# Revisiting the "Video" in Video-Language Understanding

Shyamal Buch[1], Cristóbal Eyzaguirre[1], Adrien Gaidon[2], Jiajun Wu[1], Li Fei-Fei[1], Juan Carlos Niebles[1]

1 STANFORD VISION & LEARNING    2 TOYOTA RESEARCH INSTITUTE

CVPR JUNE 19-24 2022 NEW ORLEANS • LOUISIANA

Project website: atp-video-language.stanford.edu

## Motivation

**Video understanding** offers the potential to go *beyond* image-level semantics (scenes, objects) towards event temporality + causality.
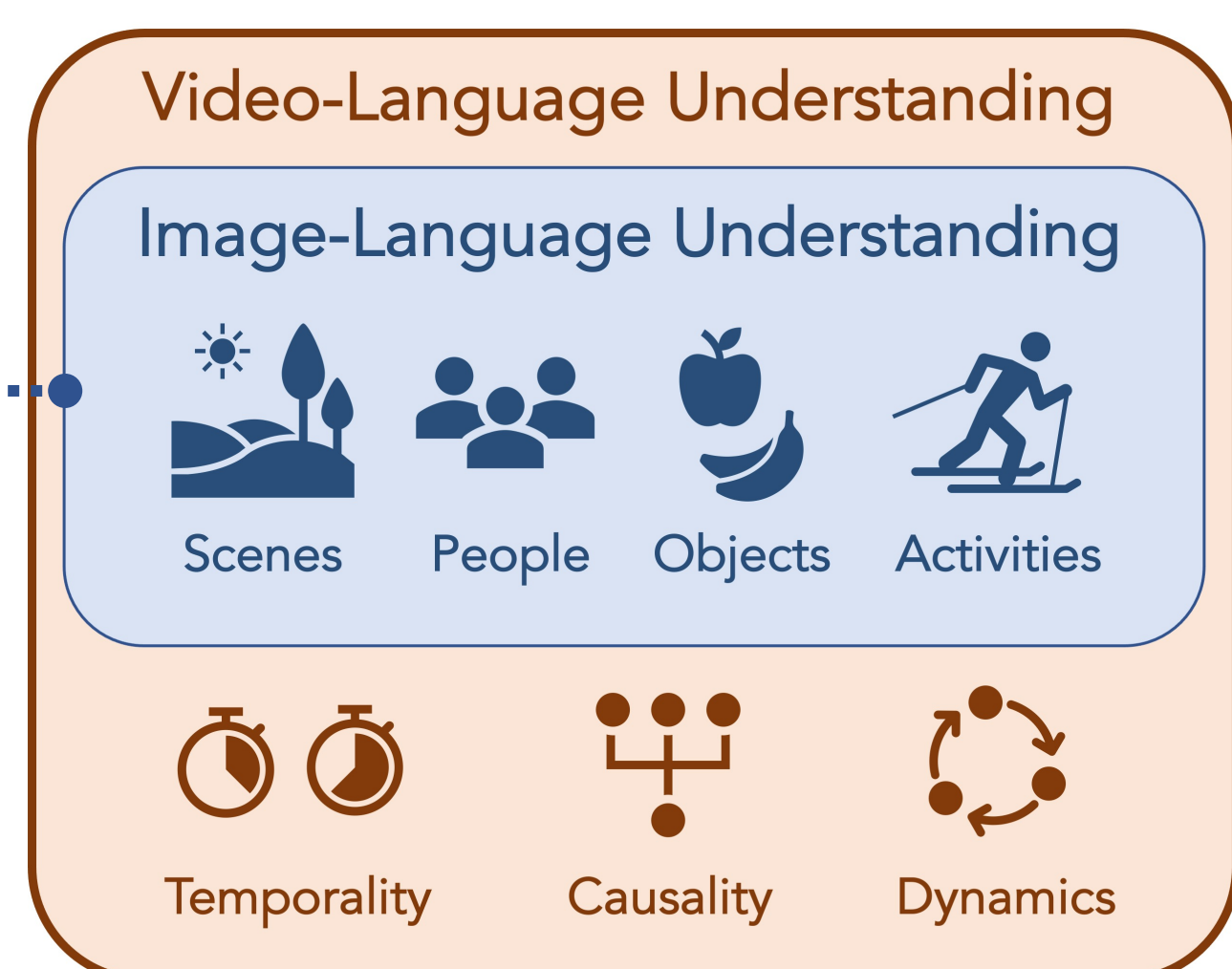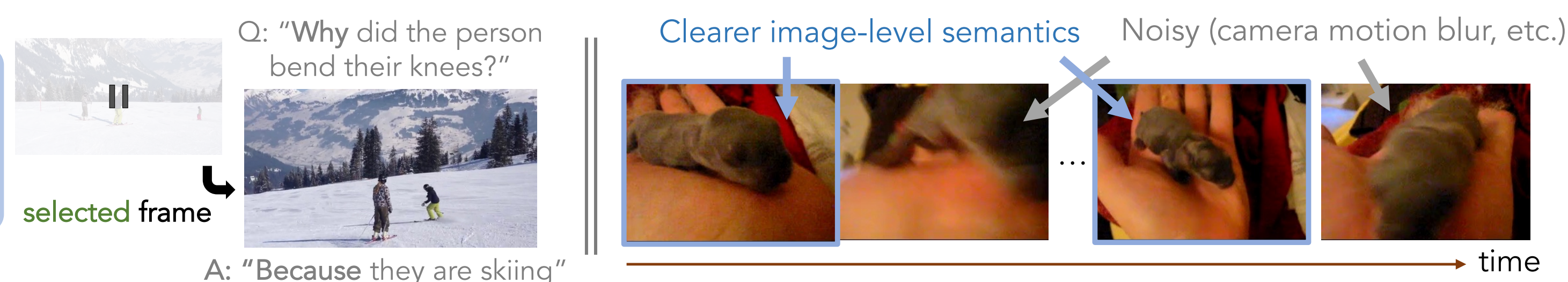
**Our work** re-examines a foundational question [5,6] in video research:

⭐ What makes a video task uniquely suited for **videos**, beyond what can be understood from a **single image**?

Our focus is on **video + language**, where *language* has the potential to describe richer event properties and relationships in *videos*.

**Challenge:** Standard approaches may **under-represent "image-centric" bound** for video-language understanding!

Standard approaches: Select a random frame? Average pool?

Q: "Why did the person bend their knees?"
selected frame
A: "Because they are skiing"

Clearer image-level semantics    Noisy (camera motion blur, etc.)

time

### Video-Language Understanding

Image-Language Understanding

Scenes    People    Objects    Activities

Temporality    Causality    Dynamics

## Atemporal Probe (ATP) for Video-Language Analysis

Standard Video-Language Benchmarks

ATP

How well does "image-centric" (from a **well-selected frame**) understanding address "video" understanding?

*Solvable with Image-Centric Understanding*

*Video-Level Understanding Required*

(a) Single Selected Encoding → Final Video-Language Task

Discrete Selection

Atemporal Probe (ATP)

*No temporal information*

Frozen Image Encoding

Frozen Language Encodings (Jointly Pre-trained with Image Model)

Pre-trained (frozen) image encoder $M_I$ , Pre-trained $M_L$

Sparse Frame Sample

Input Language

Input Video    time

Data

(b) Atemporal Probe (ATP)

*Output Selection*    (Image Encoding Set)

Selector

Selector Encoder (Atemporal)

*Unordered* set of Image Encodings    Language Encodings

Our **Atemporal Probe (ATP)** model builds on progress in self-supervised image-language understanding [1], and learns to **discretely select** a *frozen* image-level encoding (*without* using any temporal information). This encoding of a single frame is sent downstream – unmodified – to the **final video-language understanding task** (video question answering, text-to-video retrieval; e.g. [2,7]).
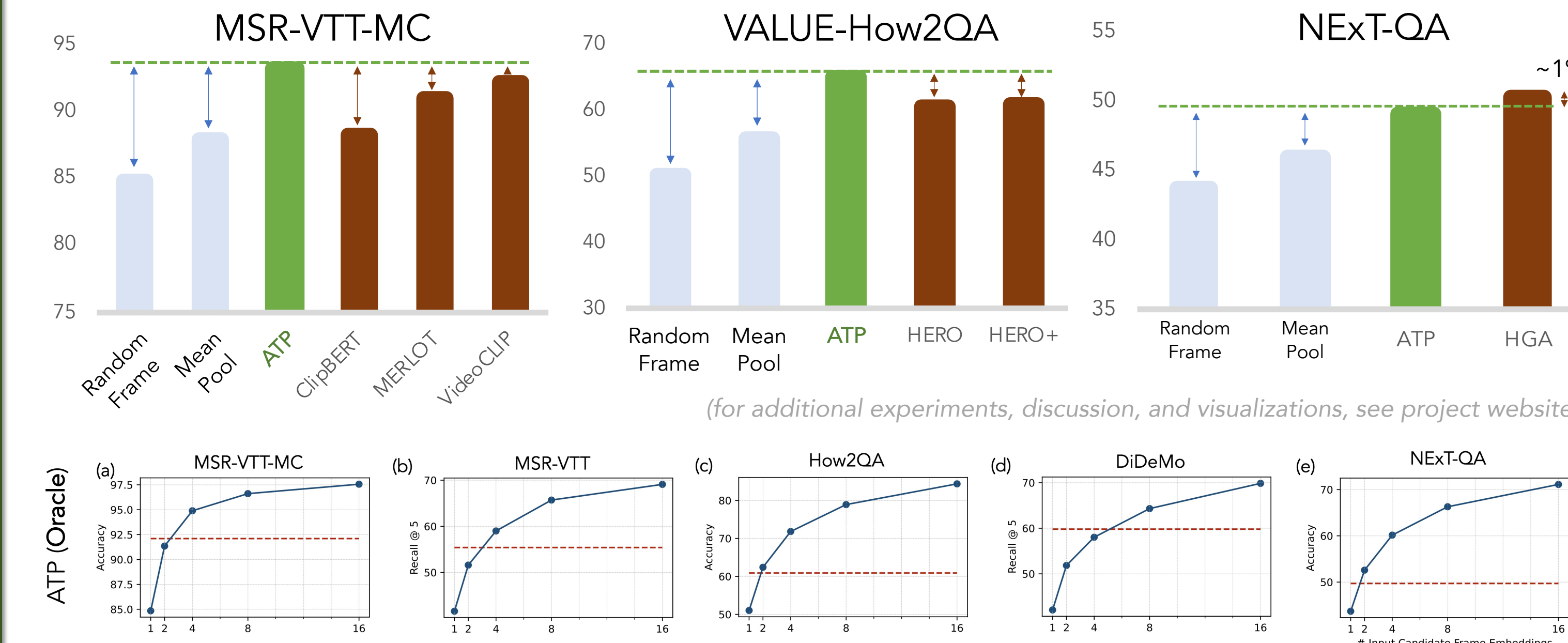
## ★ Overview

We propose the **Atemporal Probe (ATP)**:

✓ To analyze **current standard video-language** benchmarks *(stronger bound on image-centric understanding)*

... ? ATP

✓ **Improve dataset design** *(disentangling unintended biases)*

ATP

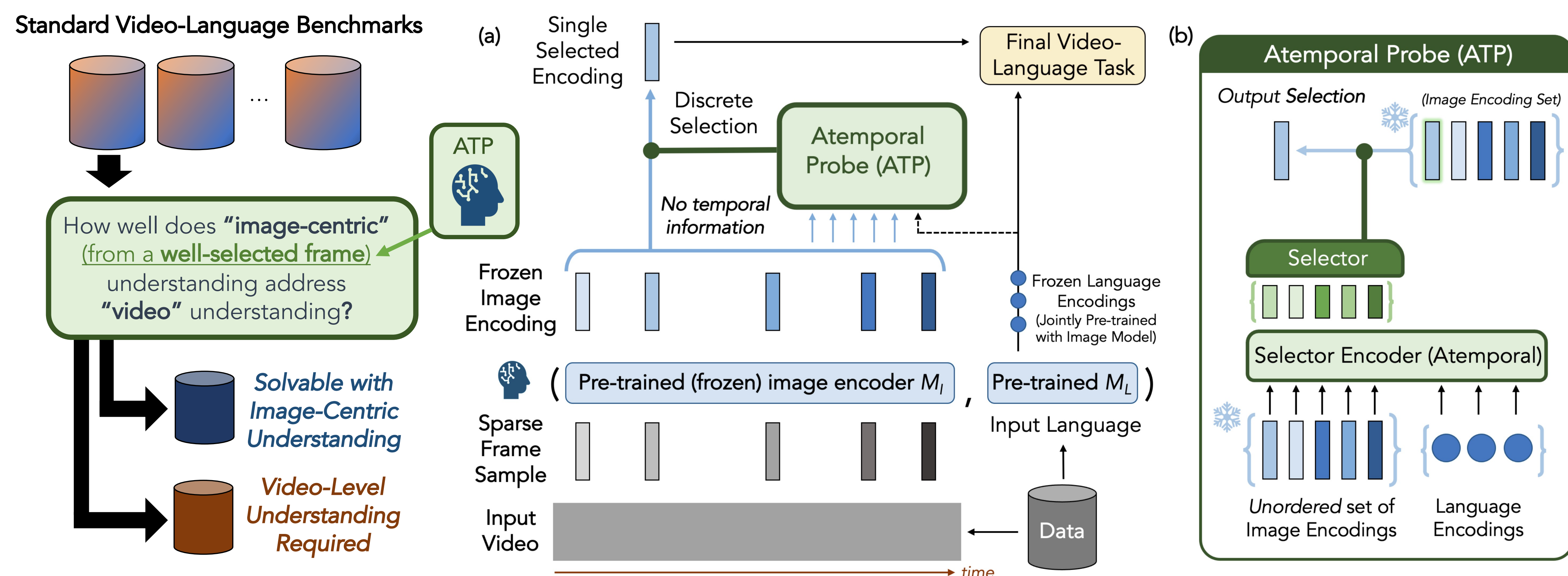✓ **Improve model design** *(better efficiency and accuracy)*

ATP → Video-level Reasoning Model

## Improving Model Design with ATP

ATP can **improve temporal (multi-frame) model design** by forwarding semantically useful candidates to reason over (reducing noise). A combined model with ATP achieves higher accuracy with fewer frames needed.

Output Video Encoding → Final Video-Language Task

Temporal Reasoning Model

Selected Image Encodings

ATP applied per-partition

Input Video (Partitioned)    time

NExT-QA

ATP    HGA    Temp[ATP]++

## Experiments and Analysis

MSR-VTT-MC

Random Frame / Mean Pool / ATP / ClipBERT / MERLOT / VideoCLIP

VALUE-How2QA

Random Frame / Mean Pool / ATP / HERO / HERO+

NExT-QA    ~1%

Random Frame / Mean Pool / ATP / HGA

*(for additional experiments, discussion, and visualizations, see project website)*

(a) MSR-VTT-MC  (b) MSR-VTT  (c) How2QA  (d) DiDeMo  (e) NExT-QA

ATP (Oracle)    # Input Candidate Frame Embeddings

**Takeaways:** *(hold even when dataset explicitly designed for temporal + causal video-language)*
1. **Datasets** can be (surprisingly) well-addressed by image-centric understanding
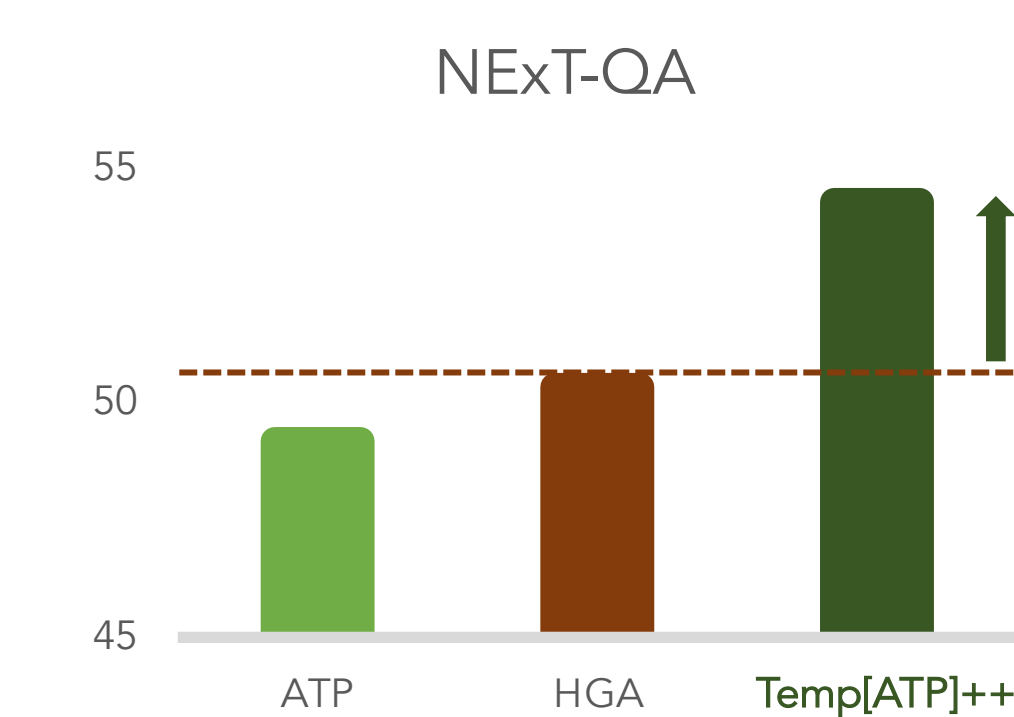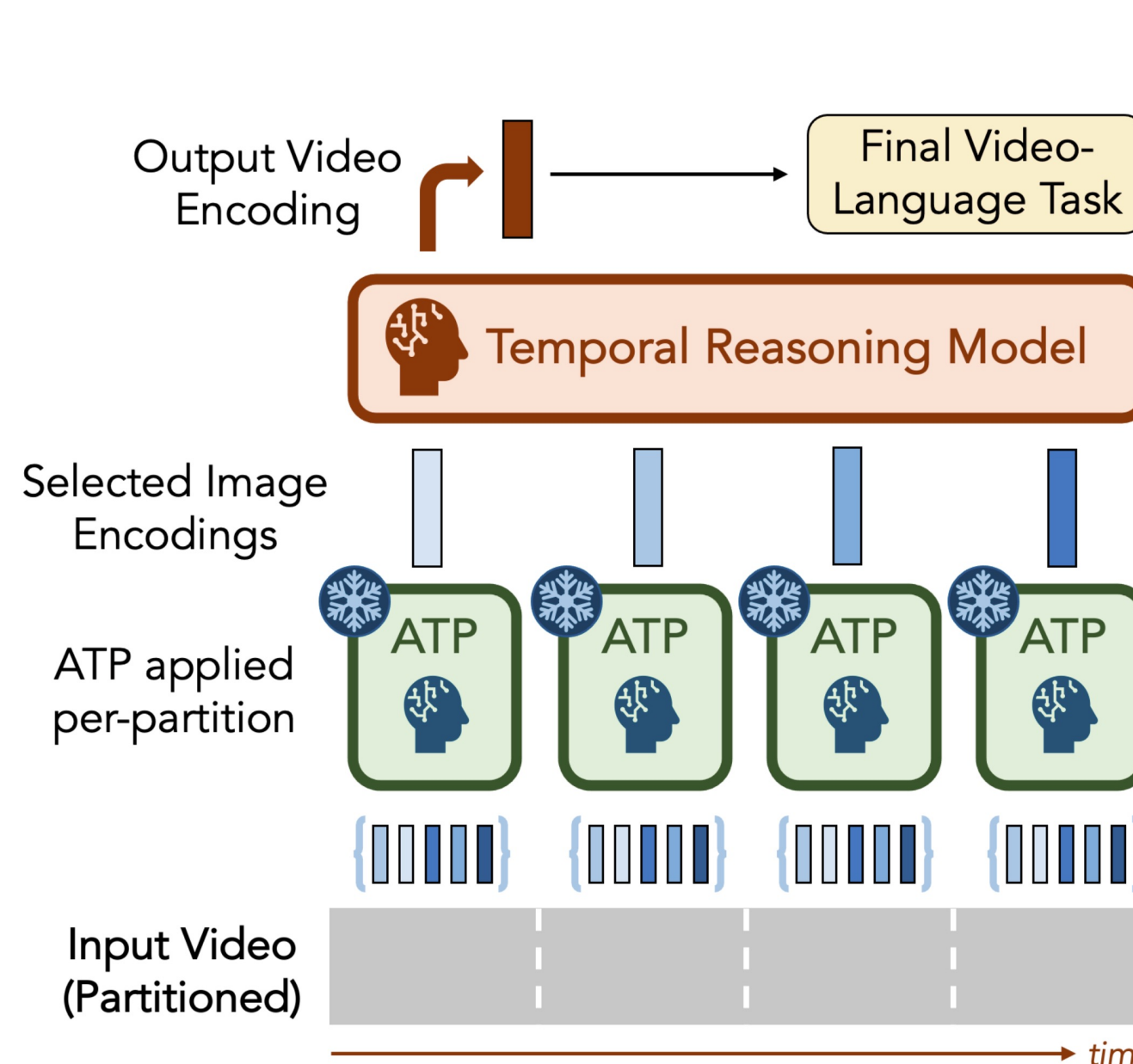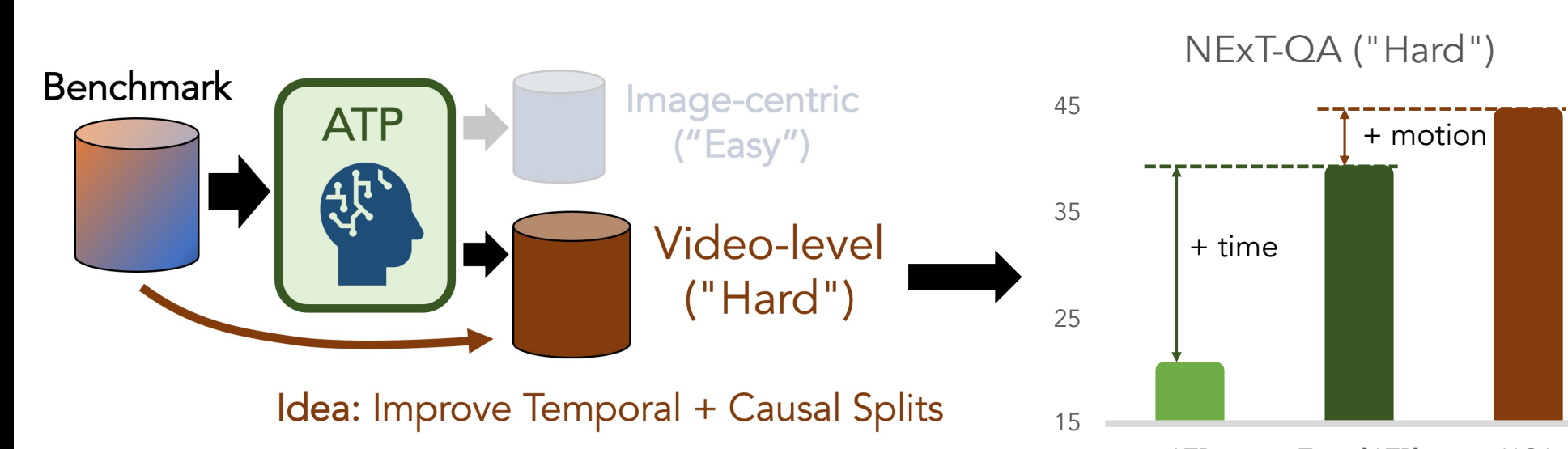2. Video-level **models** may be significantly impacted by processing noisy frames

## Improving Dataset Design with ATP

= ATP selected frame    = downstream selected answer (bold = labeled answer)

(a) Q: How is the girl moving around in the video? ✓
walking / jumping / cycling / swinging up and down

(b) Q: How do the two men play the instrument? ✗
roll the handle / strum the string / hit with sticks / pat with hand

ATP can help identify (multi-frame) temporally challenging data, and is a **promising tool** for in-the-loop **dataset design**.

Benchmark → ATP → Image-centric ("Easy") / Video-level ("Hard")

Idea: Improve Temporal + Causal Splits

NExT-QA ("Hard")

+ time  + motion

ATP    Temp[ATP]    HGA

**Selected Works Cited:** *(full list of references + acknowledgements in paper)*

[1] Radford et al. "Learning Transferable Visual Models from Natural Language Supervision." ICML 2021.
[2] Xiao et al. "NExT-QA: Next Phase of Question-Answering to Explaining Temporal Actions." CVPR 2021.
[3] Zellers*, Lu*, Hessel* et al. "MERLOT: Multimodal Neural Script Knowledge Models." NeurIPS 2021.
[4] Xu et al. "VideoCLIP: Contrastive Pre-training for Zero-Shot Video-Text Understanding." EMNLP 2021.
[5] Huang et al. "What Makes a Video a Video: Analyzing Temporal Information in Videos." CVPR 2018.
[6] Schindler and van Gool. "Action Snippets: How many frames does human action recognition require?" CVPR 2008
[7] Li*, Lei* et al. "VALUE: A Multi-task Benchmark for Video-and-Language Understanding Evaluation." NeurIPS 2021 (D&B)