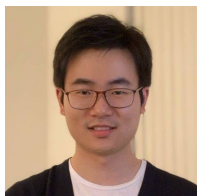Animatable version:
https://tsujuifu.github.io/slides/cvpr23_tvc.ppsx

# Tell Me What Happened: Unifying Text-guided Video Completion via Multimodal Masked Video Generation
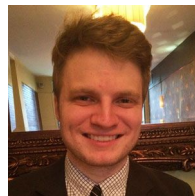
**Tsu-Jui Fu**[1]  Licheng Yu[2]  Ning Zhang[2]  Cheng-Yang Fu[2]
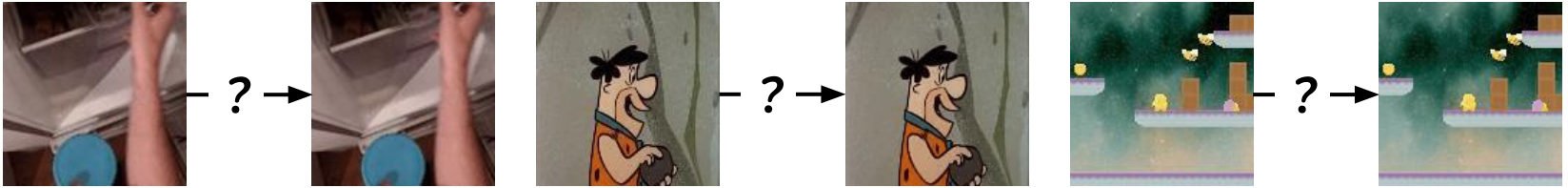
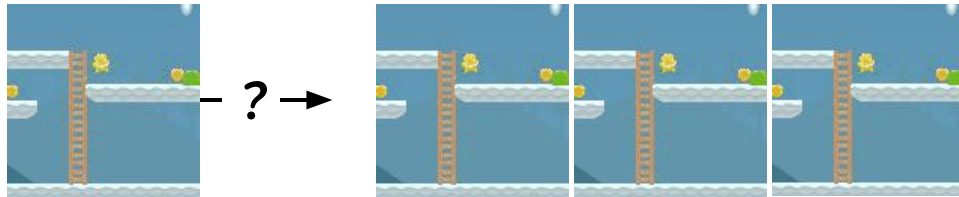Jong-Chyi Su[3]  William Wang[1]  Sean Bell[2]

[1]UC Santa Barbara, [2]Meta, [3]NEC Lab

# Video Prediction

- Generate future video frames, **given the past**
  - Maintain **reasonable continuation** and **temporal coherence**
  - Wide range of video applications (*e.g.*, **compression** / **autonomous** / **VR**)
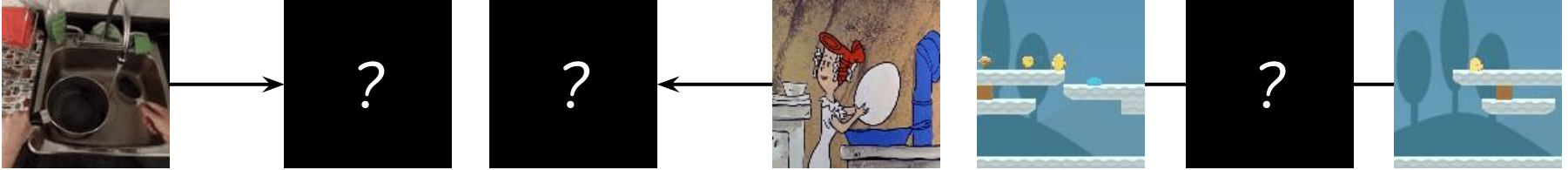


- **Uncontrollable**: there can be different outcomes

# Video Completion

- **Not just the first frame** can guide
  - Complete a video from **partial frames at arbitrary time points**
  - **Prediction** (first) / **Rewind** (last) / **Infilling** (head-tail)
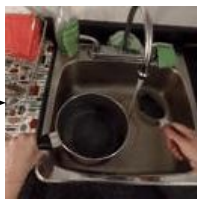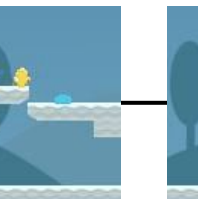
# Text-guided Video Completion

- **Not just the first frame** can guide
    - Complete a video from **partial frames at arbitrary time points**
    - **Prediction** (first) / **Rewind** (last) / **Infilling** (head-tail)

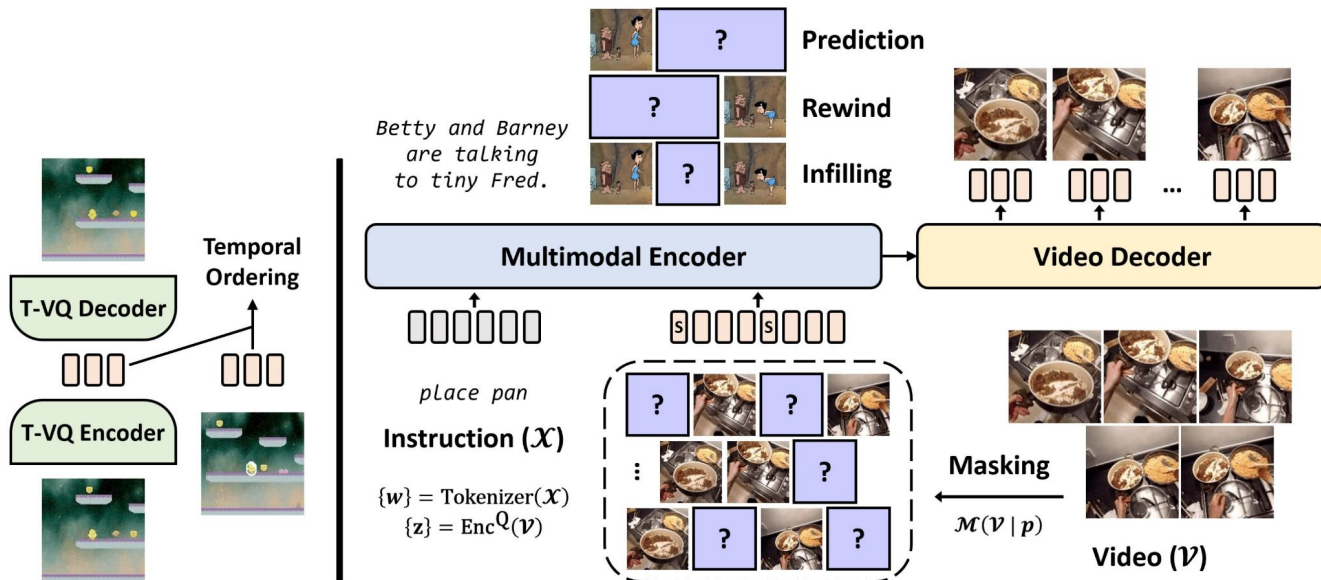- Use language to **describe missing event**



*pour water in pot*

*Wilma turns her head then she speaks.*

*Mugen runs to left. Then collects a coin and a gem.*

# Multimodal Masked Video Generation (MMVG)

- For **training**, we have full video (**V**) and caption (**X**)
- During **inference**, **only partial V and X** are provided

- **Temporal-aware** discrete video representation
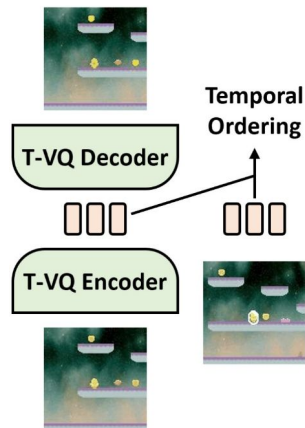- **Mask-then-fill learning** to unify video completion

# Temporal-aware (T-VQ)

- Built upon standard training of vector quantization (VQ)
  - $v' = \text{Dec}^Q(q(\text{Enc}^Q(v) \mid C))$

- **Inject temporal relationship** into discrete tokens (z)
  - Contrastive **temporal ordering**
  - Lead to **smooth video modeling**

$$z_i = q(\text{Enc}^Q(v_i) \mid C)$$

$$\hat{v}_i = \text{Dec}^Q(z_i)$$

$$\mathcal{L}_{\text{VQ}} = \underbrace{||\hat{v}_i - v_i||_1}_{\text{reconstrcution}} + \underbrace{||\text{sg}[\text{Enc}^Q(v_i)] - C_{z_i}||_2^2}_{\text{codebook}}$$

$$+ \underbrace{\beta||\text{sg}[C_{z_i}] - \text{Enc}^Q(v_i)||_2^2}_{\text{commit}} + \underbrace{||\mathcal{F}(\hat{v}_i) - \mathcal{F}(v_i)||_1}_{\text{matching}}$$

$$o_i = \text{FC}^T(z_i, z_j)$$

$$\mathcal{L}_T = \text{BCELoss}(o_i, 0 \text{ if } i > j \text{ else } 1)$$



Temporal Ordering

T-VQ Decoder

T-VQ Encoder

# Generation from Masked Video

- Masking strategy (*M*) **masks out** V with probability (*p*)
  - **Replace each fragment** as unique [SPAN]
  - For example: *M* **reserves the 3<sup>rd</sup> and 5<sup>th</sup>** for 5-length

- Multimodal Encoder (Enc$^{\text{M}}$) performs **cross-modal fusion**
  - **Tokenize** V via Enc$^{\text{Q}}$ and **X**
  - **Extract** encoding feature (*h*) via Transformer
- Video Decoder (Dec$^{\text{M}}$) produces **entire V**
  - **Auto-regressively decode** tokens upon *h*
  - **Reconstruct** all frames via Dec$^{\text{Q}}$

$$\overline{\mathcal{V}} : \{[\text{S}], v_3, [\text{S}], v_5\} = \mathcal{M}(\mathcal{V} \mid p)$$
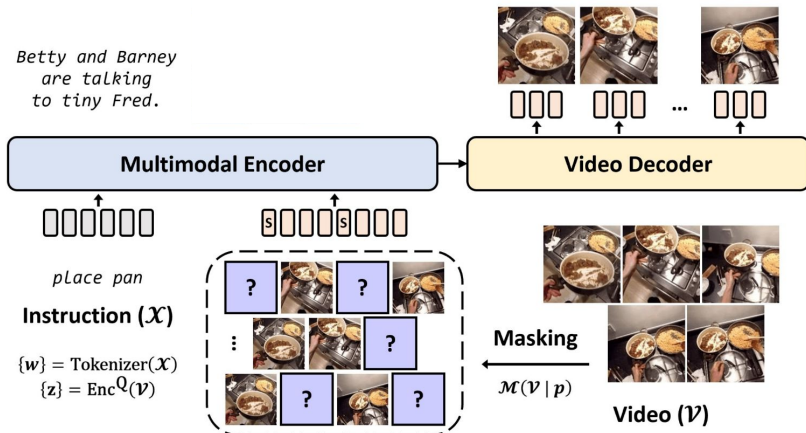
$$\{[\text{S}], \text{z}_3, [\text{S}], \text{z}_5\} \quad f_i^w, f_j^v = \text{LP}^w(w_i), \text{LP}^v(\text{z}_j)$$

$$\{h\} = \text{Enc}^{\text{M}}([\{f^w\}, \{f^v\}])$$

$$\hat{z}_t = \text{Dec}^{\text{M}}(\{\hat{z_1}, ..., \hat{z_{t-1}}\} \mid \{h\})$$

$$\mathcal{L}_t = \text{CELoss}(\hat{z}_t, \text{z}_t)$$

$$\hat{\mathcal{V}} = \text{Dec}^{\text{Q}}(\{\hat{z}\}_{t=1}^{N})$$



Betty and Barney are talking to tiny Fred.

**Multimodal Encoder**

**Video Decoder**

*place pan*

**Instruction (𝒳)**

$\{w\} = \text{Tokenizer}(\mathcal{X})$
$\{\text{z}\} = \text{Enc}^{\text{Q}}(\mathcal{V})$

**Masking**

$\mathcal{M}(\mathcal{V} \mid p)$

**Video (𝒱)**

# Generation from Masked Video

- Masking strategy (*M*) **masks out** V with probability (*p*)
  - **Replace each fragment** as unique [SPAN]
  - For example: *M* **reserves the 3$^{rd}$ and 5$^{th}$** for 5-length

- Multimodal Encoder (Enc$^M$) performs **cross-modal fusion**
  - **Tokenize** V via Enc$^Q$ and **X**
  - **Extract** encoding feature (*h*) via Transformer
- Video Decoder (Dec$^M$) produces **entire V**
  - **Auto-regressively decode** tokens upon *h*
  - **Reconstruct** all frames via Dec$^Q$

- **Unifying video completion** during inference
  - **Prediction**: [{*w*}, {z$_1$, [S]}]
  - **Rewind**: [{*w*}, {[S], z$_N$}]
  - **Infilling**: [{*w*}, {z$_1$, [S], z$_N$}]



Betty and Barney are talking to tiny Fred.

Prediction

Rewind

Infilling

**Multimodal Encoder** → **Video Decoder**

*place pan*

**Instruction ($\mathcal{X}$)**

$\{w\} = \text{Tokenizer}(\mathcal{X})$
$\{z\} = \text{Enc}^Q(\mathcal{V})$

**Masking**

$\mathcal{M}(\mathcal{V} \mid p)$

**Video ($\mathcal{V}$)**

# Experimental Setup

- **Datasets**
  - Kitchen / Flintstones / MUGEN

Resolution: $128^2$

| Dataset | Domain | #Frame | #Word | FPS |
|---|---|---|---|---|
| Kitchen | Egocentric | 8.3 | 2.8 | 6 |
| Flintstones | Animation | 15 | 16.5 | 5 |
| MUGEN | Gaming | 16 | 20.6 | 5 |



*pick up meat*    *open fridge*

*Fred is driving. Then Barney talks to Fred.*    *Barney talks to Betty in a room.*

*Mugen jumps again. Move to the right to collect a gem.*    *Mugen climbs up and jumps to collect the coin.*

# Experimental Setup

- **Datasets**
  - Kitchen / Flintstones / MUGEN

- **Evaluation Metrics**
  - **FVD (↓)**: distance of **I3D video feature** (*vs.* GT)
  - **RCS (↑)**: **relative visual-text similarity** from CLIP (*vs.* instruction)

# Experimental Setup

- **Datasets**
  - Kitchen / Flintstones / MUGEN

- **Evaluation Metrics**
  - **FVD (↓)**: distance of **I3D video feature** (*vs.* GT)
  - **RCS (↑)**: **relative visual-text similarity** from CLIP (*vs.* instruction)

- **Baselines**
  - **Auto-regressive VQ**: TATS (**requires specific training** for rewind / infilling)

# Text-guided Video Prediction



Kitchen — FVD↓: TATS w/o text 106.9, MMVG w/o text 103.8; RCS↑: TATS w/o text 64.4, MMVG w/o text 64.5

Flintstones — FVD↓: TATS w/o text 127.5, MMVG w/o text 123.8; RCS↑: TATS w/o text 60.3, MMVG w/o text 60.8

MUGEN — FVD↓: TATS w/o text 376.5, MMVG w/o text 369.4

Legend: TATS w/o text, MMVG w/o text, TATS, MMVG$^U$, MMVG$^S$

# Text-guided Video Prediction

# Text-guided Video Rewind

# Text-guided Video Infilling



*take plate*
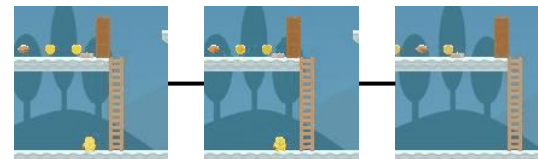
*Barney is trying to sing in a living room.*

*Mugen jumps up the stage. It runs from left to right and jumps on a worm.*

*take lid off*

*Wilma looks back at Betty who is speaking then she turns her head.*

*Mugen climbs up a ladder. It jumps onto a stack of boxes, drops down, and is killed by a worm.*
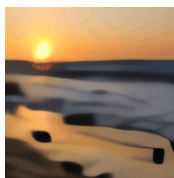
# Summary

- Text-guided Video Completion
  - **Control video generation from partial frames** via language command
- Multimodal Masked Video Generation (MMVG)
  - T-VQ for **temporal-aware** discrete frame tokens
  - **Unified masked training** to support all prediction / rewind / infilling
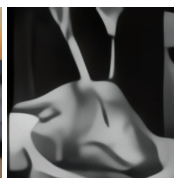  - **Benefit general video generation** as well



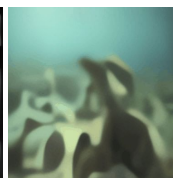*Mugen jumps from left to right to an upper platform.*

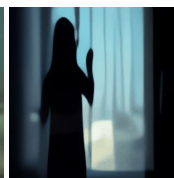*Mugen runs right to left and collects a gem.*
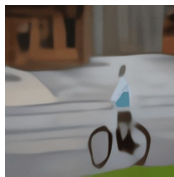
*Mugen walks to the right and jumps over a mouse to collect coins.*

*Mugen walks to the right and jumps to collect coins.*

sunset on the baltic sea

cut chicken with knife

green sea turtle swims and relaxes
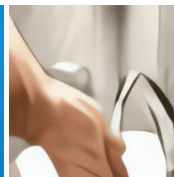
woman unveils curtain

child rides a bike

downtown city with traffic car

rotates apple lollipop

wash hand