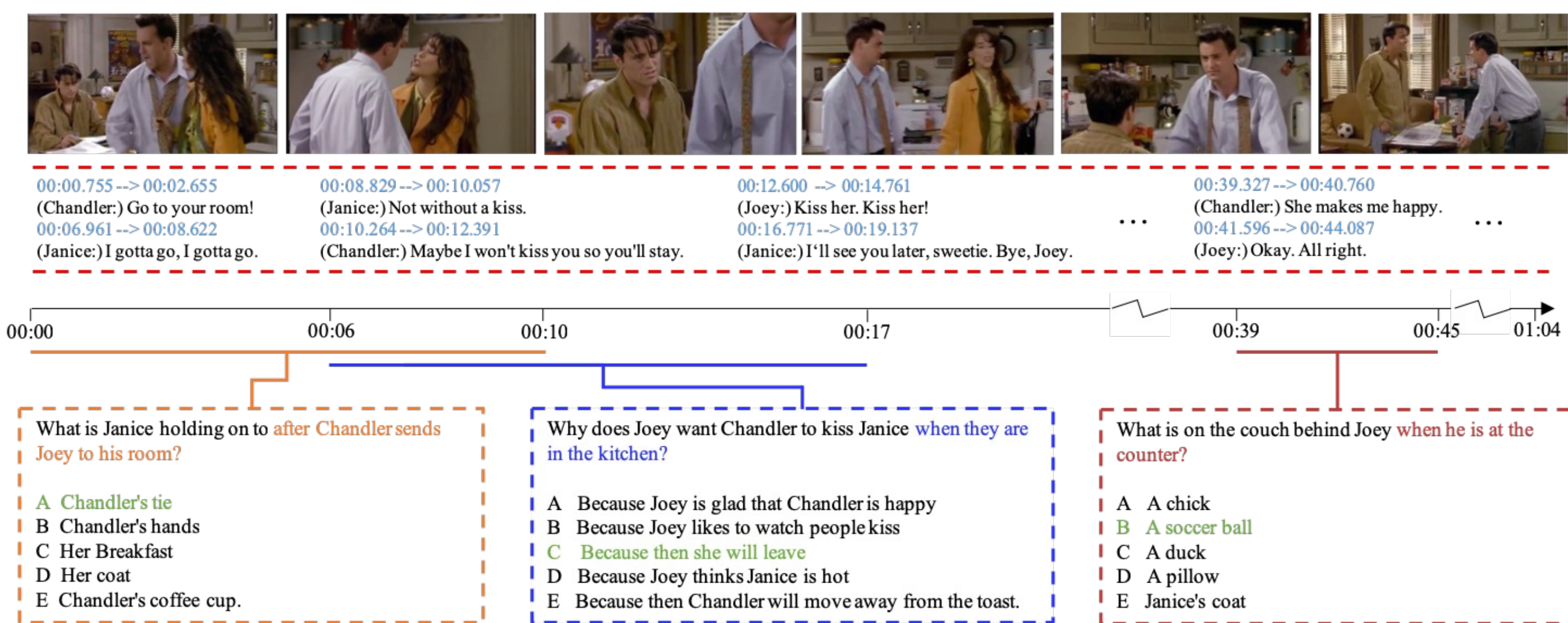# Self-Adaptive Sampling for Efficient Video Question-Answering on Image--Text Models

## Wei Han, Hui Chen, Min-Yen Kan, Soujanya Poria
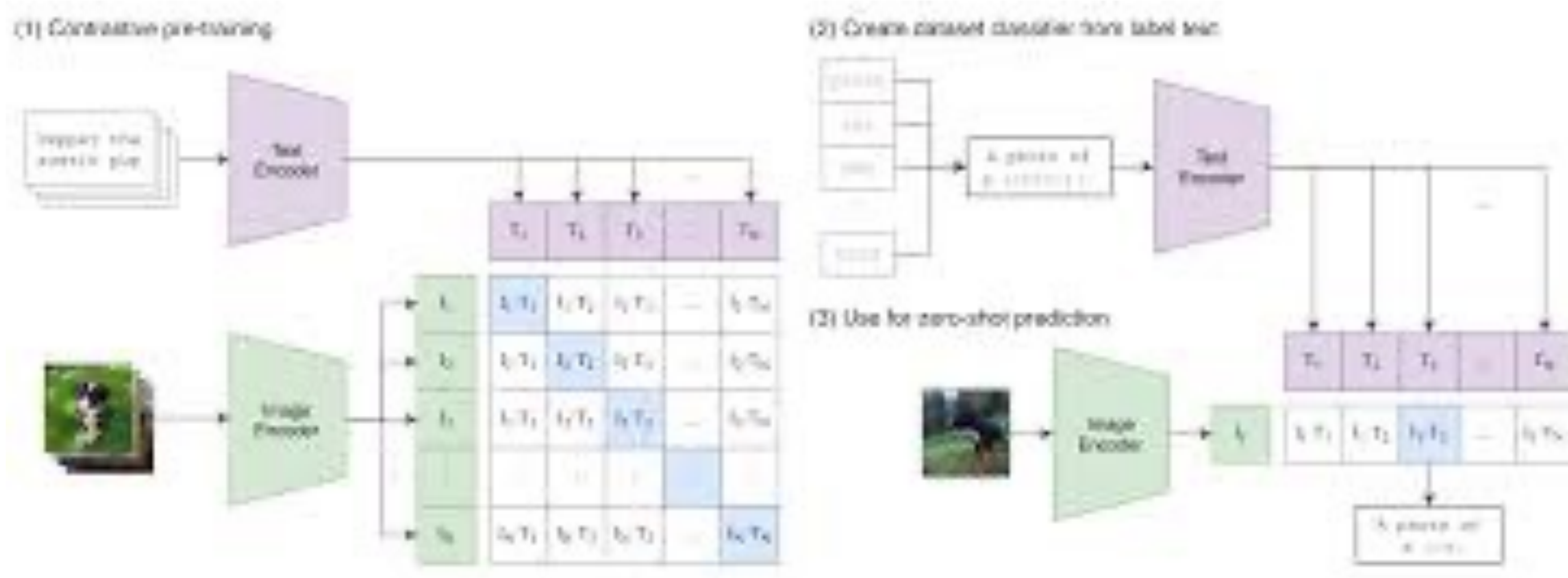
*Singapore University of Technology & Design, Singapore*
*National University of Singapore, Singapore*

## Introduction

- Video Question Answering (ViQA): Given a short video, answer the question based on the video



- Image—Text Models (ITMs): a subclass of visual language models (VLMs) that accept image sequences and text as input and generate text outputs, such as CLIP. To process video input, a series of frames must be first sampled from that video.



## Related Works

- Current Sampling Strategy
  - Learning-free sampling is cost-effective but hard to reach optimal
  - Learning-based sampling can adapt to different question input, but requires additional computational cost (huge) and difficult to converge
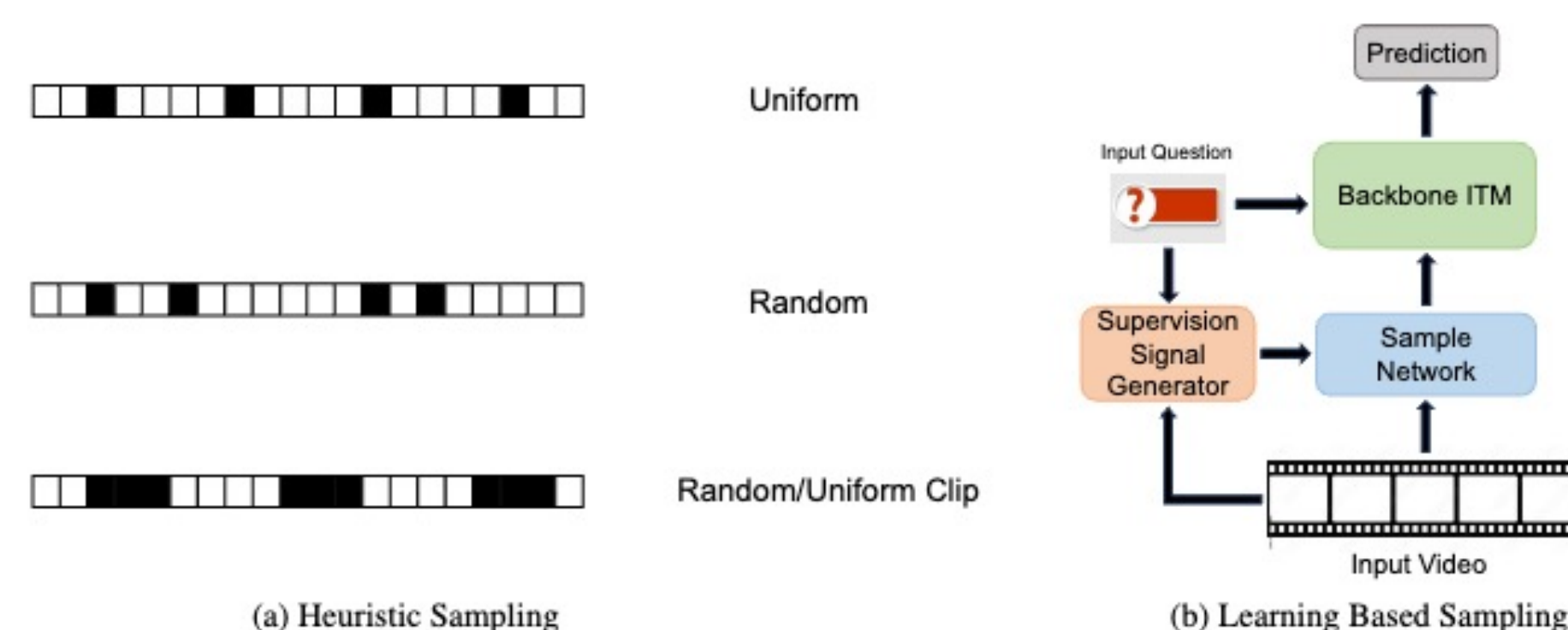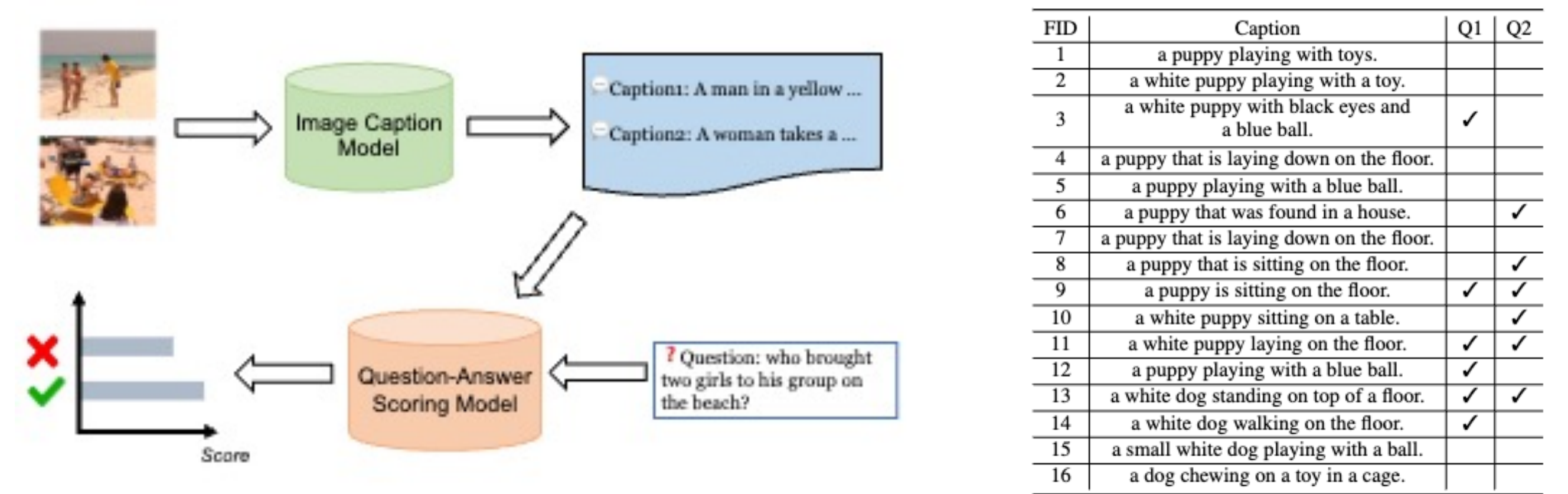


Figure 2: Existing sample strategies for video–question answering tasks. In heuristic sampling, the black boxes indicate selected frames.

## Research Question

- Can we move the sampling stage offline (decouple it from the main network)?
- Can we find a simple yet effective formulation for the offline sampling?
- Is question-aware sampling always required (can we design a question-agnostic one)?

## Method

- Most Implied Frames (MIF)
  - A captioner and a scorer to calculate scores for each frame
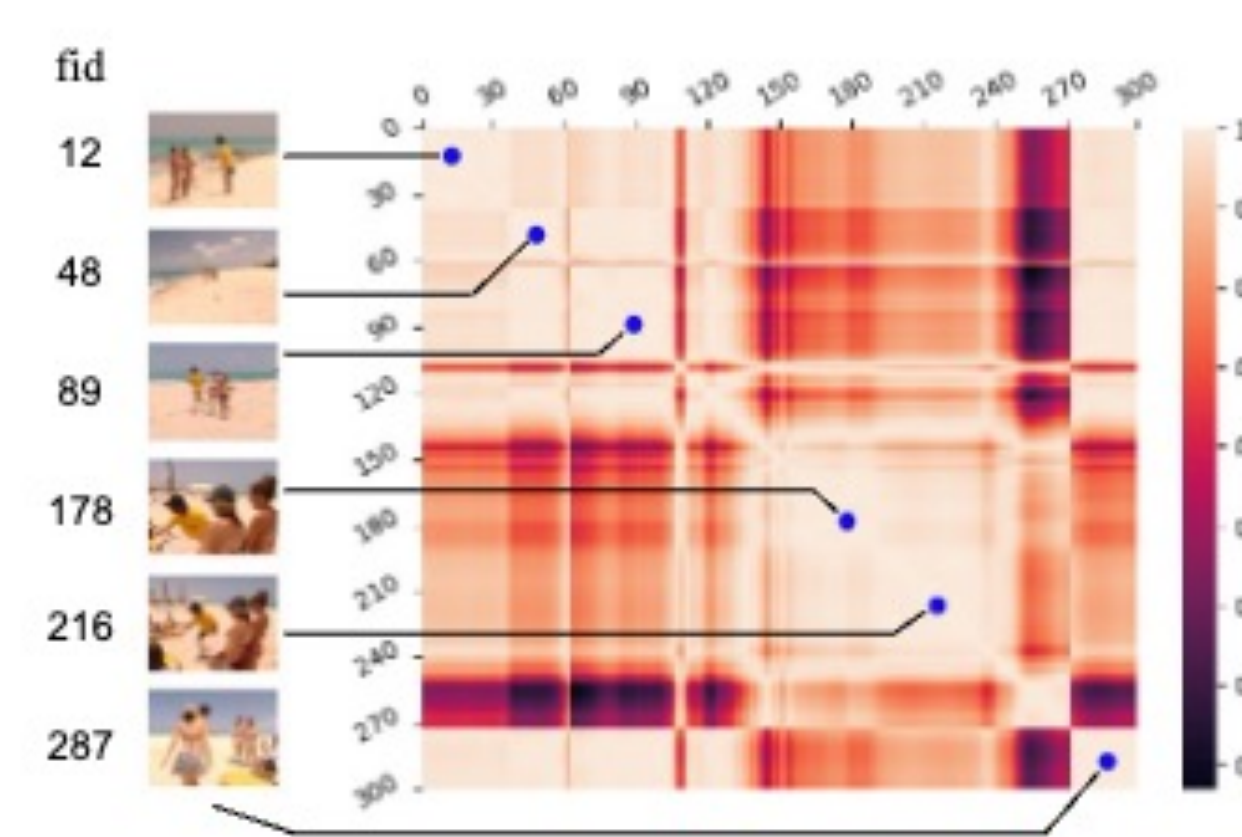  - Choose the frames of the highest scores as sampled ones



- Ablative study of MIF: Is question-aware sampling a must?
  - We change the scale (capability) of captioner and grader respectively
  - There is no obvious correlation between captioner-grader capability and accuracy

| $\mathcal{M}_c$ | $\mathcal{M}_g$ | MSVD | MSRVTT |
|---|---|---|---|
| **Separate Model** | | | |
| GIT-S | BERT-S | 46.5 | 42.3 |
| GIT-B | BERT-B | 46.7 | 42.4 |
| GIT-L | BERT-L | 46.9 | 42.1 |
| **Unified Model** | | | |
| BLIP2-T5-XL | | 46.6 | 42.0 |
| BLIP2-T5-XXL | | 46.2 | 42.2 |

- Most Dominant Frames (MDF)
  - Based on previous analysis, we can move one step forward by even discarding the question-aware component
  - Sampling scores are calculated on visual feature similarity



**Algorithm 1:** Most Dominant Frames (MDF)

**Input:** Video frames $V = \{v_1, v_2, ..., v_T\}$, vision model $\mathcal{M}$, width-adjusting rate $\lambda$
**Output:** Visual prefix $F = \{f_1, f_2, ..., f_N\}$
1 Encode frames using the vision model $E = \mathcal{M}(V) = \{e_1, e_2, ..., e_T\}$
2 Compute $dom$ score for all frames and set $W$, according to Eq. 4 and Eq. 5.
3 Init $F = \{f_{\arg\max_t dom(t)}\}$, index set $I = \{0, 1, ..., i - W, i + W, ..., T\}$
4 **while** $|F| < N$ and $I \neq \emptyset$ **do**
5     $t' \leftarrow \arg\max_t dom(t)$
6     $F \leftarrow F \cup \{f_{t'}\}$
7     $I \leftarrow I \setminus \{t''\}_{t'' - t' < W}$
8 **if** $|F| < N$ **then**
9     $\tau \leftarrow \arg\text{top}_N(\{dom(t)\}_{t \in T})$
10     return $F \cup \{f_t\}_{t \in \tau}$
11 **else**
12     return $F$

## Results

| Model | MSVD | MSRVTT | TGIF |
|---|---|---|---|
| **GIT Backbone** | | | |
| Base (Wang et al., 2022) | 52.2 | 41.1 | 67.5 |
| IGV (Li et al., 2022c) | 53.2 | 41.5 | 68.1 |
| VCSR (Wei et al., 2023) | 52.7 | 41.6 | 68.6 |
| MIF | 54.5 | **42.3** | 69.9 |
| MDF | **55.3** | 42.0 | **70.0** |
| **AIO Backbone** | | | |
| Base (Wang et al., 2023) | 46.1 | 42.7 | 64.0 |
| IGV (Li et al., 2022c) | 46.3 | 43.3 | 64.7 |
| VCSR (Wei et al., 2023) | 46.4 | 43.0 | 64.5 |
| MIF | 46.7 | **44.0** | 65.9 |
| MDF | **46.9** | 43.8 | **66.2** |

Table 3: Test set results on MSVD, MSRVTT and TGIF. Best scores are bolded.

- Both MIF and MDF achieve good performance
- MDF is competitive to MDF, showing that question-aware sampling is not necessary