

# Latency Matters: Real-Time Action Forecasting Transformer



Harshayu  
Girase\*



Nakul  
Agarwal



Chiho  
Choi



Karttikeya  
Mangalam\*

\*denotes equal technical contribution

Poster Session THU-AM-218

JUNE 18-22, 2023

**CVPR**   
VANCOUVER, CANADA



# Problem Formulation

## Input

Video consisting of past frames  
(without action labels)



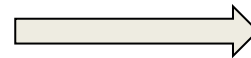
$F_0$



...



$F_t$



## Output

Predicted action at a  
predetermined time  $t_f$   
after the present

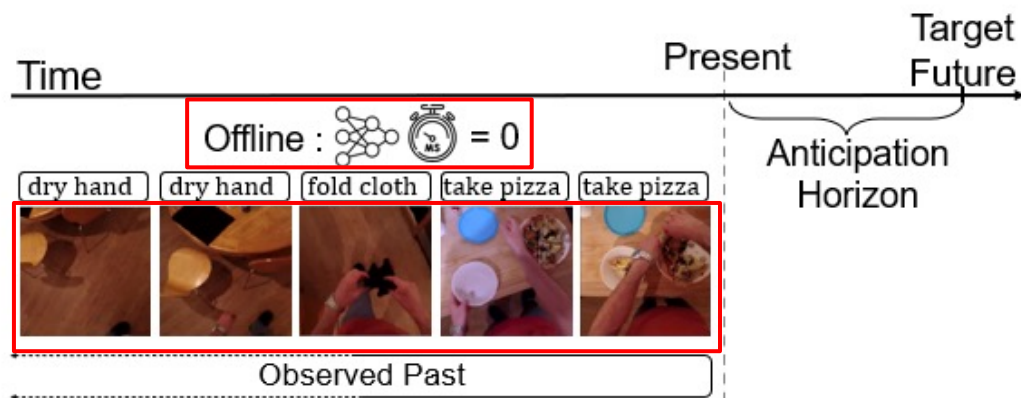
*wash knife*



# Real-time Action Forecasting Evaluation

## Offline Evaluation

- Ignores model latency
- Uses all video data up to present time T to predict the future



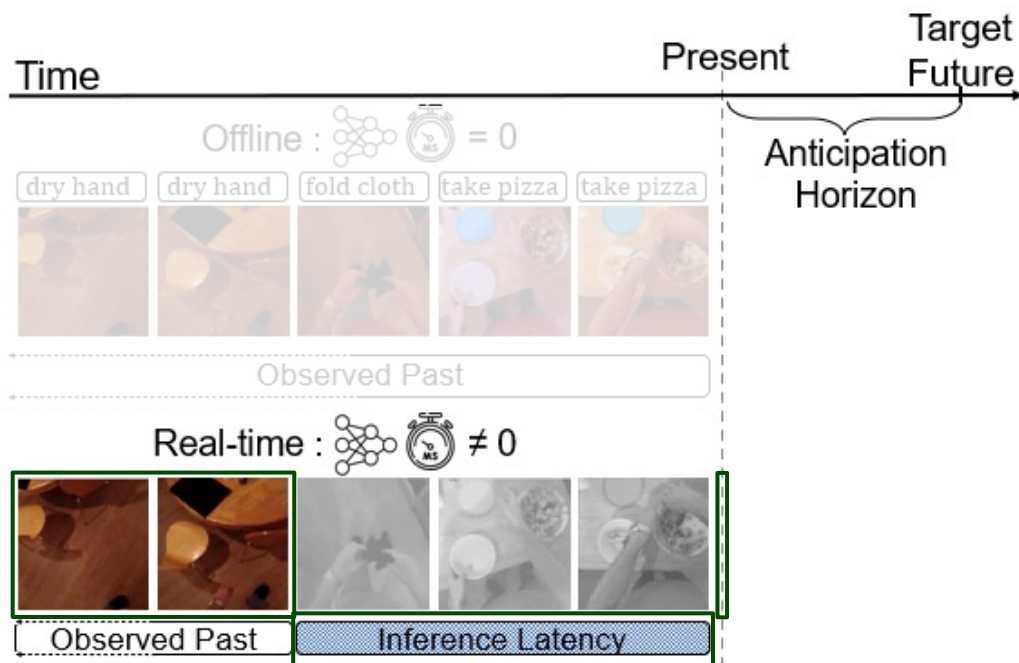
# Real-time Action Forecasting Evaluation

## Offline Evaluation

- Ignores model latency
- Uses all video data up to present time  $T$  to predict the future

## Real-time Evaluation

- Takes model latency into account
- Model can only use video data up to time  $T - t_{\text{latency}}$  to predict
- Forces prediction to arrive at time  $T$



# Real-time Action Forecasting Evaluation

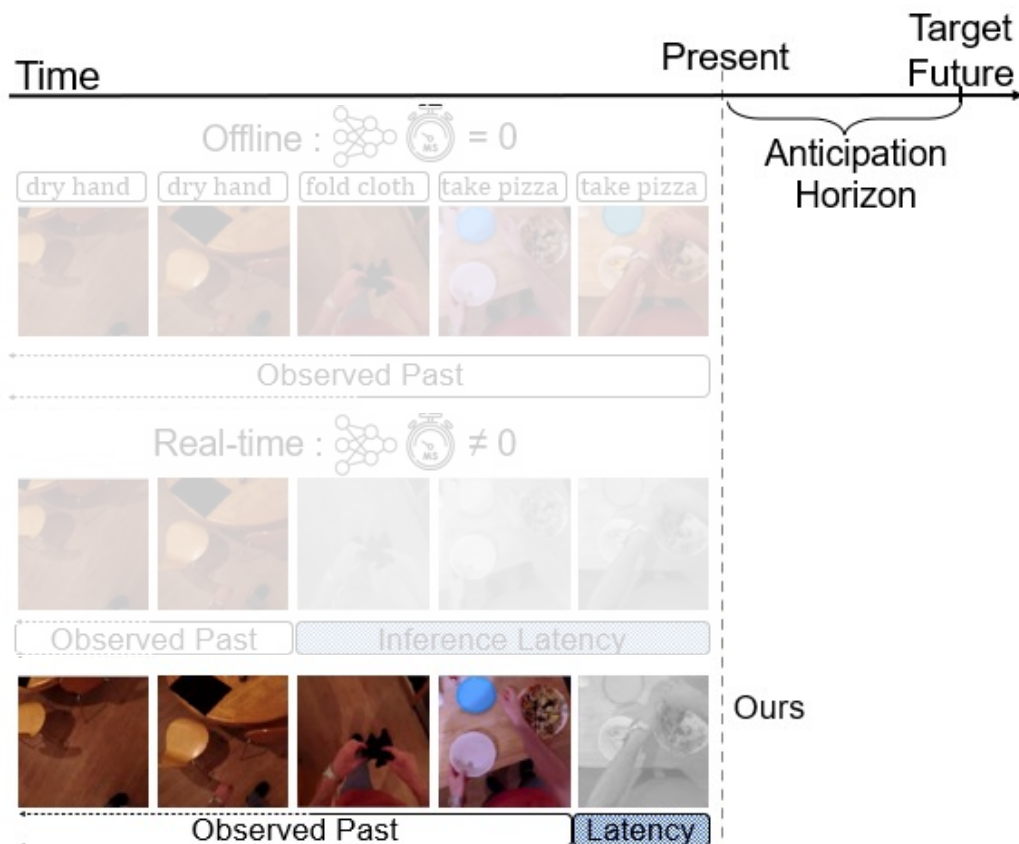
## Offline Evaluation

- Ignores model latency
- Uses all video data up to present time  $T$  to predict the future

## Real-time Evaluation

- Takes model latency into account
- Model can only use video data up to time  $T - t_{\text{latency}}$  to predict
- Forces prediction to arrive at time  $T$

We propose RAFTformer, a novel action anticipation transformer that balances high performance with low latency



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

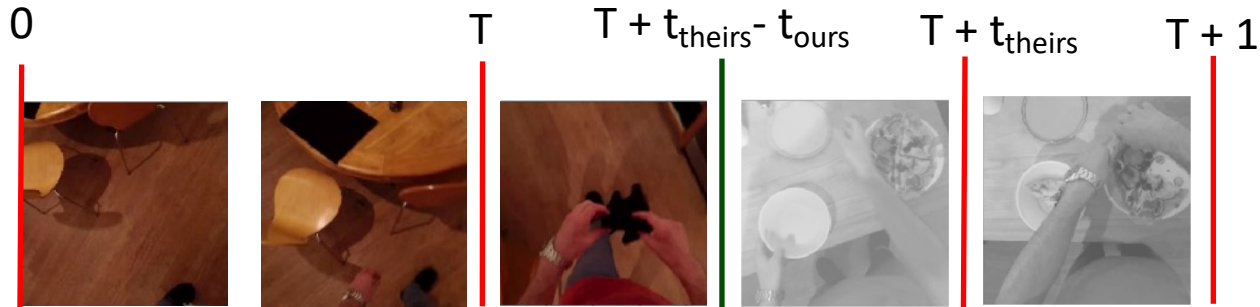
Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	26	23	40
RAFTformer	None	K700		33.7	37.1	18.0	26	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	

~9x less  
latency

State-of-the-art  
results    ~8x less  
parameters    ~94% less  
GPU hours



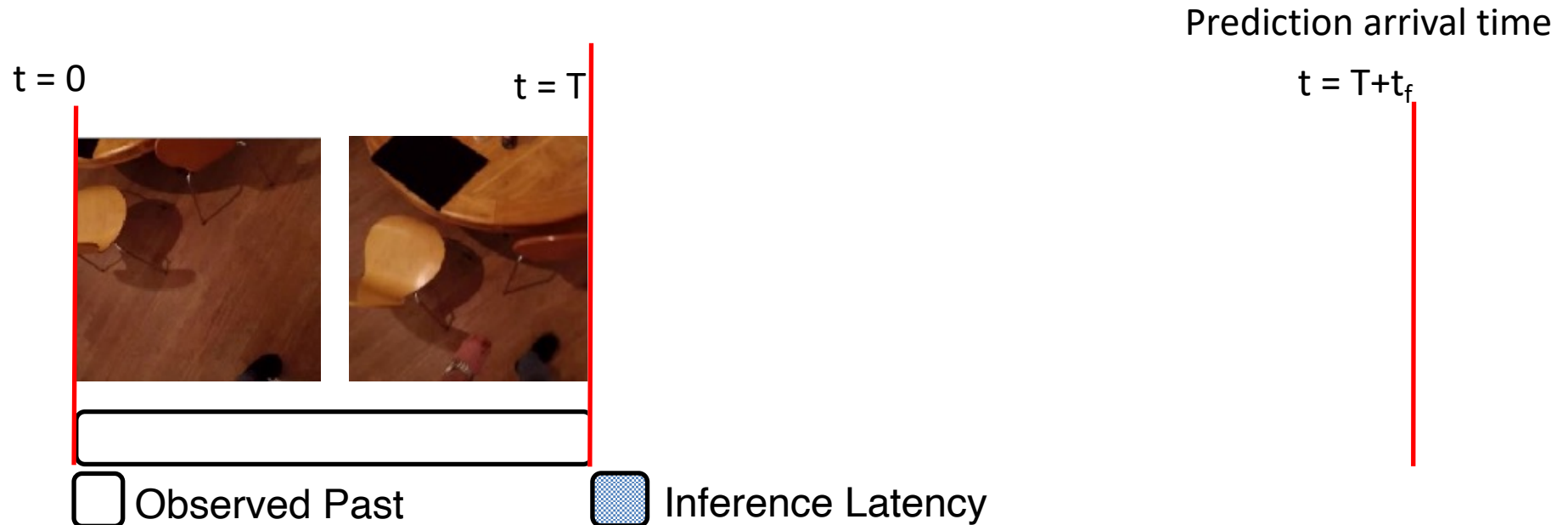
# Results: Online Setting



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action
AVT [25]	IN21K	$t_{avt} = 420$	$T$	$T + t_{avt}$	$T + 1$	30.2	31.7	14.9
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{avt} - t_{ours}$	$T + t_{avt}$	$T + 1$	34.1	38.2	<b>19.3 (+4.4)</b>
MemViT [80]	K400	$t_{vit} = 160$	$T$	$T + t_{vit}$	$T + 1$	32.8	33.2	15.1
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{vit} - t_{ours}$	$T + t_{vit}$	$T + 1$	33.8	37.1	<b>18.1 (+3.0)</b>
MemViT [80]	K700	$t_{vit} = 350$	$T$	$T + t_{vit}$	$T + 1$	32.2	37.0	17.7
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{vit} - t_{ours}$	$T + t_{vit}$	$T + 1$	33.7	37.9	<b>19.0 (+1.3)</b>

# Offline Forecasting Evaluation: Shortcomings

Forecasting models often must be done in real time.  
But current good forecasting models have high latency.



The latency can even be so high that the model furnishes predictions after the future action has already happened





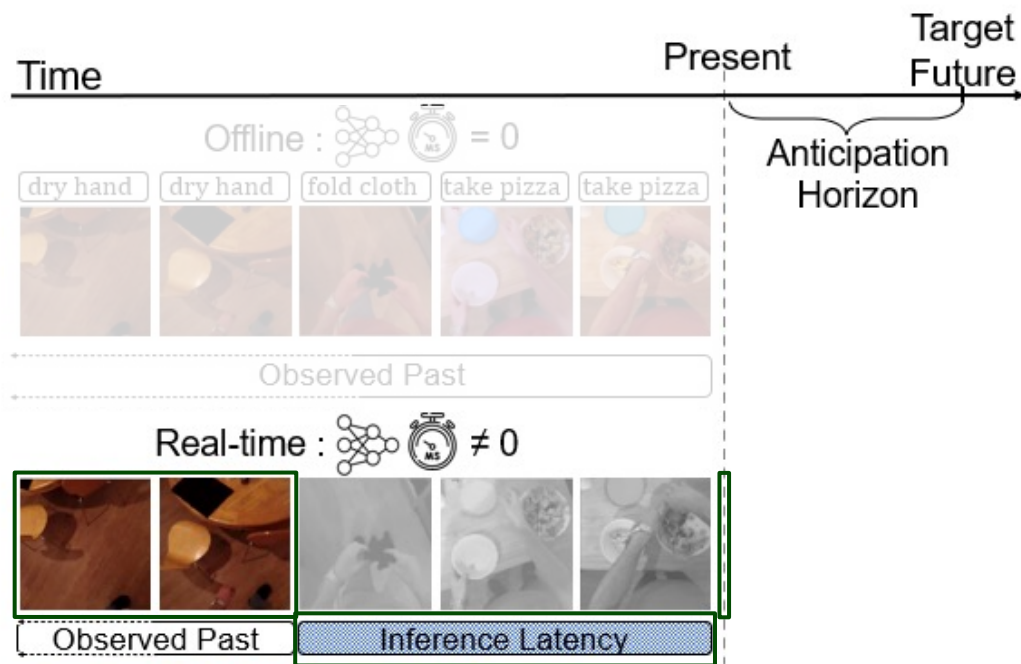
# Real-time Action Forecasting Evaluation

## Offline Evaluation

- Ignores model latency
- Uses all video data up to present time  $T$  to predict the future

## Real-time Evaluation

- Takes model latency into account
- Model can only use video data up to time  $T - t_{\text{latency}}$  to predict
- Forces prediction to arrive at time  $T$
- Tradeoff between latency and real-time performance. Larger models can lead to poorer real-time performance.



# Bigger is not necessarily better in real-time!

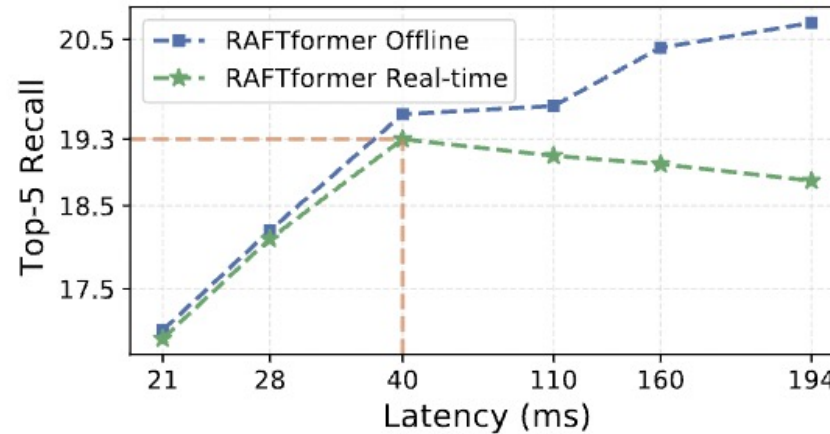
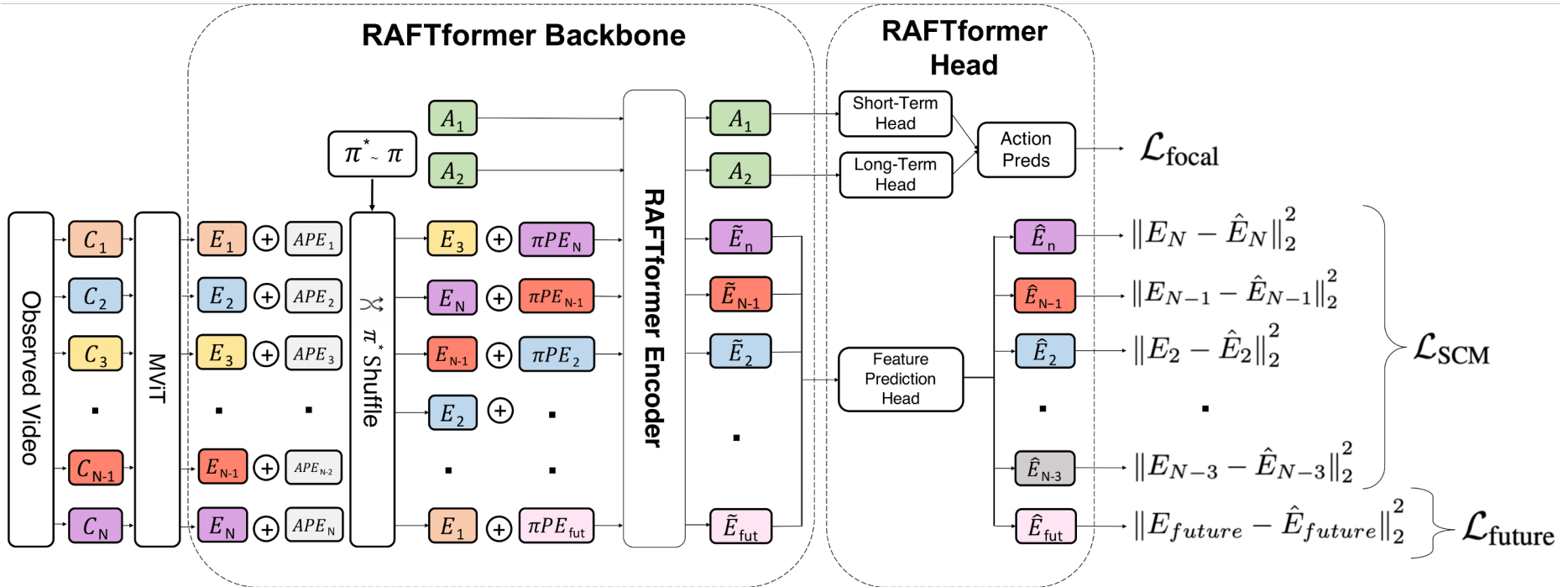
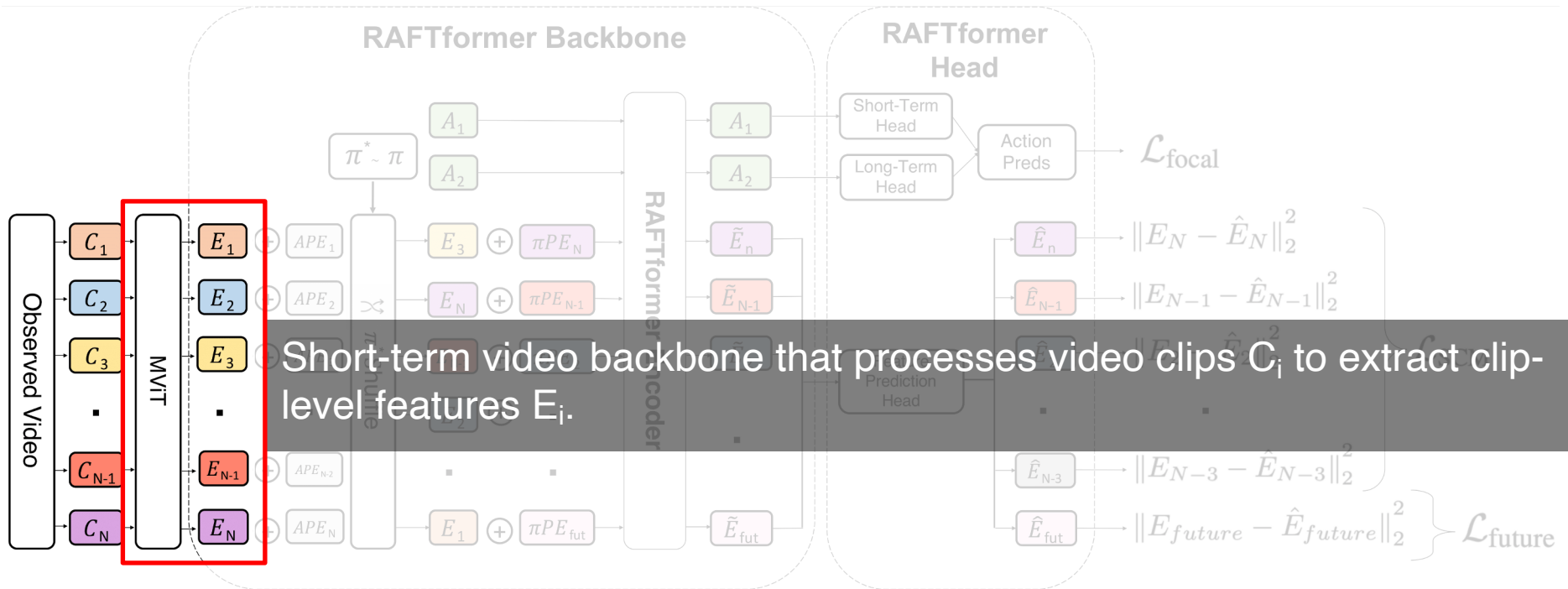


Figure 2. **Evaluation Performance vs. Latency.** Bigger models perform better in latency agnostic offline settings. In the real-time evaluation setting, we observe that, beyond a limit, bigger models with higher latency cause a drop in forecasting performance. In practical deployment, there exists a trade-off between latency and high-fidelity forecasts. See §4.3.1 for details.

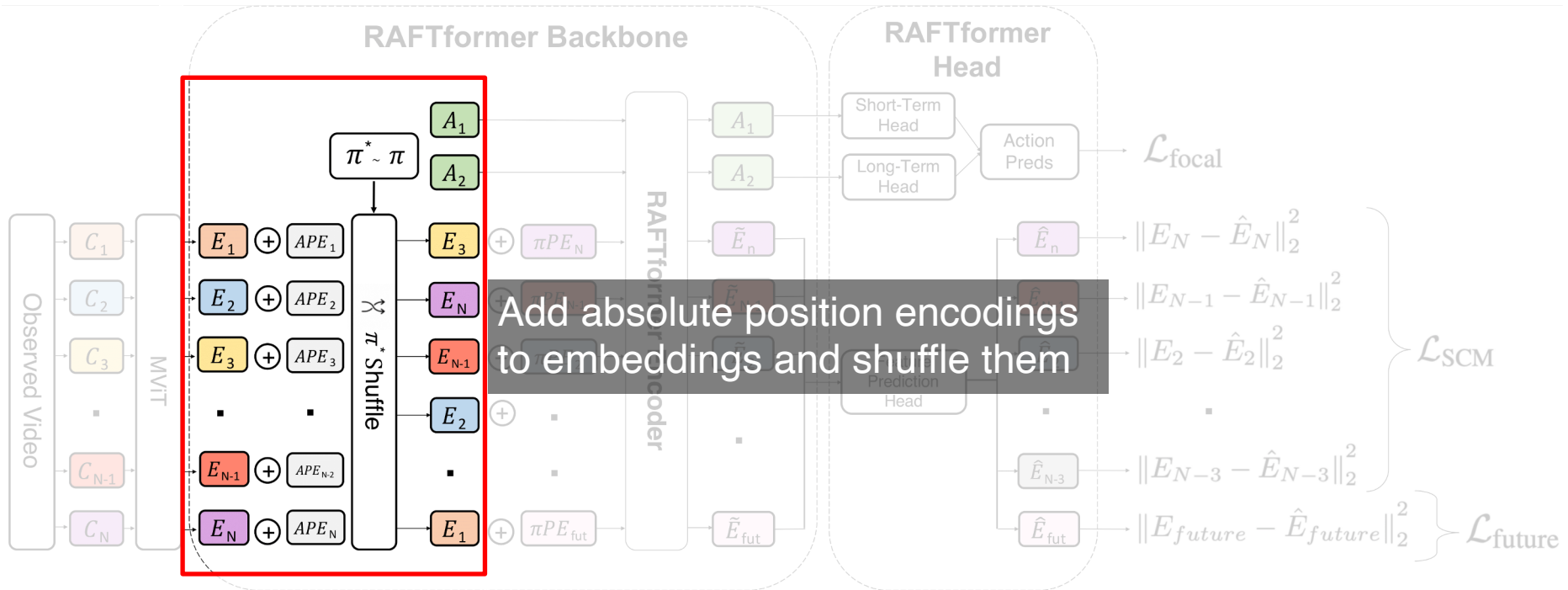
# RAFTformer Architecture



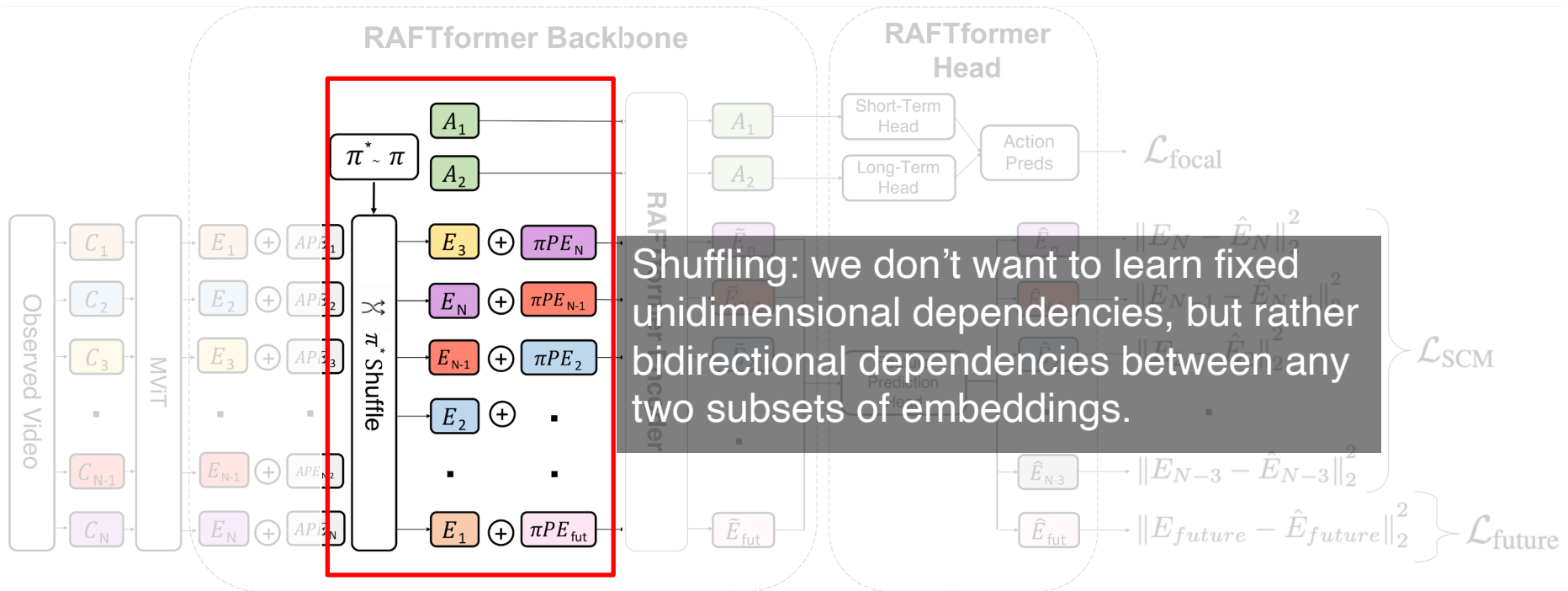
# RAFTformer Architecture



# RAFTformer Architecture

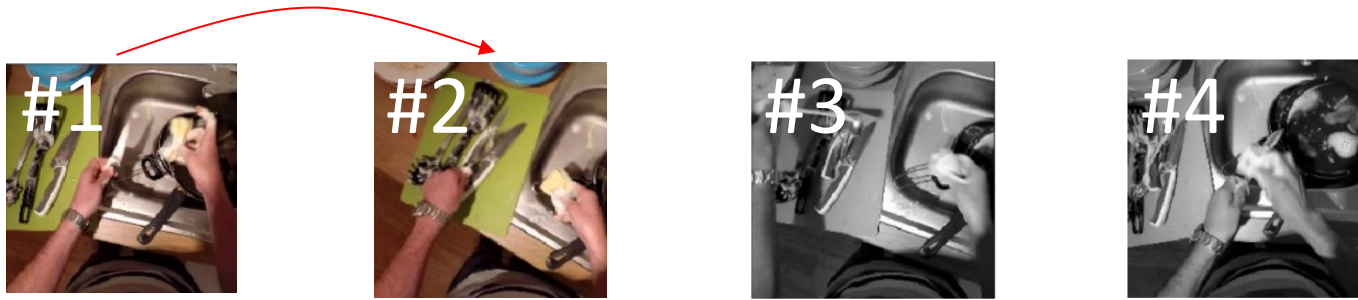


# RAFTformer Architecture



# Why Do Shuffling?

Without shuffling, a model using causal attention masking can only learn sequential unidimensional dependencies.



# Why Do Shuffling?

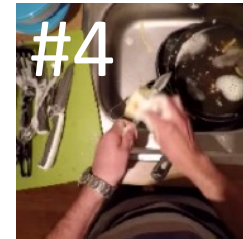
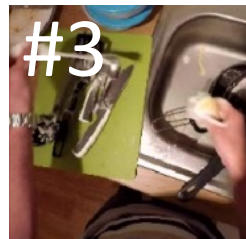
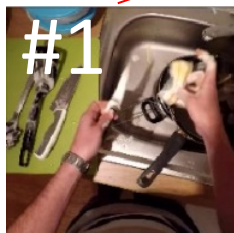
Without shuffling, a model using causal attention masking can only learn sequential unidimensional dependencies.





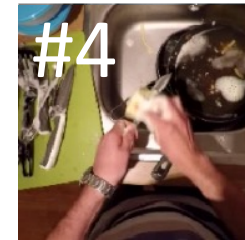
# Why Do Shuffling?

Without shuffling, a model using causal attention masking can only learn sequential unidimensional dependencies.



# Why Do Shuffling?

Shuffling tokens carefully\* allows the model to learn bidirectional dependencies between any two token subsets.

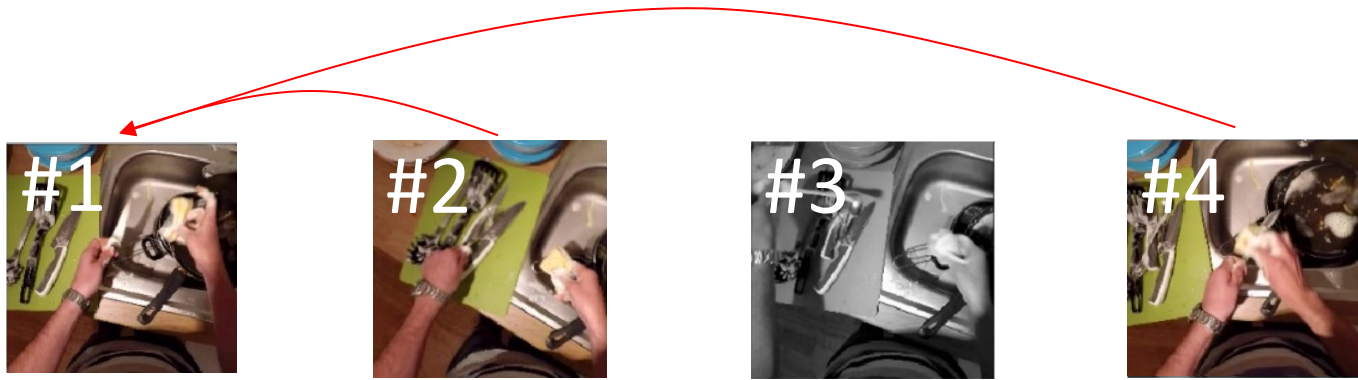


\*requires careful changes to masked attention to ensure no multi-hop information leakage through self-attention layers.  
Please see paper for details.



# Why Do Shuffling?

Shuffling tokens carefully\* allows the model to learn bidirectional dependencies between any two token subsets.



\*requires careful changes to masked attention to ensure no multi-hop information leakage through self-attention layers.  
Please see paper for details.



# Why Do Shuffling?

Shuffling tokens carefully\* allows the model to learn bidirectional dependencies between any two token subsets.



\*requires careful changes to masked attention to ensure no multi-hop information leakage through self-attention layers.  
Please see paper for details.



# Encoding Permutation $\pi$

Goal: We want RAFTformer encoder to “know” sampled permutation  $\pi^*$ , so we want to embed  $\pi^*$  vectorially for use in the encoder.

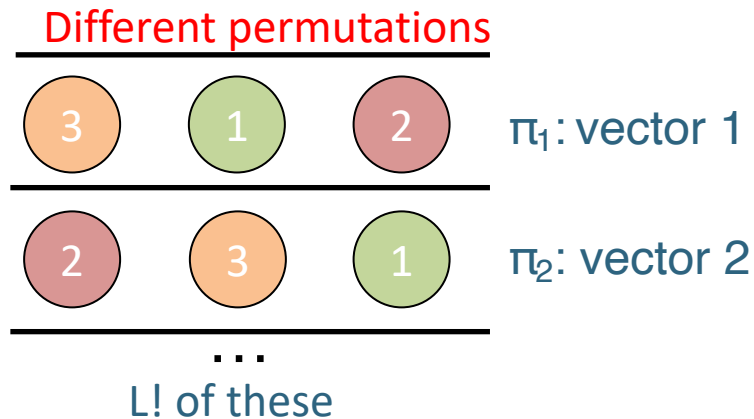


# Encoding Permutation $\pi$

## Naïve Method

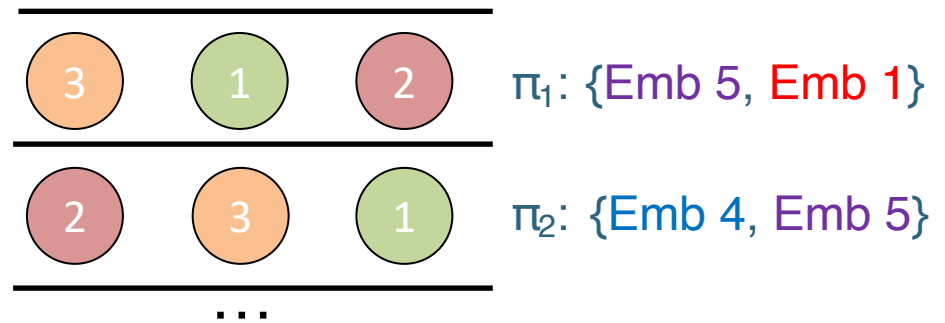
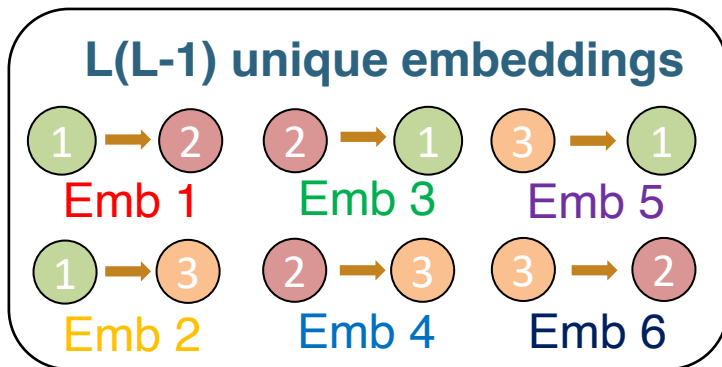
**Method:** Assign each  $\pi$  a single learnable vector

$L!$  unique embeddings needed to encode all possible permutations  $\pi$   
(one for each permutation)



# Encoding Permutation $\pi$ Predecessor Successor Method

**Method:** Encode each  $\pi$  as a set of predecessor  $\rightarrow$  successor relationships.  
 $L(L-1)$  unique embeddings needed to encode all possible permutations  $\pi$   
(one for each pair  $i \rightarrow j$ )



# Encoding Permutation $\pi$

## $\pi$ PE Method

**Method:** Encode only the successor PE in the permuted  $\pi^*$  (predecessor is already encoded through Absolute PE)

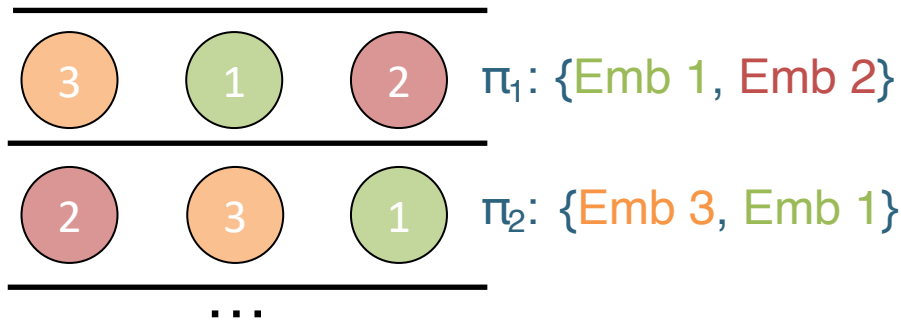
$L(L-1)$  unique embeddings needed to encode all possible permutations  $\pi$  (one for each  $i$ )

L unique embeddings

Emb 1 

Emb 2 

Emb 3 





# Illustration of $\pi$ PE

$$E_3 \oplus \pi PE_N$$

$$E_N \oplus \pi PE_{N-1}$$

$$E_{N-1} \oplus \pi PE_2$$

$$E_2 \oplus \cdot$$

·

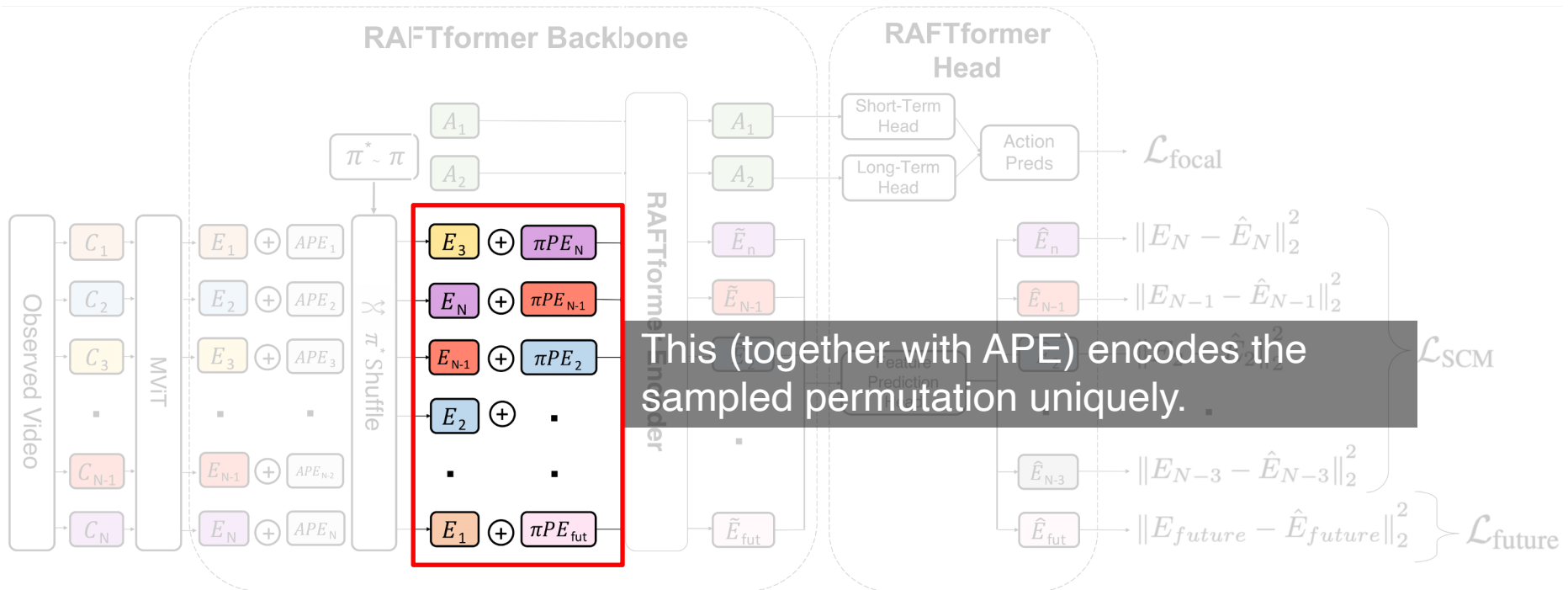
$$E_1 \oplus \pi PE_{\text{fut}}$$

$\pi$ PE is the encoding of the original temporal position of the successor in the permuted sequence

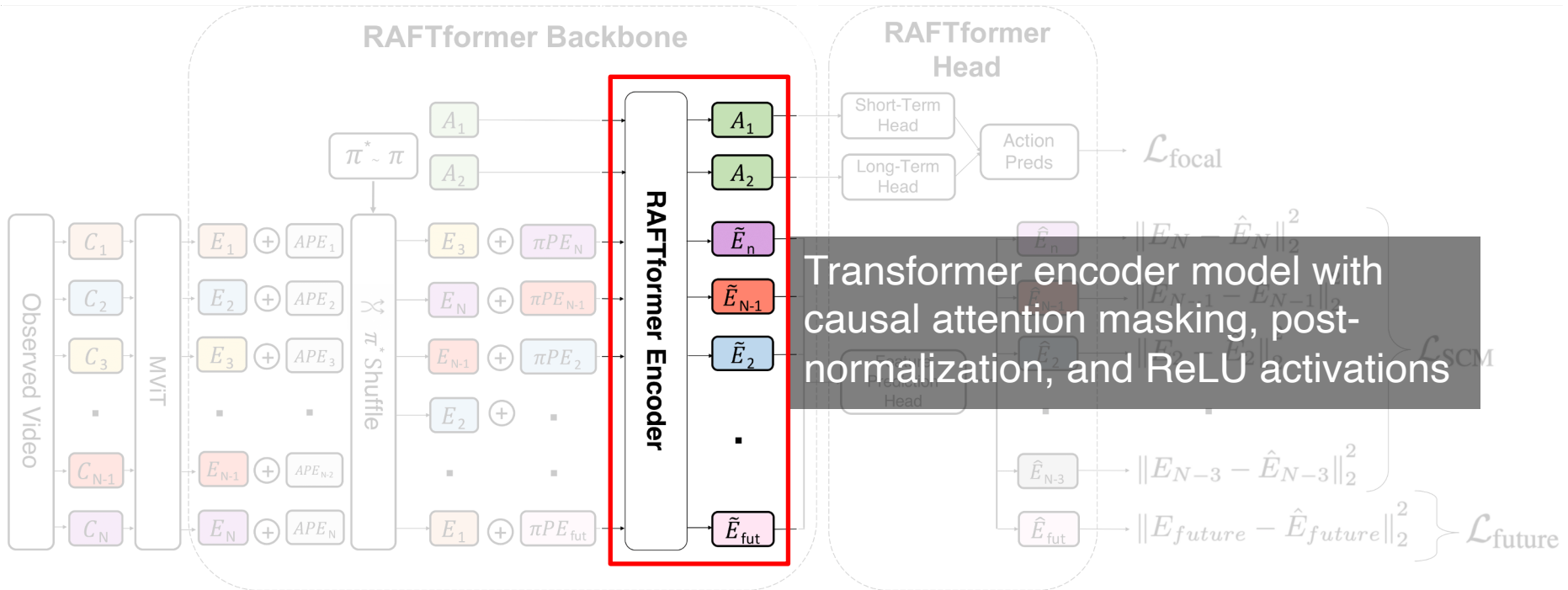
The last token adds  $\pi PE_{\text{fut}}$ , which is used to help generate the “future token”



# RAFTformer Architecture



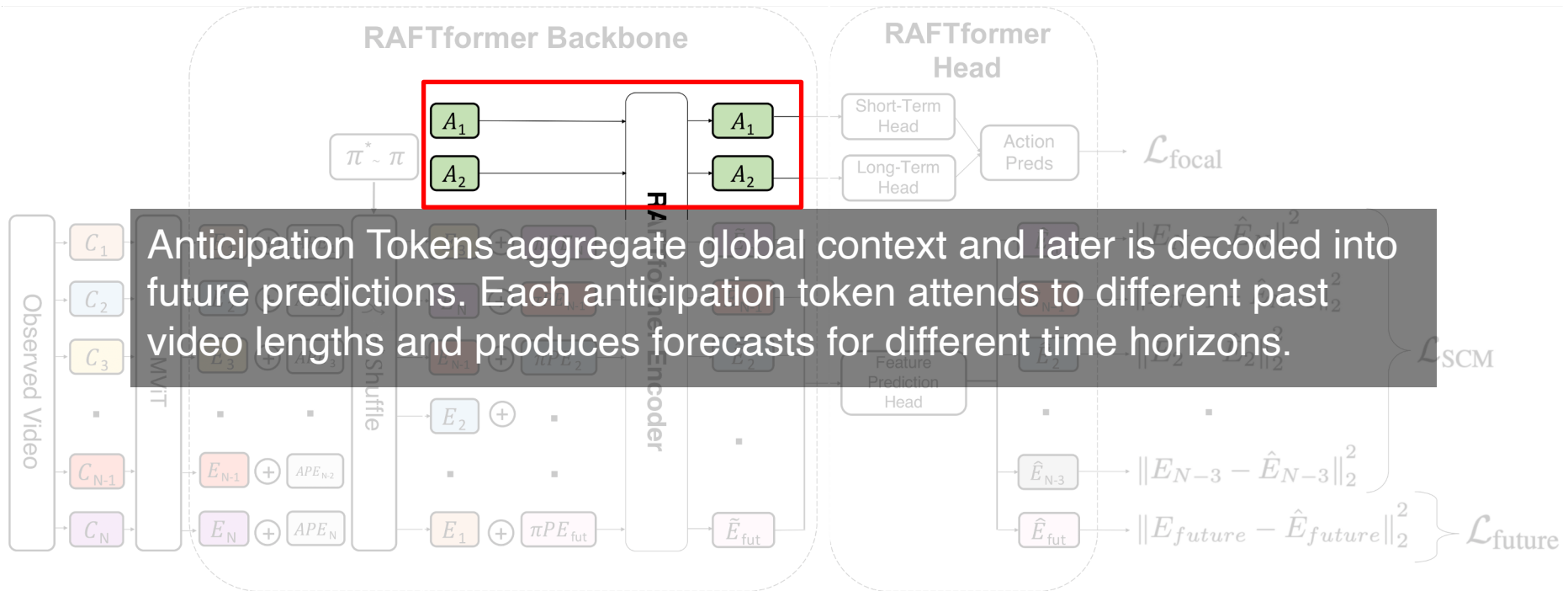
# RAFTformer Architecture



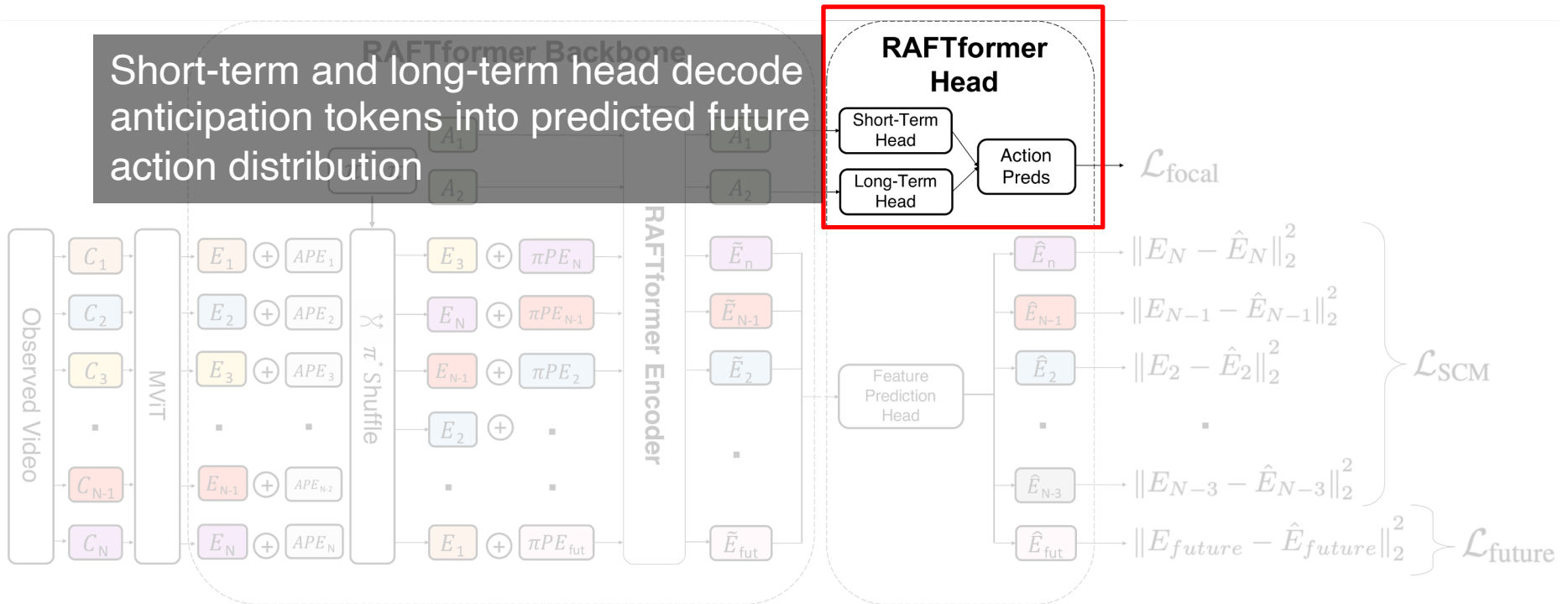
\*RAFTformer encoder has a special form of masked attention that prevents information leakage under shuffling. Please see paper for details.



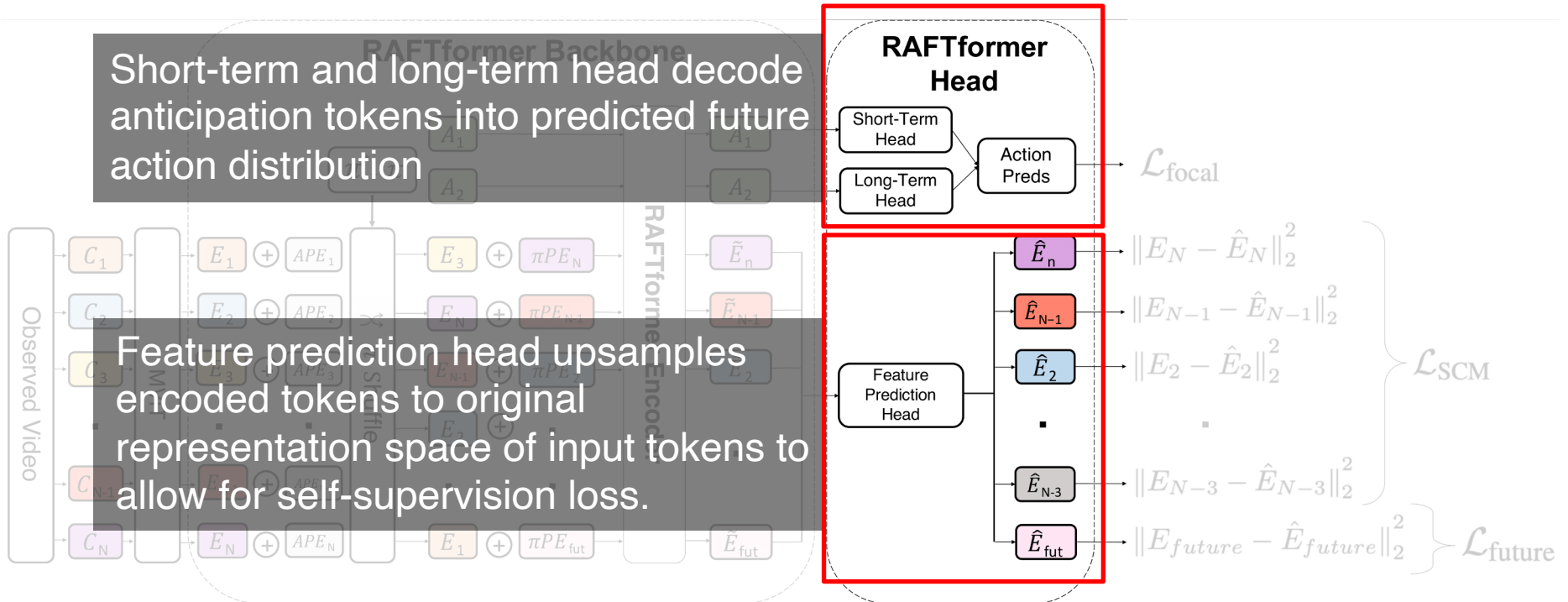
# RAFTformer Architecture



# RAFTformer Architecture



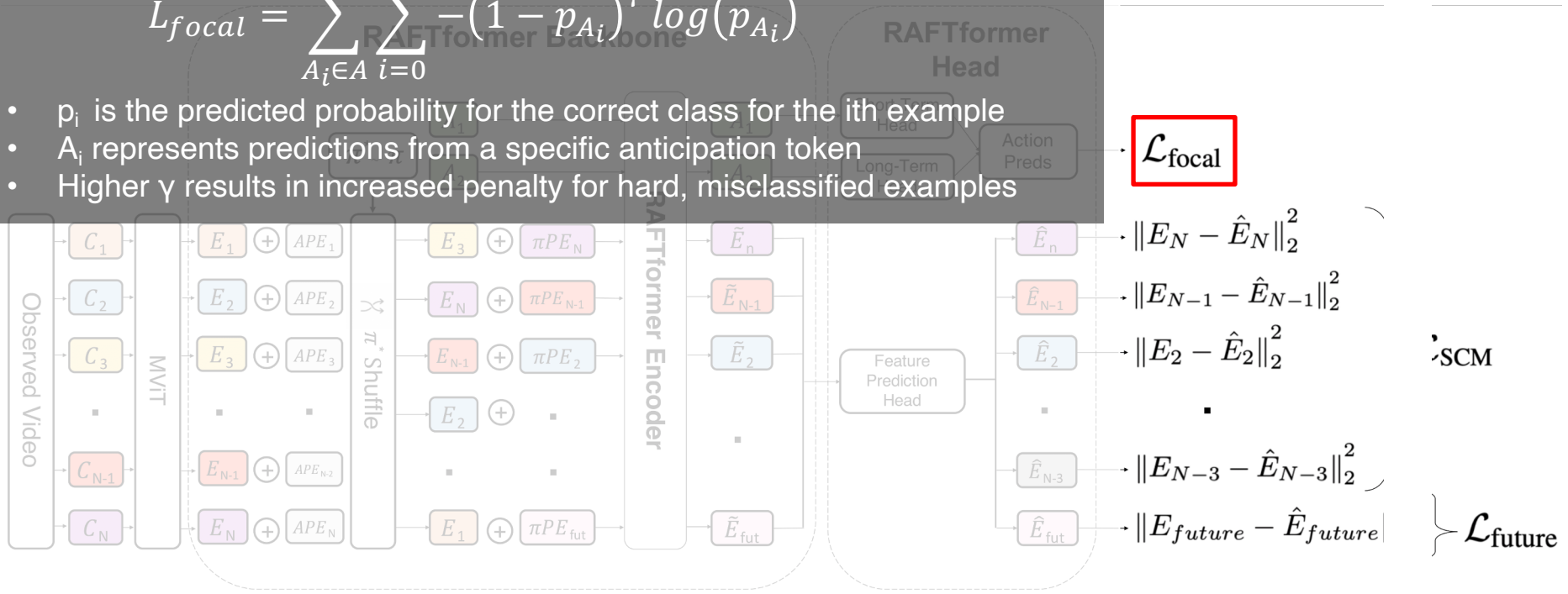
# RAFTformer Architecture



# RAFTformer Architecture

$$L_{focal} = \sum_{A_i \in A} \sum_{i=0}^n -(1 - p_{A_i})^\gamma \log(p_{A_i})$$

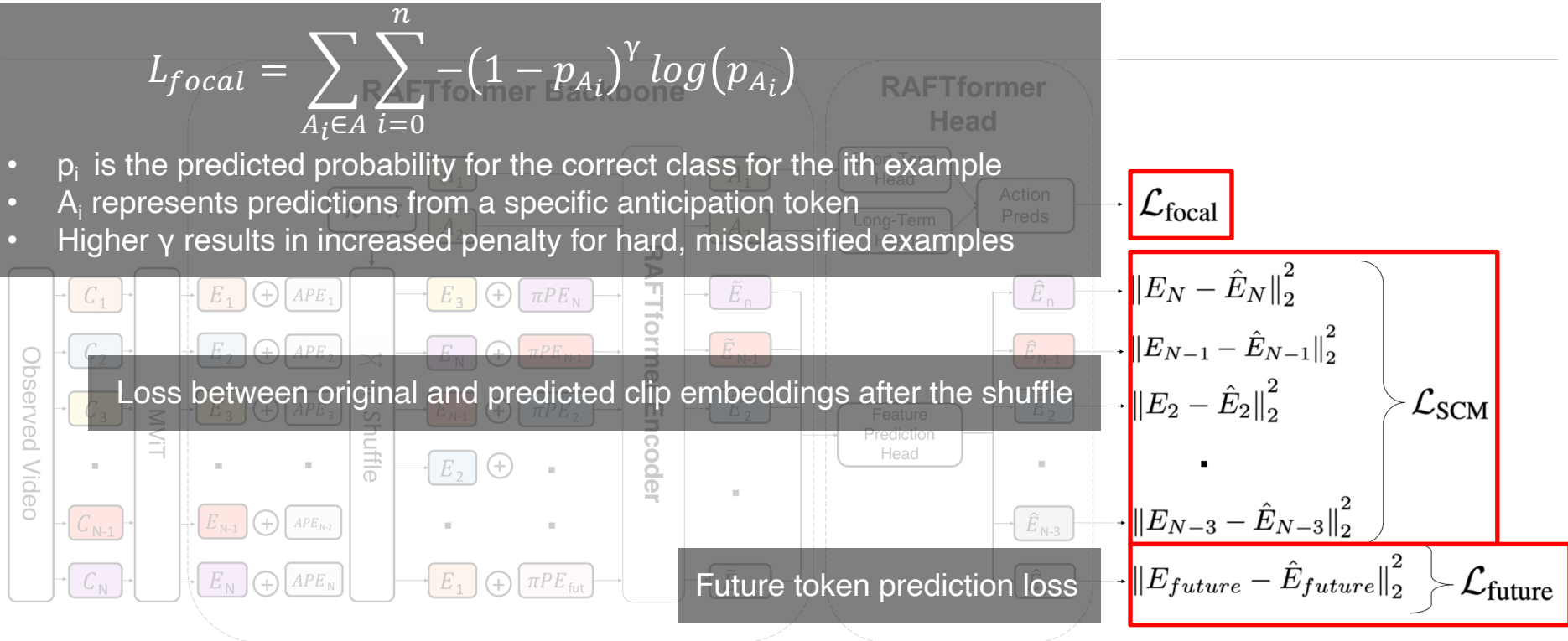
- $p_i$  is the predicted probability for the correct class for the  $i$ th example
- $A_i$  represents predictions from a specific anticipation token
- Higher  $\gamma$  results in increased penalty for hard, misclassified examples



# RAFTformer Architecture

$$L_{focal} = \sum_{A_i \in \mathcal{A}} \sum_{i=0}^n -(1 - p_{A_i})^\gamma \log(p_{A_i})$$

- $p_i$  is the predicted probability for the correct class for the  $i$ th example
- $A_i$  represents predictions from a specific anticipation token
- Higher  $\gamma$  results in increased penalty for hard, misclassified examples

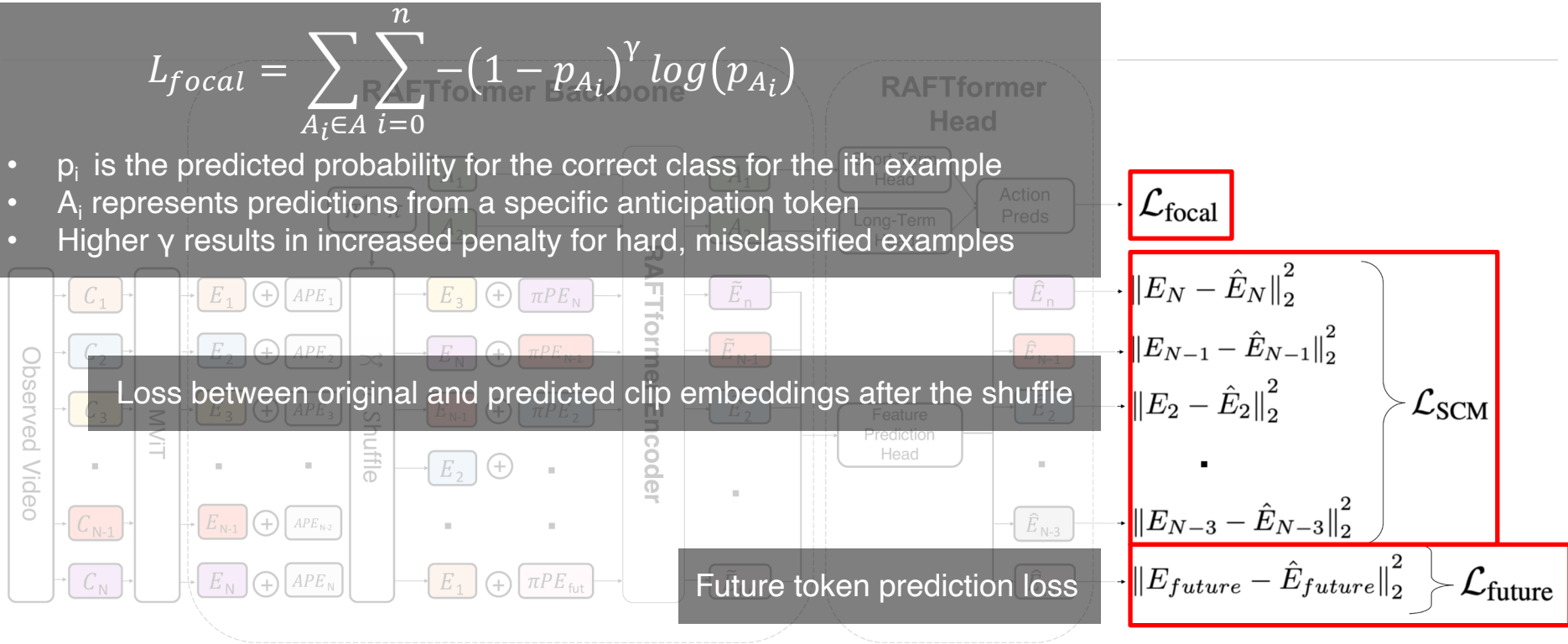




# RAFTformer Architecture

$$L_{focal} = \sum_{A_i \in \mathcal{A}} \sum_{i=0}^n -(1 - p_{A_i})^\gamma \log(p_{A_i})$$

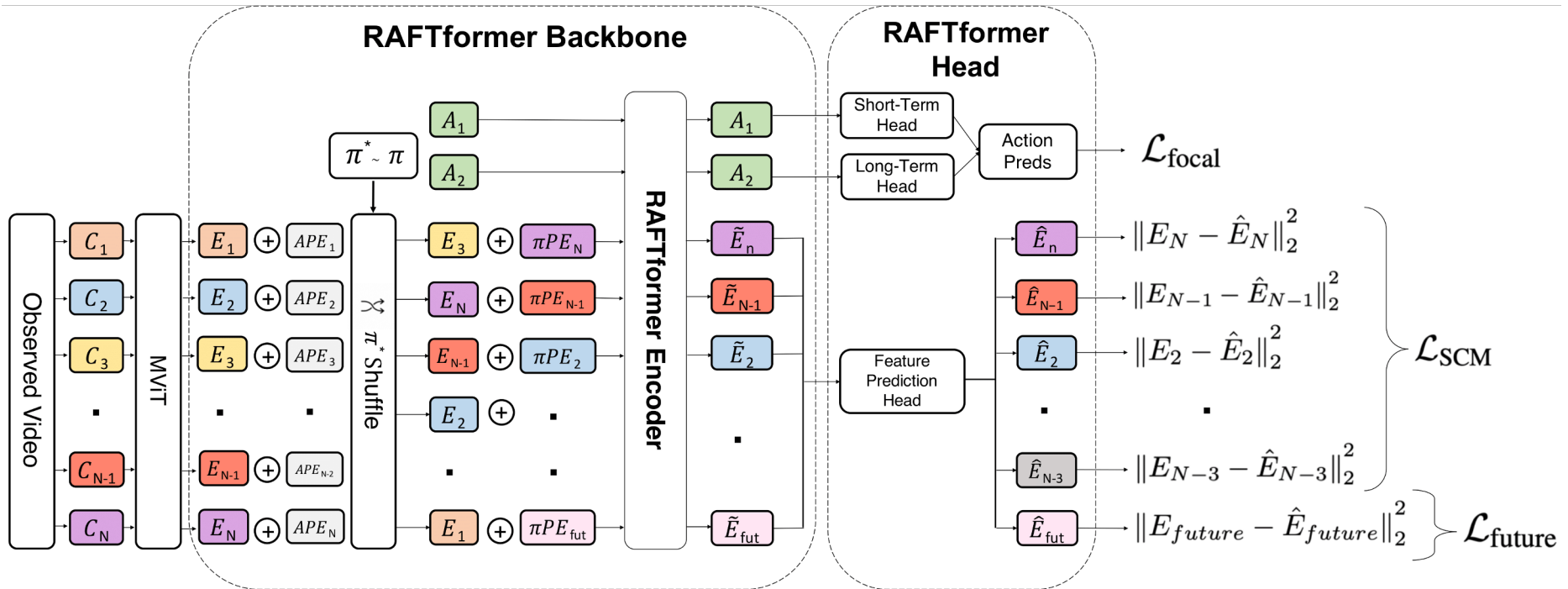
- $p_i$  is the predicted probability for the correct class for the  $i$ th example
- $A_i$  represents predictions from a specific anticipation token
- Higher  $\gamma$  results in increased penalty for hard, misclassified examples



$$L = L_{focal} + \lambda_1 L_{SCM} + \lambda_2 L_{future}$$



# RAFTformer Architecture



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	<b>26</b>	<b>23</b>	<b>40</b>
RAFTformer	None	K700		33.7	37.1	18.0	<b>26</b>	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	<b>26</b>	<b>23</b>	<b>40</b>
RAFTformer	None	K700		33.7	37.1	18.0	<b>26</b>	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	

State-of-the-art  
results



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	26	23	40
RAFTformer	None	K700		33.7	37.1	18.0	26	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	

State-of-the-art  
results      ~8x less  
parameters



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	26	23	40
RAFTformer	None	K700		33.7	37.1	18.0	26	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	

State-of-the-art  
results      ~8x less  
parameters      ~94% less  
GPU hours



# Results: Offline Setting

All past frames up to the present time  $T$   
are used to predict the action at time  $T + t_f$

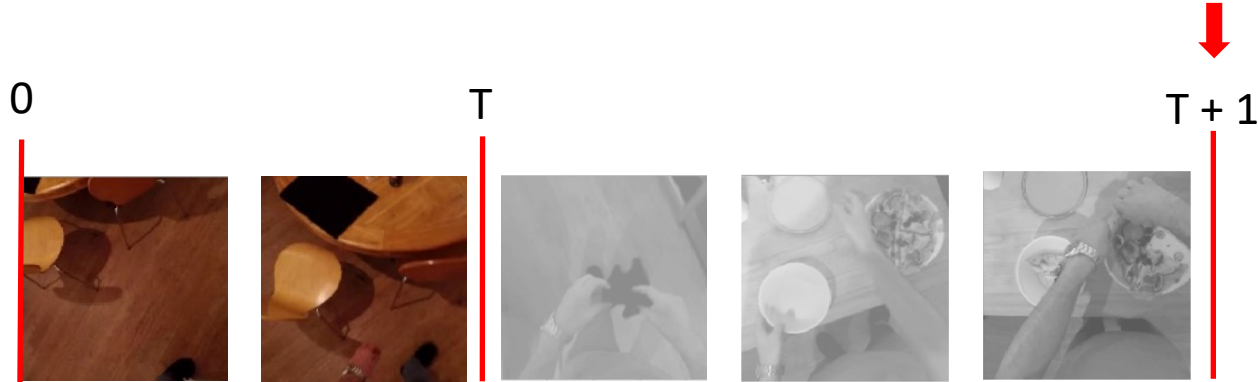
Split	Method	Addl. Modality	Init	Epic Boxes	Top-5 Recall			Parameters ( $\times 10^6$ )	GPU Hours	Inference Latency (ms)
					Verb	Noun	Action			
Val	TempAgg [70]	None	IN1K		24.2	29.8	13.0	-	-	-
	RULSTM [19]	None	IN1K		-	-	13.3	-	-	-
	RULSTM [19]	Obj+Flow	IN1K	✓	30.8	27.8	14.0	-	-	-
	TempAgg [70]	Obj+Flow+ROI	IN1K	✓	23.2	31.4	14.7	-	-	-
	AVT [25]	None	IN21K		30.2	31.7	14.9	378	-	420
	AVT+ [25]	Obj	IN21K	✓	28.2	32.0	15.9	-	-	-
	TSN-AVT+ [25]	Obj	IN21K	✓	31.8	25.5	14.8	-	-	-
	MeMVit [80]	None	K400		32.8	33.2	15.1	59	-	160
	MeMVit [80]	None	K700		32.2	37.0	17.7	212	368	350
	RAFTformer	None	K400 + IN1K		33.3	35.5	17.6	26	23	40
RAFTformer	None	K700		33.7	37.1	18.0	26	27	110	
RAFTformer-2B	None	K700 + IN1K		<b>33.8</b>	<b>37.9</b>	<b>19.1</b>	52	50	160	

~9x less  
latency

State-of-the-art  
results    ~8x less  
parameters    ~94% less  
GPU hours



# Results: Online Setting



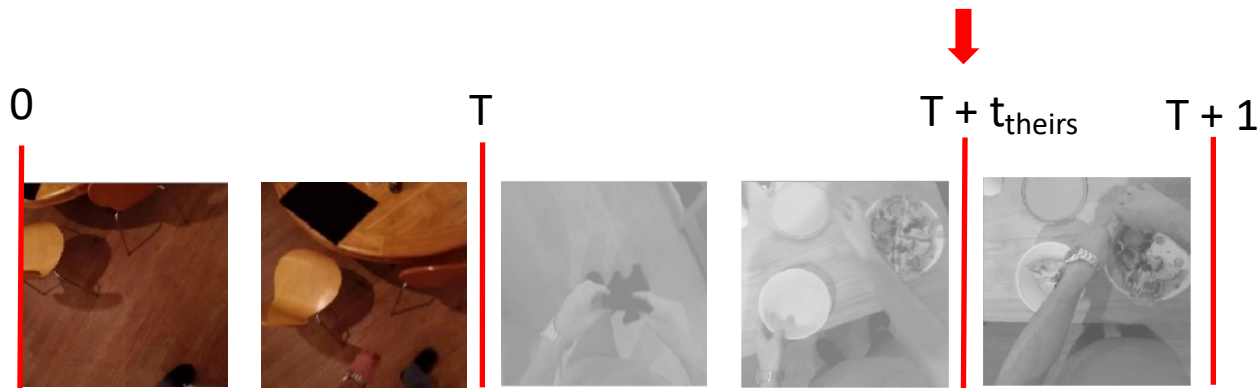
Base models are given a video up to time  $T$ , which they used to predict the action at time  $T+1$



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action



# Results: Online Setting

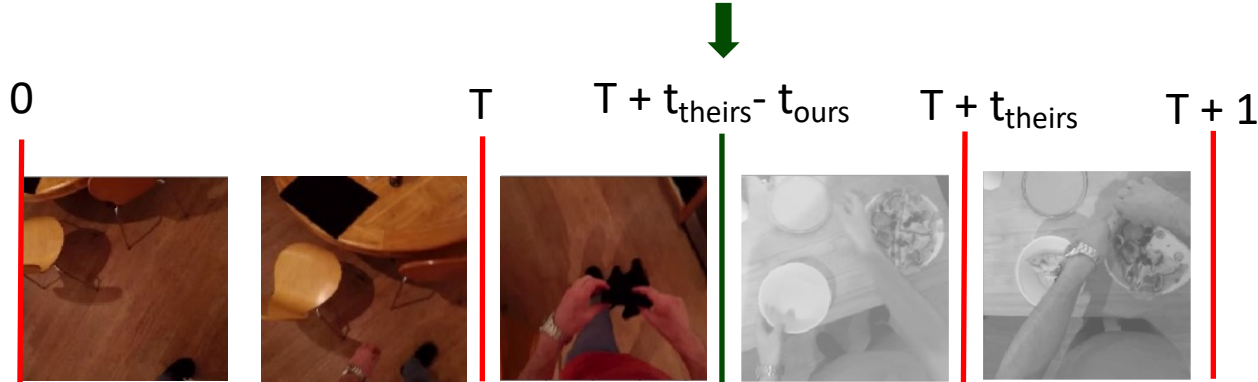


With a latency of  $t_{\text{theirs}}$ , the prediction will arrive at time  $T + t_{\text{theirs}}$



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action

# Results: Online Setting



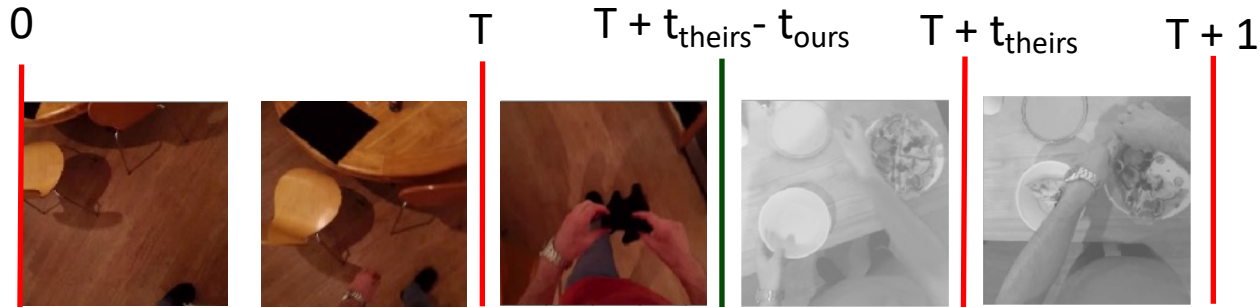
When comparing two models, the prediction arrival times should be the same. With a latency of  $t_{ours}$ , RAFTFormer must start prediction at time  $T + t_{theirs} - t_{ours}$  so that the prediction arrives at time  $T + t_{theirs}$



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action

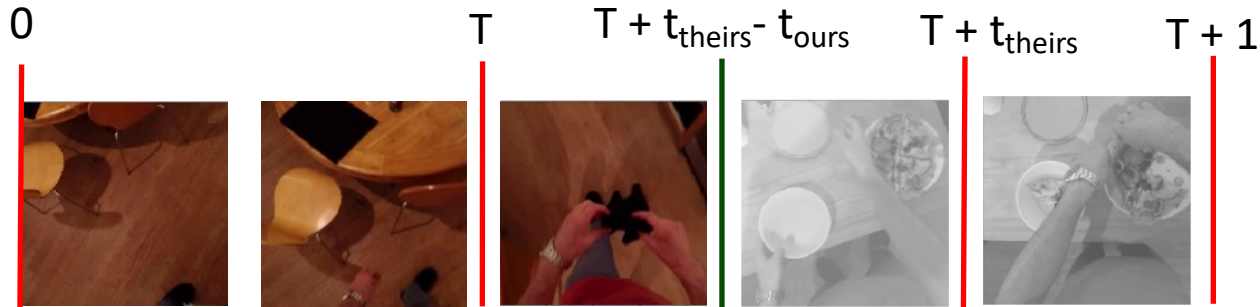


# Results: Online Setting



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action
AVT [25]	IN21K	$t_{avt} = 420$	$T$	$T + t_{avt}$	$T + 1$	30.2	31.7	14.9
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{avt} - t_{ours}$	$T + t_{avt}$	$T + 1$	34.1	38.2	<b>19.3 (+4.4)</b>
MemViT [80]	K400	$t_{vit} = 160$	$T$	$T + t_{vit}$	$T + 1$	32.8	33.2	15.1
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{vit} - t_{ours}$	$T + t_{vit}$	$T + 1$	33.8	37.1	<b>18.1 (+3.0)</b>

# Results: Online Setting



Model	Init	Latency ( $t_l$ ms)	Inference Start Time Stamp	Inference End Time Stamp	Target Time Stamp	Top-5 Recall		
						Verb	Noun	Action
AVT [25]	IN21K	$t_{avt} = 420$	$T$	$T + t_{avt}$	$T + 1$	30.2	31.7	14.9
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{avt} - t_{ours}$	$T + t_{avt}$	$T + 1$	34.1	38.2	<b>19.3 (+4.4)</b>
MemViT [80]	K400	$t_{vit} = 160$	$T$	$T + t_{vit}$	$T + 1$	32.8	33.2	15.1
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{vit} - t_{ours}$	$T + t_{vit}$	$T + 1$	33.8	37.1	<b>18.1 (+3.0)</b>
MemViT [80]	K700	$t_{vit} = 350$	$T$	$T + t_{vit}$	$T + 1$	32.2	37.0	17.7
RAFTformer	K400 + IN1k	$t_{ours} = 40$	$T + t_{vit} - t_{ours}$	$T + t_{vit}$	$T + 1$	33.7	37.9	<b>19.0 (+1.3)</b>

# Latency Matters: Real-Time Action Forecasting Transformer

# Thank You!

For full results on all datasets, [paper](#), [code](#) and further details  
please visit the project homepage



 SCAN ME

<https://karttikeya.github.io/publication/RAFTformer/>