# LazyVLM: Neuro-Symbolic Approach to Video Analytics

Xiangru Jian*
University of Waterloo
xiangru.jian@uwaterloo.ca

Wei Pang*
University of Waterloo
w3pang@uwaterloo.ca

Zhengyuan Dong*
University of Waterloo
zhengyuan.dong@uwaterloo.ca

Chao Zhang*
University of Waterloo
chao.zhang@uwaterloo.ca

M. Tamer Özsu
University of Waterloo
tamer.ozsu@uwaterloo.ca

## ABSTRACT

Current video analytics approaches face a fundamental trade-off between flexibility and efficiency. End-to-end Vision Language Models (VLMs) often struggle with long-context processing and incur high computational costs, while neural-symbolic methods depend heavily on manual labeling and rigid rule design. In this paper, we introduce LazyVLM, a neuro-symbolic video analytics system that provides a user-friendly query interface similar to VLMs, while addressing their scalability limitation. LazyVLM enables users to effortlessly drop in video data and specify complex multi-frame video queries using a semi-structured text interface for video analytics. To address the scalability limitations of VLMs, LazyVLM decomposes multi-frame video queries into fine-grained operations and offloads the bulk of the processing to efficient relational query execution and vector similarity search. We demonstrate that LazyVLM provides a robust, efficient, and user-friendly solution for querying open-domain video data at scale.

## 1 INTRODUCTION

With the massive volume of video data in real-world applications, analyzing video content has become increasingly crucial across various domains [2, 7–9, 11, 13, 14]. A common task in video analytics involves identifying a specific short video segment within a longer video. For instance, consider a surveillance video capturing road traffic: retrieving a moment where a motorcycle is positioned to the right of a bus and then moves to the left of the bus within a few seconds can be instrumental in detecting anomalous behavior in road transportation. We refer to this type of query that retrieves multi-frame events as a *video moment retrieval query* (VMRQ).

One category of existing systems for VMRQ processing relies on manual-intensive configurations and interactions. In such systems, users must define task-specific machine learning models for video content processing, which requires prior domain knowledge. These systems typically employ either a SQL-like query language [7, 8] or a query-by-example interface [2, 11, 13]. The former struggles to express complex video moments due to the need for multiple join expressions or recursive joins to define video moments, while the latter demands extensive user interaction with the video content to perform the search manually (labeling intermediate query results). These designs result in low human efficiency.

Alternatively, vision large multimodal models (VLMs) can be used for video query processing [1]. VLMs offer an intuitive, end-to-end video processing interface: users can simply load a video and input natural language queries into the context window, providing high human efficiency. However, VLMs fall short when processing video queries at scale due to several limitations. First, VLMs' inference is highly time-consuming due to the underlying self-attention mechanism, which requires $O(n^2)$ time complexity where $n$ is the total number of tokens in the context window—this overhead is primarily dominated by video length. Second, the autoregressive nature of decoder-only VLMs enforces a serial computation mechanism, significantly limiting opportunities for parallel processing. Finally, executing ad-hoc exploratory queries or video updates, such as adding new videos, necessitates repeated processing of the entire context window, further exacerbating inefficiencies. Therefore, using VLMs out of the box for video query processing leads to low system efficiency.

To overcome the above issues, we propose LazyVLM, a neuro-symbolic approach designed for scalable video analysis. LazyVLM introduces a semi-structured text interface for defining video events, allowing users to describe visual content in frames using subject-predicate-object (SPO) triples, where each element (S, P, or O) is specified in text (see Example 2.1). Video can be loaded into LazyVLM without defining task-specific models, *i.e.*, by simply dropping in the video with no additional effort required. By offering an interface similar to that of VLMs, LazyVLM improves human efficiency, compared to traditional manual-intensive systems.

LazyVLM significantly reduces the computation cost, compared to out-of-box VLMs, thereby improving system efficiency for video analytics. This is achieved by generating structured views of video data—specifically, scene graphs for frames—and embedding nodes (representing entities) within these graphs. LazyVLM then decomposes video queries into semantic search and symbolic search. Semantic search is based on embedding vectors and is dedicated to searching an entity (either a subject or an object in the SPO triples). Symbolic search is based on relational queries and focuses on verifying the existence of relationships (predicates between subjects and objects in SPO triples). Additionally, LazyVLM leverages VLMs to refine the results of relational queries. However, since most computation is offloaded to semantic search and vector search, VLMs are used in a lightweight manner. Instead of processing an entire video, LazyVLM selectively applies VLMs to individual frames that have been verified through relational queries—dramatically reducing computational overhead. Thanks to its fine-grained query decomposition, LazyVLM enables parallel computing, allowing multiple relational queries, vector searches, and VLM refinements to be executed simultaneously. Moreover, structured views and embedding vectors are precomputed once and stored. This makes LazyVLM update-friendly, supporting incremental updates by inserting new

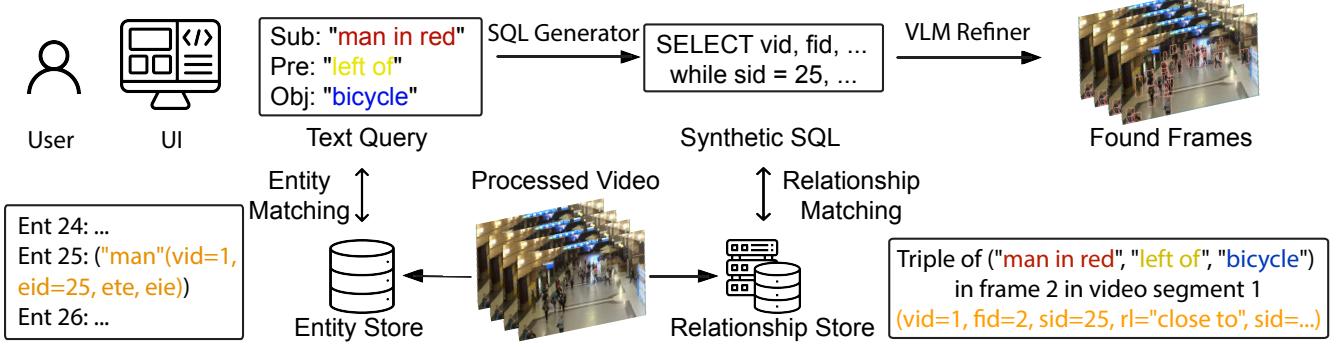---

*Equal contribution to this research.

**Figure 1: Overview of query processing in LazyVLM. The diagram illustrates the processing of a semi-structured text query, which includes entity descriptions (*e.g.*, *"man in red"*) and relationship terms (*e.g.*, *"near"*). The query is processed through a sequence of stages: entity matching via vector similarity search, SQL query processing to retrieve candidate relationships, relationship verification using a VLM to refine results, and temporal matching to identify the final set of video segments.**

structured views and vectors, eliminating the need to reprocess the entire video from scratch.

LazyVLM enhances both human efficiency and system efficiency by combining the strengths of symbolic search, semantic search and VLMs. This integration enables a powerful and efficient framework for querying open-domain video data. The remainder of this paper provides a detailed overview of LazyVLM's architecture and interaction mechanisms.

## 2 SYSTEM OVERVIEW

Figure 1 illustrates the processing pipeline of LazyVLM. Video data is first preprocessed to generate structured views, which are stored in the **Entity Store** and **Relationship Store** (see Section 2.2 for details). LazyVLM provides a semi-structured text interface that allows users to query the loaded video data. At query time, LazyVLM employs a series of components to retrieve relevant information from the two stores, followed by the use of a VLM to refine the pruned set of candidate results. The remainder of this section describes the key components of LazyVLM in detail.

## 2.1 Query Interface and Functionality

The query interface of LazyVLM is based using SPO triples to describe the visual content in video frames. We use the following example to illustrate how users can describe multi-frame events in a natural and intuitive manner.

*Example 2.1.* Consider the following query that defines a complex event: a man with a backpack is near a bicycle, and another man in red clothes moves from the left of the bicycle to the right of the bicycle after more than 2 seconds. This query is formally specified in the following steps in LazyVLM. **(1) Entity Description:** $E = \{e_1, e_2, e_3\}$, with $e_1$.text = *"man with backpack"*, $e_2$.text = *"bicycle"*, and $e_3$.text = *"man in red"*. These defined entities can serve as either a subject or an object in SPO triples. **(2) Relationship Description:** $R = \{r_1, r_2, r_3\}$, where $r_1$.text = *"is near"*, $r_2$.text = *"leftOf"*, and $r_3$.text = *"rightOf"*. These relationships function as predicates in SPO triples. **(3) Frame Description:** $F = (f_0, f_1)$, where $f_0 = \{(e_1, r_1, e_2), (e_3, r_2, e_2)\}$ represents *"man*

*with backpack is near bicycle; man in red is on the left of bicycle"*, and $f_1 = \{(e_1, r_1, e_2), (e_3, r_3, e_2)\}$ represents *"man with backpack is near bicycle; man in red is on the right of bicycle"*. **(4) Temporal Constraint:** $f_1 - f_0 > 4$, ensuring that the second frame occurs more than 2 seconds after the first, assuming a frame rate of 2 frames per second.

In LazyVLM, input video is automatically divided into non-overlapping video clips, *i.e.*, short segments lasting a few seconds or minutes, with the length defined by the user. LazyVLM then processes an input query, *e.g.*, the one in Example 2.1, by searching for events specified in the query, retrieving, and returning the clips that contain the detected event.

*Functionality Supported.* LazyVLM supports a comprehensive range of core video analytics functionalities. At the object level, LazyVLM supports (1) *object detection*, identifying and localizing entities such as cars, people, and bicycles within video frames, and (2) *object tracking*, which follows these objects across multiple frames, facilitating trajectory-based analysis. In addition, LazyVLM includes (3) *attribute recognition*, allowing queries based on object properties like color or size. At the relationship level, LazyVLM includes (4) *relationship detection*, which captures spatial and interaction-based relationships (*e.g.*, car near pedestrian or person holding an object). Beyond basic object and relationship analytics, LazyVLM supports advanced query operations: (5) *conjunction queries* allow users to specify multiple conditions (*e.g.*, detecting a person and a vehicle in the same frame), while (6) *sequencing queries* enforce temporal order between events (*e.g.*, detecting a person walking before entering a car). Additionally, (7) *window queries* constrain events within a defined time duration (*e.g.*, detecting a car stopping within 10 seconds after a pedestrian appearing).

*User Interaction Flow.* User interaction is straightforward with LazyVLM. **(1) Upload Video Dataset:** Users begin by uploading their videos. LazyVLM automatically segments the videos, extracts visual features, and builds the corresponding structured views. **(2) Compose a Query:** Through the intuitive text interface, users describe the event of interest. For example, to retrieve clips where *"a man with a backpack is near a bicycle, and a man in red moves from*

*being on the left of the bicycle to the right of the bicycle after more than 2 seconds"*, users enter the corresponding entity descriptions, relationship descriptions, frame information, and temporal constraints, as shown in Example 2.1. The query is then submitted to the system, and users can view the results after the query execution in the interface. Detailed steps are presented in Section 3.

## 2.2 Video Preprocessing

Video preprocessing converts raw video content into a structured, searchable format through the following integrated stages.

*Video Segmentation.* When a video dataset is uploaded, each video $V$ is automatically partitioned into a sequence of non-overlapping segments, $V = (v_1, v_2, \ldots, v_n)$. Each segment $v_i$ comprises a fixed number of frames, $v_i = (f_1, f_2, \ldots, f_m)$, facilitating parallel processing and management of long videos.

*Content Extraction.* Each video segment is processed to extract detailed visual information. To represent the visual content of each frame, we construct a scene graph that encodes the frame as a set of subject–predicate–object (SPO) triples in the form (subject, predicate, object). For example, a frame may include triples such as (*"man with backpack", "near", "bicycle"*) and (*"bicycle", "on", "sidewalk"*). We employ the IETrans [12] model to generate these scene graphs. For extracted triples, each subject or object—internally referred to as an entity—is represented as a tuple ($vid, eid, ete, eie$), and all entities are stored in the **Entity Store**. $vid$ denotes the identifier of the video segment containing the entity; $eid$ is a unique entity identifier within the segment, obtained through entity tracking using YOLOv8 [6]; $ete$ is a text embedding derived from the entity's description using the e5-mistral-7b [10] model; and $eie$ is an image embedding that captures the entity's visual appearance, generated by the VLM2Vec [5] model. In addition to entity representations, inter-object relationships are captured as tuples ($vid, fid, sid, rl, oid$), which are stored in the **Relationship Store**. $vid$ and $fid$ refer to the video segment and frame identifiers, respectively; $sid$ and $oid$ are the unique entity identifiers of the subject and object involved in the relationship; and $rl$ denotes the relationship label, *e.g.*, *"near"*.

## 2.3 Query Processing

Users express an event query as a four-part specification, including entity descriptions, relationship descriptions, frame-level specifications, and temporal constraints, as illustrated in Example 2.1. The query engine of LazyVLM processes the query through a series of stages: *Entity Matching*, *SQL Query Generation*, *Relationship Matching and Refinement*, and *Temporal Matching*.

*Entity Matching.* For each entity defined in the query, a vector similarity search is performed to match the textual description of the entity against the embeddings stored in the Entity Store. The result is a set of candidate entities for each query entity, represented as ($vid, eid$) pairs, *i.e.*, matched entities in specific video segments.

*SQL Query Generation.* Based on the entity matching results, the query engine automatically generates SQL queries to retrieve candidate frames from the Relationship Store. These queries filter rows by matching entity identifiers, *i.e.*, ($vid, eid$) pairs, and they return frames that potentially include the specified entities. The

output is a set of rows from the Relationship Store corresponding to each query entity, denoted as candidate frames.

*Relationship Matching and Refinement.* For each SPO triple in the query, the query engine performs a join between the candidate frames for the subject and object, based on shared $vid$ and $fid$ values. This identifies potential relationships between subject and object candidates. Following this coarse-grained match, a refinement step is applied to further verify relationships. Specifically, a lightweight local VLM (*e.g.*, Qwen-2.5-VL 7B [1]) is used for the verification. Optionally, users may apply cost-efficient closed-source VLMs (*e.g.*, GPT-4o-mini) for deployment-friendly setups. This refinement ensures that the relationships specified in the query are visually and semantically grounded in the candidate frames. Finally, the engine joins the refined relationship results to identify frames where all specified SPO triples in a query frame co-occur, producing a set of candidate video frames for each query frame.

*Temporal Matching.* In the final stage, the engine checks whether the candidate video frames satisfy the temporal constraints defined in the query. This involves join over frame identifiers to ensure compliance with temporal logic, *e.g.*, $f_1 - f_0 > 4$. Ultimately, video identifiers including candidate video frames that satisfy the temporal conditions are returned as the final query results.

A major advantage of the proposed query processing pipeline is that each step is inherently parallelizable. For instance, entity matching tasks can be executed in parallel, as they are independent of one another. In addition, the query engine applies VLMs only for fine-grained relationship verification on a pruned set of candidate frames, rather than processing all video frames, making it significantly more lightweight compared to end-to-end VLM-based approaches for video query processing.

## 3 DEMONSTRATION

We demonstrate LazyVLM using the MOT20 [4] and TAO [3] datasets. Specifically, we use Example 2.1 to perform a query on the MOT20-02 dataset. We guide the users through the following detailed steps, illustrating the interactive and intuitive interface of LazyVLM.

**Step ❶: Load Dataset and Enter Hyperparameters.** Users select the MOT20-02 dataset path from a dropdown menu. The interface shows dataset metadata, including the total number of video segments (*e.g.*, 11 segments for MOT20-02), as well as paths for preprocessed and raw data. Users then configure several hyperparameters for query execution: the *top-k* parameter (*e.g.*, top 3 results), the *temperature* for controlling search strictness, and the thresholds for *vision embedding* and *textual embedding* similarity searches, which affect entity and relationship matching accuracy.

**Step ❷: Enter Entities.** Users input descriptive text labels for the entities involved in their query via a dedicated input field. In the provided example query, users define entities like *"man with backpack," "bicycle,"* and *"man in red"*. These entities are then listed on the interface and can be reviewed or removed if necessary.

**Step ❸: Enter Relationships.** Users specify relationships that describe interactions or spatial positions between entities, such as *"is near," "is on the left of,"* or *"is talking to"*. Each entered relationship appears on the interface for adjusting before proceeding.

**Step ❹: Enter Triples.** Users construct SPO triples to precisely define the interactions between entities. The interface supports
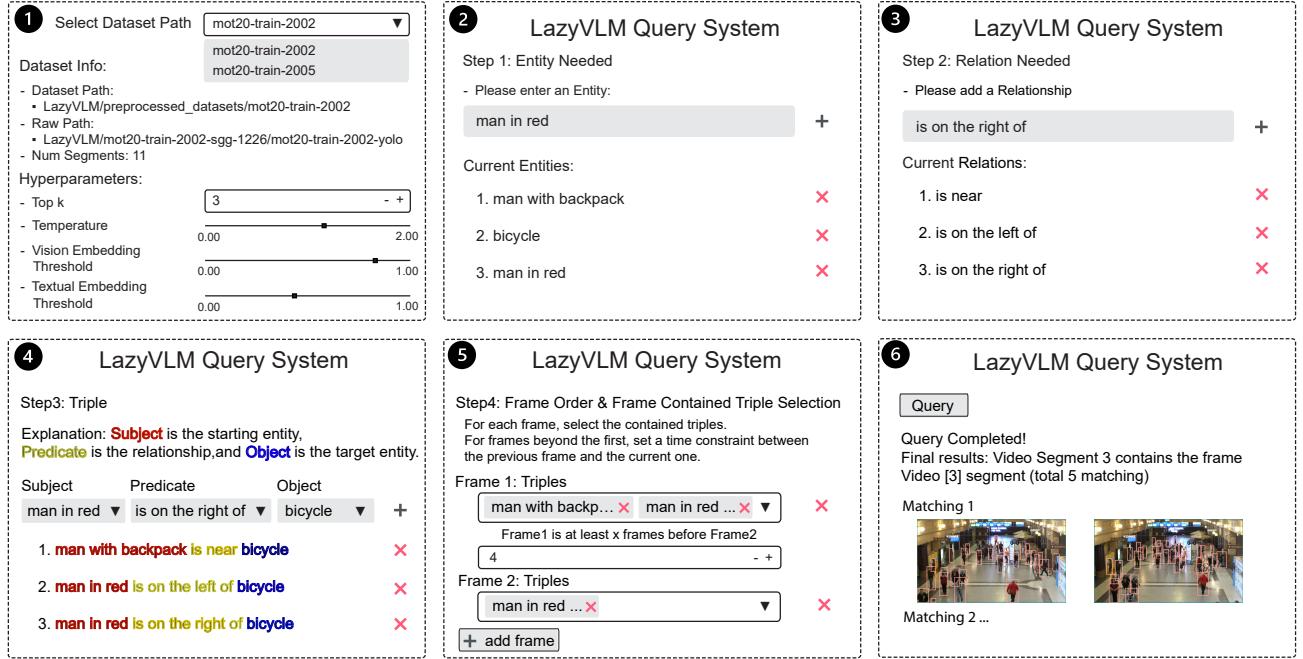
**Figure 2: Pipeline of user interactions in LazyVLM for specifying and executing a video query: Step ❶: Load Dataset and Enter Hyperparameters; Step ❷: Enter Entities; Step ❸: Enter Relationships; Step ❹: Enter Triples; Step ❺: Enter Frames and Temporal Constraints; and Step ❻: Query Execution and Presentation of Results.**

triple formation by allowing users to select from previously entered entities and relationships. For instance, users define triples like *"man with backpack is near bicycle," "man in red is on the left of bicycle,"* and *"man in red is on the right of bicycle"*. Each formed triple is displayed for user verification.

**Step ❺: Enter Frames and Temporal Constraints.** Users organize the defined triples into specific frames according to their query's temporal constraint. For each frame, users explicitly select which triples it contains from a dropdown list. For the example provided, Frame 1 is set to contain triples *"man with backpack is near bicycle"* and *"man in red is on the left of bicycle,"* while Frame 2 includes *"man with backpack is near bicycle"* and *"man in red is on the right of bicycle."* Users also define temporal constraints, such as specifying Frame 1 to occur at least 4 frames before Frame 2 since we have 2 frames per second in the video.

**Step ❻: Query Execution and Presentation of Results.** Upon completing the query setup, users initiate query execution by clicking the *"Query"* button. LazyVLM processes the query and displays matching results. The results detail the precise video segments and exact frame identifiers corresponding to each user-defined frame.

## 4 CONCLUSION

We present LazyVLM, a neuro-symbolic video analytics system that efficiently supports complex, multi-frame queries through intuitive, text-based interaction. Our demonstration highlights LazyVLM's scalability, accuracy, and usability in practical, open-domain video analytics scenarios.

## REFERENCES

[1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-VL Technical Report. arXiv:2502.13923 [cs.CV] https://arxiv.org/abs/2502.13923

[2] Maureen Daum, Enhao Zhang, Dong He, Stephen Mussmann, Brandon Haynes, Ranjay Krishna, and Magdalena Balazinska. 2023. VOCALExplore: Pay-as-You-Go Video Data Exploration and Model Building. *Proc. VLDB Endow.* 16, 13 (Sept. 2023), 4188–4201. https://doi.org/10.14778/3625054.3625057

[3] Achal Dave, Tarasha Khurana, Pavel Tokmakov, Cordelia Schmid, and Deva Ramanan. 2020. TAO: A Large-Scale Benchmark for Tracking Any Object. In *European Conference on Computer Vision.* https://arxiv.org/abs/2005.10356

[4] P. Dendorfer, H. Rezatofighi, A. Milan, J. Shi, D. Cremers, I. Reid, S. Roth, K. Schindler, and L. Leal-Taixé. 2020. MOT20: A benchmark for multi object tracking in crowded scenes. *arXiv:2003.09003[cs]* (March 2020). http://arxiv.org/abs/1906.04567 arXiv: 2003.09003.

[5] Ziyan Jiang, Rui Meng, Xinyi Yang, Semih Yavuz, Yingbo Zhou, and Wenhu Chen. 2025. VLM2Vec: Training Vision-Language Models for Massive Multimodal Embedding Tasks. In *The Thirteenth International Conference on Learning Representations.* https://openreview.net/forum?id=TE0KOzWYAF

[6] Glenn Jocher et al. 2023. *Ultralytics YOLOv8.* https://github.com/ultralytics/ultralytics

[7] Gaurav Tarlok Kakkar, Jiashen Cao, Pramod Chunduri, Zhuangdi Xu, Suryatej Reddy Vyalla, Prashanth Dintyala, Anirudh Prabakaran, Jaeho Bang, Aubhro Sengupta, Kaushik Ravichandran, Ishwarya Sivakumar, Aryan Rajoria, Ashmita Raju, Tushar Aggarwal, Abdullah Shah, Sanjana Garg, Shashank Suman, Myna Prasanna Kalluraya, Subrata Mitra, Ali Payani, Yao Lu, Umakishore Ramachandran, and Joy Arulraj. 2023. EVA: An End-to-End Exploratory Video Analytics System. In *Proceedings of the Seventh Workshop on Data Management for End-to-End Machine Learning* (Seattle, WA, USA) *(DEEM '23).* Association for Computing Machinery, New York, NY, USA, Article 8, 5 pages. https://doi.org/10.1145/3595360.3595858

[8] Daniel Kang, Peter Bailis, and Matei Zaharia. 2019. BlazeIt: optimizing declarative aggregation and limit queries for neural network-based video analytics. *Proc. VLDB Endow.* 13, 4 (dec 2019), 533–546. https://doi.org/10.14778/3372716.3372725

[9] Francisco Romero, Caleb Winston, Johann Hauswald, Matei Zaharia, and Christos Kozyrakis. 2023. Zelda: Video Analytics using Vision-Language Models. arXiv:2305.03785 [cs.DB] https://arxiv.org/abs/2305.03785

[10] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Improving Text Embeddings with Large Language Models. arXiv:2401.00368 [cs.CL] https://arxiv.org/abs/2401.00368

[11] Renzhi Wu, Pramod Chunduri, Ali Payani, Xu Chu, Joy Arulraj, and Kexin Rong. 2024. SketchQL: Video Moment Querying with a Visual Query Interface. *Proc. ACM Manag. Data* 2, 4, Article 204 (Sept. 2024), 27 pages. https://doi.org/10.1145/3677140

[12] Ao Zhang, Yuan Yao, Qianyu Chen, Wei Ji, Zhiyuan Liu, Maosong Sun, and Tat-Seng Chua. 2022. Fine-Grained Scene Graph Generation with Data Transfer. In *Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer Nature Switzerland, Cham, 409–424.

[13] Enhao Zhang, Maureen Daum, Dong He, Manasi Ganti, Brandon Haynes, Ranjay Krishna, and Magdalena Balazinska. 2023. EQUI-VOCAL Demonstration: Synthesizing Video Queries from User Interactions. *Proc. VLDB Endow.* 16, 12 (Aug. 2023), 3978–3981. https://doi.org/10.14778/3611540.3611600

[14] Enhao Zhang, Maureen Daum, Dong He, Brandon Haynes, Ranjay Krishna, and Magdalena Balazinska. 2023. EQUI-VOCAL: Synthesizing Queries for Compositional Video Events from Limited User Interactions. *Proc. VLDB Endow.* 16, 11 (jul 2023), 2714–2727. https://doi.org/10.14778/3611479.3611482