# Video Foundation Model

# Text-to-Video Generation

**PROMPT:**

*"A man gives a small dog a bath in a shower."*

# Playable Video Generation - 2021

**Unlabeled videos**

Generative Model

Dynamics Model

right

stay

forward

right

forward

[Menapace et al. CVPR 2021]

# Playable Video Generation - 2021

# Controllable Video Synthesis

➢ Control the generation of (coarse) actions

➢ Control the camera

➢ Control the style

# Text-Driven Video Synthesis



The player sidesteps to the left and stops almost in the middle behind the baseline

The player moves to the left and hits with another forehand to the right side of the field

# Method

# Efficient Video Generation: Snap Video



"[...] a movie set where an otter serves as a film director. [...] with furrowed brows and raised paws shouting 'Action!'"



"Two unicorns in armors are playing a game of chess, in a medieval castle, high definition, photo-realistic style."

[Menapace et al. CVPR 2024]

# Efficient Video Generation: Snap Video

Transformer–based video diffusion models



(a) Computational Paradigms for Videos

(b) Snap Video FIT Architecture

[Menapace et al. CVPR 2024]

## Object-Centric Video Manipulation

We edit an object of a real video modifying the **shape** through a keyframe and the **appearance** through a driving image.



[Peruzzo et al. *under submission*]

# Object-Centric Video Manipulation

# Video-Language Alignment

Lin, Zhiqiu, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. "Evaluating text-to-visual generation with image-to-text generation." In *ECCV* 2024.
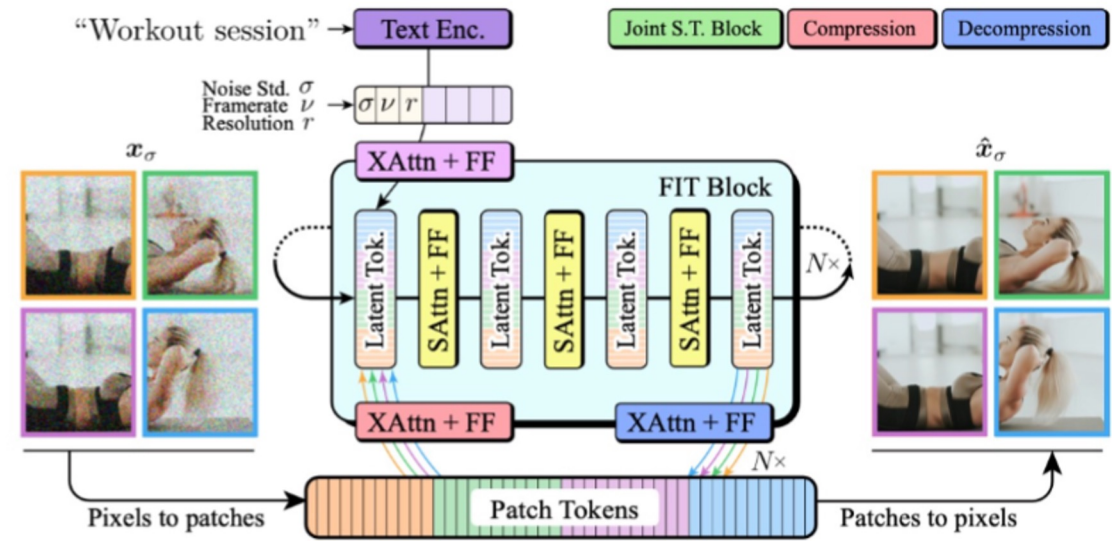Bansal, Hritik, Yonatan Bitton, Idan Szpektor, Kai-Wei Chang, and Aditya Grover. "Videocon: Robust video-language alignment via contrast captions." In *CVPR* 2024.

# Video-Language Alignment

# Synthetic Videos for Video-Language Alignment

# Method

# Results

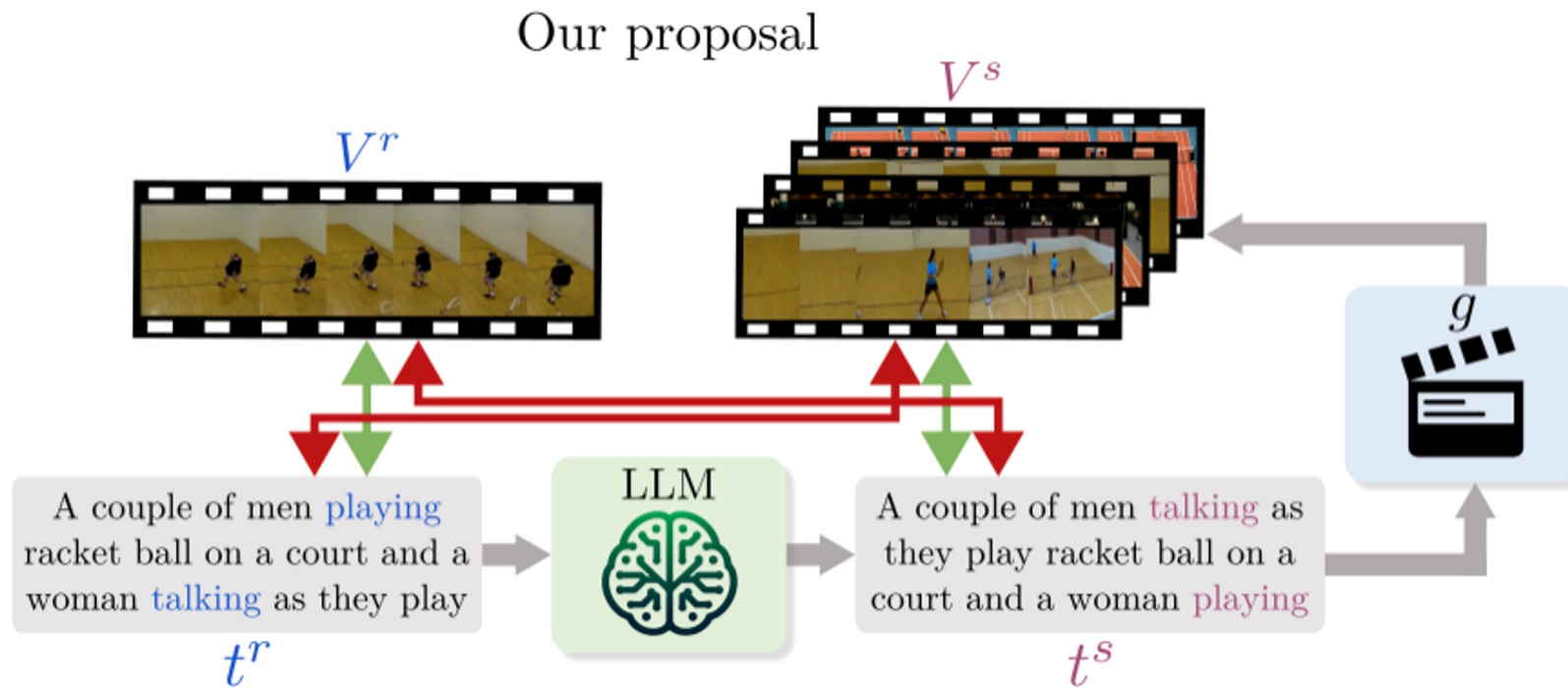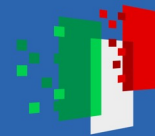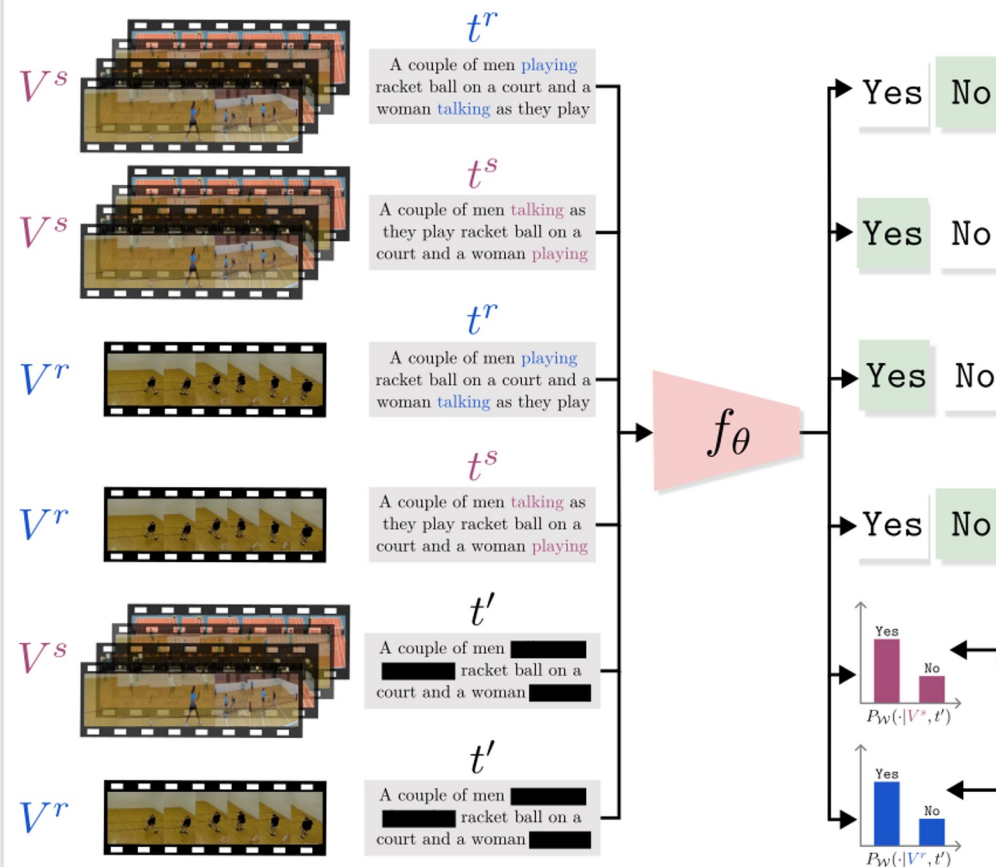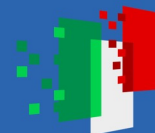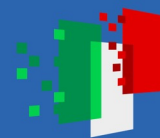| | VIDEO-LANGUAGE ENTAILMENT (VIDEOCON) | | | TEXT-TO-VIDEO RETRIEVAL | | VIDEO QA |
|---|---|---|---|---|---|---|
| | LLM | Human | Human-Hard | SSv2-Temporal | SSv2-Events | ATP-Hard |
| VIDEO-LLAVA (Lin et al., 2023) | 62.96 | 70.21 | 65.88 | 11.51 | 7.60 | **39.11** |
| VIDEO-LLAVA (VIDEOCON) (Lin et al., 2023) | 80.90 | **78.41** | 73.73 | 16.15 | **9.91** | 39.03 |
| SYNVITA (VIDEO-LLAVA) | **81.19** | 78.32 | **75.46** | **17.14** | 8.46 | 37.70 |
| MPLUG-OWL 7B (Ye et al., 2023) | 57.24 | 67.02 | 64.39 | 11.08 | 6.75 | 37.96 |
| MPLUG-OWL 7B (VIDEOCON)* (Bansal et al., 2023) | **88.39** | 77.16 | 74.76 | 13.00 | 10.37 | 35.46 |
| SYNVITA (MPLUG-OWL 7B) | 83.61 | **78.95** | **76.38** | **14.40** | **11.15** | **38.12** |

[Zanella et al. Under submission]

# Conclusions

➢ Numerous models for video generation (with many more on the horizon)

➢ Synthetic videos offer immense value, also for improving video understanding systems.

➢ Current challenges: misalignment between generated videos and prompts, scalability issues, appearance gap between synthetic and real videos.

➢ Wide range of applications: video surveillance, multimedia analysis, robotics, etc.