## Let's think Frame by Frame with VIP: A Video Infilling and Prediction Dataset for Evaluating Video Chain-of-Thought

Vaishnavi Himakunthala\*, Andy Ouyang\*, Daniel Rose\*, Ryan He\*, Alex Mei, Yujie Lu, Chinmay Sonar, Michael Saxon, William Wang University of California, Santa Barbara



Motivation

UCSB NLP GROUP

- Investigate video reasoning (VR) abilities of vision-language (VL) systems, the next logical step after text and image
- Evaluate multi-hop, multi-frame reasoning on general real-life videos, missing in VR datasets
- Automate data creation to minimize



## **VIP** at a **Glance**

- An **inference time dataset** to assess models' video reasoning abilities
- Contains video keyframes from a wide variety of realistic videos
- Uses two textual representations for keyframes, including a new format of keyframe captioning – FAMOuS
- Introduces two new tasks, <u>Video</u>









				VIP D	ataset			
VIP's Video Categories	<u>VIP vs. Video Reasoning Datasets</u>							<b>Task Definition</b>
Hobbies & Leisure	Dataset	Frame	Structured	Domain	Vid. Len.	Cap. Len.	Test Samples	Video Infilling – Given n previous
Food & Drink	MSR-VTT	X	X	Open	$20.7\mathrm{s}$	9.6	3K	and next keyframes' descriptions
Business &	YouCook2	X	X	Cooking	$5.26\mathrm{m}$	8.8	$2\mathrm{K}$	and next keynames descriptions,
Industrial	ActyNet-Cap	X	×	Open	$2\mathrm{m}$	13.5	$5\mathrm{K}$	predict the p masked keyframes'
Pets &	HowTo100M	X	X	Instructional	18s	4	$24\mathrm{K}$	
Animals	VATEX	X	X	Open	10s	15.2	$6\mathrm{K}$	descriptions
Travel	VideoStory	$\checkmark$	X	Events	$12.6\mathrm{m}$	12.1	16	-
Autos &	WebVid-2M	X	X	Open	18s	12	$5\mathrm{K}$	Goal: infill <b>p</b> keyframes
Vehicles Home &	VIP	$\checkmark$	$\checkmark$	Open	$3.6\mathrm{m}$	114.2	$1.5\mathrm{K}$	





- Low scores demonstrate difficulty of our tasks on existing LLMs
- Marginal boost in performance with additional context frames
- Filler words/verbosity in dense captions lead similarity metrics to score them higher than FAMOuS descriptions
- Models perform better on the infilling task compared to the prediction task, as bidirectional context is a less complex task



- Low scores on dynamic components (Action, Focus), motivating other strategies for dynamic video reasoning
- Slightly stronger performance for basic components (Objects, Mood), likely due to a consistent background