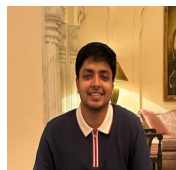# Aligning Moment in Time using Video Queries

Yogesh Kumar[1*], Uday Agarwal[1*], Manish Gupta[2], Anand Mishra[1]

*Equal Contribution
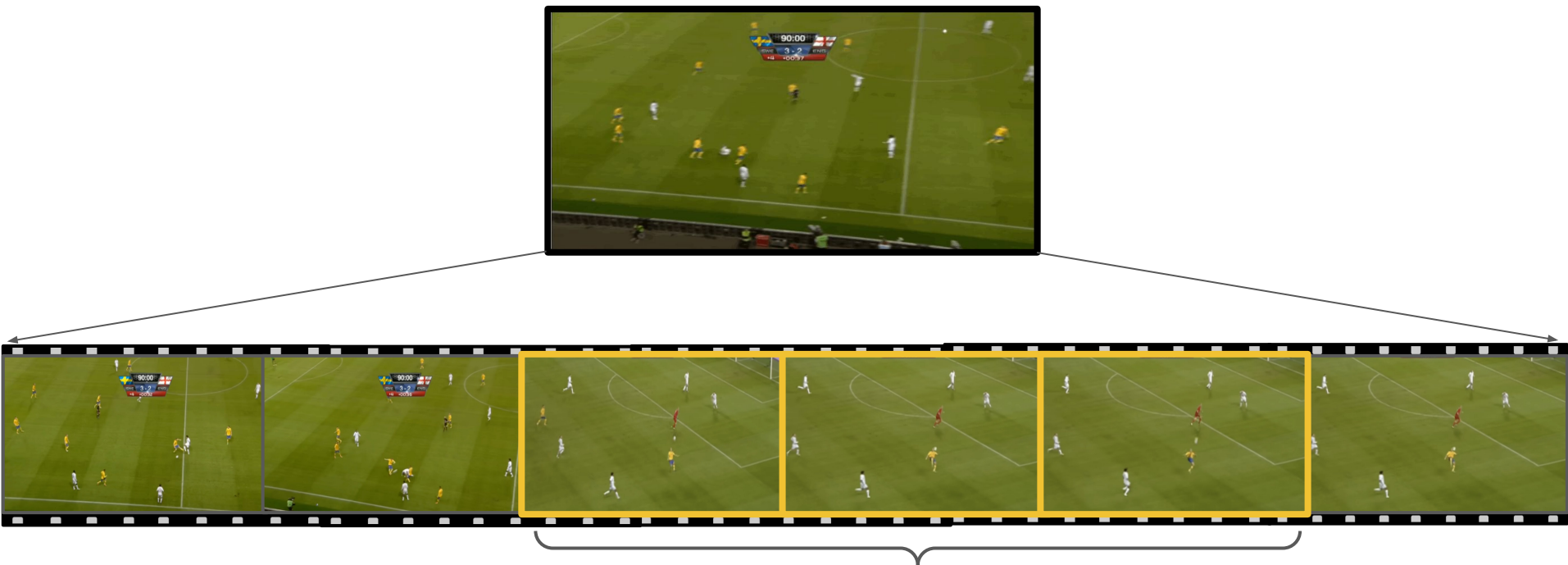
भारतीय प्रौद्योगिकी संस्थान जोधपुर
**Indian Institute of Technology Jodhpur**
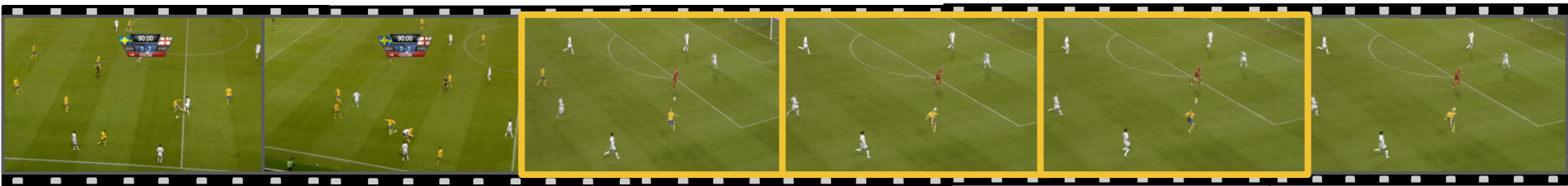
Microsoft

1

# Video Moment Retrieval



*Segment of user Interest*

# How to Represent a Query?

*Option 1: Text Query*

"A man positions himself beneath the ball, leaping into the air, bending his knees and arching his back. As he flips backward, he connects with the ball using his foot."



*Segment of user Interest*

# How to Represent a Query?

**Option 1***: Text Query*

"A man positions himself beneath the ball, leaping into the air, bending his knees and arching his back. As he flips backward, he connects with the ball using his foot."
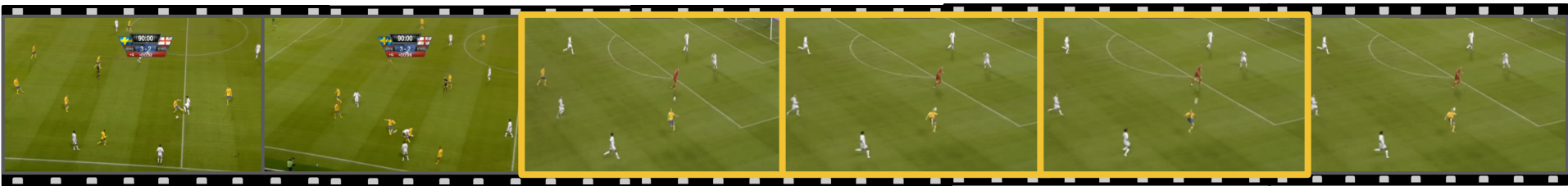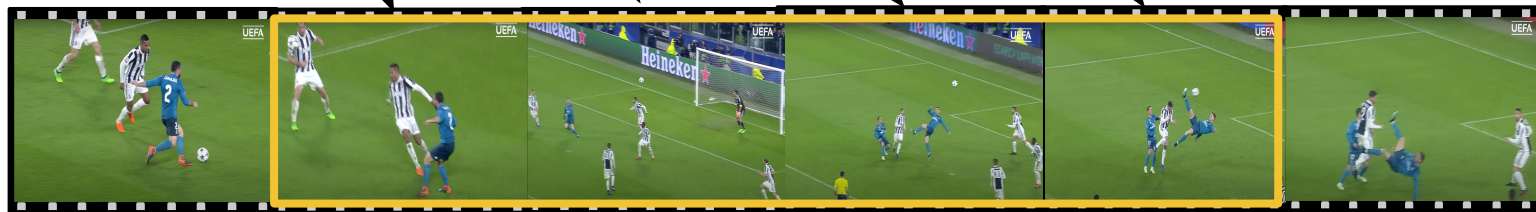


**Option 2***: Video Query*



*Segment of user Interest*
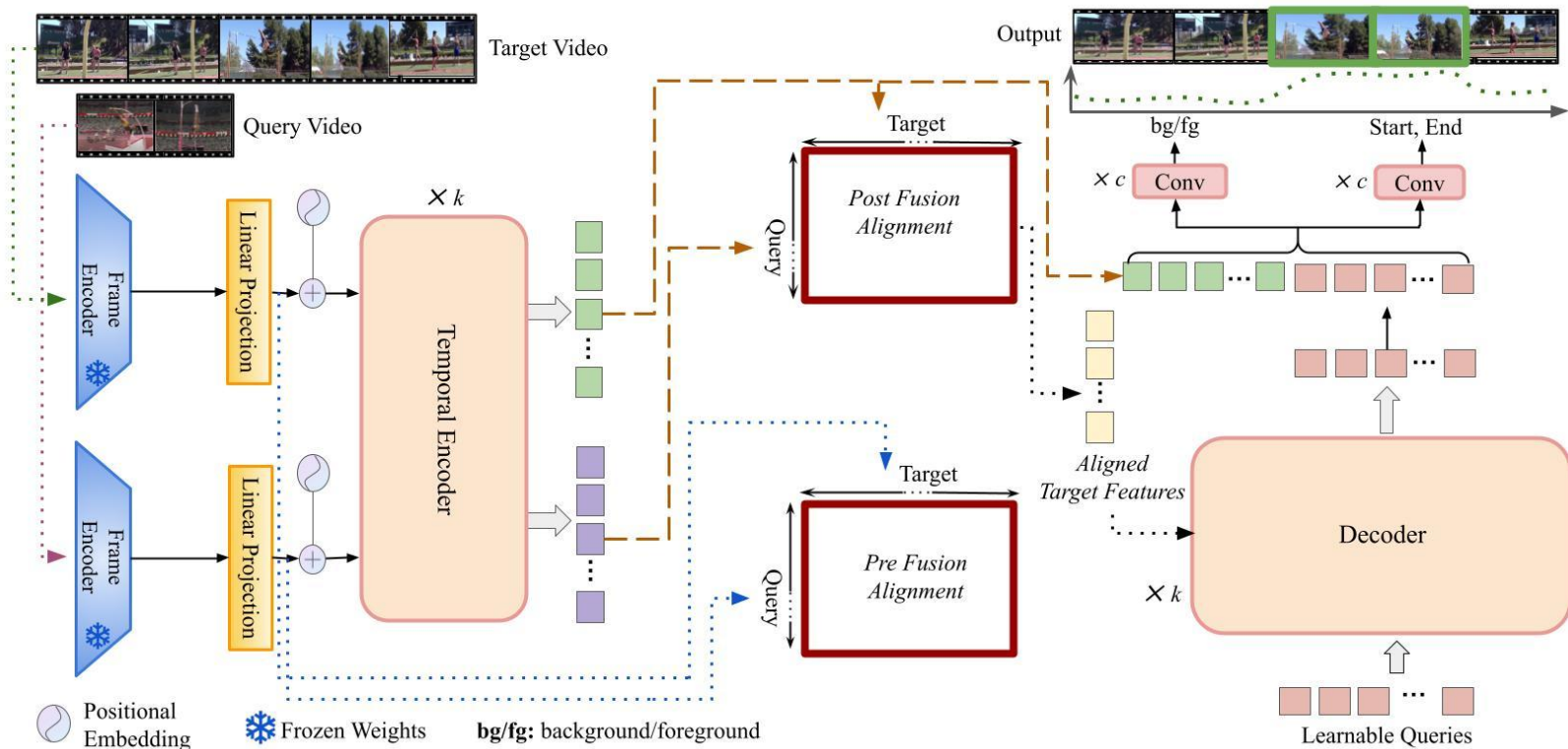
# Advantage of Video Query
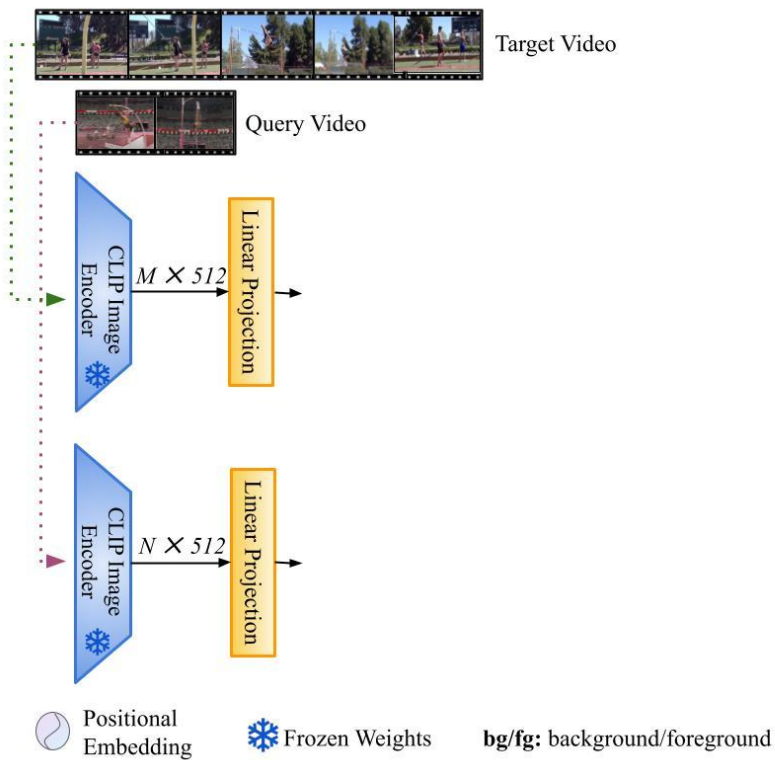
**Query Video**



*Frame to Frame Alignment*

**Target Video**

- *Enables fine-grain alignment and correlation Learning*
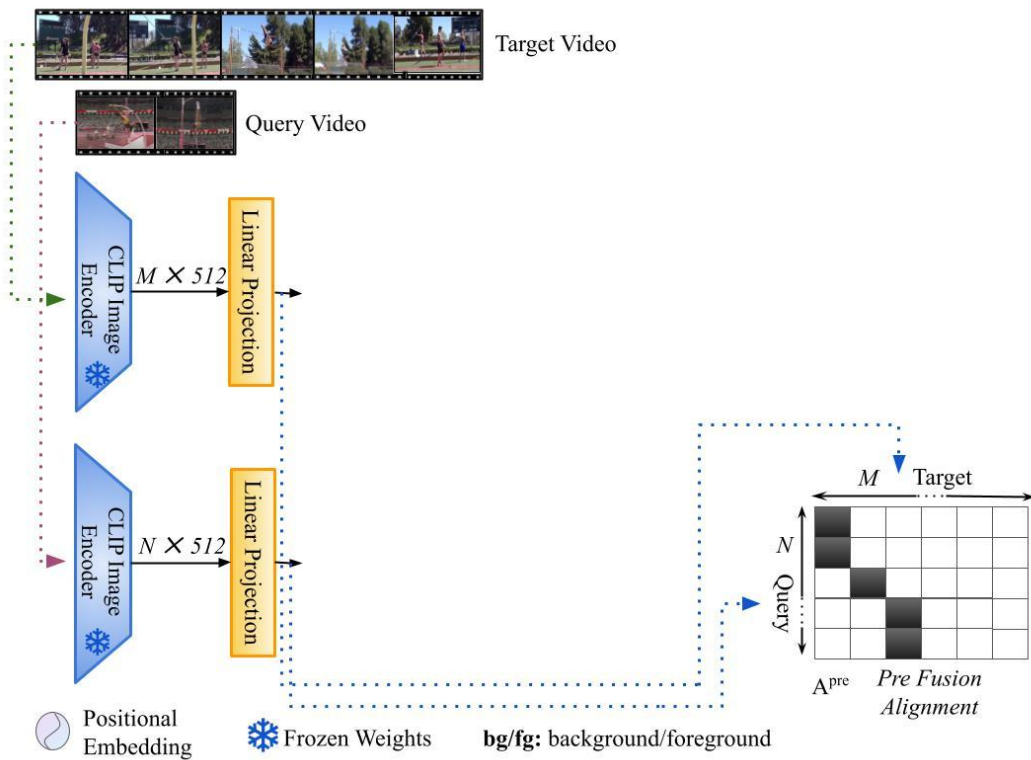- *Capture rich spatial-temporal cues directly from video query*
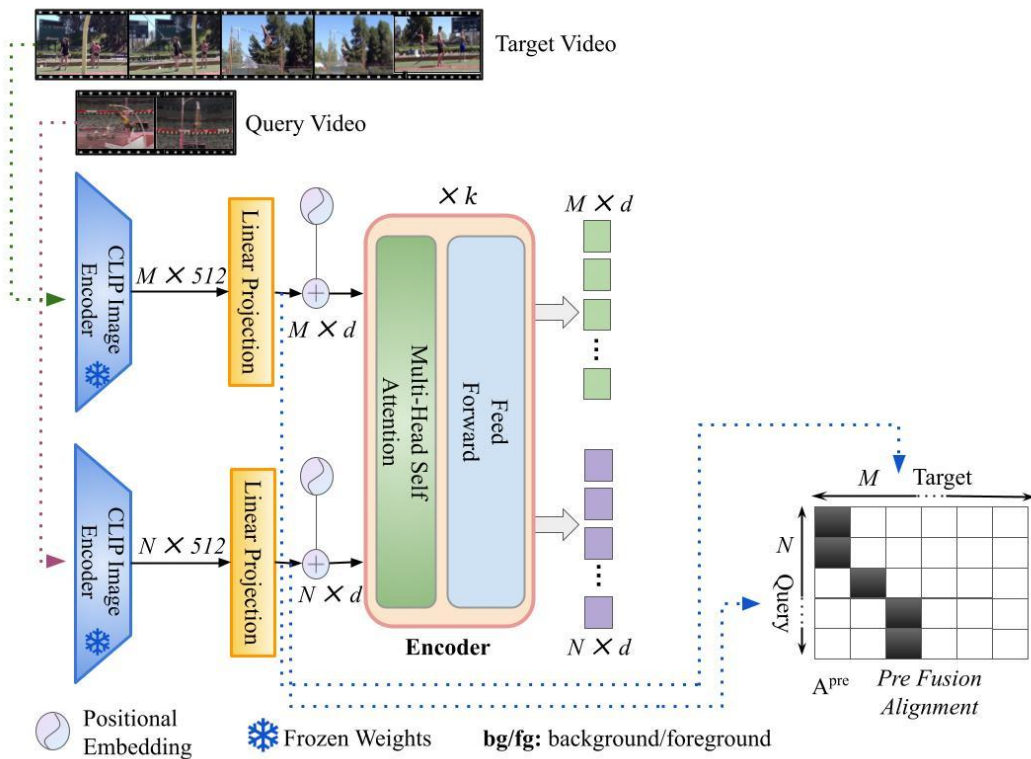
# Proposed Method: MATR (**M**oment **A**lignment **TR**ansformer)



6

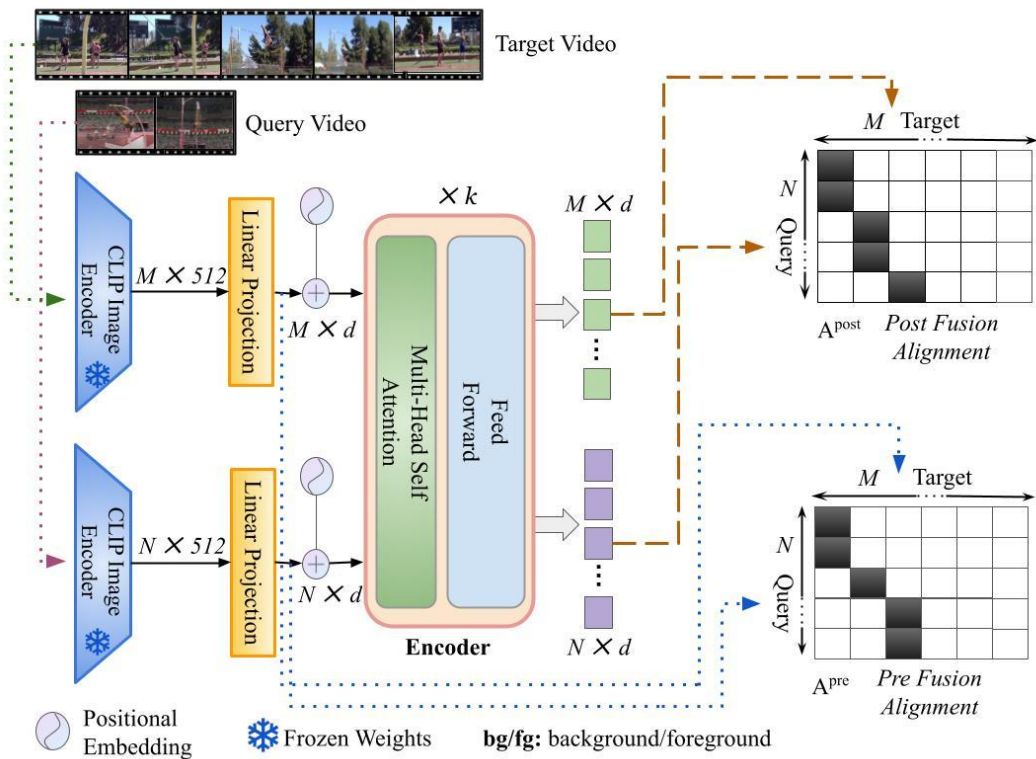# **MATR**: *Video Feature Extraction*



Target Video

Query Video

CLIP Image Encoder — $M \times 512$ — Linear Projection

CLIP Image Encoder — $N \times 512$ — Linear Projection

Positional Embedding    ❄ Frozen Weights    **bg/fg:** background/foreground

# **MATR**: *Fre-fusion Alignment*



Target Video

Query Video

CLIP Image Encoder — $M \times 512$ — Linear Projection

CLIP Image Encoder — $N \times 512$ — Linear Projection

$M$ Target

$N$ Query

$A^{pre}$  *Pre Fusion Alignment*

Positional Embedding

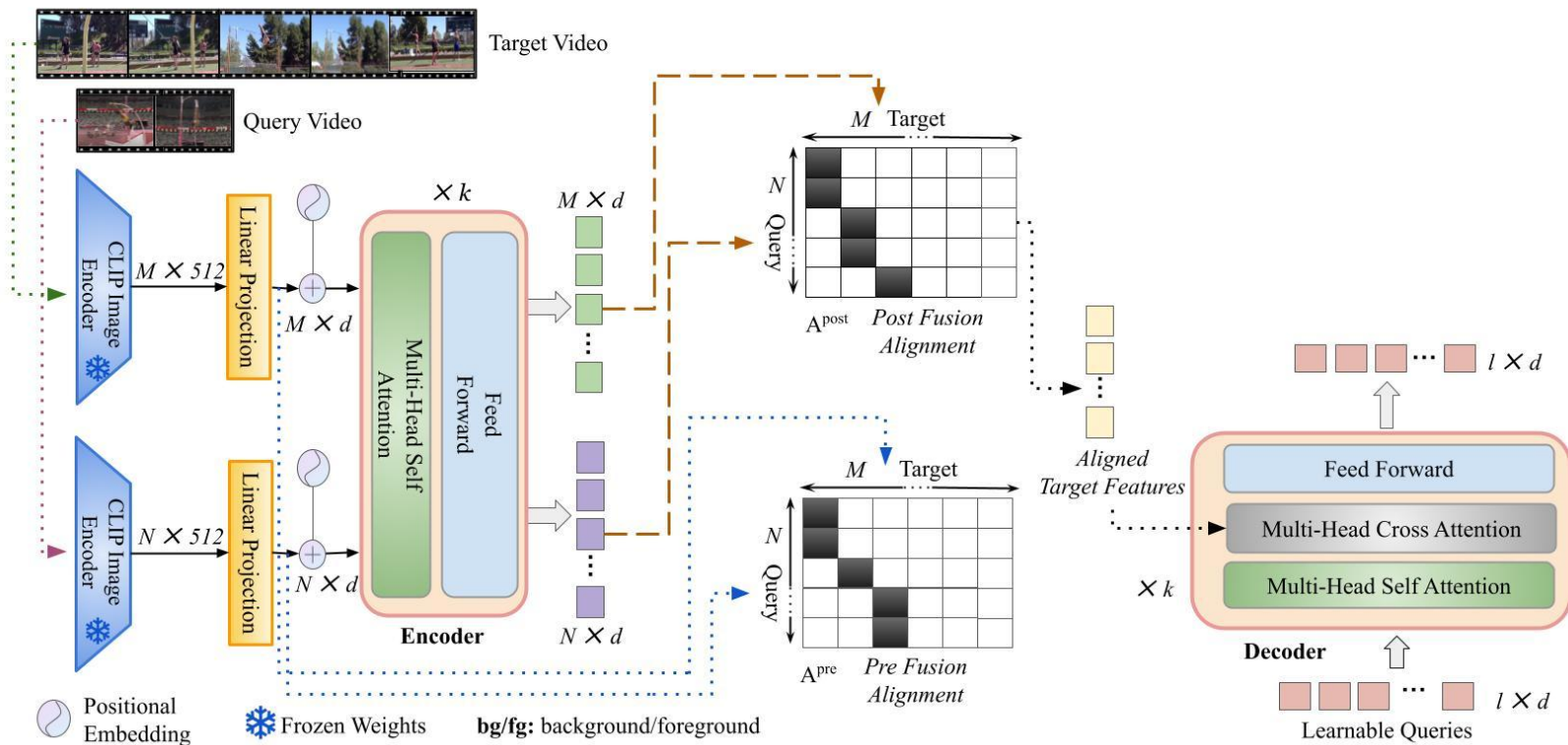Frozen Weights

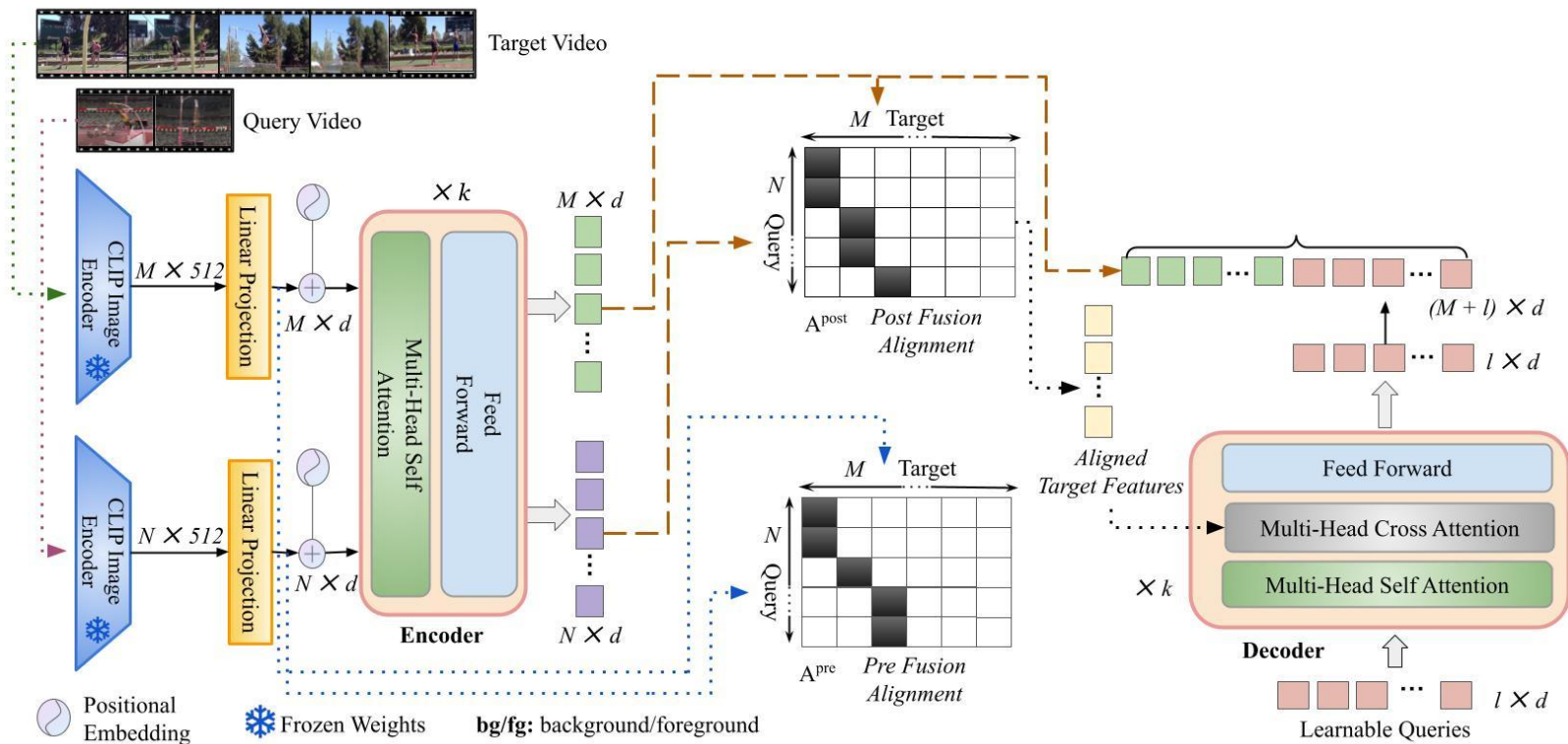**bg/fg:** background/foreground

# MATR: *Fre-fusion Alignment*
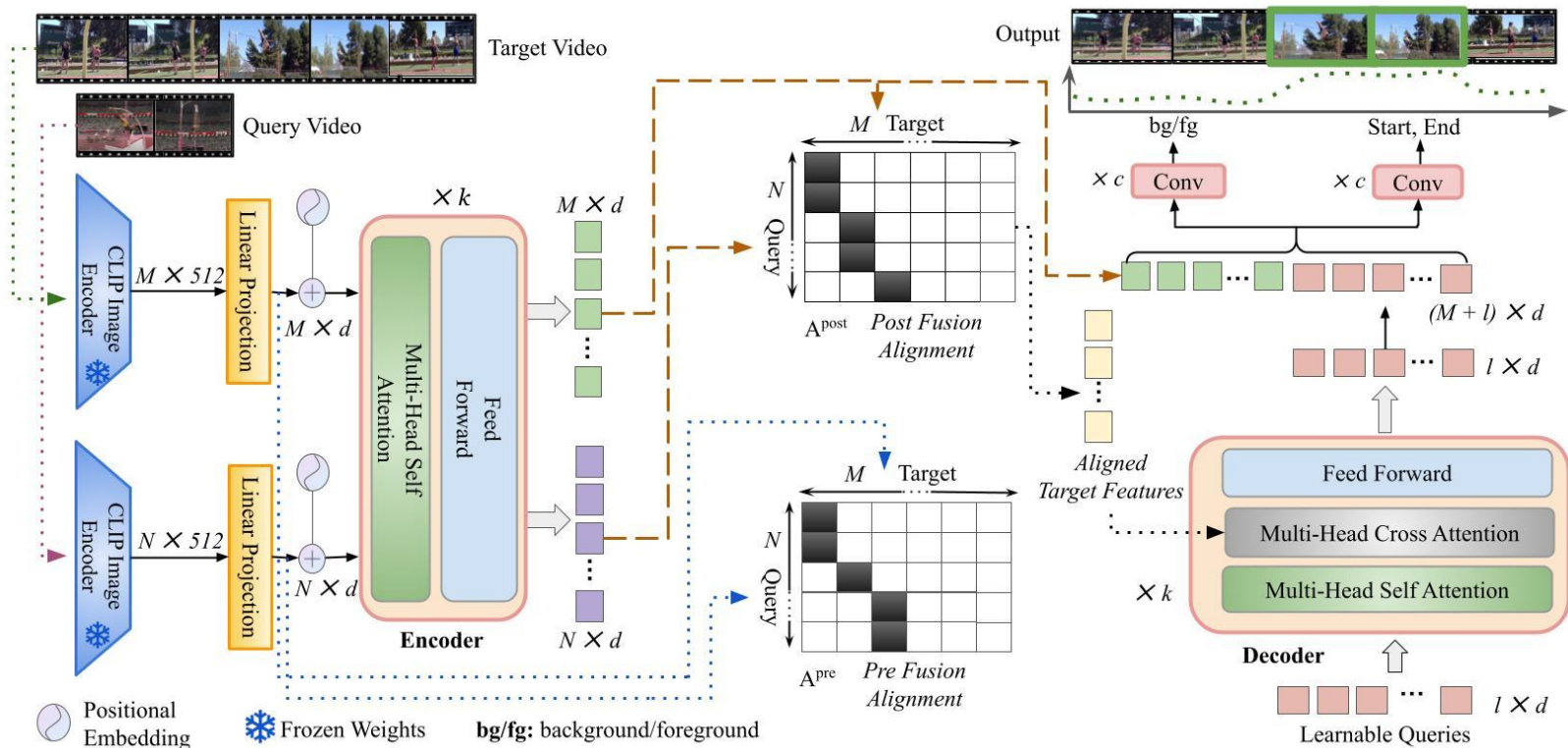
# MATR: *Post-fusion Alignment*

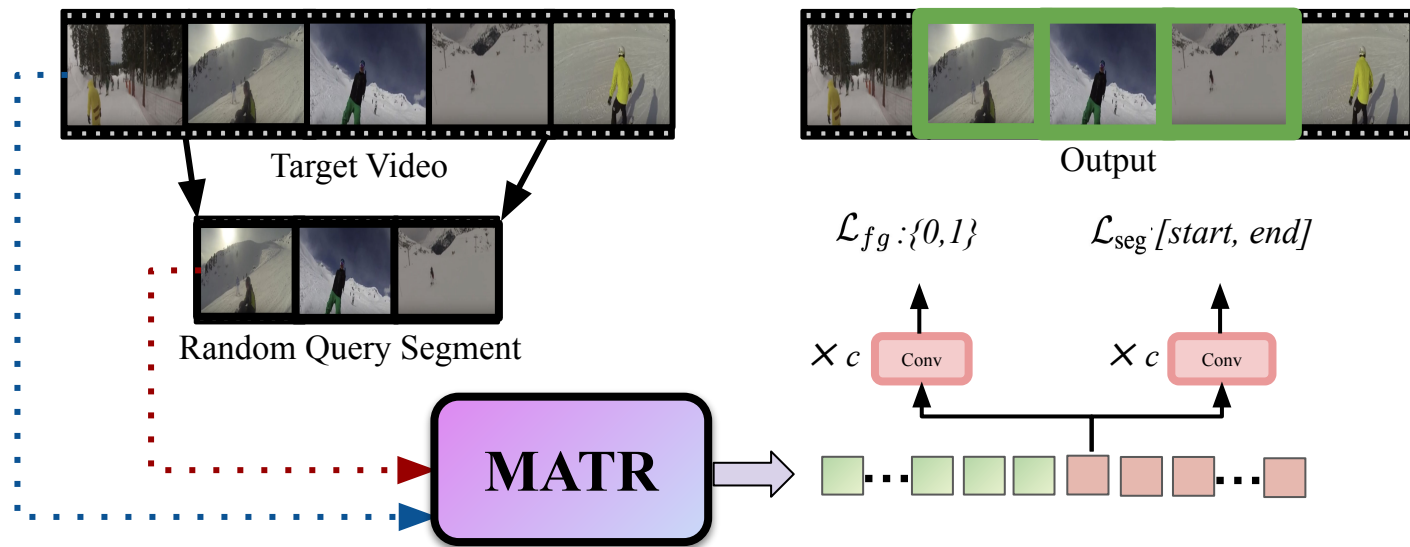# MATR: *Passing Aligned Target Feature to Decoder*

# **MATR**: *Combining Query fused and Query aligned representations of target Video*

# **MATR**: *Predicting moments using heads on combined target representation*
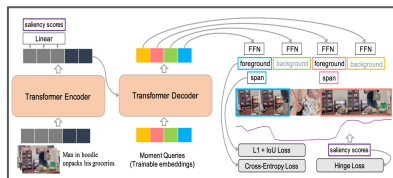
# **Pre-training:** *How to better initialize MATR?*



Target Video

Random Query Segment

**MATR**

Output

$\mathcal{L}_{fg} : \{0,1\}$        $\mathcal{L}_{seg} \cdot [start, end]$

$\times c$  Conv        $\times c$  Conv

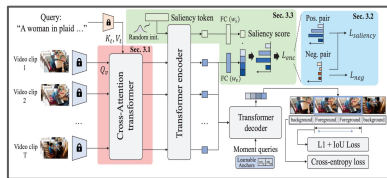*Random clip localization*

# Competitive Approaches
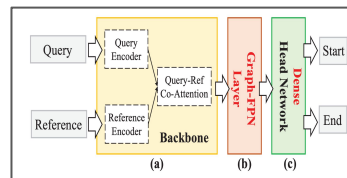


**Text-VMR Methods**

Moment-DETR        QD-DETR

[Lei et al., NeurIPS'21;Moon et al., CVPR'23]
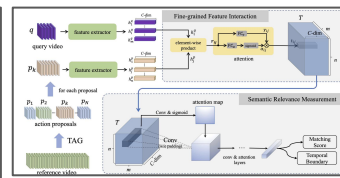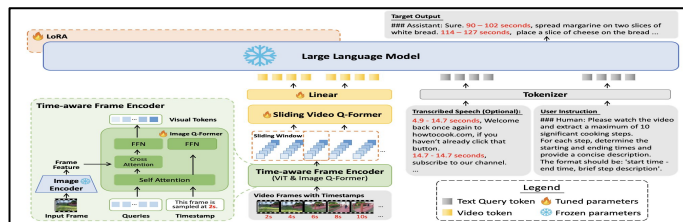
**Video to Video VMR Methods**

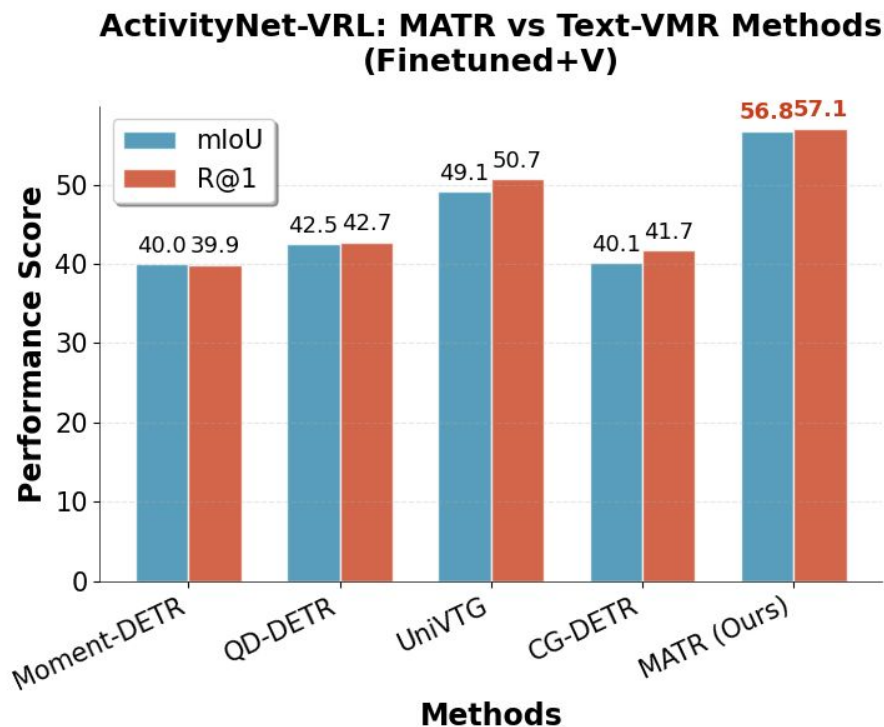GDP        FFI_SRM

[Chen et al., AAAI'20;Huo et al., TMM'23]

**Vision Language Models**

TimeChat

[Ren et al., CVPR 2023]

# **Results:** *Comparison with Text-VMR Methods*



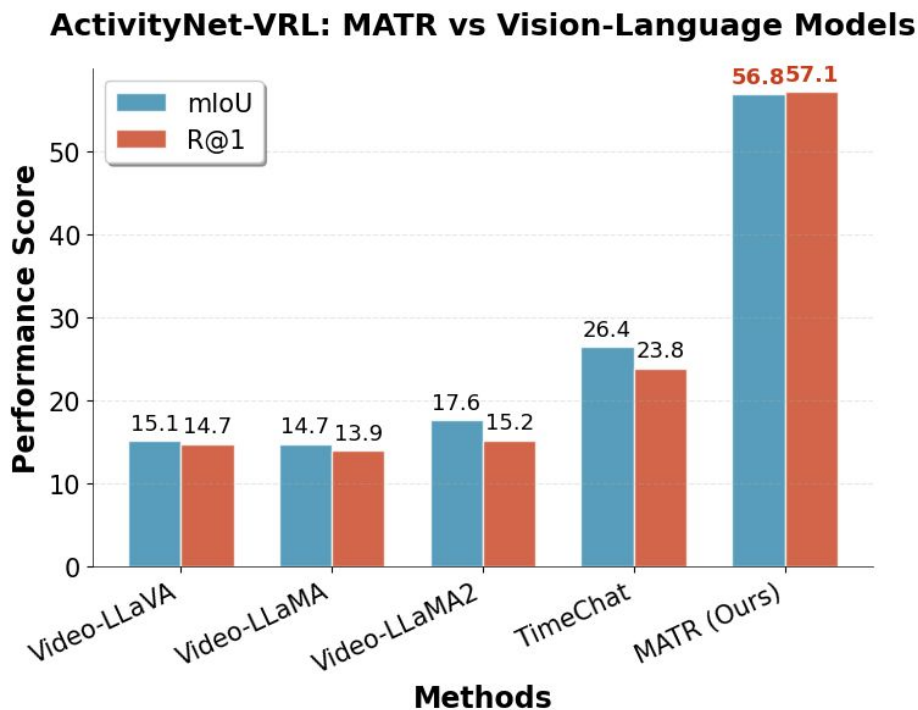ActivityNet-VRL: MATR vs Text-VMR Methods (Finetuned+V)

# **Results:** *Comparison with Video-to-Video Methods*



ActivityNet-VRL: MATR vs Fully-Supervised Methods

# **Results:** *Comparison with Vision-Language Models*



ActivityNet-VRL: MATR vs Vision-Language Models

# **Ablations:** *Pre/Post Fusion Alignment*



ActivityNet-VRL: Ablation Study
Dual-Stage Alignment Components

# Advantage of Pre-training



ActivityNet-VRL

# Qualitative Results (1/2)



**Query Video**



**Target Video**

**Ground truth**: [2.0, 12.7], **Prediction**: [1.7, 12.7]

# Qualitative Results (2/2)



**Query Video**

**Target Video**

**Ground truth**: [13.3, 25.6], **Prediction**: [13.2, 25.8].

# Conclusion

- MATR advances Video to Video moment retrieval via:
    - *Dual-stage alignment within transformer*
    - *Self-supervised pre-training*
    - *Strong performance across benchmarks*
- Future Directions
    - *Multi-Moment Extension*
    - *Multimodal Queries (Video + Text)*

https://github.com/vl2g/MATR

भारतीय प्रौद्योगिकी संस्थान जोधपुर
**Indian Institute of Technology Jodhpur**

Microsoft