

# NaQ: Leveraging Narrations as Queries to Supervise Episodic Memory



Santhosh Kumar Ramakrishnan<sup>1</sup>



Ziad Al-Halah<sup>2</sup>



Kristen Grauman<sup>1,3</sup>

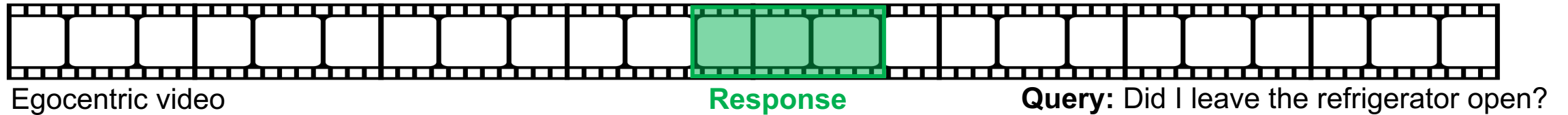
**Poster session:** TUE-PM-245

**Project page:** <https://vision.cs.utexas.edu/projects/naq/>

**Code:** <https://github.com/srama2512/NaQ>

# Overview

**Task:** Episodic Memory via Natural Language Queries (NLQ)



**Challenge:** Limited training data (e.g., 11k queries over 130 hours of video)

**Our idea:** Augment NLQ training by learning to localize “narrations”

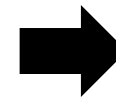
Example narration text: `C rinses hand; C closes tap`



Convert to NLQ  
annotation



**Narrations-as-query data**  
(video, query, response window)



Train NLQ models

# Episodic Memory (EM)

**Goal:** Enable AR assistants for super-human memory



**Query:** Who did I interact with when I played with the dog for the second time in the living room?



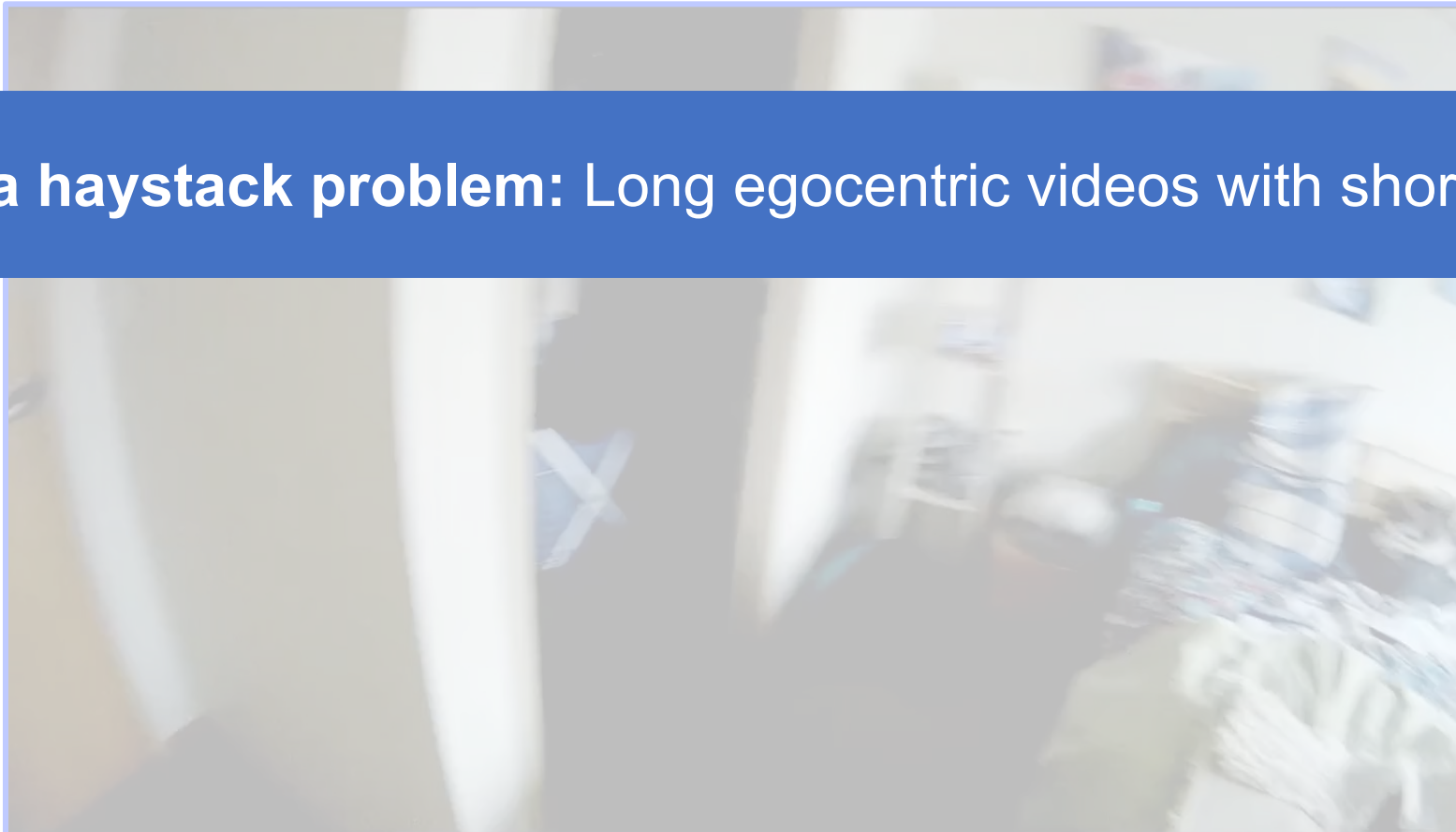
# Episodic Memory (EM)

**Goal:** Enable AR assistants for super-human memory



**Query:** Who did I interact with when I played with the dog for the second time in the living room?

**Needle in a haystack problem:** Long egocentric videos with short responses



# Episodic Memory benchmark on Ego4D

Temporally localize responses to **Natural language queries (NLQ)**

## Queries formulated based on templates

Category	Template
Objects	Where is object X before / after event Y?
	Where is object X?
	What did I put in X?
	How many X's? (quantity question)
	What X did I Y?
	In what location did I see object X ?
	What X is Y?
	State of an object
	Where is my object X?
	Place
People	Who did I interact with when I did activity X?
	Who did I talk to in location X?
	When did I interact with person with role X?

## NLQ dataset statistics

Split	Train	Val	Test
# video hours	136	45	46
# clips	1.0k	0.3k	0.3k
# queries	11.3k	3.9k	4.0k

**Average clip duration:** 8.2 mins  
**Average response duration:** 10.5 sec  
*Needle in a haystack problem*

**Key challenge:** Limited annotation quantity and sparsity



# NLQ annotation procedure

## Step 1: Preview long video

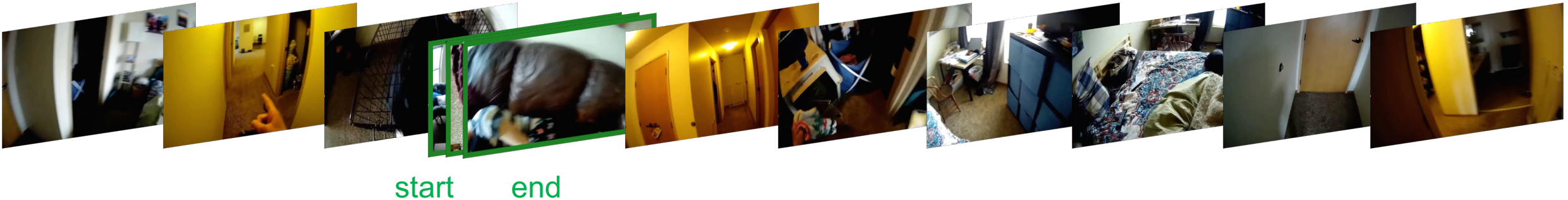


## Step 2: Formulate creative question

*Who did I interact with when I played with the dog for the second time in the living room?*

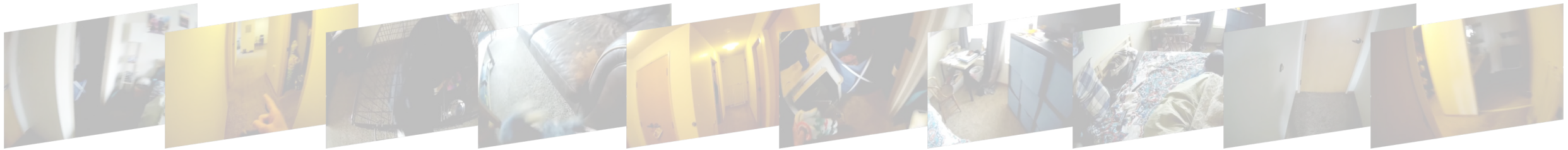
- Template-based
- Unambiguous
- Precise response localization

## Step 3: Annotate response (start, end) times



# NLQ annotation procedure

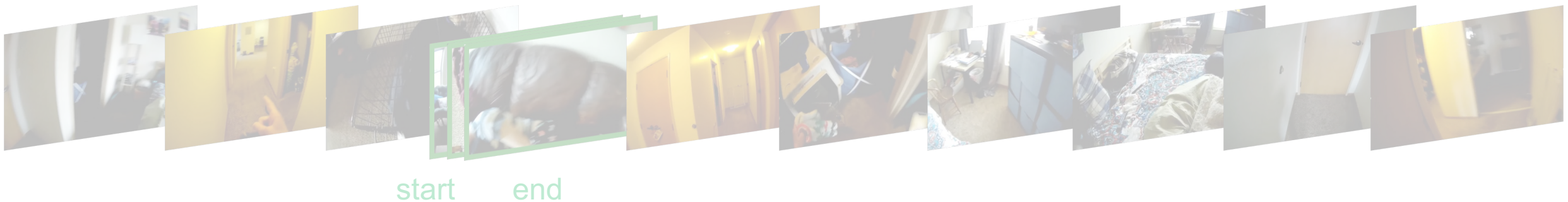
## Step 1: Preview long video



Expensive and slow process limits scalability of annotations

- Template-based
- Unambiguous
- Precise response localization

## Step 3: Annotate response (start, end) times



# NaQ: Narrations-as-Queries

**Key insight:** Augment NLQ training by learning to localize *narrations*  
Timestamped play-by-play descriptions of camera-wearer's activities.



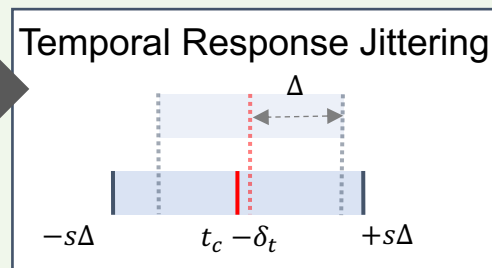
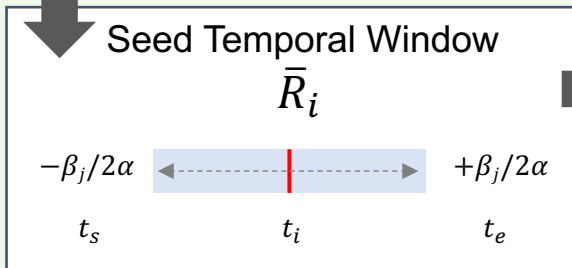
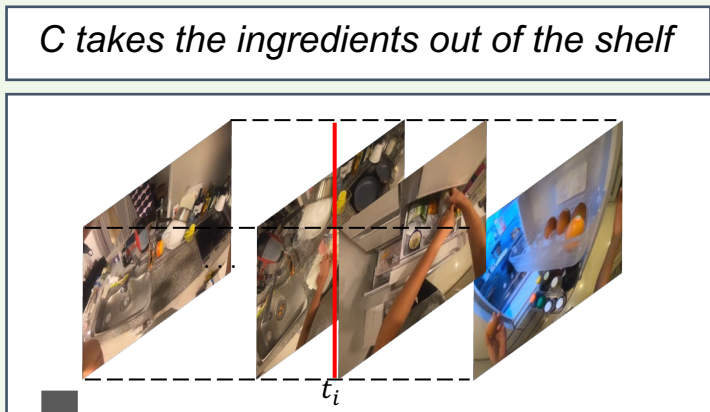
- ✓ **Easier to annotate**
  - Describe as you watch the video
- ✓ **Available on a large scale**
  - 200x more narrations than NLQ annotations
- ✓ **Multi-purpose annotations**
  - Not annotated specifically for NLQ
  - Applications across several benchmarks
  - Likely to be expanded over time



# NaQ data-augmentation for scaling NLQ

**Simple-yet-effective approach:** Augment NLQ dataset using NaQ and perform large-scale training

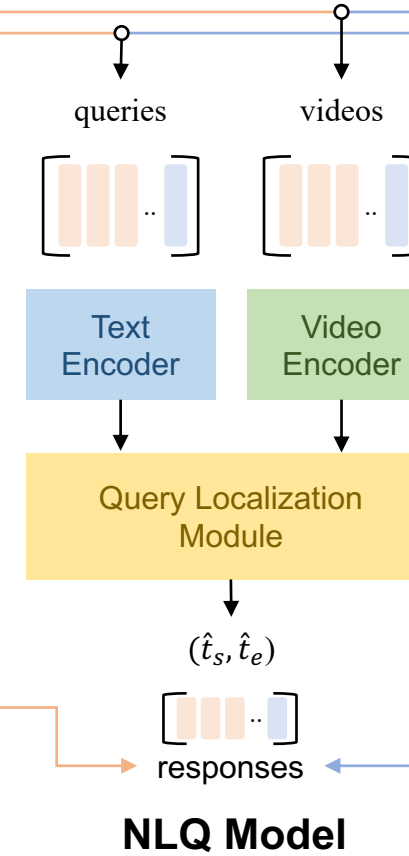
## Narrations-as-Queries (NaQ)



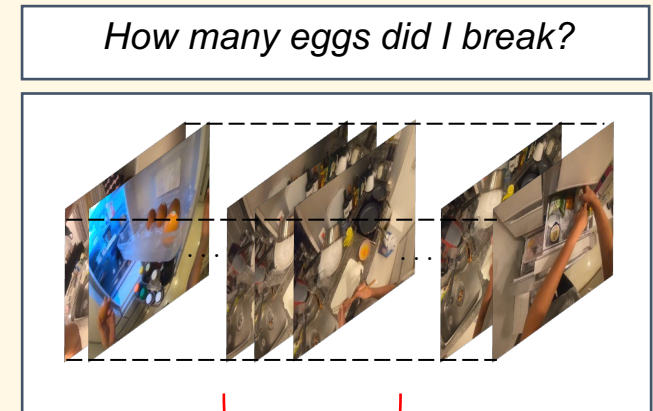
$R_i$

$T_i$

$V_j$



## NLQ Dataset



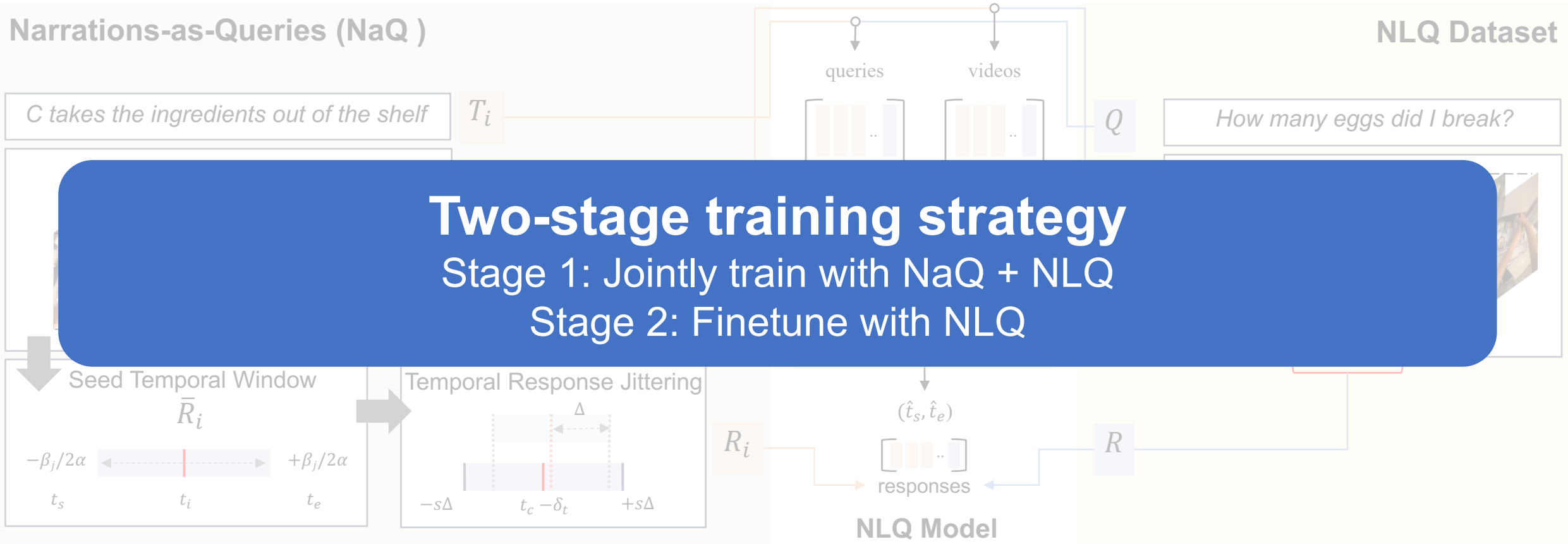
$Q$

$V$

$R$

# NaQ data-augmentation for scaling NLQ

**Simple-yet-effective approach:** Augment NLQ dataset using NaQ and perform large-scale training



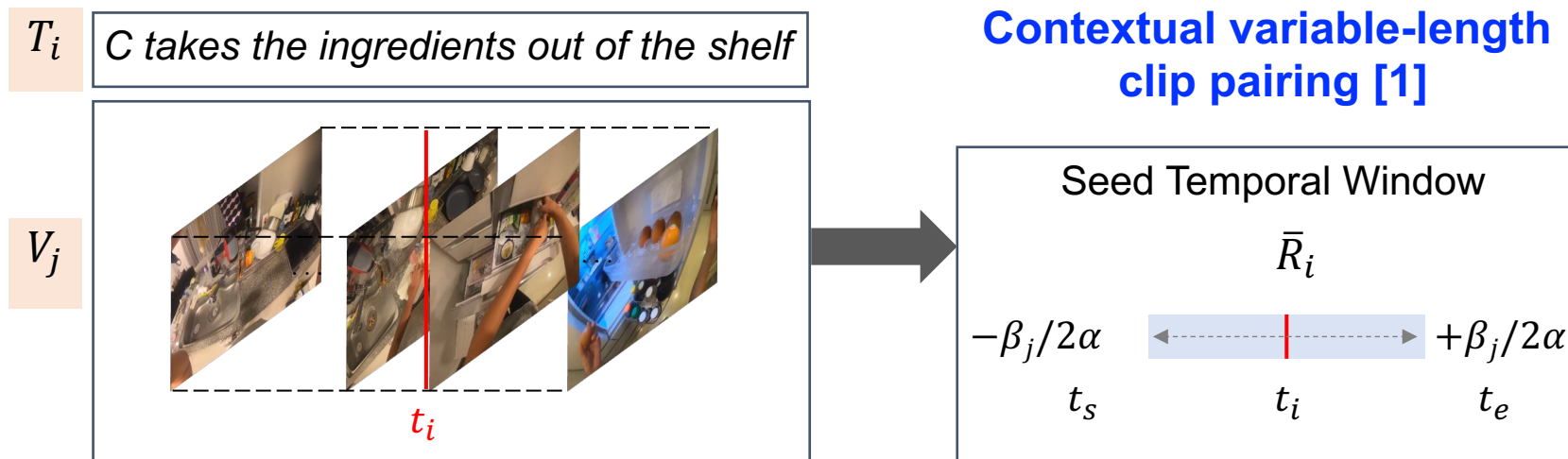
# Converting narrations $\rightarrow$ NLQ queries

**Narration annotation:**  $\langle V_j, T_i, t_i \rangle$

$V_j$  : Video  
 $T_i$  : Narration text  
 $t_i$  : Time-stamp

**NaQ annotation for NLQ:**  $\langle V_j, T_i, R_i \rangle$

$V_j$  : Video  
 $T_i$  : Narration text as query  
 $R_i$  :  $(t_s, t_e)$  response window



$\beta_j$  = average separation between consecutive narrations in video  $j$   
 $\alpha$  = average of  $\beta_j$  over all videos

# Converting narrations $\rightarrow$ NLQ queries

Narration annotation:  $\langle V_j, T_i, t_i \rangle$

$V_j$  : Video  
 $T_i$  : Narration text  
 $t_i$  : Time-stamp



NaQ annotation for NLQ:  $\langle V_j, T_i, R_i \rangle$

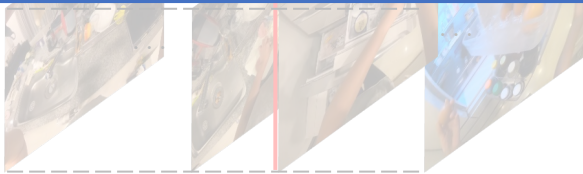
$V_j$  : Video  
 $T_i$  : Narration text as query  
 $R_i$  :  $(t_s, t_e)$  response window

NaQ augmentation significantly expands the training data

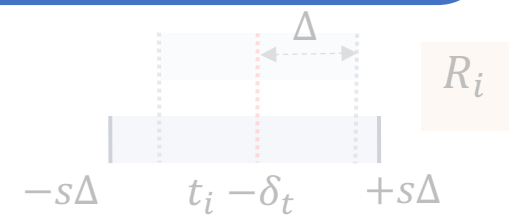
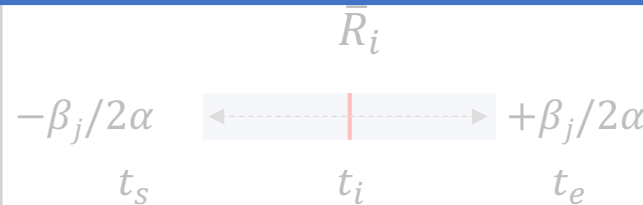
- 11k  $\rightarrow$  860k queries
- 1k  $\rightarrow$  5k video clips

$T_i$

$V_j$



$t_i$



$s$  = randomly sampled scaling factor  
 $\delta_t$  = random translation factor  
 $\Delta$  = half-width of original temporal window



# Experimental setup

## Dataset

Ego4D NLQ dataset [1]

## Evaluation metrics

**Mean Recall @ k:** Recall @ top k retrieval averaged over IoU=[0.3, 0.5]

## Baselines

**VSLNet [1,2]:** Span-based localization approach to vision-language grounding

**EgoVLP [3] :** Enhances VSLNet with clip features learned through egocentric video-language pretraining

**ReLER\* [4] :** Improves over VSLNet architecture + uses video-level data augmentation

\*we further improve the ReLER baseline using EgoVLP features

[1] Grauman, Kristen, et al. "Ego4d: Around the world in 3,000 hours of egocentric video." *CVPR 2022*

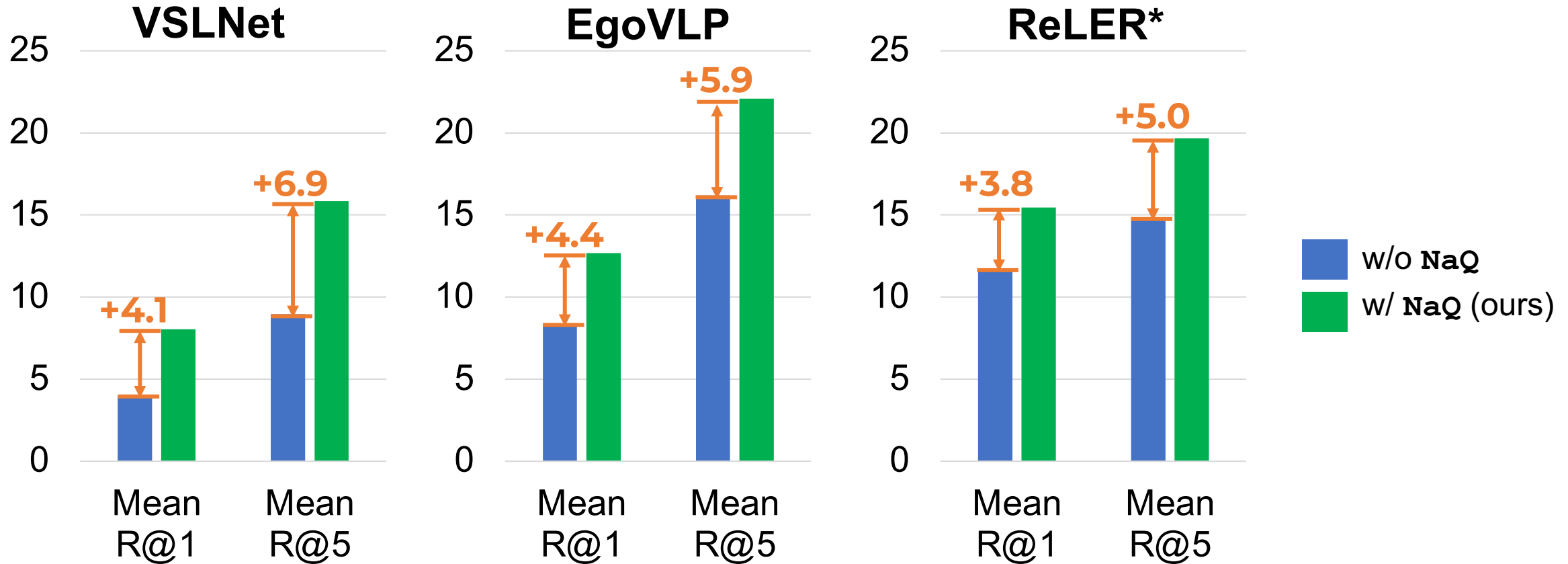
[2] Zhang, Hao, et al. "Span-based Localizing Network for Natural Language Video Localization." *ACL 2020*

[3] Lin, Kevin Qinghong, et al. "Egocentric video-language pretraining." *NeurIPS 2022*

[4] Shao, Jiayi, Xiaohan Wang, and Yi Yang. "ReLER@ ZJU Submission to the Ego4D Moment Queries Challenge 2022." *arXiv 2022*

# Experimental results

NaQ augmentation *consistently* and *significantly* enhances all baselines



Our approach improves NLQ performance by **up to 7% absolute mean recall**

\*we further improve the ReLER baseline using EgoVLP features

# Experimental results

NaQ sets the state-of-the-art results on the public Ego4D NLQ leaderboard

Method	R@1	R@1	Mean	R@5	R@5
	IoU=0.3	IoU=0.5	R@1 <sup>†</sup>	IoU=0.3	IoU=0.5
<b>NaQ++ (ours)<sup>‡</sup></b>	<b>21.70</b>	<b>13.64</b>	<b>17.67</b>	25.12	16.33
<b>NaQ (ours)</b>	<b>18.46</b>	<b>10.74</b>	<b>14.59</b>	21.50	13.74
InternVideo [5]	16.46	10.06	13.26	22.95	16.11
Badgers@UW-Mad. [27]	15.71	9.57	12.64	<b>28.45</b>	<b>18.03</b>
CONE [18]	15.26	9.24	12.25	26.42	16.51
ReLER [24]	12.89	8.14	10.51	15.41	9.94
EgoVLP [23]	10.46	6.24	8.35	16.76	11.29
VSLNet [38]	5.42	2.75	4.08	8.79	5.07

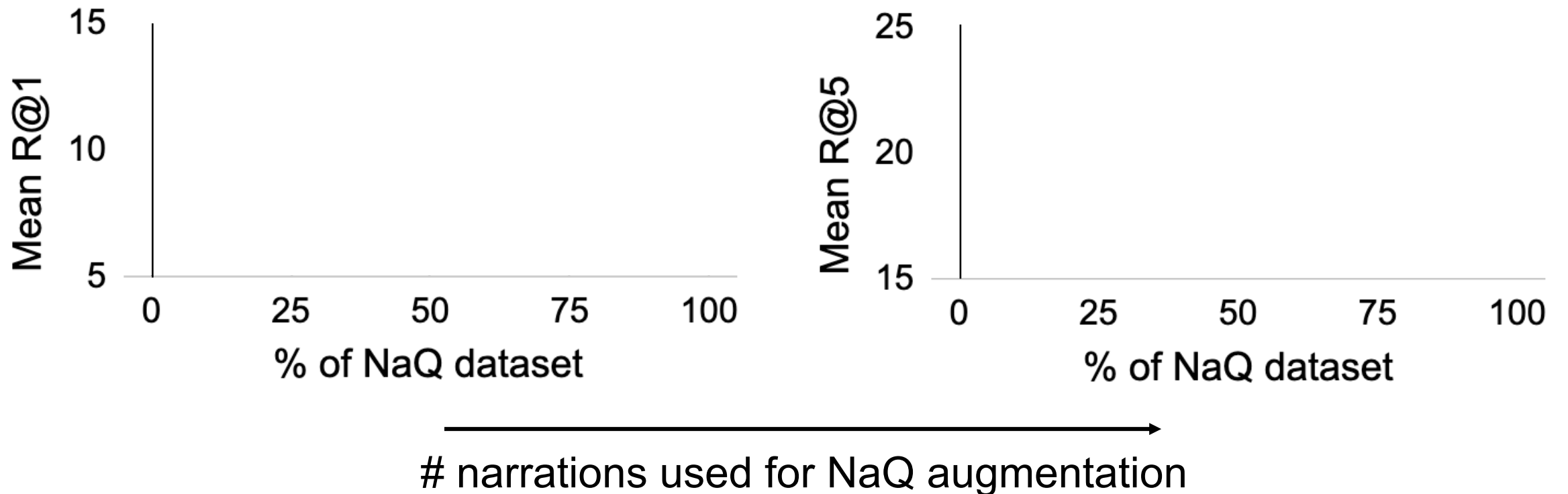
<sup>†</sup> Mean R@1 is the primary metric for deciding challenge winners

<sup>‡</sup> NaQ++ combines winning entries from prior challenges and NaQ to achieve SoTA

Our approach improves NLQ SotA by 4.5% absolute mean recall @ 1

# Experimental results

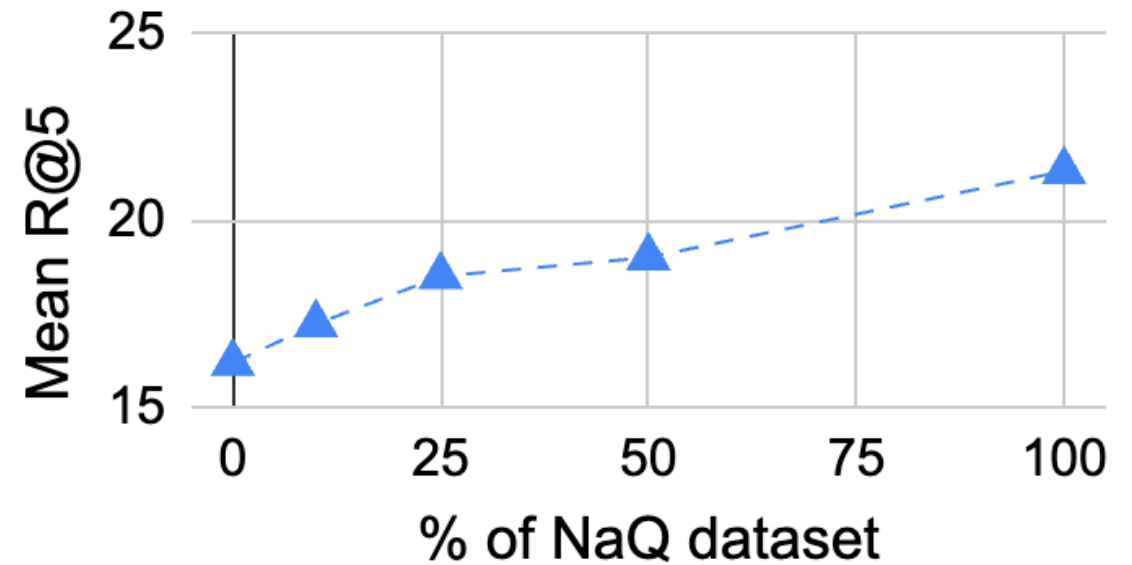
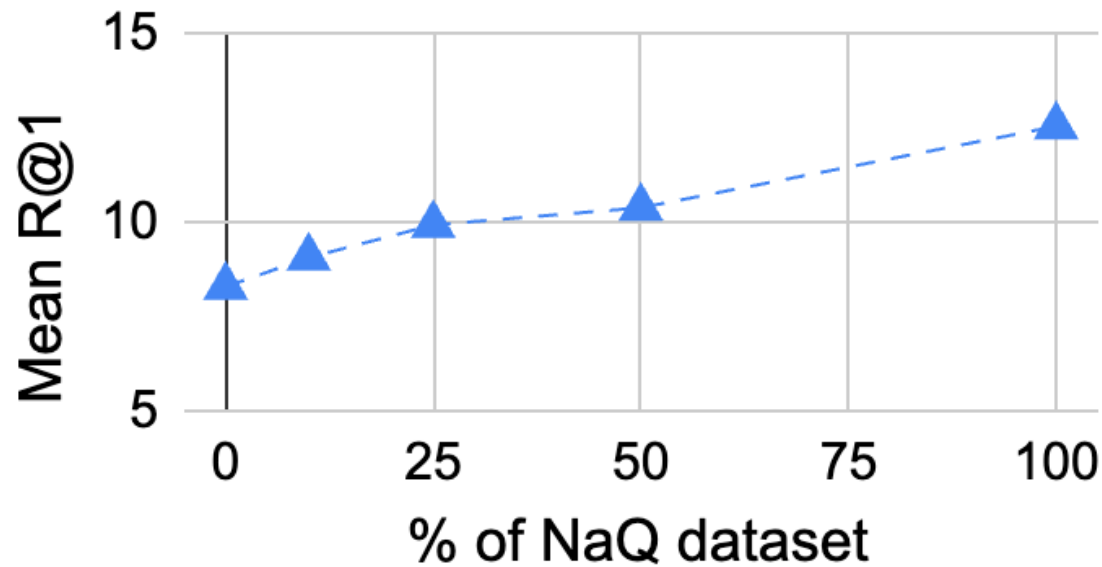
NaQ performance scales with the number of narrations used for training





# Experimental results

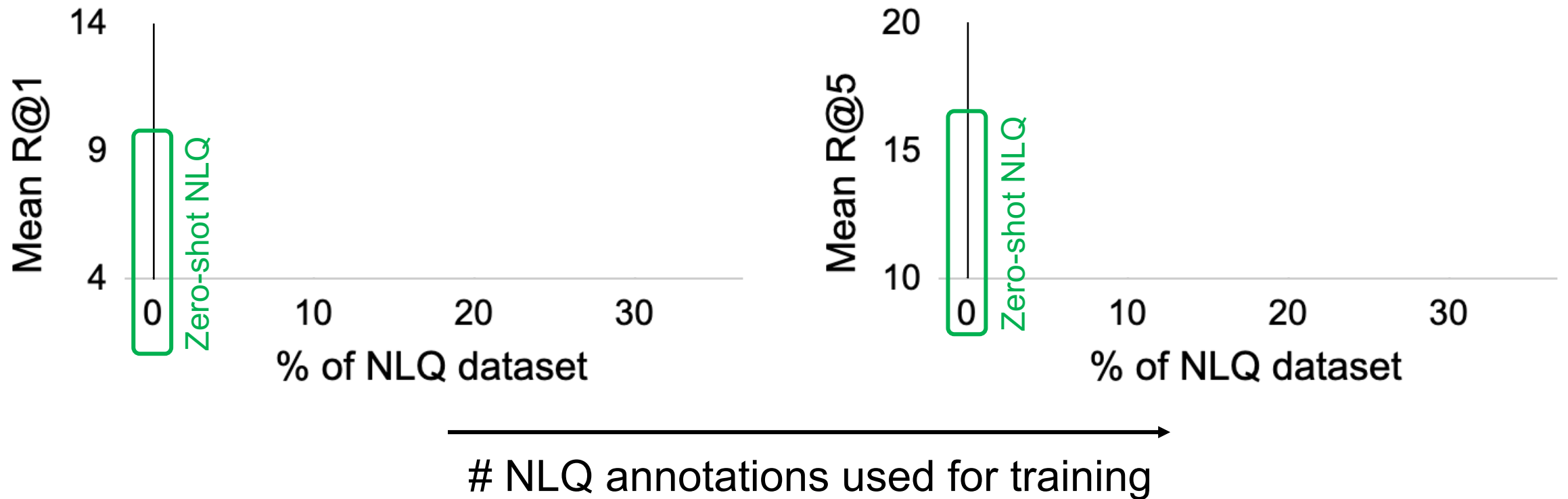
NaQ performance scales with the number of narrations used for training



→  
# narrations used for NaQ augmentation

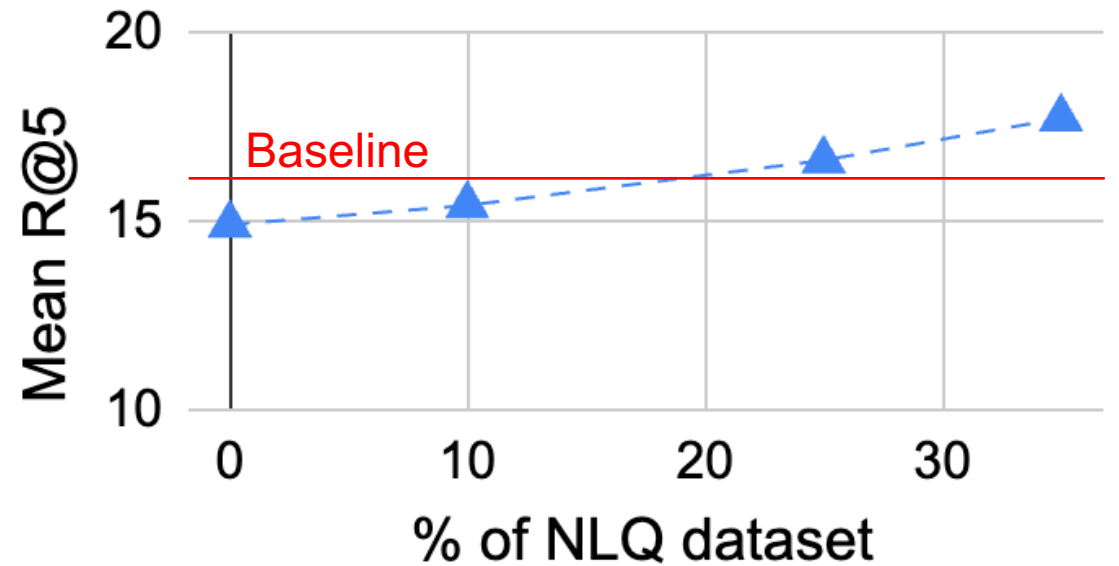
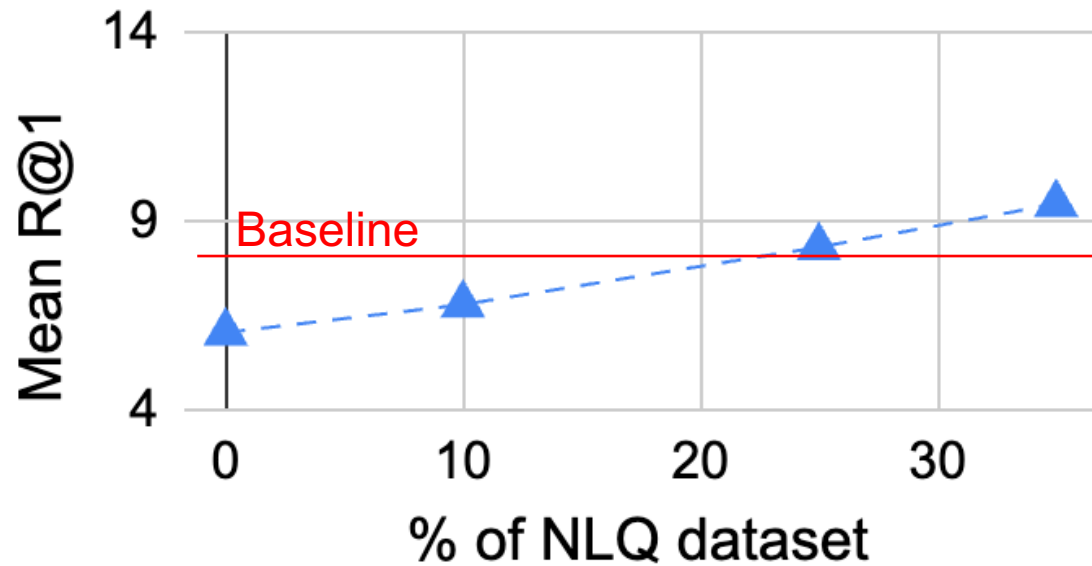
# Experimental results

NaQ facilitates zero-/few-shot NLQ



# Experimental results

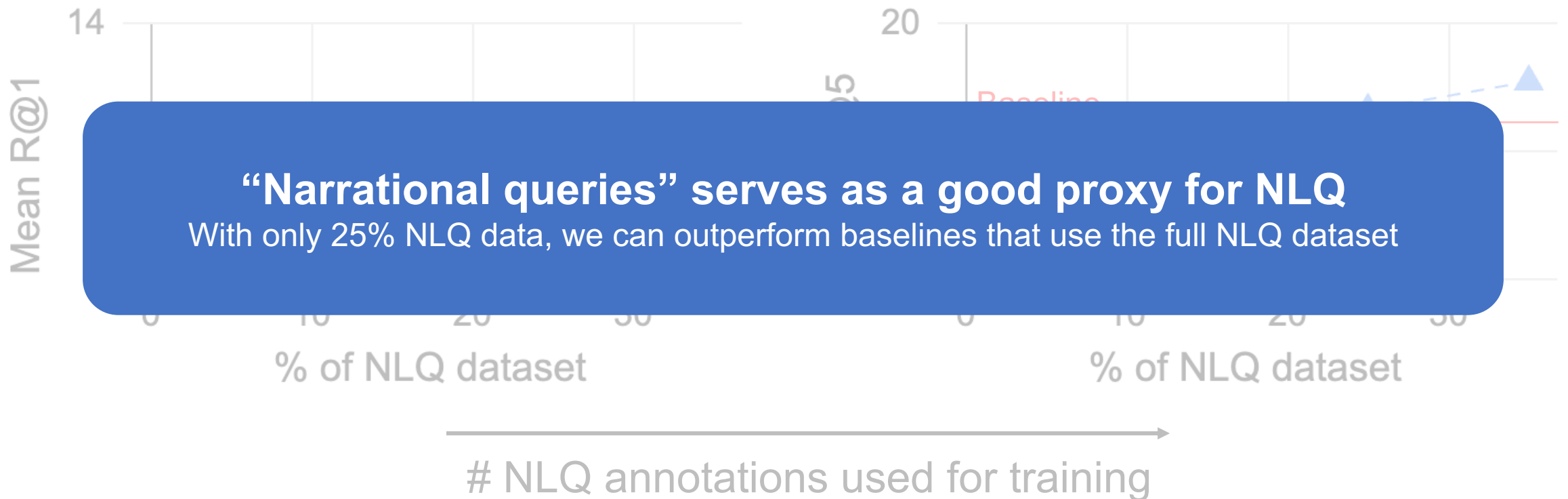
NaQ facilitates zero-/few-shot NLQ



→  
# NLQ annotations used for training

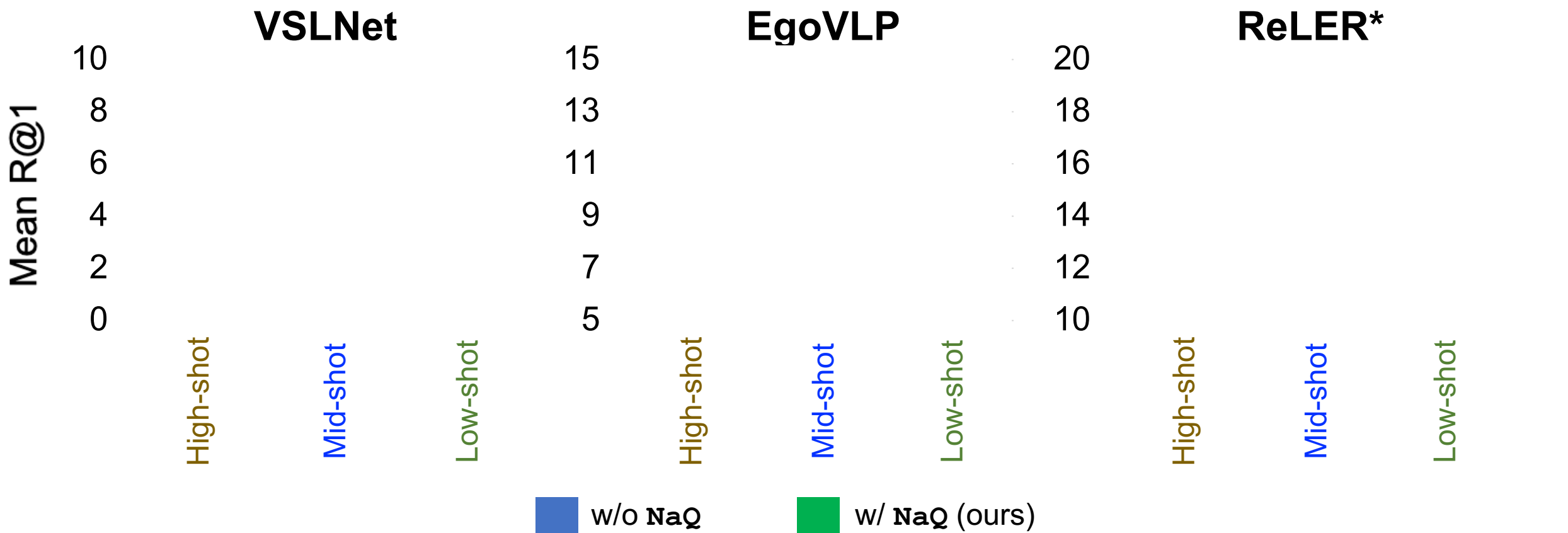
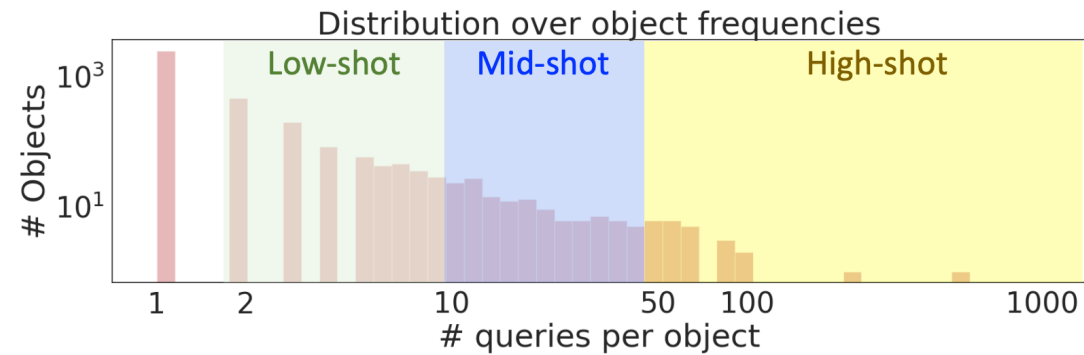
# Experimental results

NaQ facilitates zero-/few-shot NLQ



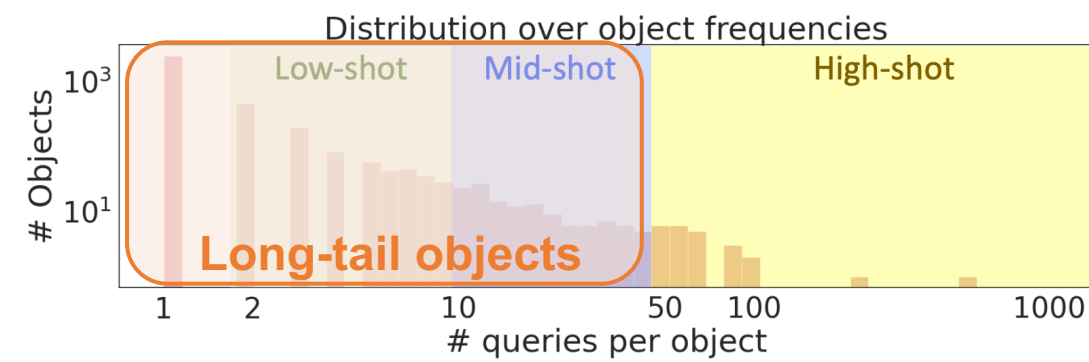


# Experimental results

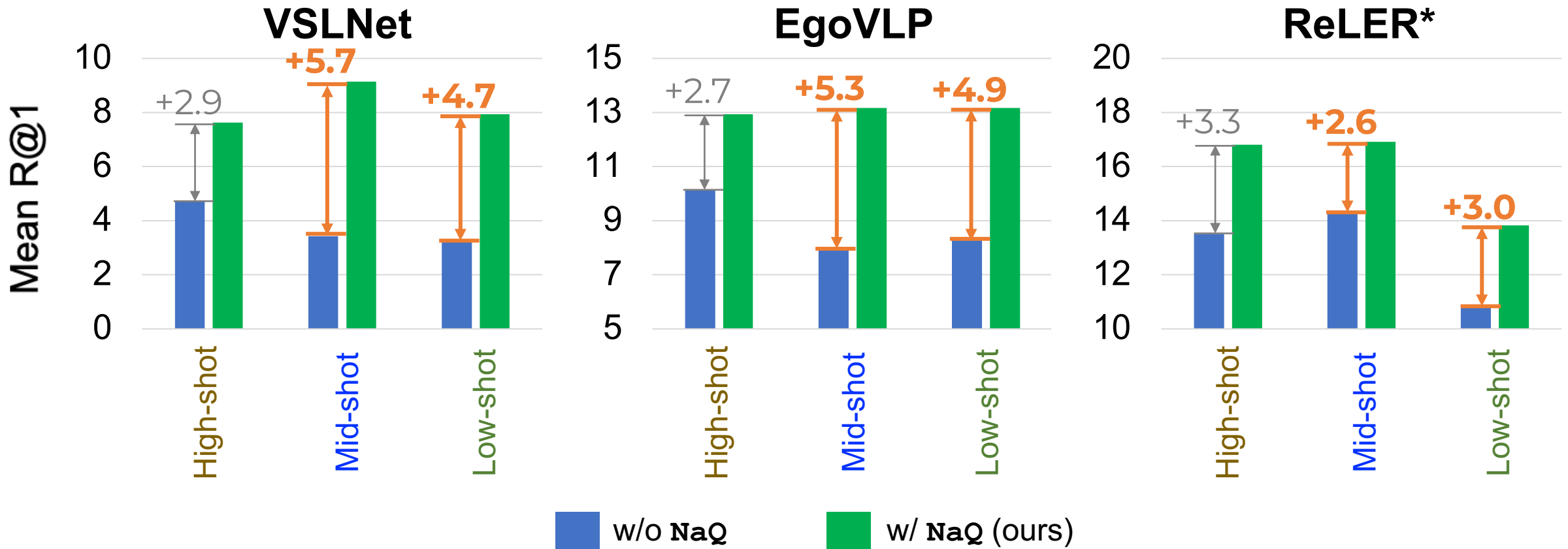


\*we further improve the ReLER baseline using EgoVLP features

# Experimental results

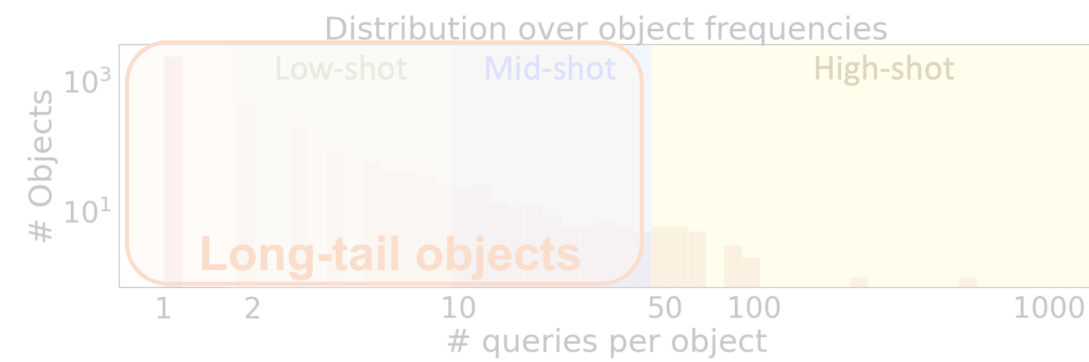


NaQ significantly improves responding to queries about **long-tail objects**

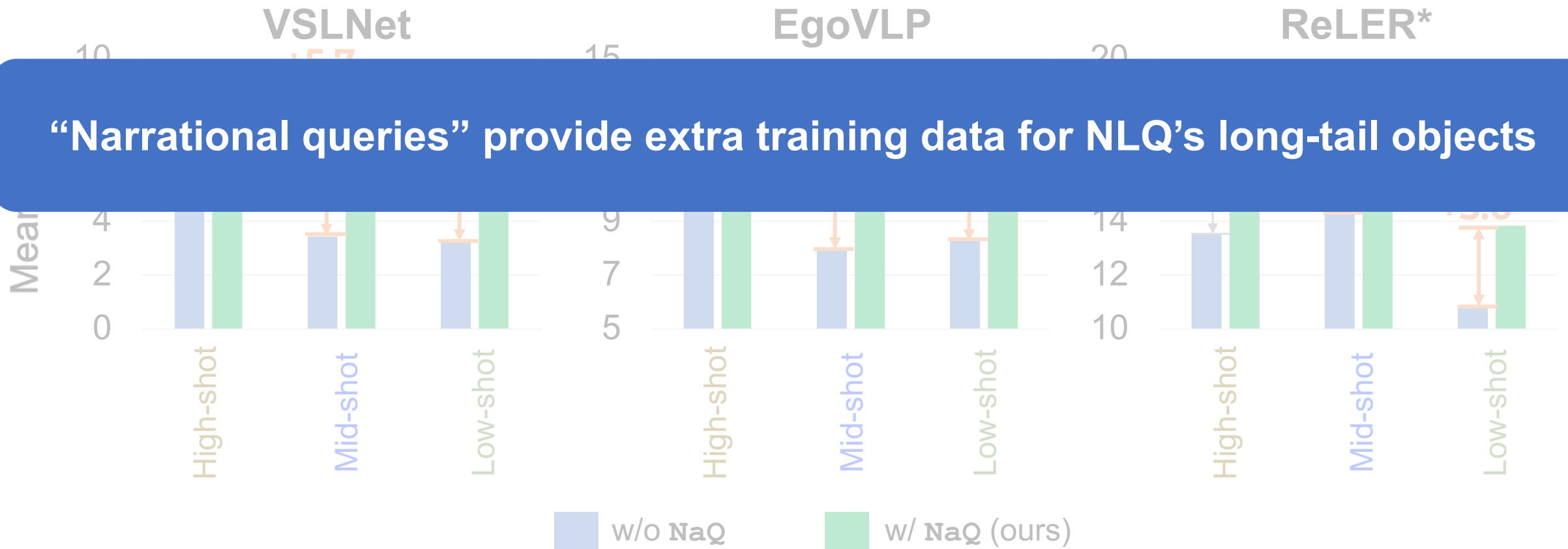


\*we further improve the ReLER baseline using EgoVLP features

# Experimental results



NaQ significantly improves responding to queries about **long-tail objects**



\*we further improve the ReLER baseline using EgoVLP features



# Qualitative results

NaQ succeeds, while baseline fails, to reason about the long-tail object “sieve”

Query: *Where did I last put the sieve?*

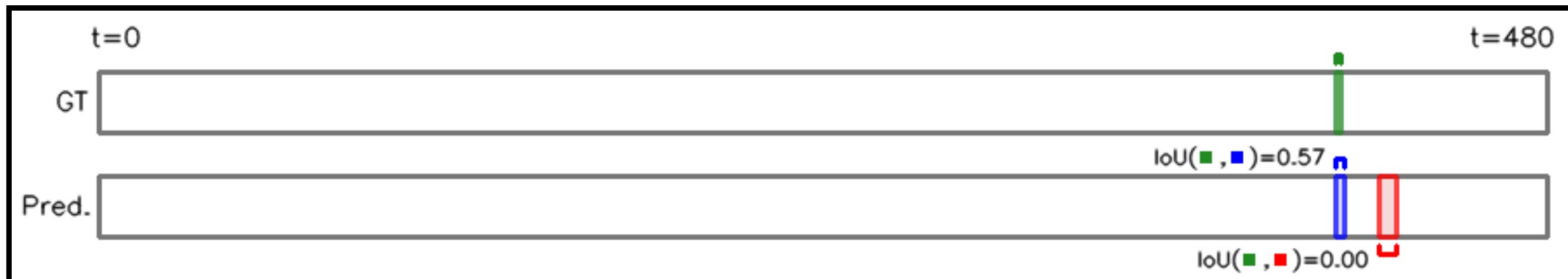
Ground-truth



ReLER\* + NaQ

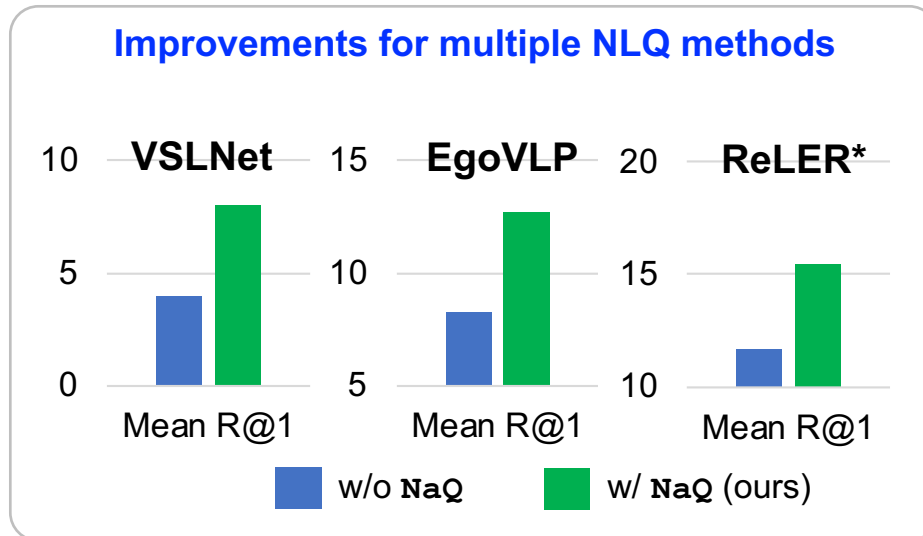


ReLER\*



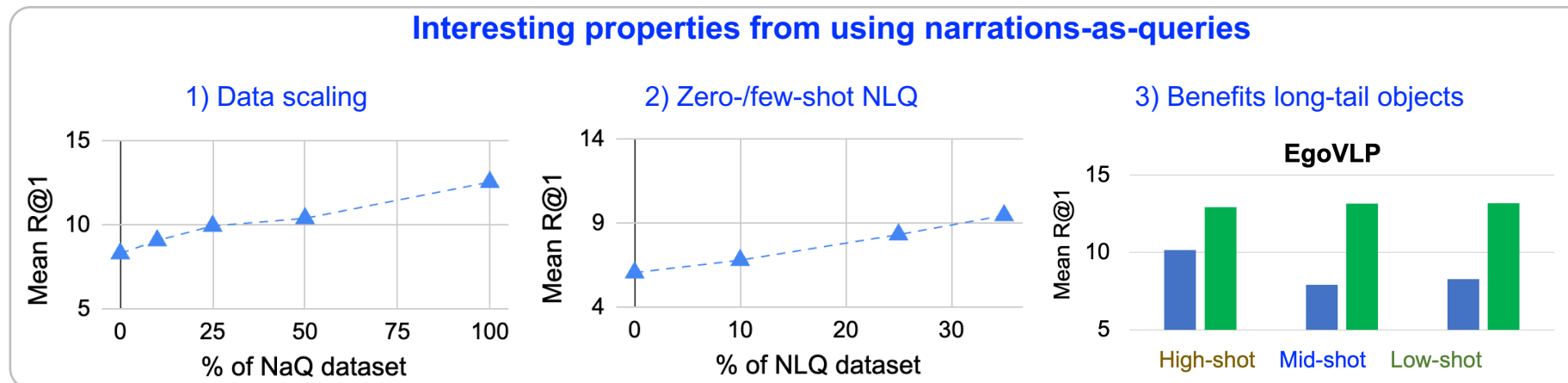
# Conclusion

## NaQ: Simple-yet-effective augmentation strategy for Episodic Memory NLQ



### Obtains SotA on Ego4D NLQ

Method	R@1 IoU=0.3	R@1 IoU=0.5	Mean R@1 <sup>†</sup>
<b>NaQ++ (ours)<sup>‡</sup></b>	<b>21.70</b>	<b>13.64</b>	<b>17.67</b>
<b>NaQ (ours)</b>	<b>18.46</b>	<b>10.74</b>	<b>14.59</b>
InternVideo [5]	16.46	10.06	13.26
Badgers@UW-Mad. [27]	15.71	9.57	12.64
CONE [18]	15.26	9.24	12.25
VSLNet [38]	5.42	2.75	4.08





# NaQ: Leveraging Narrations as Queries to Supervise Episodic Memory



Santhosh Kumar Ramakrishnan<sup>1</sup>



Ziad Al-Halah<sup>2</sup>



Kristen Grauman<sup>1,3</sup>

**Poster session:** TUE-PM-245

**Project page:** <https://vision.cs.utexas.edu/projects/naq/>

**Code:** <https://github.com/srama2512/NaQ>