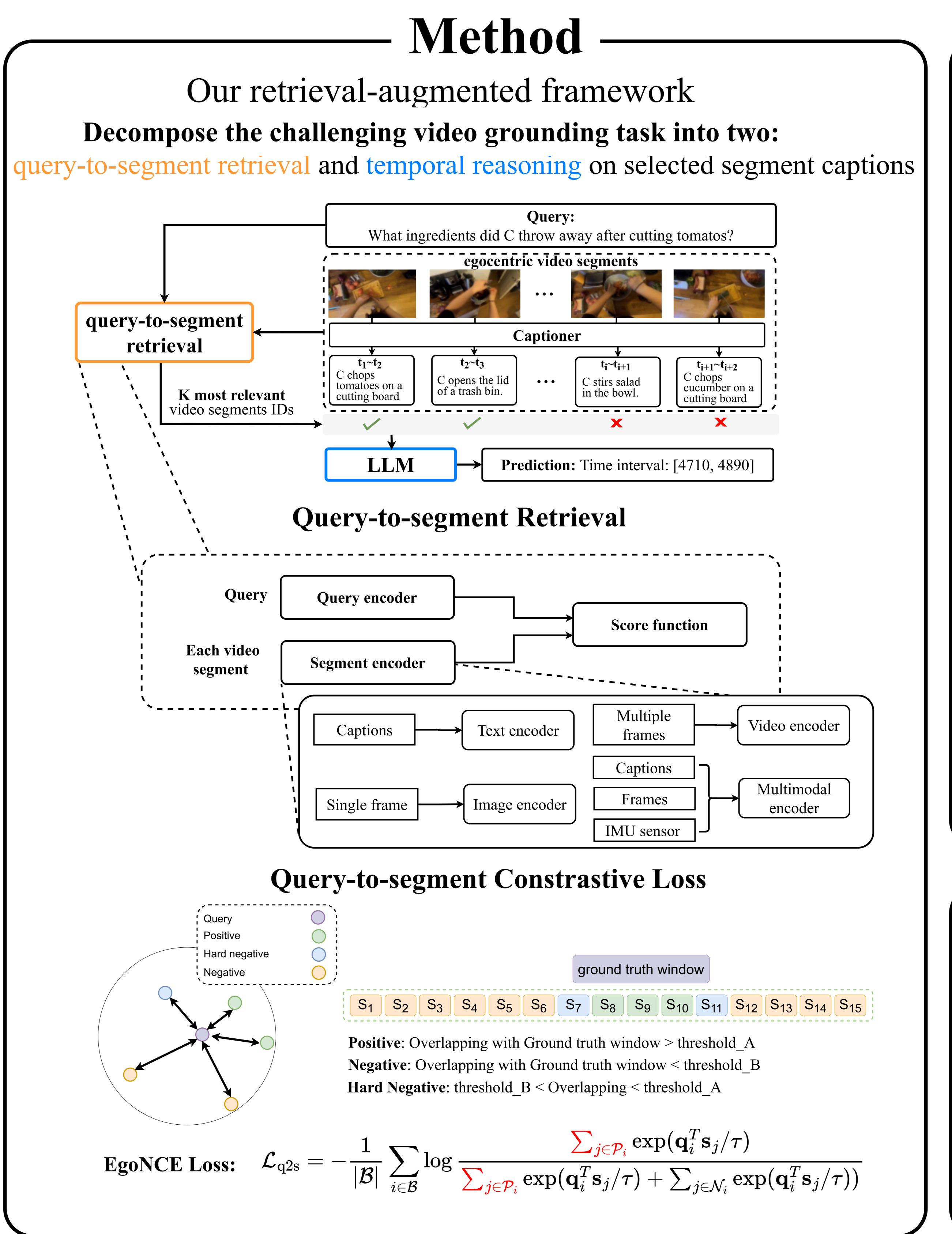


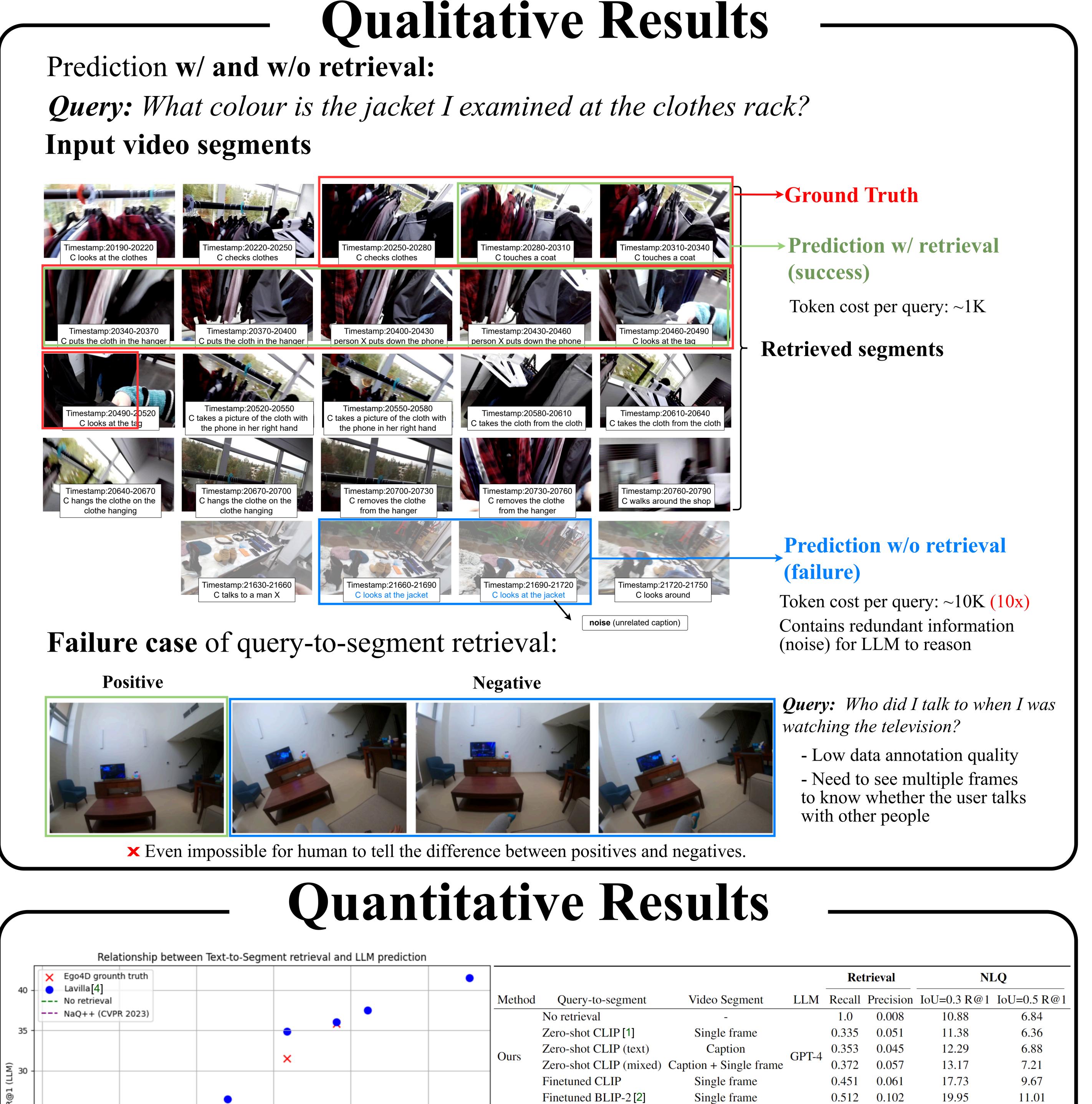
- Contributions

- Propose a retrieval augmented generation (RAG) framework that leverages LLM's temporal reasoning for NLQ
- Demonstrate that the effective query-to-segment retrieval can improve the LLM's prediction

Supervise Episodic Memory, in CVPR 2023. [4] Yue Zhao et al., Learning Video Representations from Large Language Models, in CVPR 2023.

• Establish a connection between RAG and Video understanding





*: The state-of-the-art method involves training the model on both NaQ and NLQ

train-set, where NaQ contains 80x more examples than the NLQ train-set.

[1] Alec Radford et al., Learning Transferable Visual Models From Natural Language Supervision, in ICML 2021. [2] Junnan Li et al., BLIP-2: Bootstrapping Language Models, in ICML 2023. [3] Santhosh Ramakrishnan et al., NaQ: Leveraging Narrations as Queries to