

# DeepQAMVS: Query-Aware Hierarchical Pointer Networks for Multi-Video Summarization



Safa Messaoud



Ismini Lourentzou



Assma Boughoula\*



Mona Zehni\*



Zhizhen Zhao



Chengxiang Zhai



Alex Schwing

sigir21



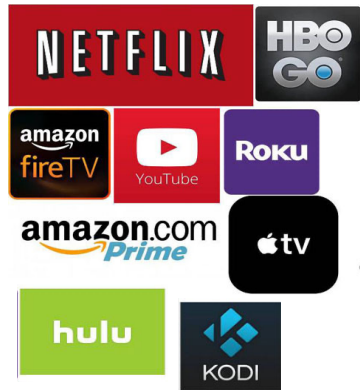
# Growing Popularity of Video and Explosion of Video Sharing Platforms



**82% of internet traffic** will be from video streaming and downloads in 2022  
(Cisco, 2019)



Video watch time:  
ca. **10 hours per week**  
per internet user in 2020  
(Statistica, 2020)



Avg. U.S. consumer pays  
**4 different streaming video subscriptions**  
(Deloitte, 2021)



**87% of the marketing professionals** use video as a marketing tool  
(Wyzowl, 2019)

<https://cedcommerce.com/video-marketing>

**How can we make this exponentially growing video content easier to consume?**

# Query-Aware Multi-Video Summarization (QAMVS)

**Query:** Prince William Wedding



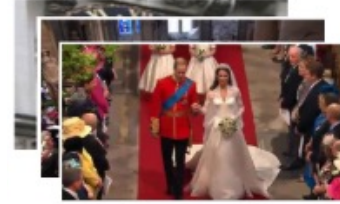
## Retrieved Videos



**Video 1**

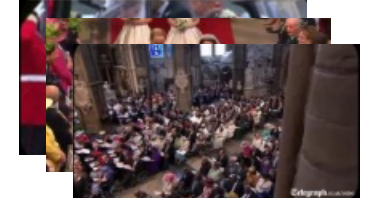


**Video 2**



**Video 3**

...



**Video N**



**Summary:** a subset of the input videos frames

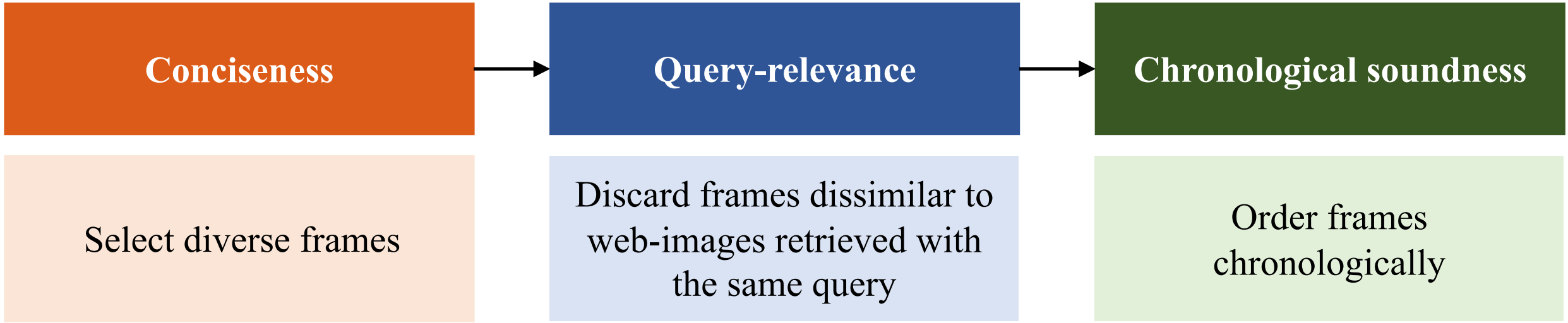


## Summary Criteria

1. Conciseness
2. Representativeness of query-relevant events
3. Chronological soundness

# QAMVS Models in the Literature

Multi-staged pipelines optimizing for summary criteria **sequentially**



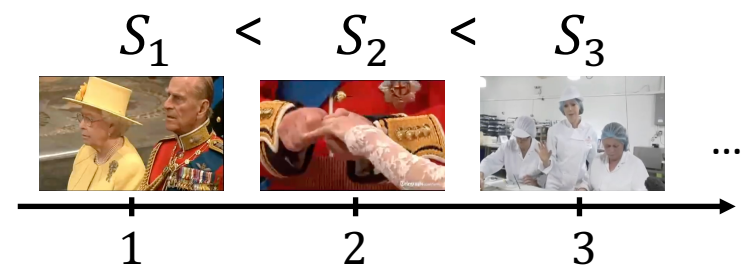
Sparse Coding (*e.g.* QUASC)  
Clustering (*e.g.* HDS)



**Textual Query** ↔ **Image space (Web-images)**

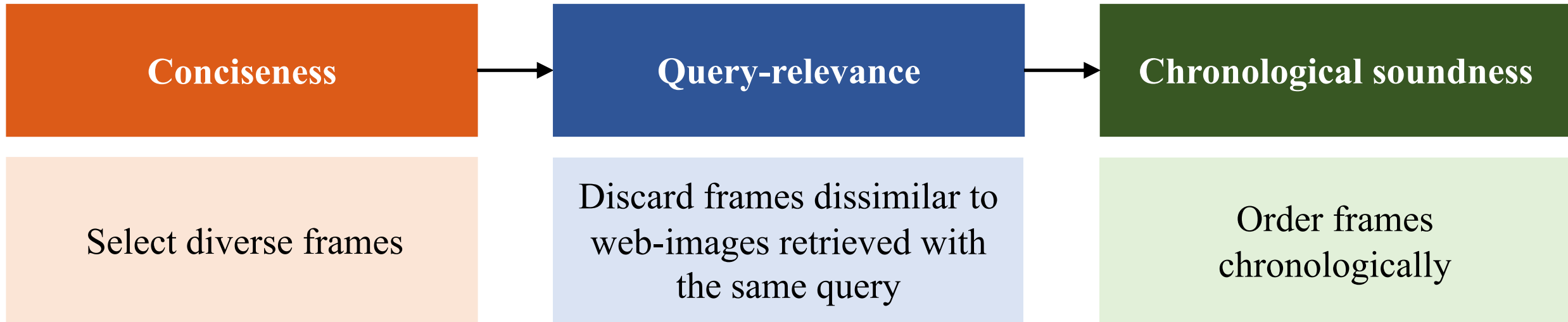


Scores based on topic closeness/videos time tags



# QAMVS Models in the Literature

Multi-staged pipelines optimizing for summary criteria **sequentially**



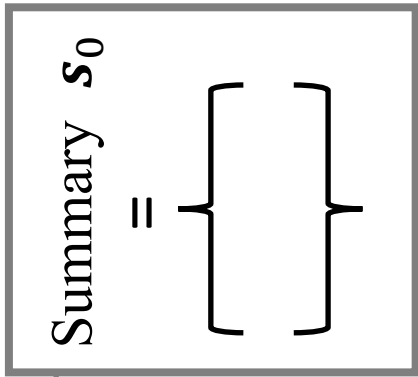
## Limitations:

- ✗ Error propagation through sequential stages
- ✗ Polynomial complexity w.r.t. number of input frames

Can we **learn** efficient end-to-end trainable models for QAMVS?

- ✗ Scarcity of annotated data
- ✗ Subjectivity of ground-truth summaries

# Reinforcement Learning for DeepQAMVS



Agent (Policy Network)



$F_1$   $F_2$   $F_3$

Video 1



$F_1$   $F_2$   $F_3$

Video 2



$F_1$   $F_2$   $F_3$

Video 3

Videos ( $\mathcal{V}$ )

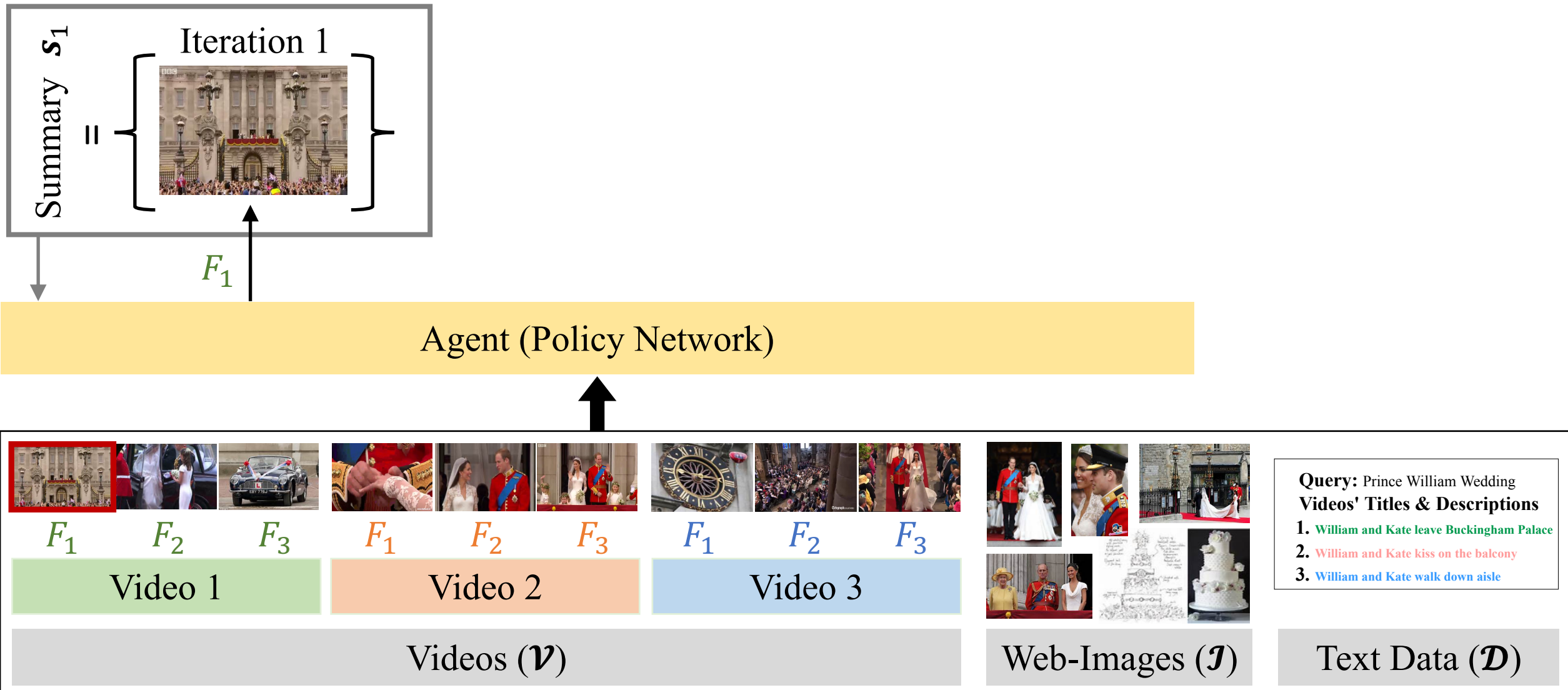


Web-Images ( $\mathcal{I}$ )

- Query:** Prince William Wedding  
**Videos' Titles & Descriptions**
1. William and Kate leave Buckingham Palace
  2. William and Kate kiss on the balcony
  3. William and Kate walk down aisle

Text Data ( $\mathcal{D}$ )

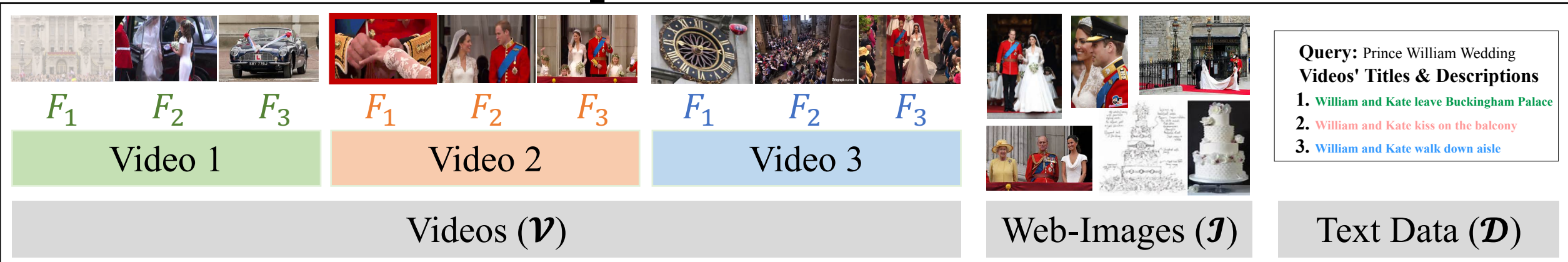
# Reinforcement Learning for DeepQAMVS



# Reinforcement Learning for DeepQAMVS

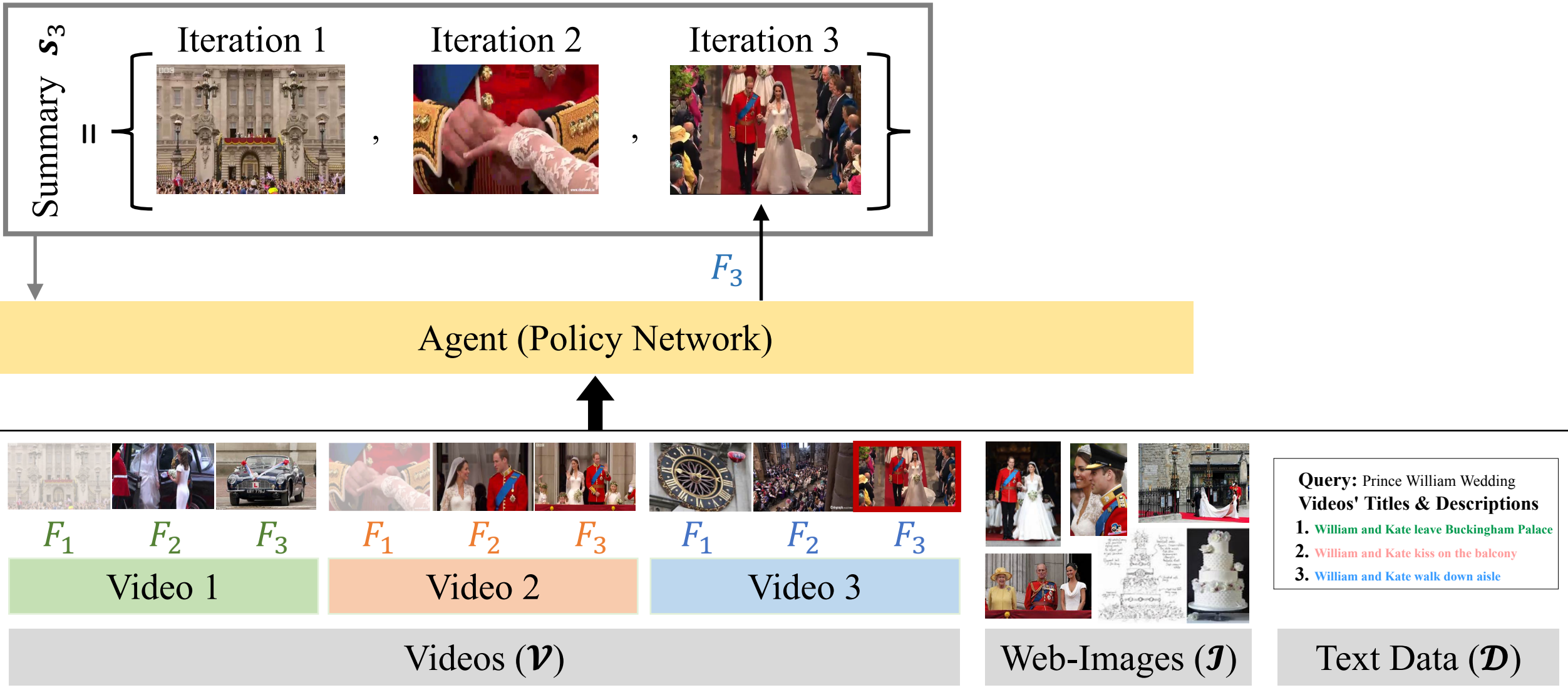


Agent (Policy Network)

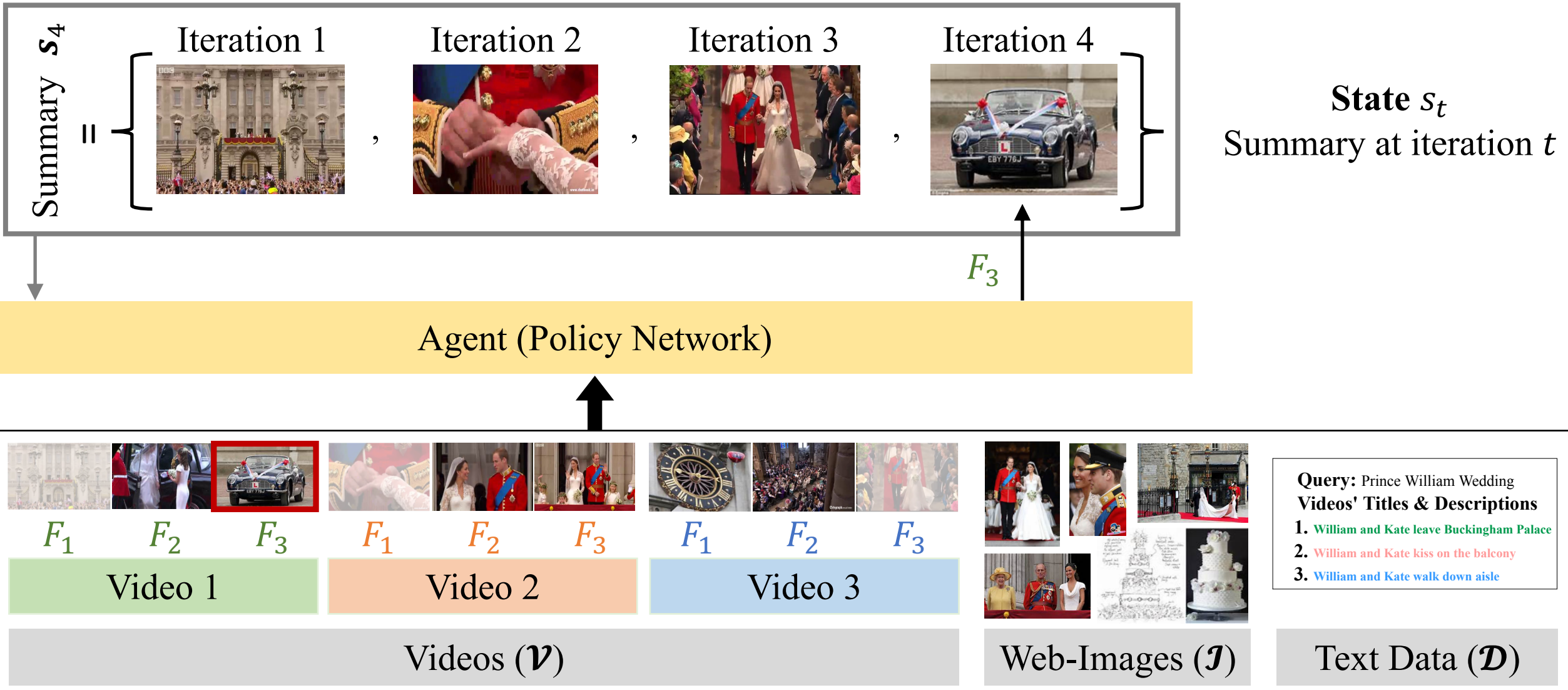




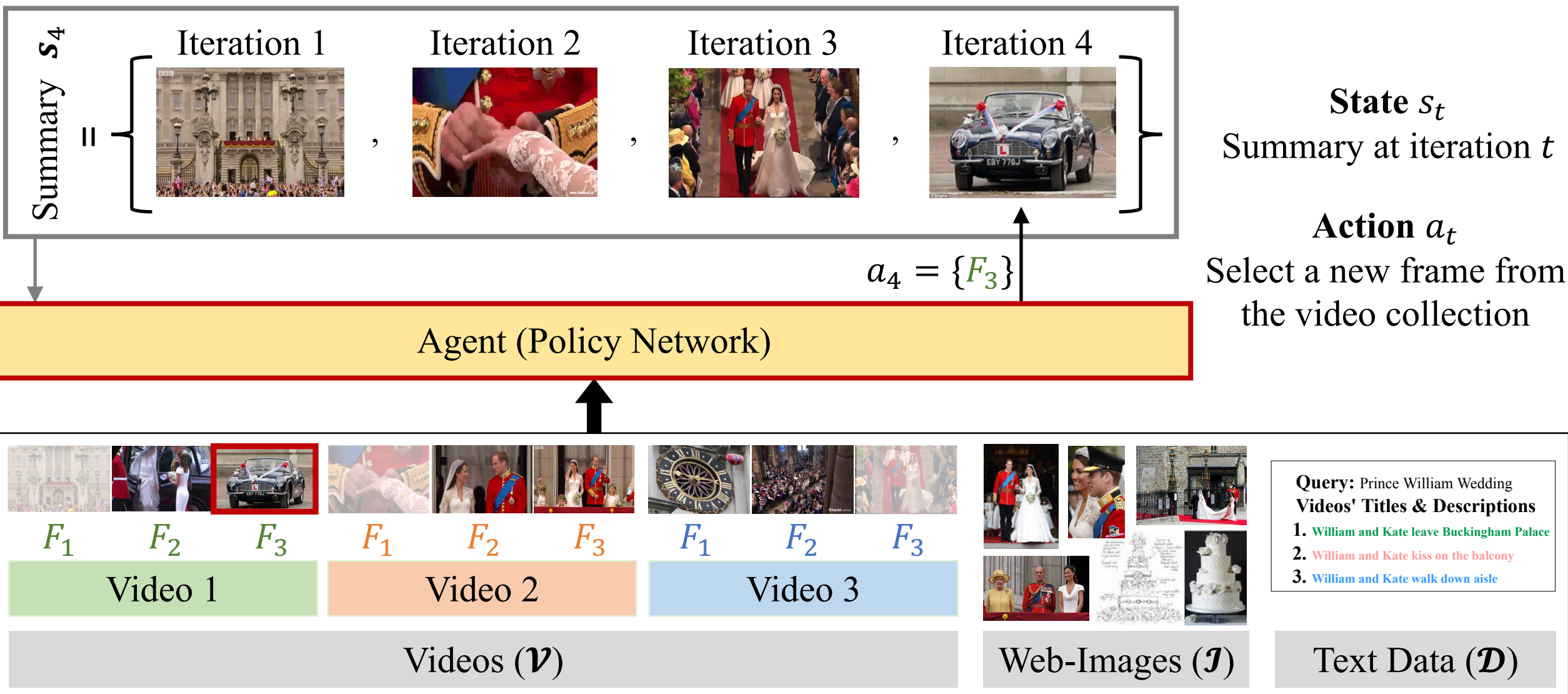
# Reinforcement Learning for DeepQAMVS



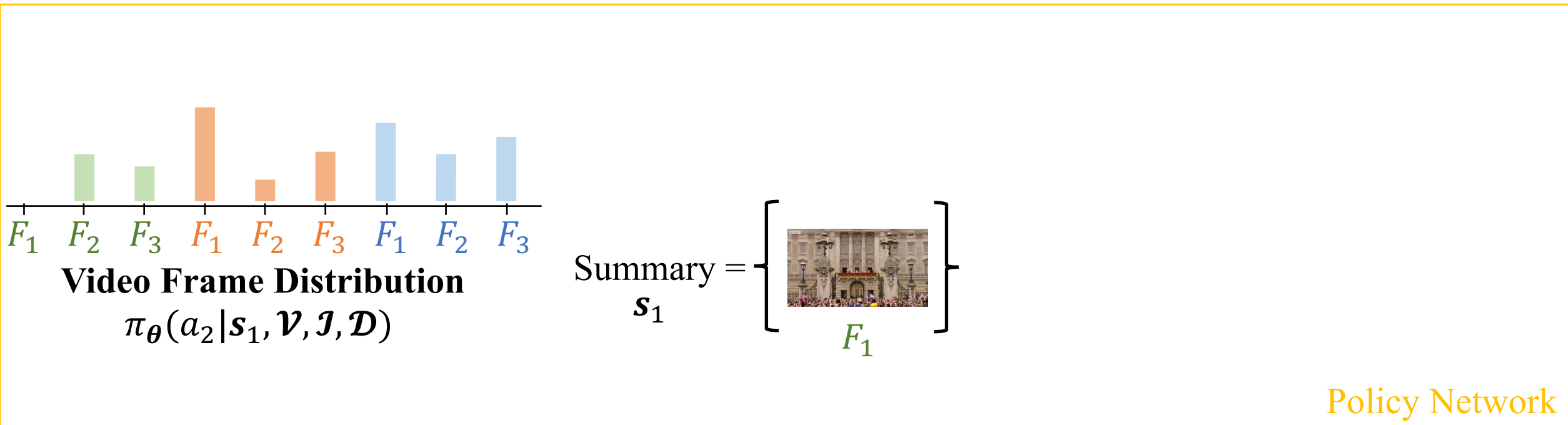
# Reinforcement Learning for DeepQAMVS



# Reinforcement Learning for DeepQAMVS



# Policy Network: Pointer Network



Policy Network



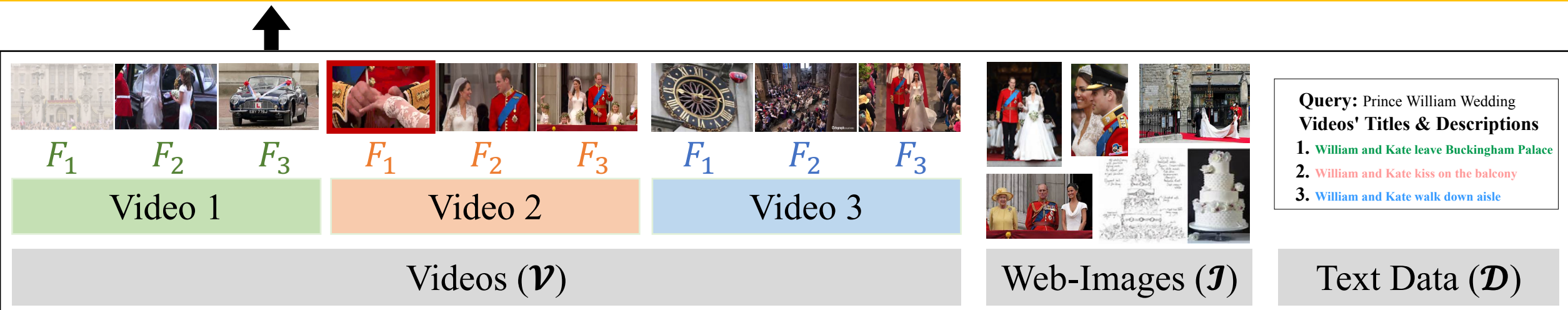
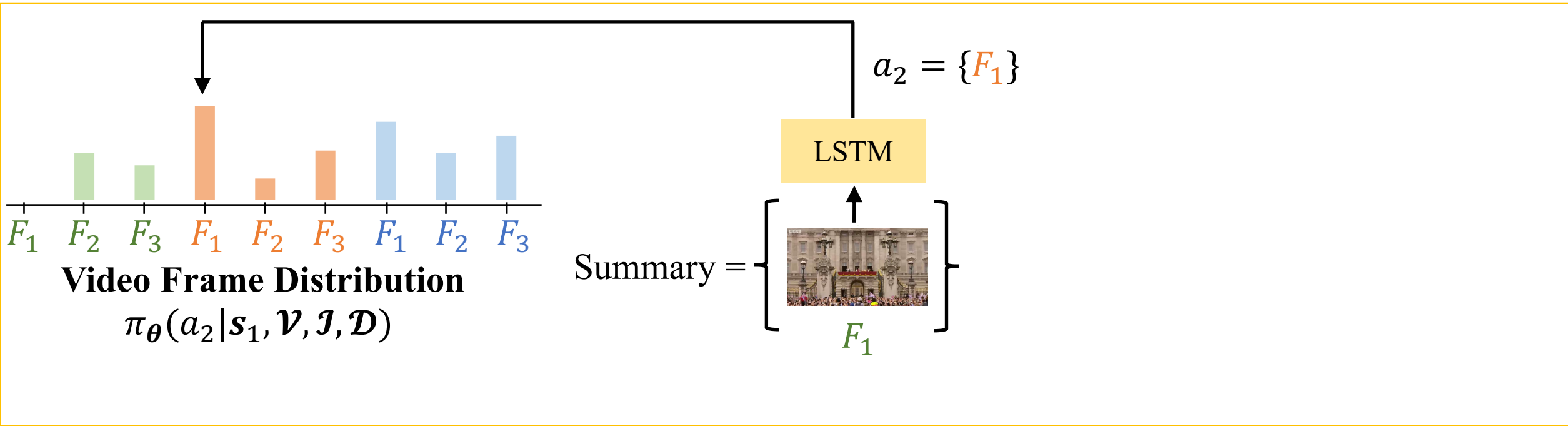
									<p><b>Query:</b> Prince William Wedding  <b>Videos' Titles &amp; Descriptions</b></p> <ol style="list-style-type: none"> <li>William and Kate leave Buckingham Palace</li> <li>William and Kate kiss on the balcony</li> <li>William and Kate walk down aisle</li> </ol>	
Video 1			Video 2			Video 3				<p>Web-Images (<math>\mathcal{J}</math>)</p>

Videos ( $\mathcal{V}$ )

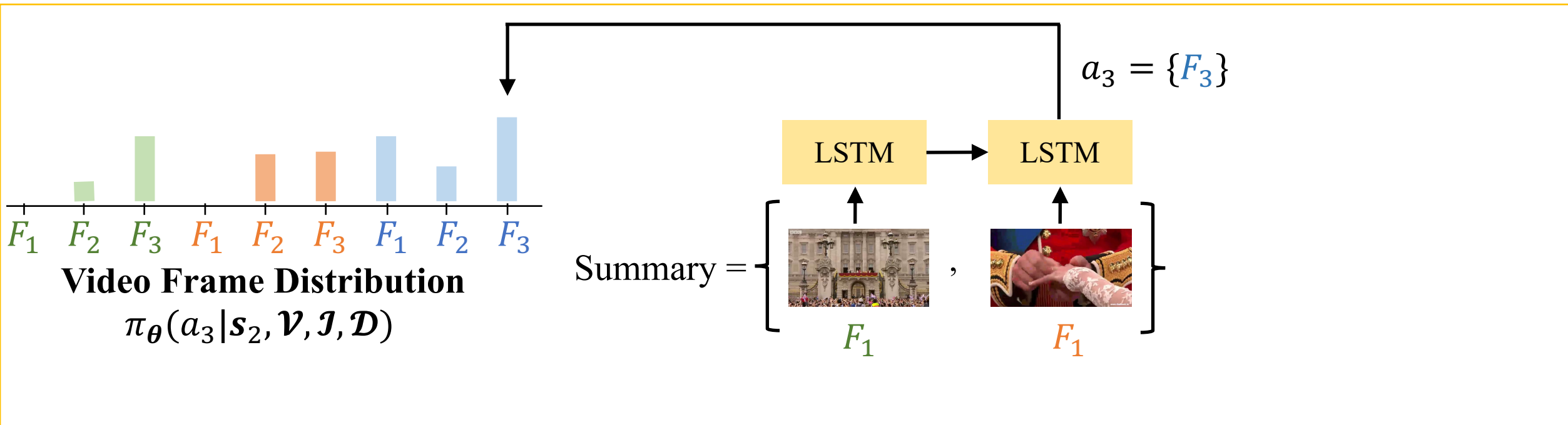
Web-Images ( $\mathcal{J}$ )

Text Data ( $\mathcal{D}$ )

# Policy Network: Pointer Network



# Policy Network: Pointer Network

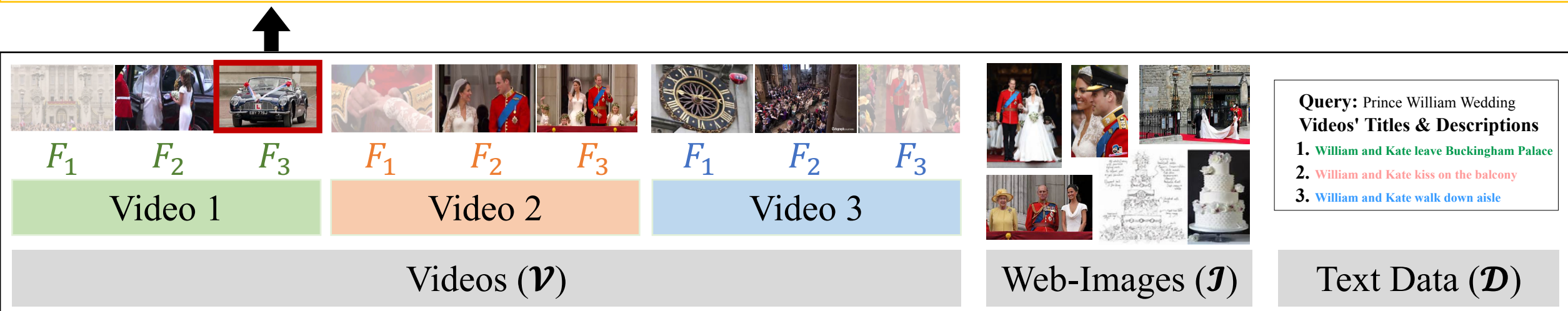
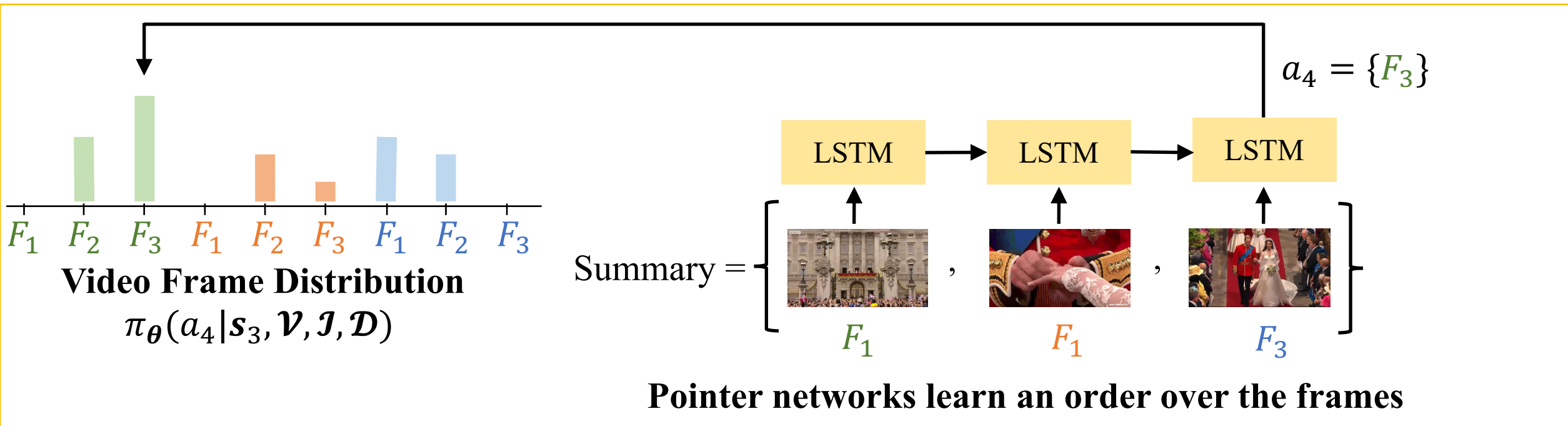


The input data is organized into three main categories:

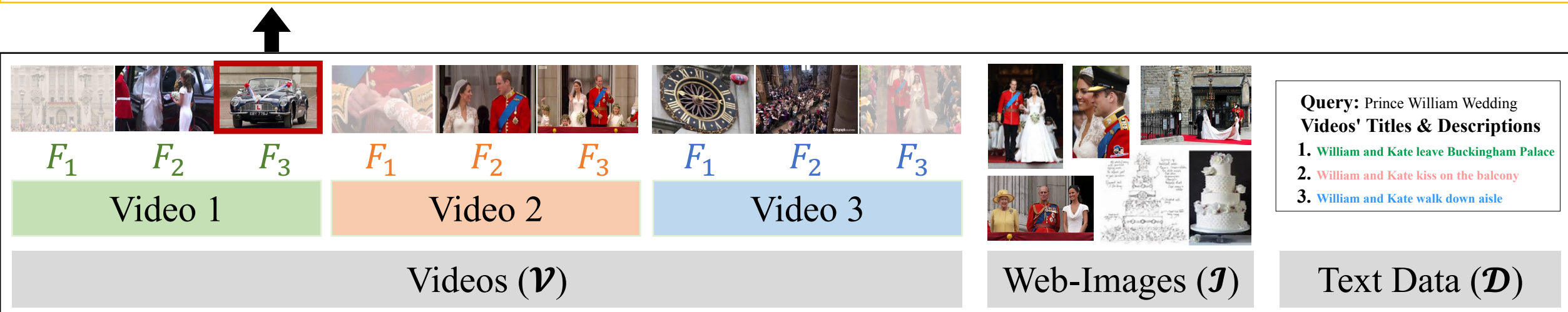
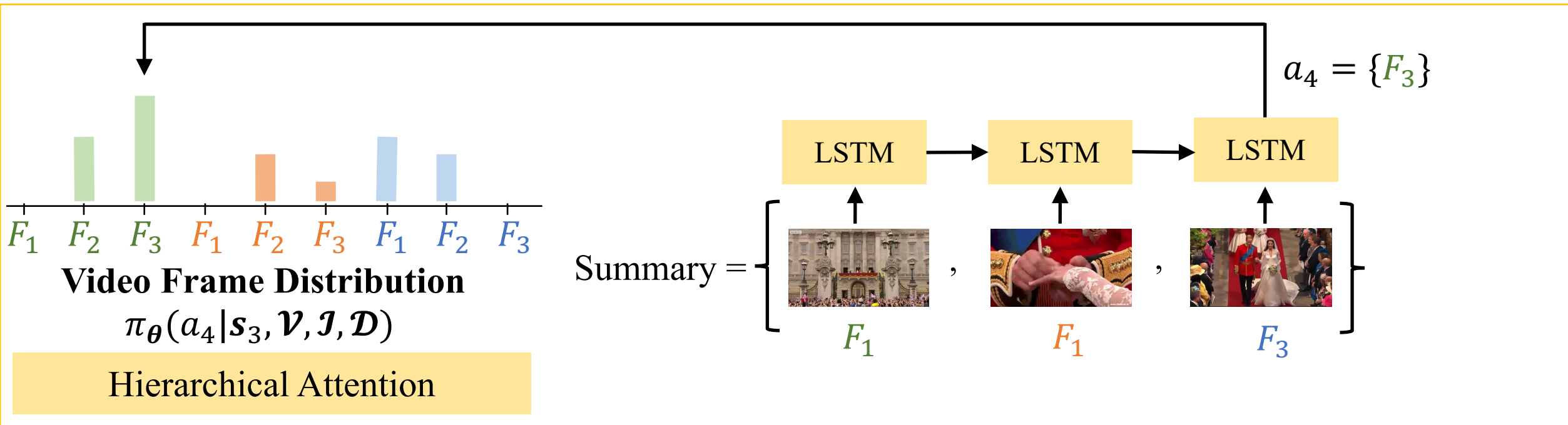
- Videos ( $\mathcal{V}$ ):** Three videos are shown, each with three frames:
  - Video 1 (Green):** Frames  $F_1, F_2, F_3$  showing a wedding scene.
  - Video 2 (Orange):** Frames  $F_1, F_2, F_3$  showing a couple kissing.
  - Video 3 (Blue):** Frames  $F_1, F_2, F_3$  showing a crowd and a couple.
- Web-Images ( $\mathcal{J}$ ):** A collection of images related to the wedding, including the couple, the cake, and the ceremony.
- Text Data ( $\mathcal{D}$ ):** A list of video titles and descriptions:
  - William and Kate leave Buckingham Palace
  - William and Kate kiss on the balcony
  - William and Kate walk down aisle

An arrow points from the input data to the Policy Network diagram above.

# Policy Network: Pointer Network



# Policy Network: Pointer Network

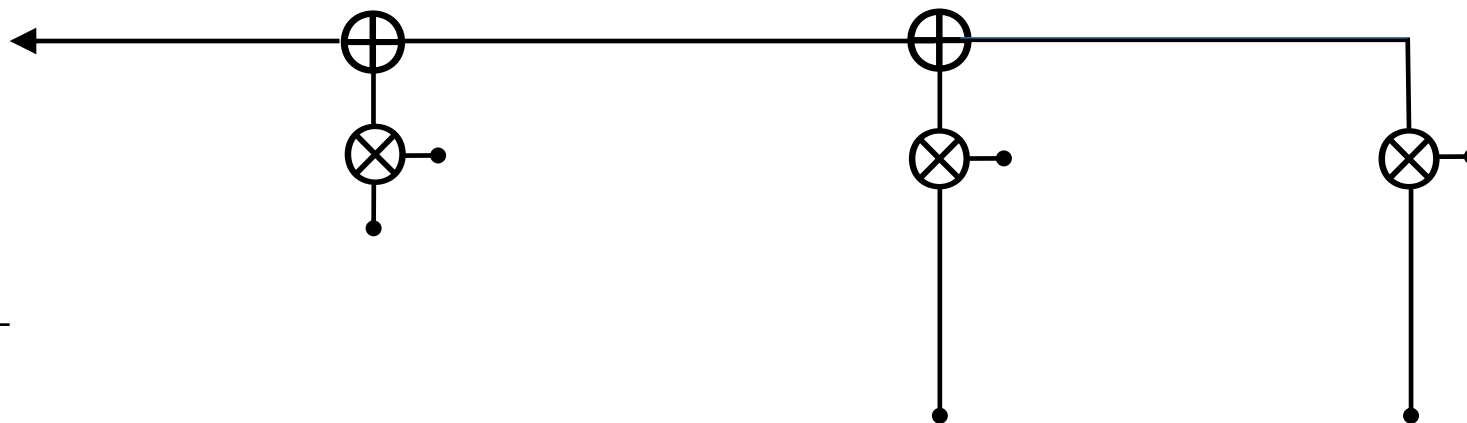
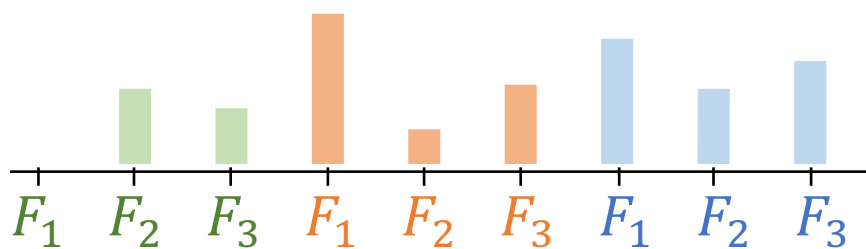




# Policy Network: Hierarchical Attention

**Video Frame Distribution**

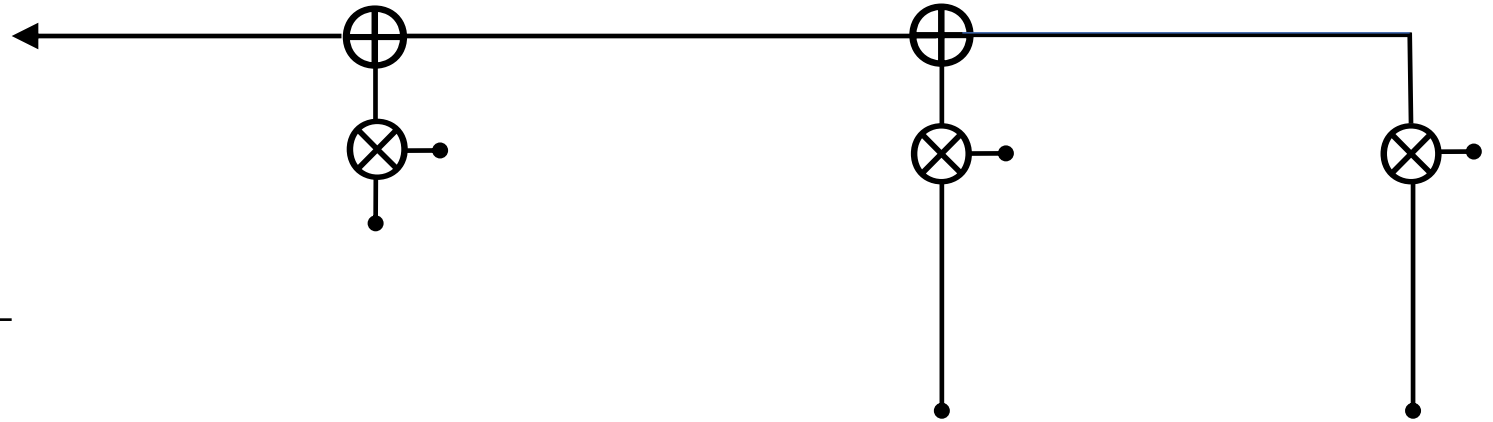
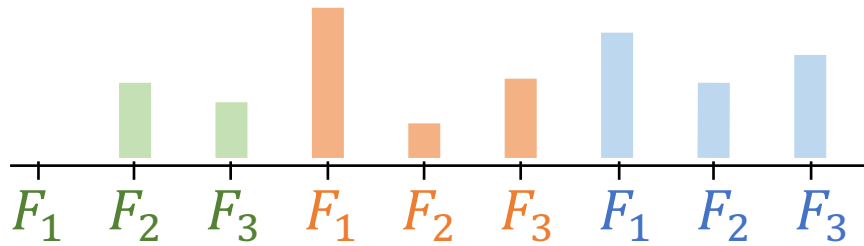
$$\pi_{\theta}(a_t | \mathbf{s}_{t-1}, \mathcal{V}, \mathcal{J}, \mathcal{D})$$



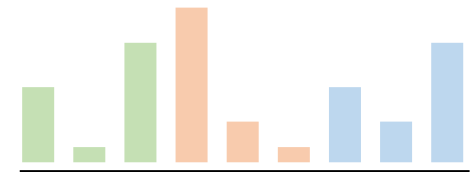
# Policy Network: Hierarchical Attention

## Video Frame Distribution

$$\pi_{\theta}(a_t | \mathbf{s}_{t-1}, \mathcal{V}, \mathcal{J}, \mathcal{D})$$



$$\pi_{\theta}(a_t | \mathbf{s}_{t-1}, \mathcal{D})$$



Query Attention

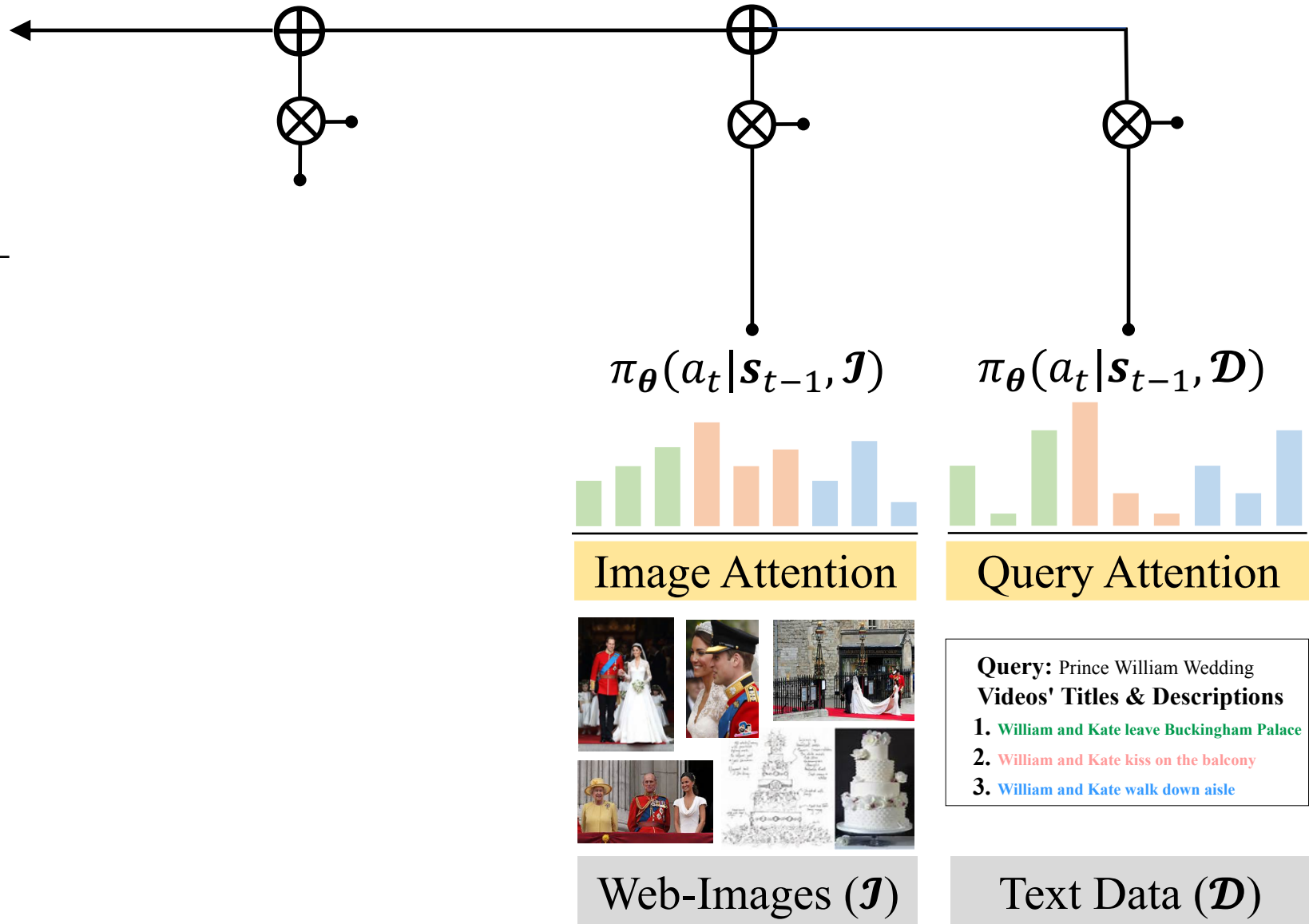
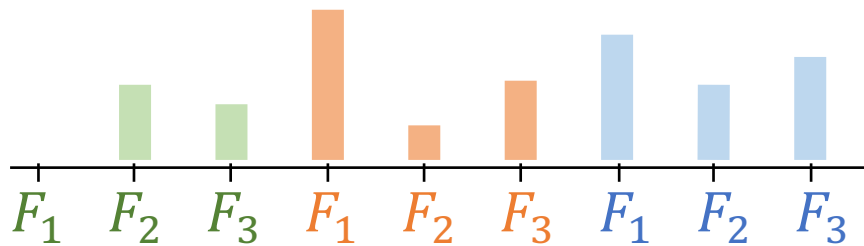
- Query:** Prince William Wedding  
**Videos' Titles & Descriptions**
1. William and Kate leave Buckingham Palace
  2. William and Kate kiss on the balcony
  3. William and Kate walk down aisle

Text Data ( $\mathcal{D}$ )

# Policy Network: Hierarchical Attention

## Video Frame Distribution

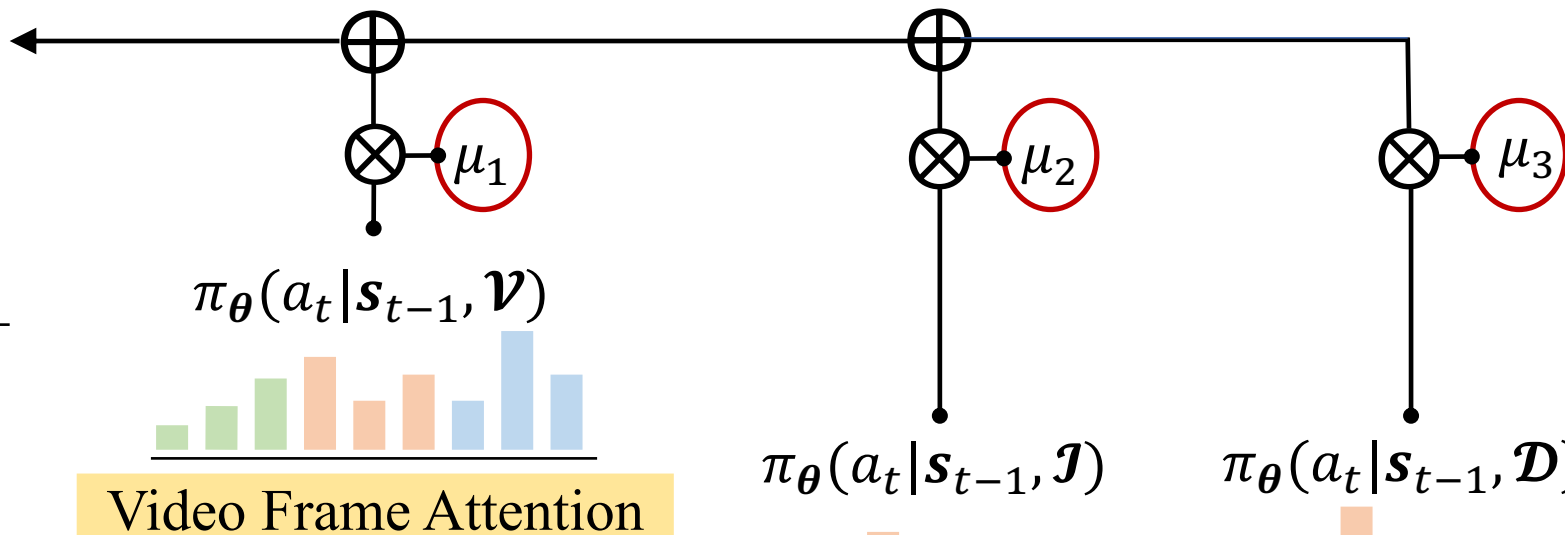
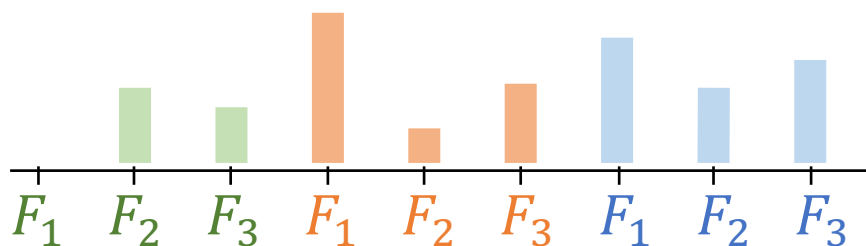
$$\pi_{\theta}(a_t | \mathbf{s}_{t-1}, \mathcal{V}, \mathcal{J}, \mathcal{D})$$



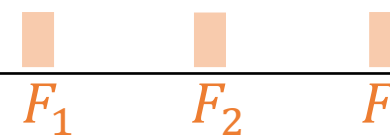
# Policy Network: Hierarchical Attention

## Video Frame Distribution

$$\pi_{\theta}(a_t | \mathbf{s}_{t-1}, \mathcal{V}, \mathcal{J}, \mathcal{D})$$



## Video Attention



## Frame Attention



Video 1

## Frame Attention



Video 2

## Frame Attention



Video 3

## Image Attention



Web-Images ( $\mathcal{J}$ )

## Query Attention

- Query:** Prince William Wedding  
Videos' Titles & Descriptions
1. William and Kate leave Buckingham Palace
  2. William and Kate kiss on the balcony
  3. William and Kate walk down aisle

Text Data ( $\mathcal{D}$ )

# Rewards

1. **Diversity Reward ( $R_{\text{div}}$ ):** Selected frames are diverse

$$R_{\text{div}} \left( \left\{ \left[ \text{img1}, \text{img2}, \text{img3} \right] \right\} \right) > R_{\text{div}} \left( \left\{ \left[ \text{img1}, \text{img1}, \text{img1} \right] \right\} \right)$$

The diagram illustrates the Diversity Reward ( $R_{\text{div}}$ ) by comparing two sets of selected frames. The left set, enclosed in a large right-facing curly bracket, contains three distinct images: a wide shot of a grand building facade, a close-up of hands in ornate, colorful attire, and a bride in a white dress walking down a red carpet. The right set, also enclosed in a large right-facing curly bracket, contains three identical copies of the first image (the building facade). An orange greater-than sign (>) is placed between the two sets, indicating that the first set of diverse frames receives a higher reward than the second set of non-diverse frames.

# Rewards

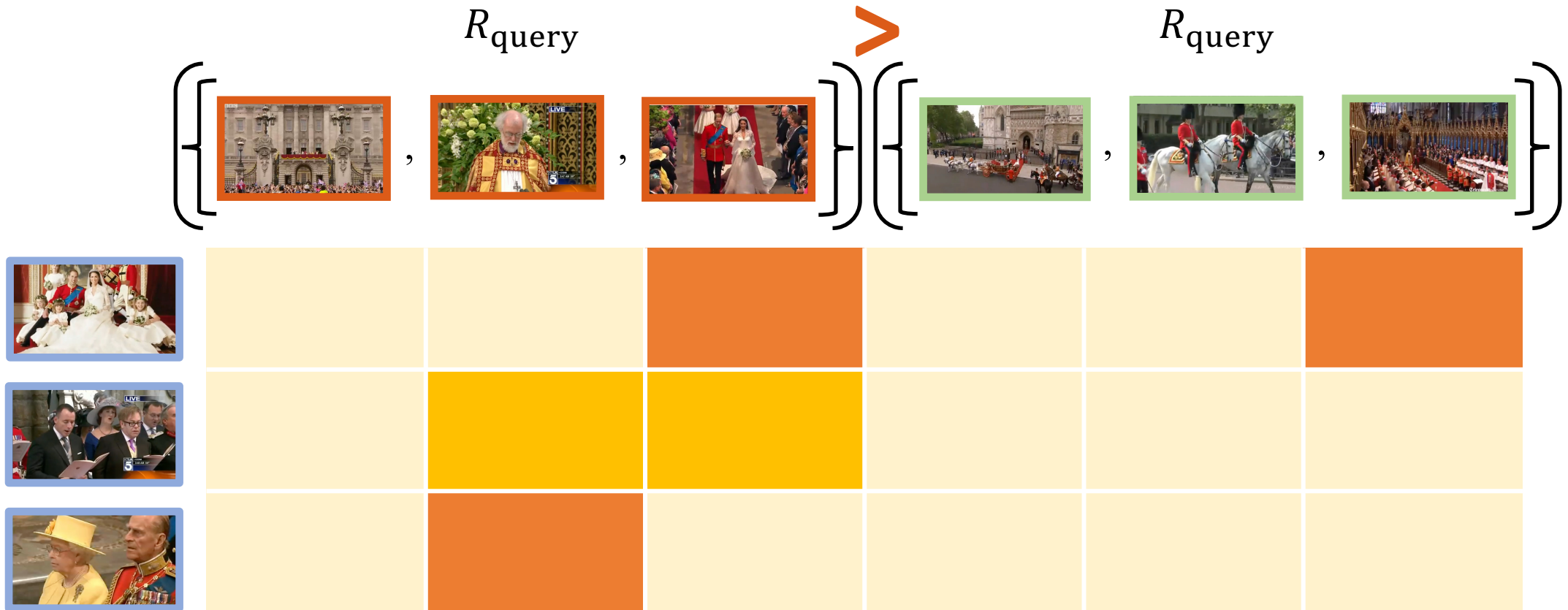
1. **Diversity Reward ( $R_{div}$ ):** Selected frames are diverse
2. **Representativeness Reward ( $R_{rep}$ ):** Selected frames are **cluster centers** of the input frames



$$R_{rep} \left( \left\{ \left[ \text{Frame 1} \right], \left[ \text{Frame 2} \right], \left[ \text{Frame 3} \right] \right\} \right) > R_{rep} \left( \left\{ \left[ \text{Frame 4} \right], \left[ \text{Frame 5} \right], \left[ \text{Frame 6} \right] \right\} \right)$$

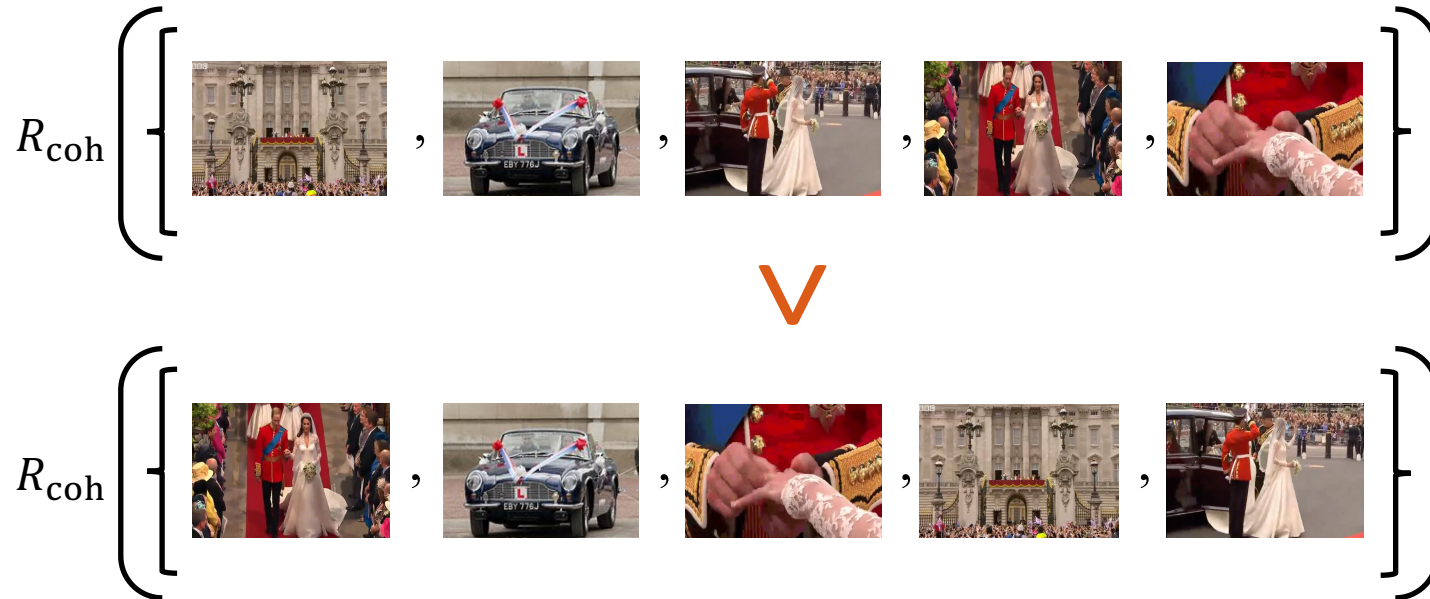
# Rewards

1. **Diversity Reward ( $R_{div}$ ):** Selected frames are diverse
2. **Representativeness Reward ( $R_{rep}$ ):** Selected frames are cluster centers of the input frames
3. **Query-adaptability Reward ( $R_{query}$ ):** Selected frames are similar to [retrieved web-images](#)



# Rewards

1. **Diversity Reward ( $R_{\text{div}}$ ):** Selected frames are diverse
2. **Representativeness Reward ( $R_{\text{rep}}$ ):** Selected frames are cluster centers of the input frames
3. **Query-adaptability Reward ( $R_{\text{query}}$ ):** Selected frames are similar to retrieved web-images
4. **Temporal Coherence Reward ( $R_{\text{coh}}$ ):** Summaries are visually coherent



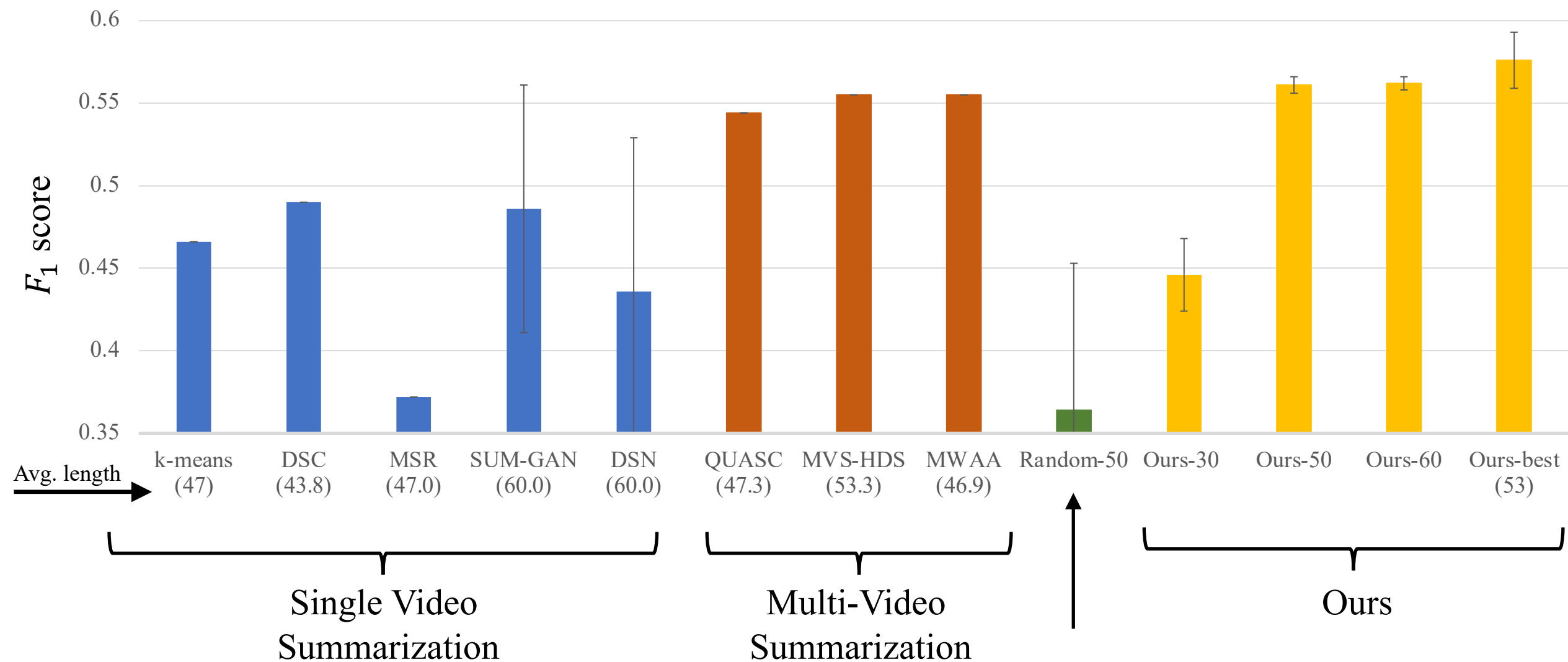


# Experiments

**Dataset:** MVS1K Dataset (4 ground-truth summaries per Query ID)

Query ID	Query	# Videos	# Frames	# Images
1	Britains Prince William wedding 2011	90	1124	324
2	Prince death 2016	104	1549	142
3	NASA discovers Earth-like planet	100	1349	226
4	American government shut-down 2013	82	962	177
5	Malaysia Airline MH370	109	1330	435
6	FIFA corruption scandal 2015	90	785	177
7	Obama re-election 2012	85	1263	207
8	Alpha go vs Lee Sedo	84	976	118
9	Kobe Bryant retirement	109	1140	221
10	Paris terror attacks	83	857	651
<b>Total</b>	-	936	-	2678

# Quantitative Results



# Qualitative Results

K-Means (6|7)



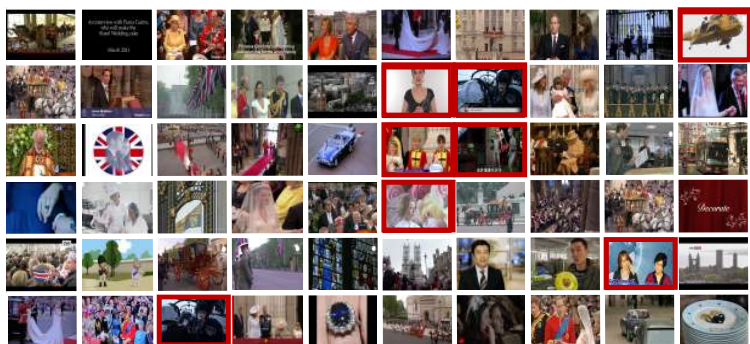
DSC (0|4)



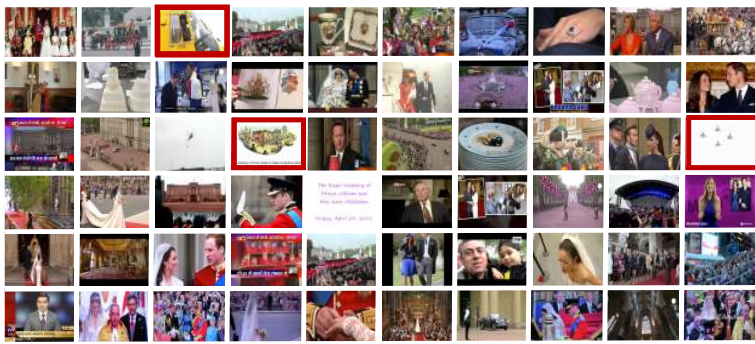
MSR (0|7)



SUM-GAN (0|8)



DSN (0|3)



HDS (2|4)



QUASC (0|3)



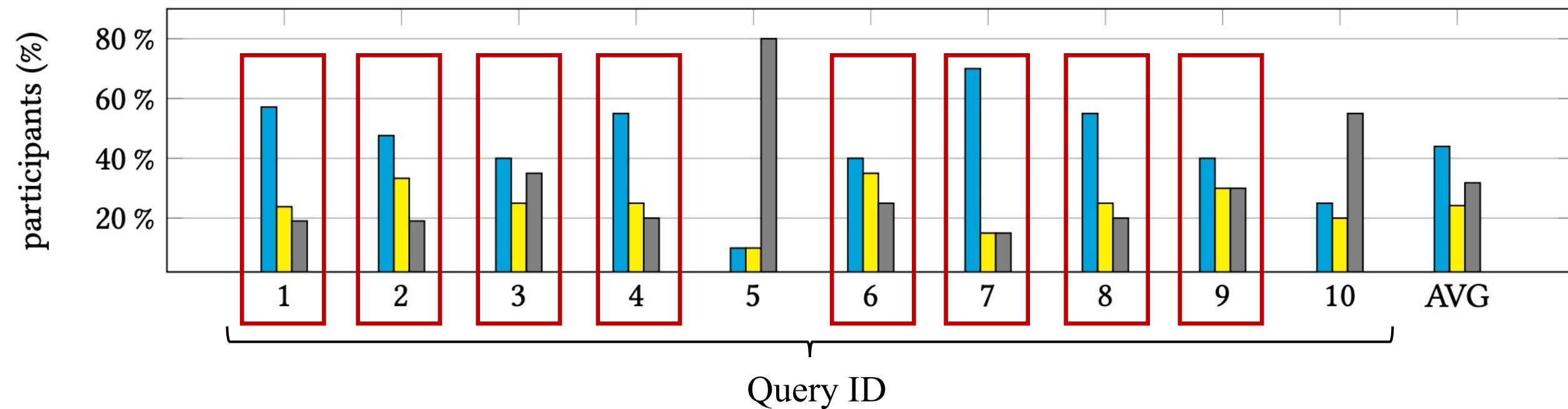
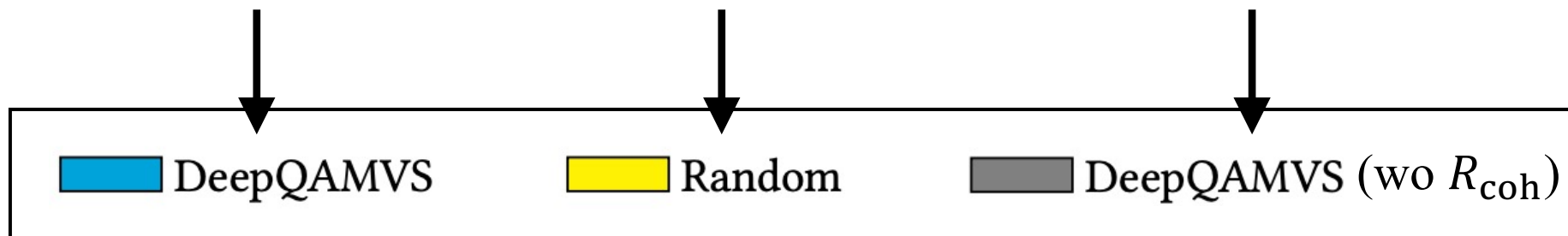
→ OURS (0|1)



■ Redundant  
■ Unimportant

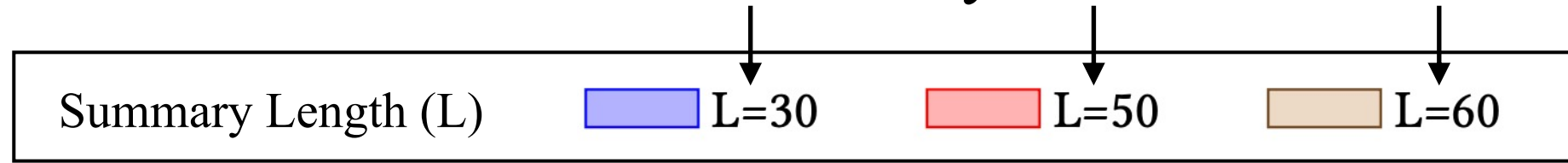
# User Study for Temporal Coherence Assessment

21 participants

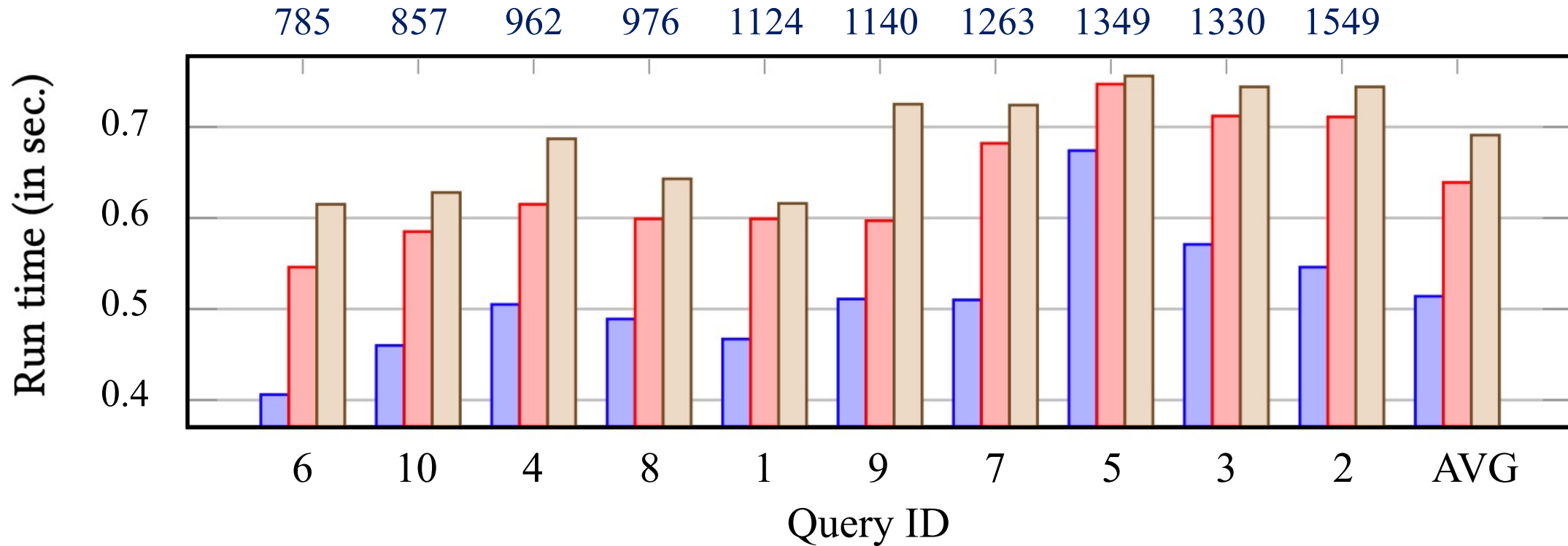


**Users preferred our DeepQAMVS summary in 8 out of 10 queries**

# Run-Time Analysis



Number of input frames



Run-time scales **linearly** with the number of input frames and summary length

# Conclusion

## Advantages

- ✓ First end-to-end trainable model for QAMVS
- ✓ SOTA results on MVS1K dataset
- ✓ Scales linearly with the number of input video frames and summary length

## Future QAMVS would benefit from:

Better evaluation metric combining visual, textual and temporal order overlap

Better rewards for temporal coherence

New dataset with segment based textual annotation

