# LVD-2M:
# A Long-take Video Dataset with Temporally Dense Captions

Tianwei Xiong[*], Yuqing Wang[*], Daquan Zhou, Zhijie Lin, Jiashi Feng, Xihui Liu[✉]

Project page: https://silentview.github.io/LVD-2M/

# LVD-2M Features

- Long videos ≥ 10 seconds
- Long-take videos without cuts
- Large motion and diverse contents
- Temporally dense captions

# LVD-2M Caption



The video depicts a person riding a mountain bike through a mountainous landscape with lush greenery and rocky terrain. The rider, wearing a black outfit, navigates a dirt and rocky trail, maneuvering the bike with control and skill. The camera perspective provides a first-person view of the ride, capturing the rider's descent down the mountain. The video appears to document an outdoor adventure, showcasing the natural beauty of the surroundings.
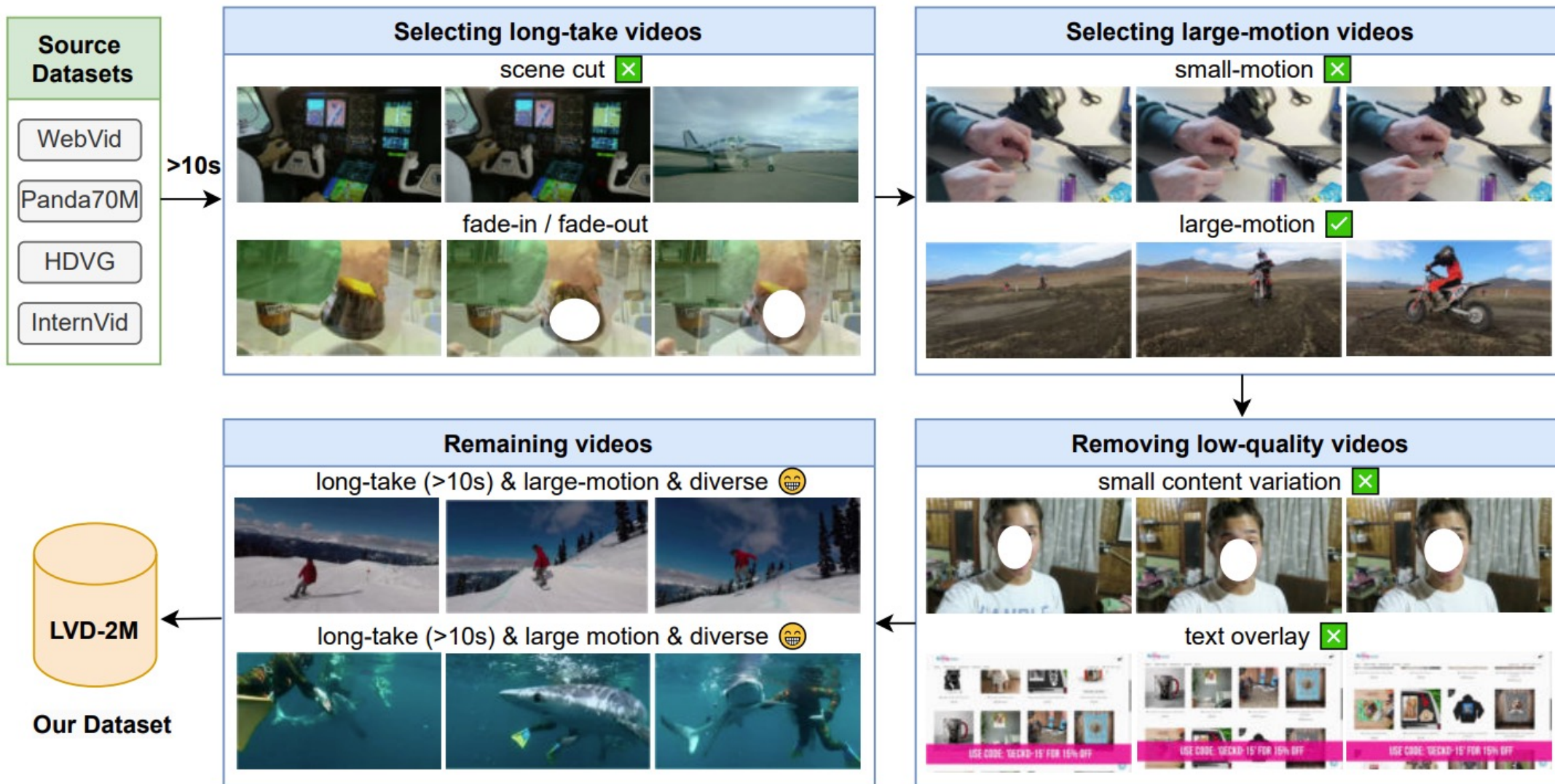
# LVD-2M Caption



The video depicts a person learning to snowboard at an indoor snowboarding facility. The learner is wearing a blue jacket and is being assisted by another individual, both wearing snowboarding gear including helmets and goggles. The video shows the learner's progress, starting with the person holding the instructor's hand for support, then leaning forward as they begin to fall, and eventually falling onto the snow-covered ground, with the instructor remaining nearby to offer guidance.
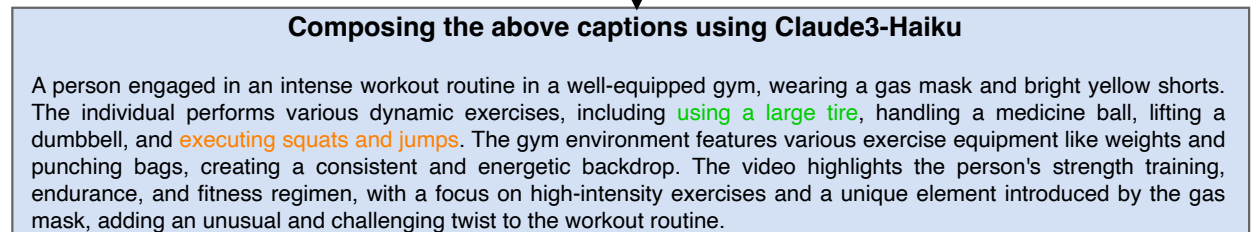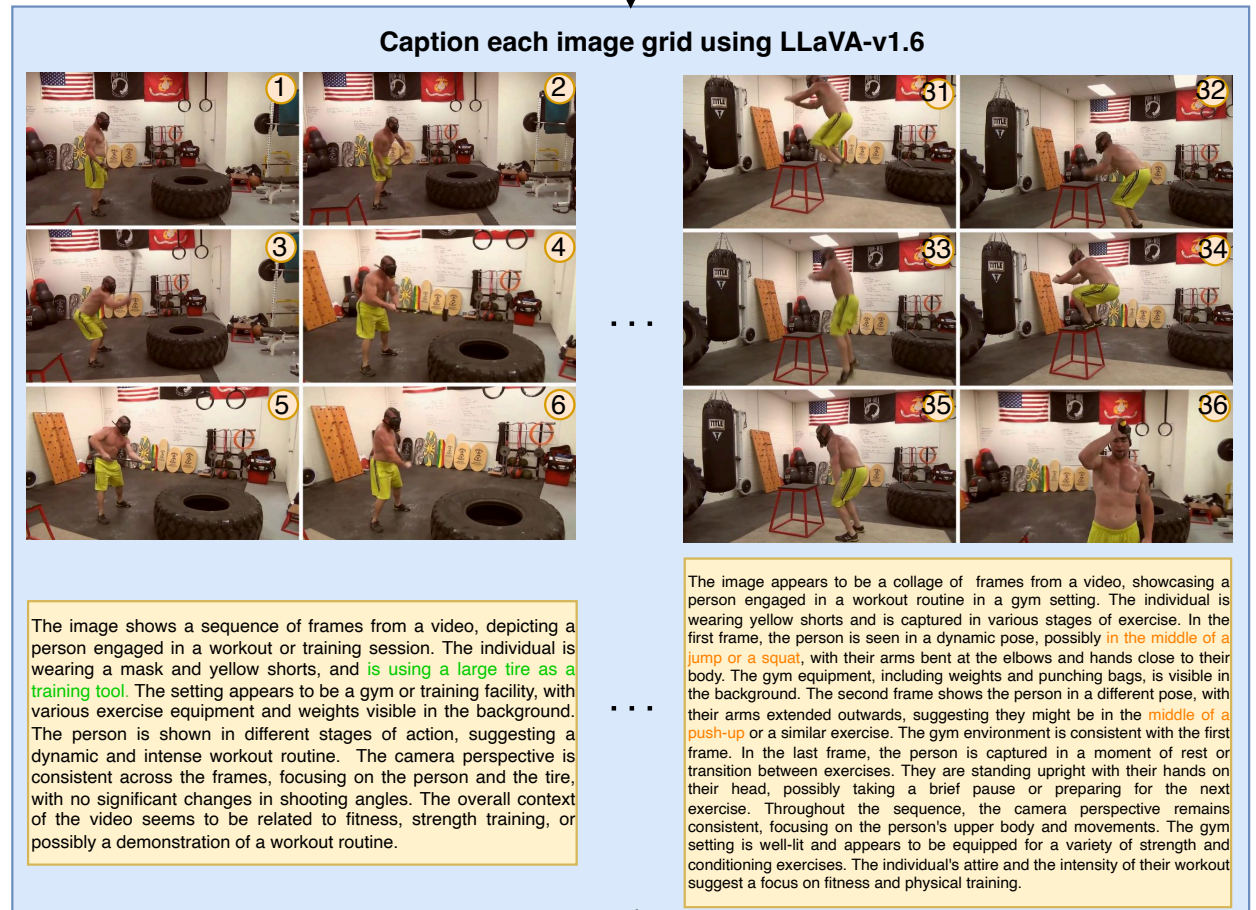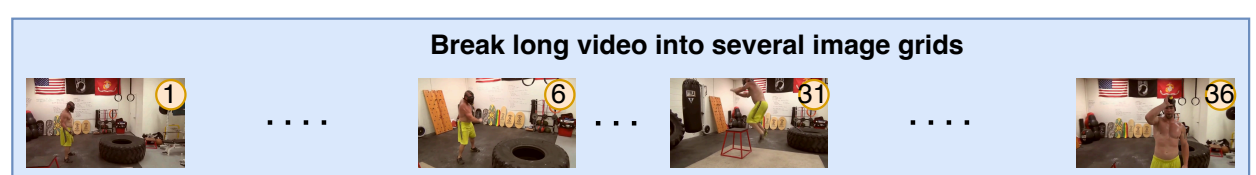
# LVD-2M Caption



The video captures an intense basketball game or practice session on an indoor court with distinctive orange, black, purple, and white color schemes. The players, dressed in athletic attire, engage in various actions such as dribbling, shooting, passing, defending, and interacting with each other. The court, equipped with multiple hoops, serves as a dynamic backdrop for the ongoing game, highlighting the competitive and active nature of the sport.

# LVD-2M Video Data Pipeline
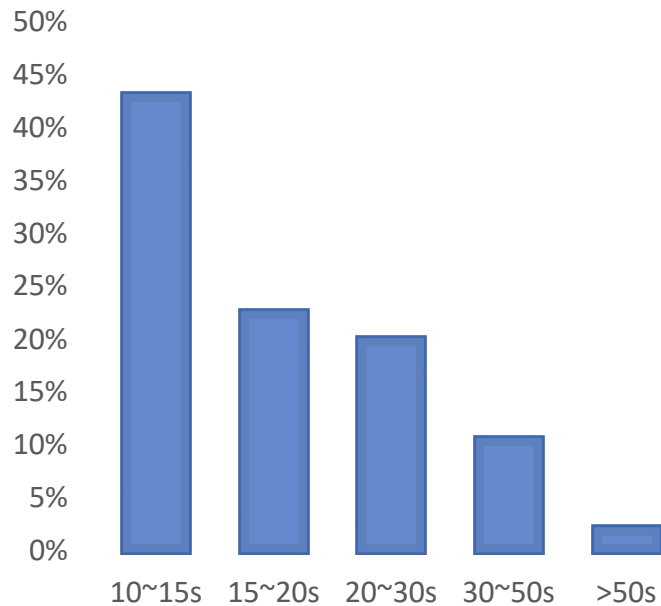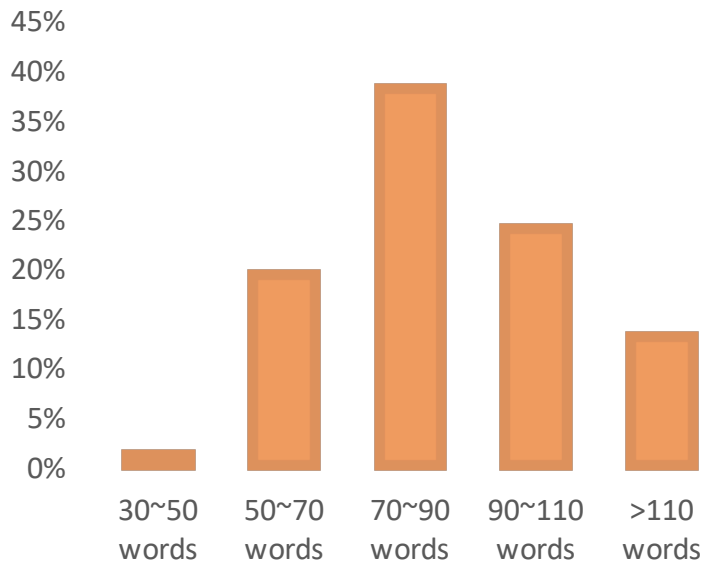
# LVD-2M for Model Finetuning

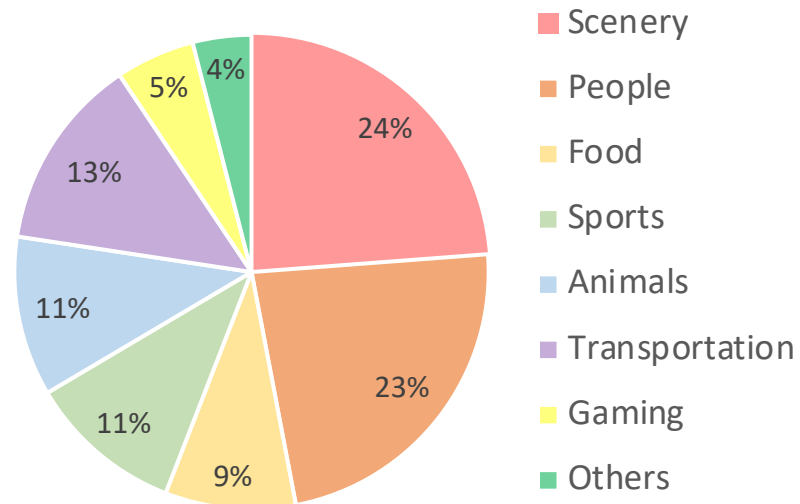Long-frame-length Fine-tuning on WebVid-10M[1]



Long-frame-length Fine-tuning on LVD-2M



[1] Bain, Max, et al. "Frozen in time: A joint video and image encoder for end-to-end retrieval." ICCV. 2021.

# LVD-2M for Model Finetuning

Vbench[1] evaluation for the two finetuned diffusion-based T2V models on LVD-2M and WebVid-10M [2] separately. Metrics exhibiting an absolute difference greater than 8% between the two models are underlined for emphasis.
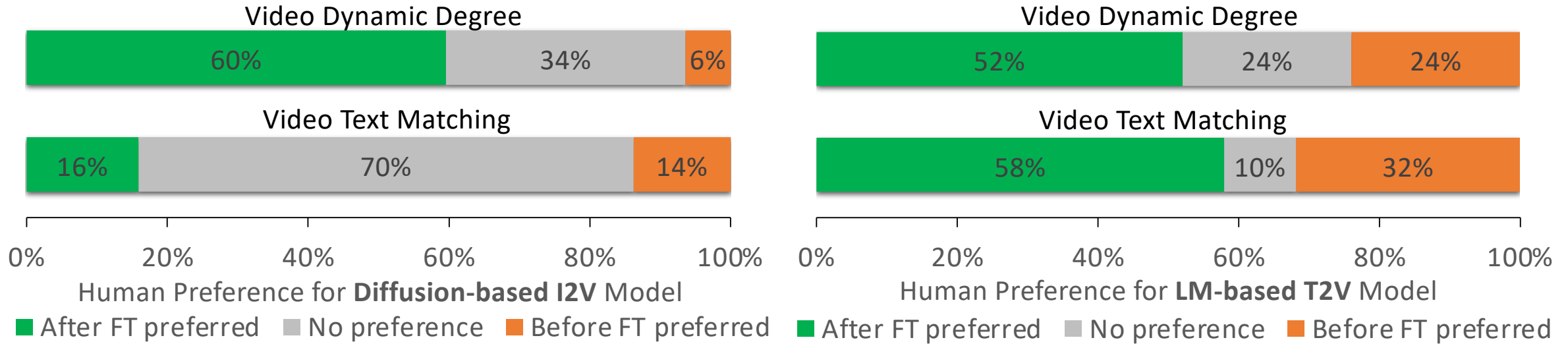
| Finetuning Dataset | Subject Consistency | Background Consistency | Temporal Flickering | Motion Smoothness | Dynamic Degree | Aesthetic Quality | Imaging Quality | Object Class |
|---|---|---|---|---|---|---|---|---|
| WebVid-10M | 95.81% | **98.02%** | **98.00%** | 97.87% | 20.00% | **58.02%** | **72.63%** | 76.95% |
| LVD-2M | **96.12%** | 96.92% | 97.44% | **98.43%** | **28.06%** | 57.56% | 70.72% | **86.93%** |

| Finetuning Dataset | Multiple Objects | Human Action | Color | Spatial Relationship | Scene | Appearance Style | Temporal Style | Overall Consistency |
|---|---|---|---|---|---|---|---|---|
| WebVid-10M | **26.02%** | 61.40% | 75.51% | 51.06% | 29.19% | 20.12% | 19.34% | **21.43%** |
| LVD-2M | 22.76% | **76.20%** | **79.32%** | **51.40%** | **32.95%** | **20.60%** | **20.25%** | 21.29% |

[1] Huang, Ziqi, et al. "Vbench: Comprehensive benchmark suite for video generative models." CVPR. 2024.
[2] Bain, Max, et al. "Frozen in time: A joint video and image encoder for end-to-end retrieval." ICCV. 2021.

# LVD-2M for Model Finetuning

**Video Dynamic Degree** (Diffusion-based I2V): 60% After FT preferred, 34% No preference, 6% Before FT preferred

**Video Text Matching** (Diffusion-based I2V): 16% After FT preferred, 70% No preference, 14% Before FT preferred

**Human Preference for Diffusion-based I2V Model**

Legend: After FT preferred · No preference · Before FT preferred

**Video Dynamic Degree** (LM-based T2V): 52% After FT preferred, 24% No preference, 24% Before FT preferred

**Video Text Matching** (LM-based T2V): 58% After FT preferred, 10% No preference, 32% Before FT preferred

**Human Preference for LM-based T2V Model**

Legend: After FT preferred · No preference · Before FT preferred

Project Page: https://silentview.github.io/LVD-2M/