# ETVA: Evaluation of Text-to-Video Alignment
# via Fine-grained Question Generation and Answering

Kaisi Guan [1,2*]   Zhengfeng Lai[2]   Yuchong Sun[1]   Peng Zhang[2]

Wei Liu[2]   Kieran Liu[2]   Meng Cao[2]   Ruihua Song[1†]

[1] Renmin University of China    [2]Apple
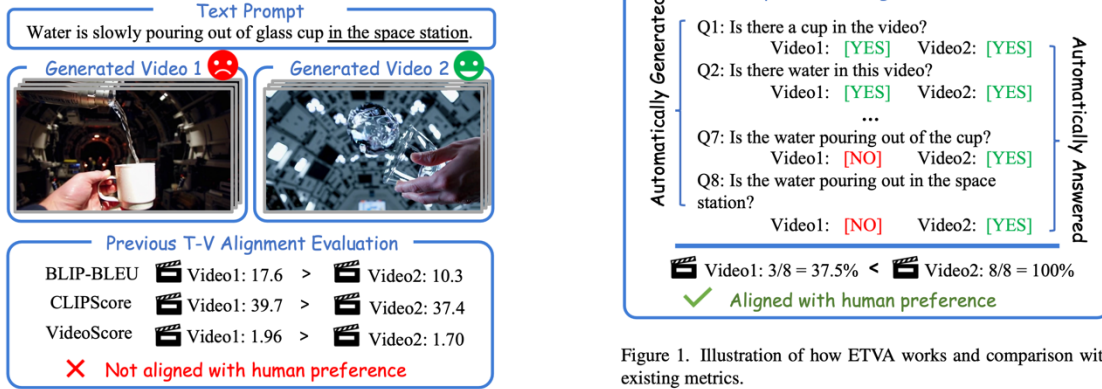
https://eftv-eval.github.io/etva-eval

Figure 1. Illustration of how ETVA works and comparison with existing metrics.
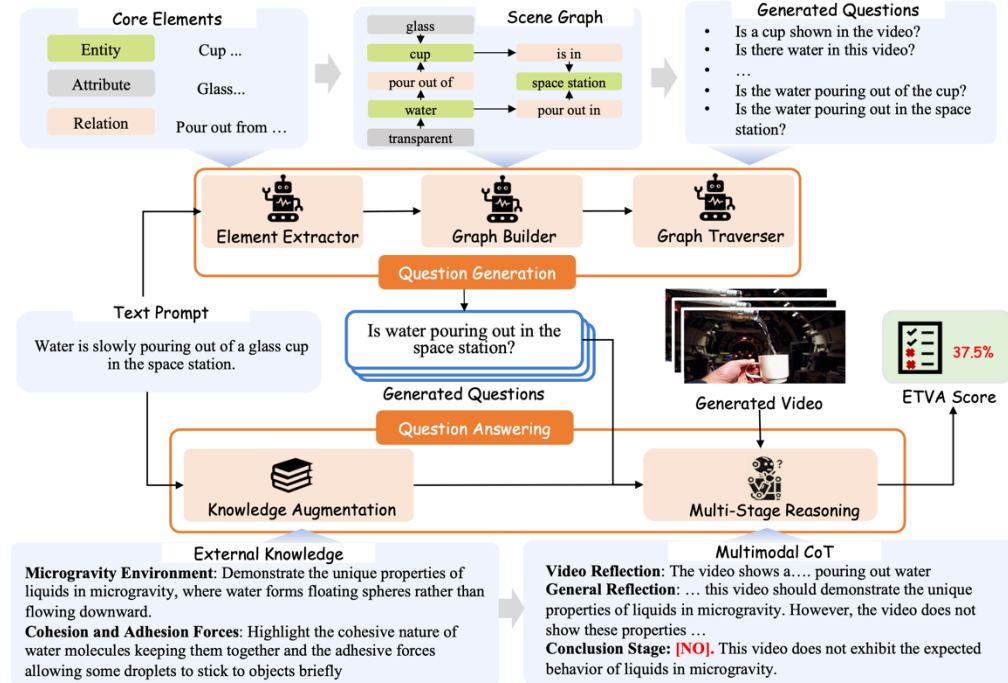


Figure 2. Overall pipeline of ETVA. ETVA contains a multi-agent framework for generating atomic questions and a knowledge-augmented multi-stage reasoning framework for question answering.