

Seq2Time: Sequential Knowledge Transfer for Video LLM

Temporal Grounding

(Supplementary Materials)

Andong Deng^{1,2,*}, Zhongpai Gao^{2,†}, Anwesa Choudhuri², Benjamin Planche², Meng Zheng²,
Bin Wang^{2,3}, Terrence Chen², Chen Chen¹, Ziyang Wu²

¹University of Central Florida, ²United Imaging Intelligence, ³Northwestern University

1. Data Preparation

1.1. Image Sequence

Data Source. To construct an image sequence dataset that provides high-quality visual-text alignment, we select LLaVA-ReCap [4] as the primary source. Built on the robust vision-language model LLaVA-NeXT-34B [5], LLaVA-ReCap contains over 3.7 million image-text pairs. Its captions are detailed, providing a rich description of the visual content, as exemplified in Figure 1. These captions enable strong correspondence between visual elements and language, serving as a reliable basis for generating high-quality sequential supervision in our training. This alignment is critical for training models to understand and reason over temporally structured visual data.

Instruction Template Design. To maximize the diversity and effectiveness of the training data, we design multiple instruction templates for each pretext task using GPT-4o [6]. These templates include variations in both questions and answers to ensure robustness in model learning. Figures 2, 3, 4, 5, and 6 illustrate examples of these templates for different tasks, such as Image Index Grounding (IIG), Indexed Image Captioning (IIC), and Adjacent Location Reasoning (ALR). This variety enhances the model’s ability to generalize across different types of instructions and strengthens its visual-language understanding capabilities.

1.2. Clip Sequence

Data Source. We utilize the Kinetics-700 dataset [1], a large-scale video understanding dataset comprising approximately 650,000 ten-second video clips annotated with 700 distinct human action classes. The dataset offers diverse and naturalistic videos capturing a wide range of everyday activities. For our Seq2Time, we focus on the training split, which contains about 550,000 clips, ensuring a broad spectrum of visual and temporal information for model training.

Clip Captioning with LongVA. To generate textual descriptions for the video clips, we employ LongVA [7],

a state-of-the-art video captioning model. LongVA generates captions based on the visual content and the associated action label, as demonstrated in Figure 7. For instance, it accurately captures detailed descriptions like “A player in an orange jersey is leaping towards the hoop...” for basketball action clips. These captions provide contextual insights into the actions and environments depicted in the videos, enriching the clip sequence data with semantic annotations.

Despite its strengths, LongVA occasionally produces hallucinated descriptions, as illustrated in Figure 8. For example, a video showing an adult and a baby reading by a table might be misinterpreted as “a group of adults and babies...” Such hallucinations can potentially introduce noise into the dataset, affecting the reliability of the generated sequences. This highlights a trade-off between the scalability of automatic captioning and the quality of annotations.

2. Visualization

In this section, we present qualitative results to demonstrate the effectiveness of our Seq2Time. Compared with the baseline model TimeChat, Seq2Time significantly improves upon repetitive text patterns, temporal localization accuracy, and the precision of event descriptions.

In Figure 9, we analyze a video illustrating the instructions for making pancakes. TimeChat struggles to recognize ingredients in certain steps, providing incomplete descriptions such as “a bowl of eggs” and “a bowl of flour,” and its temporal predictions are imprecise, resulting in a fragmented understanding of the cooking process. In contrast, Seq2Time generates accurate and comprehensive event descriptions with precise temporal annotations, such as predicting “The next step is to mix the ingredients together in a bowl” as occurring between 25.0 to 61.4 seconds, closely aligning with the ground truth annotation: “mix flour,

sugar, baking powder, and salt together in a bowl” from 25 to 61 seconds. This comparison highlights Seq2Time’s ability to capture finer-grained details, maintain alignment with the video’s temporal structure, and robustly capture the sequential flow of actions, effectively reducing redundancies and hallucinations often present in TimeChat’s outputs. These improvements underscore Seq2Time’s effectiveness in enhancing temporal grounding and understanding for long videos.

In Figure 10, a video illustrating the steps for cooking eggs is analyzed. TimeChat exhibits several issues, including generating repetitive text such as “She pours some into the bowl and stirs it” across multiple timestamps and misrecognizing objects, as seen in the caption “She sets the bowl on the stove,” where the “bowl” should be the “pan.” In contrast, Seq2Time minimizes these repetitive patterns and provides a more coherent description of the cooking process. Although Seq2Time incorrectly predicts the timestamp for the event “She stirs the mixture again,” it successfully captures the overall sequence of actions with detailed and contextually appropriate descriptions. This example highlights Seq2Time’s enhanced capability to handle complex sequences with improved temporal and semantic accuracy compared to TimeChat.

In Figure 11, the video illustrates the steps for making tacos. TimeChat fails to recognize any of the correct cooking steps, producing entirely incorrect predictions and hallucinating that the video involves baking in an oven. In contrast, our Seq2Time provides general but accurate descriptions of the steps, despite omitting some specific ingredient details. For example, Seq2Time predicts “The man adds more ingredients to the dish and continues to stir it together” between 148.1 and 259.2 seconds. While this aligns with the overall sequence, the ground truth includes finer-grained annotations, such as adding specific ingredients like cumin powder and beef.

Figure 12 presents a video of two girls preparing tofu soup. Our Seq2Time successfully captures key cooking steps, such as describing “They mix the ingredients in a pot and stir it” from 136.1 to 211.4 seconds. This aligns closely with the ground truth annotation, “add the tofu chunks and dissolve miso paste in the soup,” which spans 142 to 184 seconds. In contrast, TimeChat only captures the initial scene of the video, which is less relevant to the cooking process, failing to identify any meaningful steps.

Figure 13 shows a video depicting a woman making cakes. TimeChat fails to generate relevant captions, producing incorrect descriptions, such as “a young man baking something.” While Seq2Time accurately

identifies basic events within the video, providing correct and relevant captions that align with the overall context.

Finally, Figure 14 shows a video of a woman preparing a salad in the kitchen. Both TimeChat and Seq2Time manage to produce correct captions for this video. However, neither model provides detailed step-by-step instructions. Instead, they capture only general actions, such as “woman talks” or “show the different ingredients,” overlooking specific details of the cooking steps. This indicates an area where further refinement of both models could enhance their performance in capturing detailed procedural actions.

These examples collectively demonstrate the superiority of our Seq2Time in reducing hallucinations, improving temporal alignment, and providing more accurate descriptions compared to TimeChat. However, they also highlight opportunities for improvement, such as better capturing fine-grained details and enhancing object recognition in complex cooking scenarios.

3. Additional Results

We also evaluate our methods on VTG-LLM [3] to demonstrate its general effectiveness. Differently, for the clip sequence data, we directly leverage video data from ShareGPT4Video [2]. Following the same training strategy, we first train VTG-LLM using our Seq2Time with 1 epoch and the rest training follows the original setting of VTG-LLM. As shown in Table 1, Seq2Time brings a 5.23% improvement on the F1 score and 2.37% improvement on the R@0.5 metric, indicating consistent effectiveness for different video LLMs.

	YouCook2				Charades-STA	
	S	C	M	F1	R@0.5	R@0.7
VTG-LLM*	1.5	5.1	1.8	17.2	33.8	15.7
+ S2T	1.6	5.1	2.0	18.1	34.6	16.1

Table 1. The performance of VTG-LLM with Seq2Time as additional training data. *Reproduced VTG-LLM results.

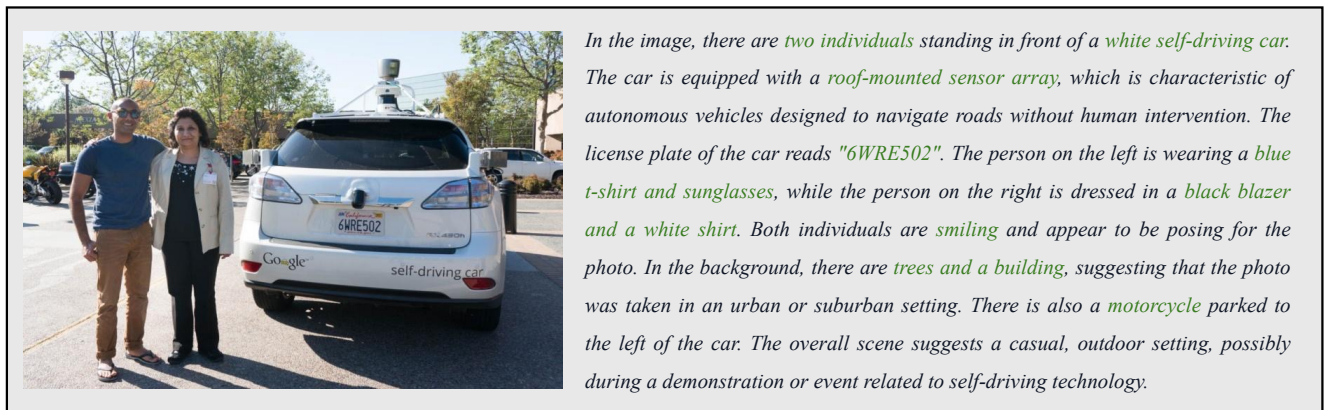
4. Limitation Analysis

From Figure 9 to Figure 14, we observe the overall effectiveness of our Seq2Time in capturing key events and providing accurate temporal annotations. However, some limitations remain, which could be addressed in future work.

First, while the model generates correct captions with precise timestamps, the described events are occasionally less critical within the context of the video. For example, in Figure 12, “0.0 – 40.4 seconds, A pair of girls are preparing a meal in the kitchen,” and in Figure 13, “158.7 – 166.4 seconds, The woman takes a bite of the

cake and smiles at the camera,” both describe valid events, but these are not essential steps in the respective cooking procedures.

Second, the model sometimes overlooks fine-grained steps that are crucial for understanding detailed processes. For instance, in Figure 11, the model predicts a general event: “33.6 – 103.1 seconds, The man prepares the ingredients for the dish and stirs them together in a pan.” While this is accurate at a high level, it misses finer details such as “put olive oil,” “drain off the fat,” and “add chopped vegetables,” which are key preparatory steps in making tacos. Addressing these limitations would enhance the granularity and relevance of the captions generated by Seq2Time.



In the image, there are two individuals standing in front of a white self-driving car. The car is equipped with a roof-mounted sensor array, which is characteristic of autonomous vehicles designed to navigate roads without human intervention. The license plate of the car reads "6WRE502". The person on the left is wearing a blue t-shirt and sunglasses, while the person on the right is dressed in a black blazer and a white shirt. Both individuals are smiling and appear to be posing for the photo. In the background, there are trees and a building, suggesting that the photo was taken in an urban or suburban setting. There is also a motorcycle parked to the left of the car. The overall scene suggests a casual, outdoor setting, possibly during a demonstration or event related to self-driving technology.

Figure 1. Illustration of an example from LLaVA-ReCap-CC3M. The caption provide details in every aspects of the image from foreground to background. For instance, even the motorcycle in the background is captured. The important descriptions are indicated in green texts.

Question Templates for Image Index Grounding with Single Target:	Answer Templates for Image Index Grounding with Single Target:
<i>"Which image matches the description: {}. Please output the image index."</i>	<i>"The image index is {}."</i>
<i>"Identify the image that corresponds to this caption: {}. Please provide the index."</i>	<i>"You should look at image number {}."</i>
<i>"What is the index of the image that describes {}?"</i>	<i>"The correct index is {}."</i>
<i>"Can you find the image that matches this description: {}. Please give the index."</i>	<i>"Check image of index {} for the answer."</i>
<i>"Which image is described by {}? Please output the index."</i>	<i>"It corresponds to image of index {}."</i>
<i>"Find the image that matches the following description: {}. Please output the index."</i>	<i>"Image number {} is the correct one."</i>
<i>"What is the index of the image that corresponds to {}?"</i>	<i>"The image matching this description is at index {}."</i>
<i>"Which image aligns with the caption: {}? Please provide the index."</i>	<i>"You will find it at image {}."</i>
<i>"Identify the image matching this description: {}. Please give the index."</i>	<i>"The right image is at index {}."</i>
<i>"What image index corresponds to the description: {}?"</i>	<i>"Refer to image {}."</i>

Figure 2. Instruction templates of image index grounding with single target image.

Question Templates for Image Index Grounding with Multiple Targets:	Answer Templates for Image Index Grounding with Multiple Targets:
<i>"Which images match the following descriptions: ",</i>	<i>"The image indices are {}."</i>
<i>"Identify the images corresponding to these captions: ",</i>	<i>"You should look at image numbers {}."</i>
<i>"What are the indices of the images that describe the following: ",</i>	<i>"The correct indices are {}."</i>
<i>"Can you find the images that match these descriptions: ",</i>	<i>"Check images with indices {} for the answer."</i>
<i>"Which images are described by the following: ",</i>	<i>"They correspond to images of indices {}."</i>
<i>"Find the images matching these descriptions: ",</i>	<i>"Image numbers {} are the correct ones."</i>
<i>"What are the image indices that correspond to these captions: ",</i>	<i>"The images matching these descriptions are at indices {}."</i>
<i>"Which images align with the captions: ",</i>	<i>"You will find them at images {}."</i>
<i>"Identify the images matching these descriptions: ",</i>	<i>"The right images are at indices {}."</i>
<i>"What image indices correspond to these descriptions: "</i>	<i>"Refer to images of indices {}."</i>
Following Prompt after Question:	
<i>The 1st caption is {}; the 2nd caption is {}; ... <omitted></i>	
<i>Please output the image indices in the same order as the captions.</i>	

Figure 3. Instruction templates of image index grounding with multiple target images.

Question Templates for Indexed Image Captioning with Single Target:	Answer Templates for Indexed Image Captioning with Single Target:
<i>"Please describe the image with index {}."</i>	<i>"The image with index {} describes {}."</i>
<i>"What does the image with index {} depict?"</i>	<i>"At index {}, the image depicts {}."</i>
<i>"Can you provide the description for the image at index {}?"</i>	<i>"The image at index {} shows {}."</i>
<i>"Describe the image that is indexed as {}."</i>	<i>"Index {} corresponds to an image that describes {}."</i>
<i>"What is the caption for the image with index {}?"</i>	<i>"The description for the image at index {} is {}."</i>
<i>"Tell me what the image at index {} shows."</i>	<i>"For the image indexed at {}, the caption is {}."</i>
<i>"Could you describe the content of the image at index {}?"</i>	<i>"The image numbered {} describes {}."</i>
<i>"What does the image numbered {} describe?"</i>	<i>"At index {}, the image is about {}."</i>
<i>"Provide the description for the image at index {}."</i>	<i>"The caption for the image at index {} is {}."</i>
<i>"What is depicted in the image with index {}?"</i>	<i>"The image at index {} depicts {}."</i>

Figure 4. Instruction templates of indexed image captioning with single target image.

Question Templates for Indexed Image Captioning with Multiple Targets :	Answer Templates for Indexed Image Captioning with Multiple Targets:
<i>"Please describe the images with the following indices: "</i>	<i>"The image with index {} describes: {}."</i>
<i>"What do the images at the following indices depict: "</i>	<i>"At index {}, the image depicts: {}."</i>
<i>"Can you provide descriptions for the images indexed as: "</i>	<i>"The image at index {} shows: {}."</i>
<i>"Describe the images corresponding to these indices: "</i>	<i>"Index {} corresponds to an image that describes: {}."</i>
<i>"What are the captions for the images with indices: "</i>	<i>"The description for the image at index {} is: {}."</i>
<i>"Tell me what the images at the following indices show: "</i>	<i>"For the image indexed at {}, the caption is: {}."</i>
<i>"Could you describe the content of the images at these indices: "</i>	<i>"The image numbered {} describes: {}."</i>
<i>"Provide the descriptions for the images at these indices: "</i>	<i>"At index {}, the image is about: {}."</i>
<i>"What is depicted in the images with the following indices: "</i>	<i>"The caption for the image at index {} is: {}."</i>
	<i>"The image at index {} depicts: {}."</i>

Figure 5. Instruction templates of indexed image captioning with multiple target images.

Question Templates for Adjacent Location Reasoning:	Answer Templates for Adjacent Location Reasoning:
<i>"Which image comes right before the one depicting {}? Please output the index and describe this image."</i>	<i>"Refer to image {}. It shows {}."</i>
<i>"Which image comes right after the one depicting {}? Please output the index and describe this image."</i>	<i>"Look at image {}. It depicts {}."</i>
<i>"What is the image right before the one showing {}? Please output the index and describe the target image."</i>	<i>"Image {} describes {}."</i>
<i>"What is the image right after the one showing {}? Please output the index and describe the target image."</i>	<i>"You should check image {}. It illustrates {}."</i>
<i>"Which image is right before the one describing {}? Please output the index and describe this image."</i>	<i>"The correct image is {}. It shows {}."</i>
<i>"Which image is right after the one describing {}? Please output the index and describe this image."</i>	<i>"Image {} presents {}."</i>
<i>"What is the image right before the one showing {}? Please provide the index and describe the image."</i>	<i>"The target image is {}. It depicts {}."</i>
<i>"What is the image right after the one showing {}? Please provide the index and describe the image."</i>	<i>"Refer to image {}. It describes {}."</i>
<i>"Identify the image that comes immediately before the one featuring {}. Please output the index and describe it."</i>	<i>"You should see image {}. It illustrates {}."</i>
<i>"Identify the image that comes immediately after the one featuring {}. Please output the index and describe it."</i>	<i>"Image {} shows {}."</i>
<i>"Which image is right before the one depicting {}? Please provide the index and describe the image."</i>	
<i>"Which image is right after the one depicting {}? Please provide the index and describe the image."</i>	
<i>"What is the image immediately before the one showing {}? Please output the index and describe this image."</i>	
<i>"What is the image immediately after the one showing {}? Please output the index and describe this image."</i>	
<i>"Which image is right before the one describing {}? Please output the index and describe the target image."</i>	
<i>"Which image is right after the one describing {}? Please output the index and describe the target image."</i>	
<i>"What is the image right before the one showing {}? Please provide the index and describe it."</i>	
<i>"What is the image right after the one showing {}? Please provide the index and describe it."</i>	
<i>"Identify the image that comes right before the one featuring {}. Please output the index and describe the image."</i>	
<i>"Identify the image that comes right after the one featuring {}. Please output the index and describe the image."</i>	

Figure 6. Instruction templates of adjacent location reasoning.



The video captures an intense moment during a basketball game. A player in an orange jersey is leaping towards the hoop, attempting to dunk or block the ball with a powerful leap, while other players from both teams are positioned strategically around the court, ready for the play's outcome. The crowd in the stands is visible, their faces a blur of anticipation and excitement as they watch the action unfold. The lighting inside the arena highlights the drama of the game, casting dramatic shadows and emphasizing the athleticism of the players on the court.

Figure 7. Instruction of an example in K700-LongVA-Captions.



*In this video, **a group of adults and babies** are gathered around a table, engrossed in reading a colorful newspaper. The adults and children are **all wearing green** hats with a distinctive leaf design on them, suggesting a theme or celebration related to spring or nature. The atmosphere is warm and familial, with **each person holding a baby** securely as they share the experience of reading together. The room is filled with homely warmth, and the focus is clearly on the shared joy of learning and spending time together.*

Figure 8. Instruction of a low-quality example in K700-LongVA-Captions. The incorrect descriptions are indicated by **bold red** text.

References

- [1] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. *arXiv preprint arXiv:1907.06987*, 2019. [1](#)
- [2] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Bin Lin, Zhenyu Tang, et al. Sharegpt4video: Improving video understanding and generation with better captions. *arXiv preprint arXiv:2406.04325*, 2024. [2](#)
- [3] Yongxin Guo, Jingyu Liu, Mingda Li, Xiaoying Tang, Xi Chen, and Bo Zhao. Vtg-llm: Integrating timestamp knowledge into video llms for enhanced video temporal grounding. *arXiv preprint arXiv:2405.13382*, 2024. [2](#)
- [4] Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. Llava-next: What else influences visual instruction tuning beyond data?, 2024. [1](#)
- [5] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024. [1](#)
- [6] OpenAI. Gpt-4 technical report. <https://openai.com/research/gpt-4>, 2023. Accessed: 2024-11-07. [1](#)
- [7] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *arXiv preprint arXiv:2406.16852*, 2024. [1](#)

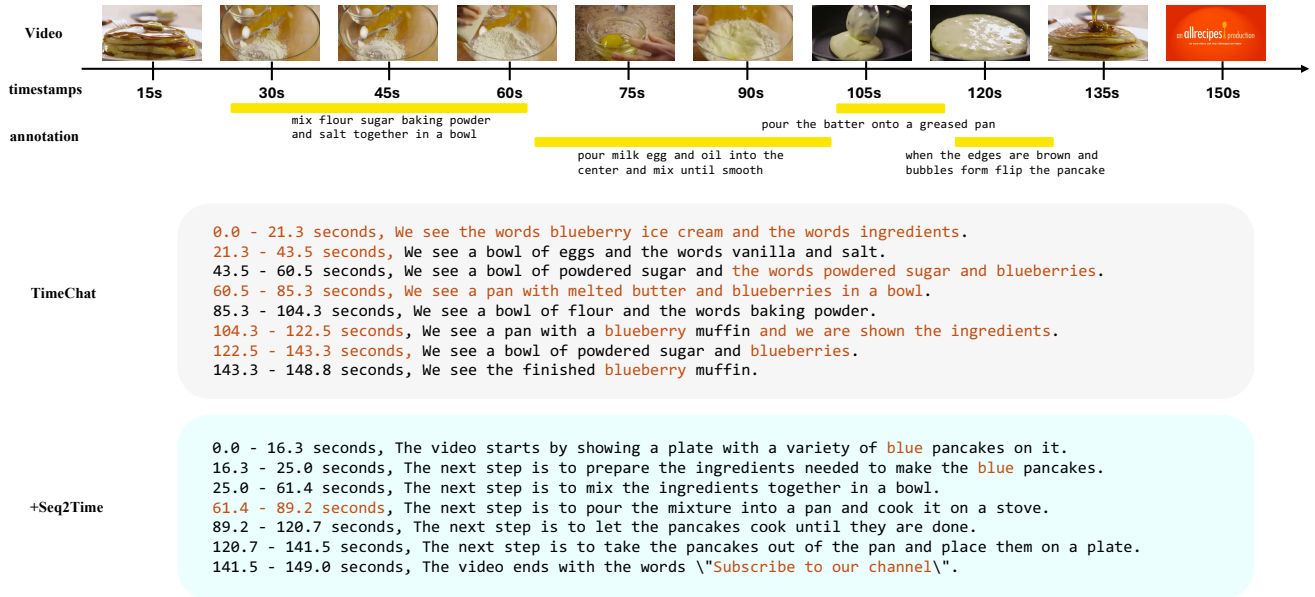


Figure 9. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text.

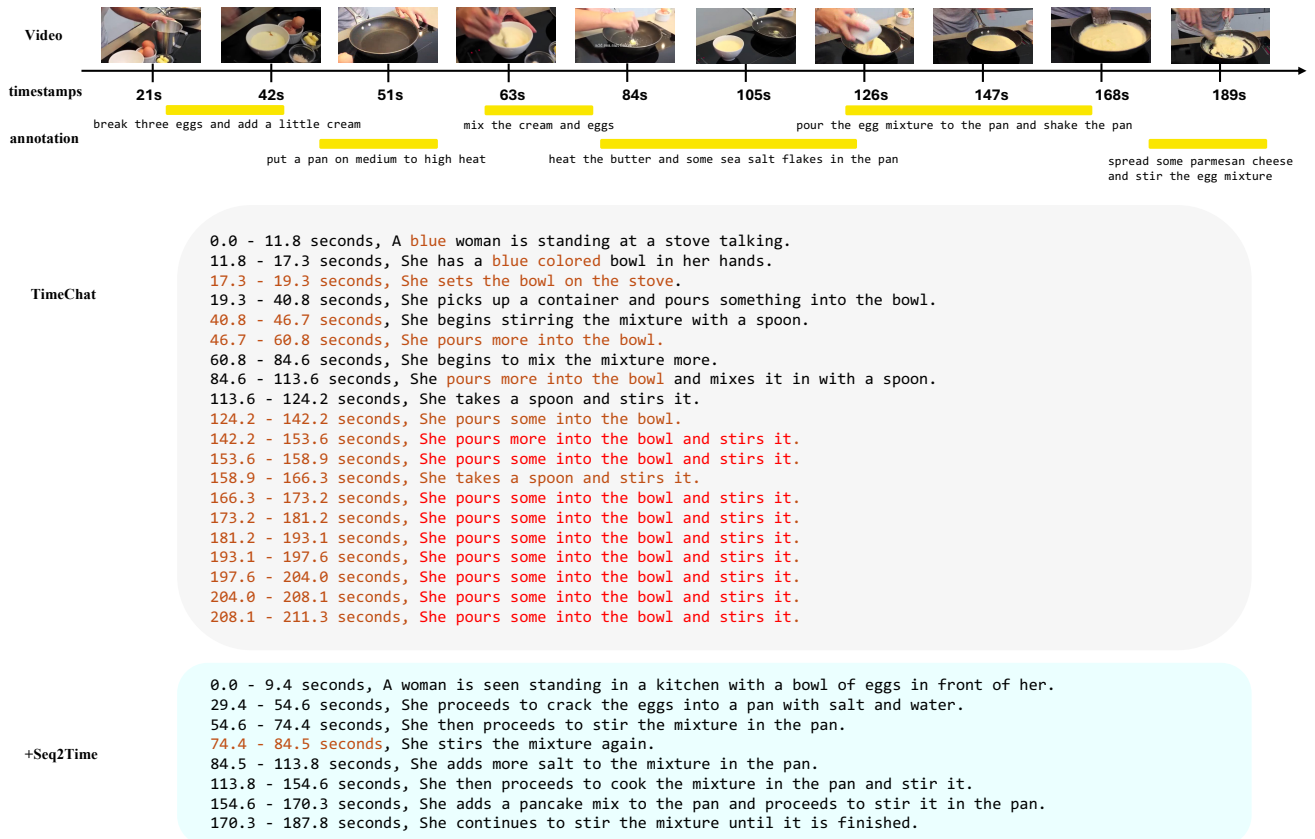


Figure 10. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text and repetitive text is highlighted in red.



Figure 11. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text.

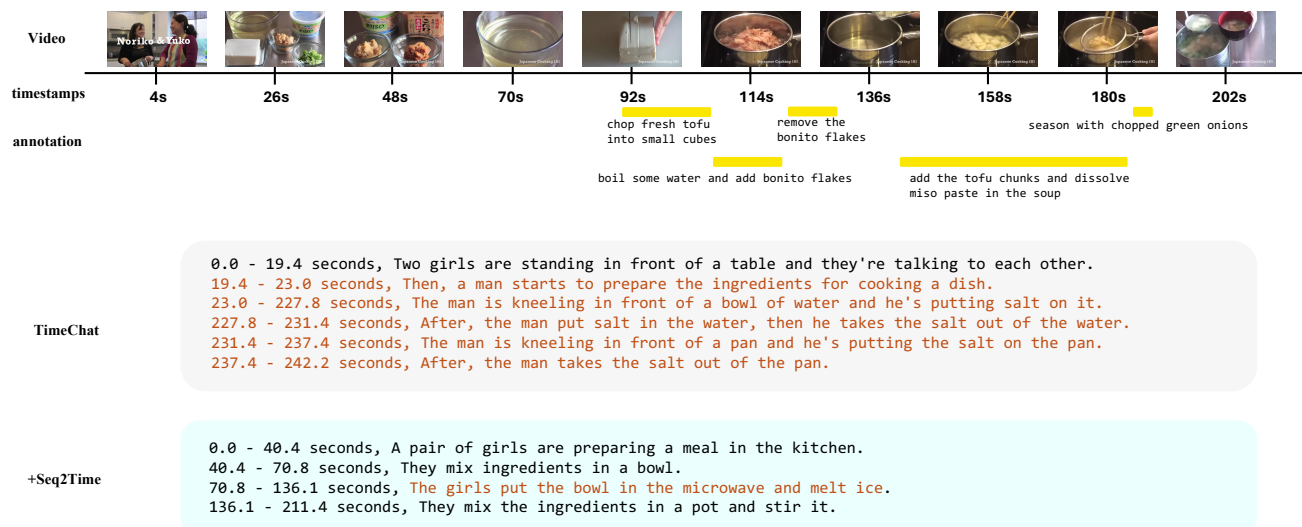


Figure 12. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text.

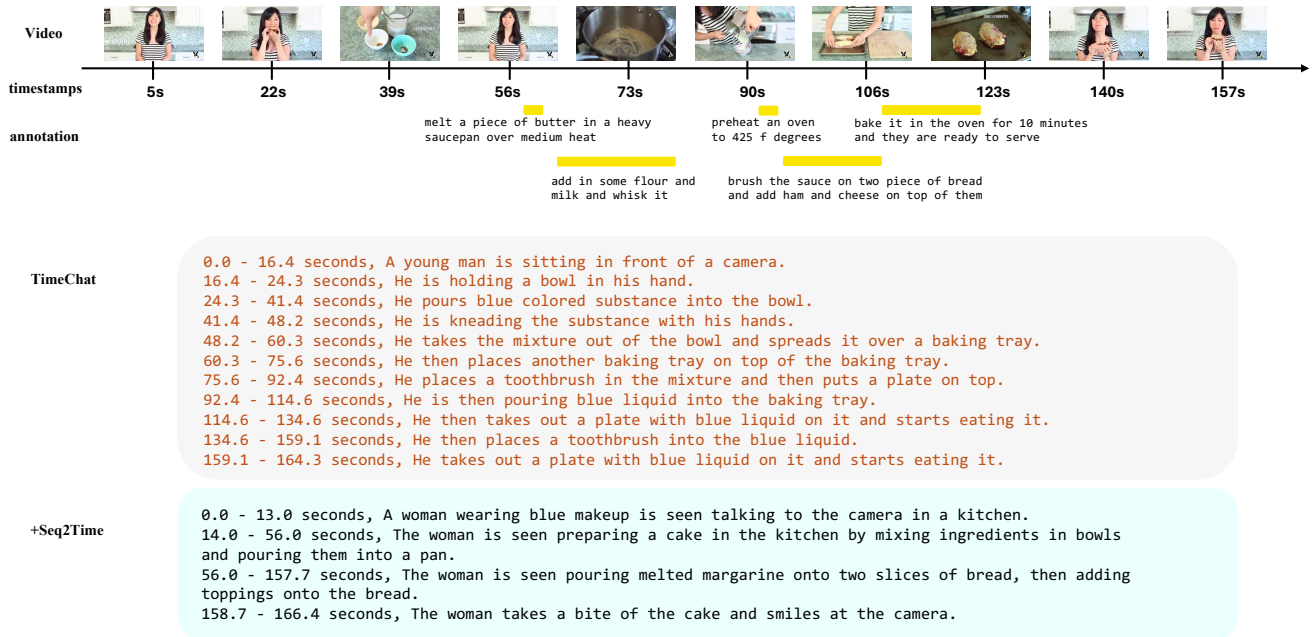


Figure 13. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text.

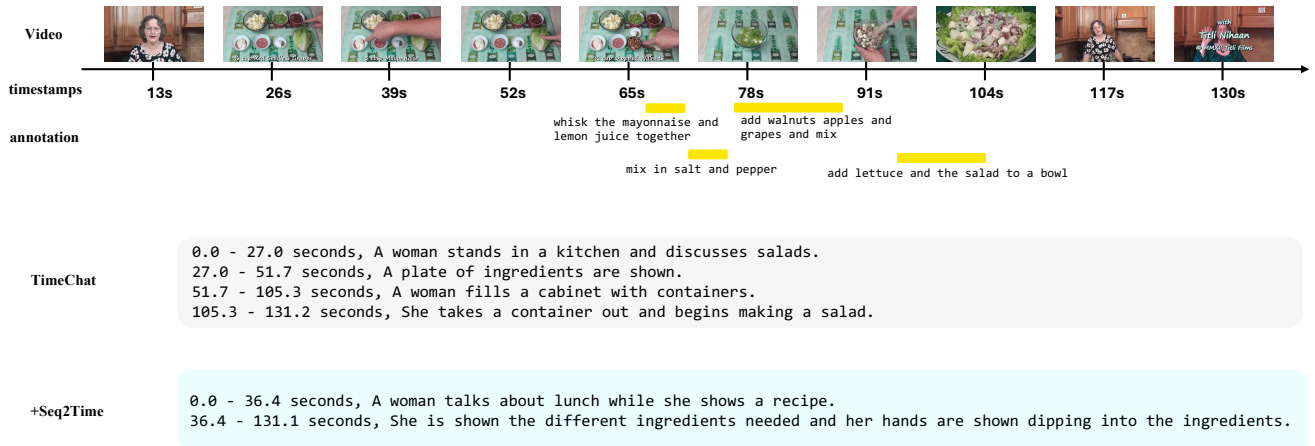


Figure 14. Qualitative example of our Seq2Time on TimeChat. The incorrect descriptions and timestamps are indicated by brown text.