

Thinking with Drafts: Speculative Temporal Reasoning for Efficient Long Video Understanding

Supplementary Material

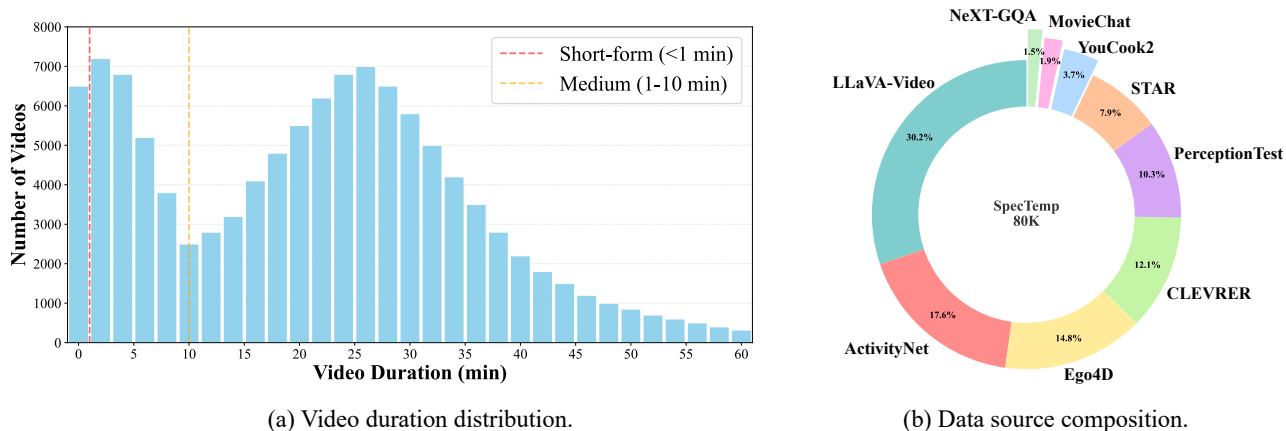


Figure 1. SpecTemp-80K dataset statistics.

001 This supplementary material is organized as follows. We
002 begin with a detailed description of the experimental setups,
003 then present additional experimental findings, followed by
004 the visualization results and a discussion of limitations.

005 Specifically, the **experiment setups** include the follow-
006 ing aspects:

- 007 • Training configurations.
- 008 • Prompt templates for dual models.
- 009 • SpecTemp-80K dataset statistics.

010 **More experimental findings** include the following aspects:

- 011 • Ablation on different iteration numbers.
- 012 • Ablation on reward components.
- 013 • Frame allocation strategies under fixed budget.
- 014 • Needle-in-a-Haystack evaluation.

015 1. Experiment Setup

016 **Training configurations.** We employ a two-stage training
017 strategy. First, we conduct supervised fine-tuning (SFT) for
018 2 epochs on the SpecTemp-80K dataset with a learning rate
019 of 1×10^{-6} . Subsequently, we perform reinforcement fine-
020 tuning (RFT) using GRPO [2] for 1 epoch with a learning
021 rate of 1×10^{-6} . The KL penalty coefficient β is set to
022 0.04. All experiments are conducted on $16 \times L20$ GPUs with
023 batch sizes of 32 for SFT and 16 for RFT.

024 **Prompt templates for dual models.** The prompt templates
025 for the target model (temporal segment prediction and answer
026 generation) and the draft model (salient frame selec-

tion from dense samples) are illustrated in Figure 4 and Figure 5, respectively.

SpecTemp-80K dataset statistics. Figure 1(a) presents the video duration distribution, showing three distinct categories: short-form videos (<1 min, 32.4%), medium-length videos (1-10 min, 51.8%), and long-form videos (>10 min, 15.8%). Figure 1(b) illustrates the data source composition of SpecTemp-80K. The dataset comprises 80,142 video-question-answer triplets collected from 9 diverse sources, with LLaVA-Video (30.2%), ActivityNet (17.6%), and Ego4D (14.8%) constituting the primary sources. This diverse composition ensures comprehensive coverage across varying temporal complexities and video domains, enabling robust training for dual-model collaborative reasoning.

041 2. More Experimental Findings

Table 1. Ablation study on maximum iteration number T_{\max} . We report accuracy (%) on LongVideoBench and Video-Holmes, along with actual iterations executed and inference latency.

T_{\max}	Video-Holmes [1]	LongVideoBench [3]	Iterations	Latency (s)
1	42.8	54.7	0.9	1.5
2	44.6	55.1	1.7	2.0
3 (Ours)	47.0	57.5	2.3	2.3
4	46.7	57.3	3.1	2.9
5	46.2	58.1	3.7	3.5

Ablation on different iteration numbers. We investigate the impact of maximum iteration number T_{\max} on performance and efficiency. As shown in Table 1, accuracy im-

Table 2. Ablation study on reward components. We evaluate the contribution of IoU reward R_{IoU} for temporal localization in the target model and visual information gain reward R_{visual} for frame selection in the draft model.

R_{IoU} (Target)	R_{visual} (Draft)	Video-Holmes [1]	LongVideoBench [3]
-	-	42.1	53.1
✓	-	45.5	55.8
-	✓	44.7	54.6
✓	✓	47.0	57.5

Table 3. Ablation study on frame allocation strategies (Initial+Per-Iter×Max-Iter) under fixed budgets of 16 and 64 frames. We report accuracy (%) and inference latency.

Strategy	Video-Holmes [1]	LongVideoBench [3]	Latency (s)
<i>16-frames budget</i>			
4+4×3	44.3	55.8	1.5
10+2×3 (Ours)	47.0	57.5	1.8
13+1×3	45.6	56.2	2.0
<i>64-frames budget</i>			
32+8×4 (Ours)	47.8	61.4	4.7
48+4×4	46.1	59.7	5.2
56+2×4	46.9	60.2	5.5

045 proves steadily as T_{max} increases from 1 to 3, achieving
 046 57.5% on LongVideoBench (+2.8%) and 47.0% on Video-
 047 Holmes (+4.2%). However, beyond $T_{max} = 3$, we observe
 048 diminishing returns and even performance degradation on
 049 Video-Holmes (46.2% at $T_{max} = 5$), suggesting that exces-
 050 sive iterations may introduce redundant frames and impair
 051 reasoning coherence. The actual average iteration numbers
 052 (0.9 to 3.7) are consistently lower than T_{max} , demonstrat-
 053 ing effective early termination. Considering that $T_{max} = 3$
 054 achieves the best Video-Holmes performance and competi-
 055 tive LongVideoBench accuracy with only 2.3s latency, we
 056 adopt this configuration as the optimal balance between per-
 057 formance and efficiency.

058 **Ablation on reward components.** We conduct ablation
 059 studies to analyze the contribution of individual reward
 060 components in our reinforcement learning framework. As
 061 shown in Table 2, both the IoU reward R_{IoU} for the tar-
 062 get model and the visual information gain reward R_{visual}
 063 for the draft model provide substantial improvements over
 064 the SFT baseline (53.1% on LongVideoBench, 42.1% on
 065 Video-Holmes). Individually, R_{IoU} contributes larger gains
 066 (+2.7% on LongVideoBench, +3.4% on Video-Holmes)
 067 compared to R_{visual} (+1.5% on LongVideoBench, +2.6%
 068 on Video-Holmes), indicating that accurate temporal lo-
 069 calization is the primary bottleneck for long video under-
 070 standing. Importantly, combining both rewards yields the
 071 best performance (57.5% on LongVideoBench, 47.0% on
 072 Video-Holmes), with improvements of +4.4% and +4.9%
 073 over baseline respectively, demonstrating that the two re-
 074 ward signals provide complementary supervision for effec-
 075 tive dual-model collaboration.

Frame allocation strategies under fixed budget. To in-
 076 vestigate the optimal frame allocation strategy under a fixed
 077 computational budget, we conduct ablation studies with
 078 different distributions of initial frames and per-iteration
 079 frames. As shown in Table 3, we evaluate three configu-
 080 rations for both 16-frame and 64-frame budgets, where the
 081 strategy is denoted as Initial+Per-Iter×Max-Iter.
 082

083 The results demonstrate that our default configurations
 084 consistently achieve the best performance across both set-
 085 tings. For the 16-frame budget, the 10+2×3 strategy attains
 086 57.5% on LongVideoBench and 47.0% on Video-Holmes,
 087 outperforming both alternatives with only 1.8s latency. For
 088 the 64-frame budget, our 32+8×4 configuration achieves
 089 61.4% and 47.8% respectively while maintaining the low-
 090 est latency (4.7s) among all 64-frame strategies.

091 These findings reveal an important trade-off in frame
 092 allocation strategies. Aggressive early sampling (4+4×3,
 093 48+4×4) prioritizes many frames per iteration but suffers
 094 from insufficient global temporal context during initial rea-
 095 soning, leading to suboptimal segment prediction. Con-
 096 versely, conservative late sampling (13+1×3, 56+2×4) al-
 097 locates excessive frames to initial observation but lacks the
 098 granularity needed for fine-grained verification in subse-
 099 quent iterations. Our balanced approach allocates moderate
 100 initial frames for coarse temporal understanding while re-
 101 serving sufficient budget for iterative refinement, effectively
 102 bridging global perception and local examination. This val-
 103 idates that our frame allocation strategy achieves an optimal
 104 equilibrium between temporal reasoning breadth and verifi-
 105 cation depth.

Needle-in-a-Haystack evaluation. To evaluate
 106 SpecTemp’s ability to retrieve visual information from
 107 extremely long videos, we conduct a Visual Needle-
 108 In-A-Haystack (NIAH) test inspired by [4]. We insert
 109 single-frame visual questions at various positions into
 110 hour-long video haystacks sampled at 1 FPS. The needle
 111 images are designed to be counterfactual or counter-
 112 commonsense to ensure the model cannot rely on language
 113 priors alone.
 114

115 As shown in Figure 3, SpecTemp successfully retrieves
 116 needle information within videos containing up to 2000
 117 frames, maintaining accuracy above 80% across different
 118 needle positions. The dual-model collaboration (Target
 119 7B + Draft 3B) achieves superior retrieval accuracy com-
 120 pared to using the 7B model alone, demonstrating that our
 121 speculative temporal reasoning framework maintains pre-
 122 cise long-range visual information retrieval while signifi-
 123 cantly improving computational efficiency.

3. Visualizations

124 We provide qualitative examples to illustrate SpecTemp’s
 125 iterative speculation-verification process across different
 126 video understanding scenarios. Figures 6 – Figures 8
 127



Figure 2. Illustration of the Visual Needle-In-A-Haystack evaluation. A single-frame visual question is inserted at varying positions within long video sequences to assess long-range retrieval capability.

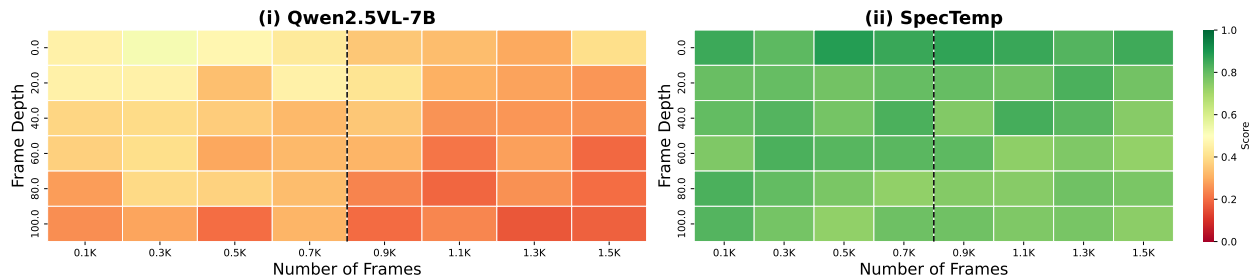


Figure 3. V-NIAH performance across frame counts and needle depths.

128 demonstrate how the target model progressively refines
 129 temporal segments while the draft model proposes salient
 130 frames, showcasing effective dual-model collaboration for
 131 long video reasoning.

132 4. Limitations

133 Despite the effectiveness of SpecTemp, we acknowledge
 134 several limitations that present opportunities for future re-
 135 search:

136 **Computational overhead during training.** Our dual-
 137 model framework requires training two separate MLLMs
 138 and performing joint reinforcement learning optimization,
 139 which increases computational costs compared to training
 140 a single model. The RL stage is particularly resource-
 141 intensive, requiring generation of multiple trajectory sam-
 142 ples per training example for group reward computation.
 143 Future work could explore more efficient training strategies,
 144 such as parameter-efficient fine-tuning methods.

145 **Generalization to extremely long videos.** While
 146 SpecTemp handles hour-long videos effectively, its perfor-
 147 mance on videos exceeding 2-3 hours has not been exten-
 148 sively evaluated. For such extremely long videos, the num-
 149 ber of potential temporal segments grows dramatically, and
 150 maintaining coherent reasoning across numerous iterations
 151 becomes challenging. Hierarchical reasoning strategies or
 152 memory mechanisms may be necessary for handling ultra-
 153 long videos.

References

- 154
- 155 [1] Junhao Cheng, Yuying Ge, Teng Wang, Yixiao Ge, Jing
 156 Liao, and Ying Shan. Video-holmes: Can mllm think
 157 like holmes for complex video reasoning? *arXiv preprint*
 158 *arXiv:2505.21374*, 2025. 1, 2
 - 159 [2] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao
 160 Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li,
 161 Y Wu, et al. Deepseekmath: Pushing the limits of mathe-
 162 matical reasoning in open language models. *arXiv preprint*
 163 *arXiv:2402.03300*, 2024. 1
 - 164 [3] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li.
 165 Longvideobench: A benchmark for long-context inter-
 166 leaved video-language understanding. *arXiv preprint*
 167 *arXiv:2407.15754*, 2024. 1, 2
 - 168 [4] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng,
 169 Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan,
 170 Chunyuan Li, and Ziwei Liu. Long context transfer from lan-
 171 guage to vision. *arXiv preprint arXiv:2406.16852*, 2024. 2

Question: {question}

Based on the {frames} sampled video frames from a video, analyze whether the provided frames contain sufficient information to answer this question.

Frame indexing: 1 to {frames}

If you need additional visual information from a specific time interval, provide your detailed reasoning of what information is missing between the <think> </think> tags, then specify the interval using <segment>(x,x+1)</segment> which indicates you need frames between frame x and frame x+1.

You will receive new video frames. You may request additional video frames up to a maximum of 3 times. Always choose an interval according to the most recent video frames.

IMPORTANT: Each round of reasoning must be based solely on the currently visible video frames. Do not rely on or repeat reasoning patterns from previous iterations.

When sufficient information is available, provide your reasoning of answering the question in <think> </think> tags and your final answer in <answer> </answer> tags.

Format: <think>your reasoning process here</think> <segment>(x,x+1)</segment> OR <think>your reasoning process here</think> <answer>your final answer here</answer>

Figure 4. **Prompt template for the target model.** The target model is guided to perform temporal reasoning and either predict evidence segments for further exploration or generate final answers.

You will be shown {num_frames} video frames. Your task is to select the two frames that are most relevant to answering the following question:

Question: {question}

Frame indexing: 1 to {num_frames}

Provide your detailed reasoning between the <think> </think> tags, and then give the two frame numbers using <frame>(x,y)</frame> tags, where $x \neq y$.

Format: <think>your reasoning process here</think> <frame>(x,y)</frame>

Figure 5. **Prompt template for the draft model.** The draft model is instructed to select the most informative and diverse frames from densely sampled temporal regions.

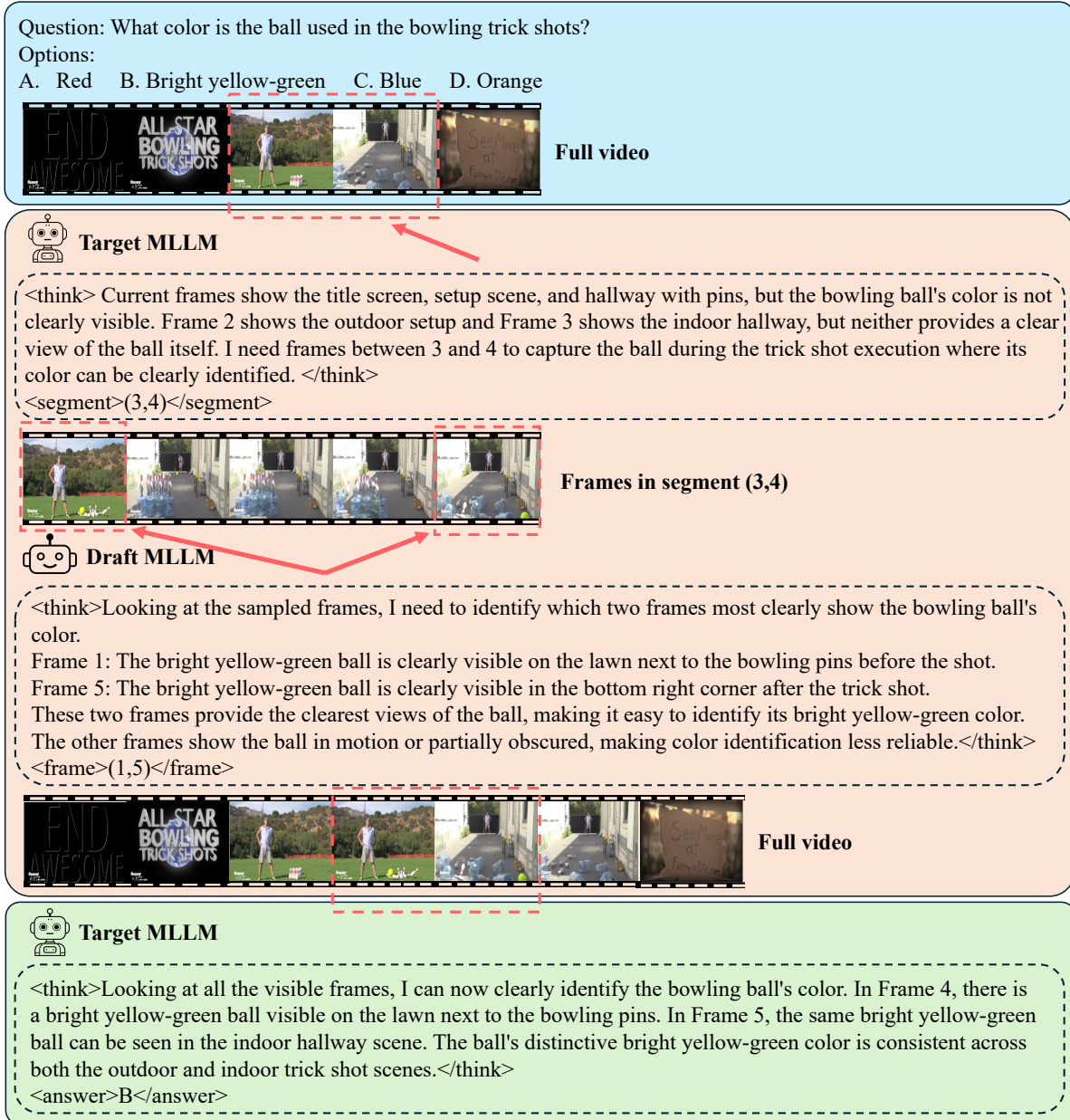


Figure 6. Qualitative case: 1 iteration for bowling ball color identification.

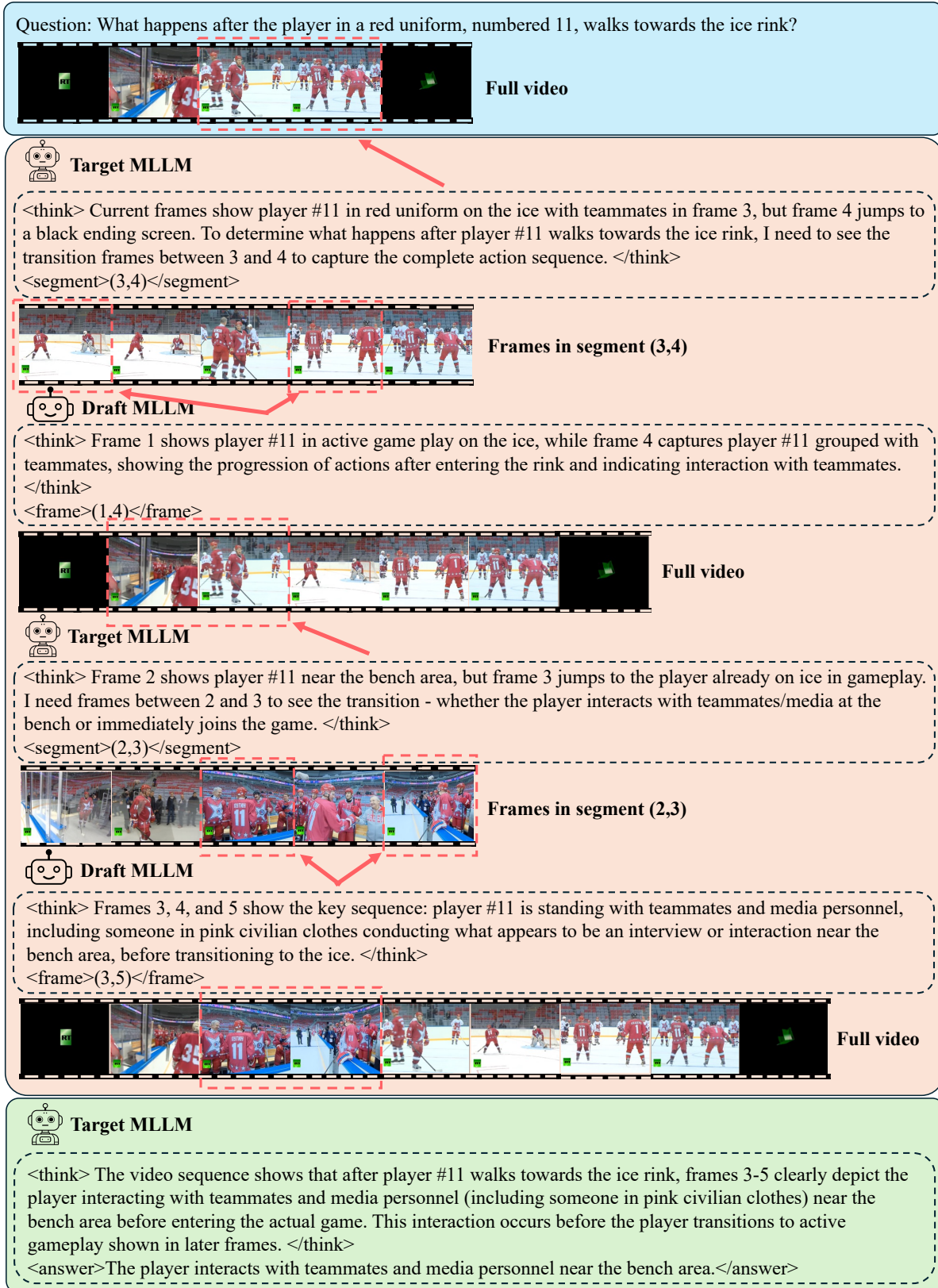


Figure 7. Qualitative case: 2 iterations for sports action sequence tracking.

Question: What is cut after the cabbage?
Options: A. Garlic B. Ginger C. Carrot D. Green onion

Full video

Target MLLM

<think>Current frames show cooking preparation with pots and the final dish, but no vegetable cutting scenes are visible. To answer what is cut after the cabbage, I need to see the actual cutting sequence. Frames between 3 and 4 should capture the food preparation stage where vegetables are being cut.</think>
<segment>(3,4)</segment>

Frames in segment (3,4)

Draft MLLM

<think> Frame 3 clearly shows cutting action on a pink cutting board, which is directly relevant to identifying what ingredient is being cut. Frame 4 shows the cooking stage, providing context for the prepared ingredients. These two frames together help establish the cutting sequence in the recipe preparation. </think>
<frame>(3,4)</frame>

Full video

Target MLLM

<think> Frame 4 shows cutting on the pink cutting board, but the sequence of what is cut after cabbage is not fully clear. I need frames between 4 and 5 to see the complete cutting sequence and identify which ingredient is cut immediately after the cabbage. </think>
<segment>(4,5)</segment>

Frames in segment (4,5)

Draft MLLM

<think> Frames 2 and 3 show the cutting sequence on the pink cutting board. Frame 2 captures cutting of small white pieces (likely garlic), and frame 3 shows cutting of cabbage, which helps establish the order of ingredient preparation. </think>
<frame>(2,3)</frame>

Full video

Target MLLM

<think> Frames 4-6 show the cutting sequence with cabbage visible in frame 6. To determine what is cut after the cabbage, I need to see the frames between 6 and 7 to capture the next ingredient in the cutting sequence. </think>
<segment>(6,7)</segment>

Frames in segment (6,7)

Draft MLLM

<think> Frame 1 shows the cabbage already cut on the board, and frame 2 shows orange carrot pieces, indicating that carrot is the ingredient cut immediately after cabbage in the preparation sequence. </think>
<frame>(1,2)</frame>

Full video

Target MLLM

<think>From the complete video sequence, I can see the cutting order: Frame 7 shows the cabbage already cut, and frame 8 shows orange carrot pieces appearing next on the cutting board. This clearly indicates that carrot is cut immediately after the cabbage.</think>
<answer>C</answer>

Figure 8. Qualitative case: 3 iterations for temporal ordering reasoning.