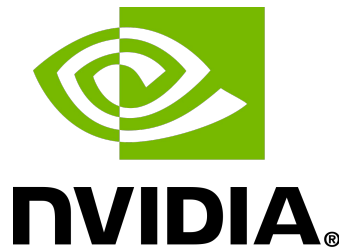




# HERMES: temporal-coHERent long-form understanding with Episodes and Semantics

*ICCV 2025*





# HERMES: temporal-coHERent long-form understanding with Episodes and Semantics

In Video Understanding, we have always traded accuracy for efficiency...

Not anymore.

Introducing **HERMES**.



# HERMES: temporal-coHERent long-forM understanding with Episodes and Semantics

Given a video such as the one below, **HERMES** does two things:



1st, it combines similar scenes into **contiguous episodes**



2nd, it **identifies key semantics** scattered throughout the video



# HERMES: temporal-coHERent long-form understanding with Episodes and Semantics

Given a video such as the one below, **HERMES** does two things:

1

1st, it combines similar scenes into **contiguous episodes**

2

2nd, it **identifies key semantics** scattered throughout the video

3

It then **combines the Episodes and the Semantics**

Episode #1

Episode #2

Episode #3

Semantics #1

Semantics #2



# HERMES: temporal-coHERent long-form understanding with Episodes and Semantics

**HERMES** achieves **SOTA** results in several benchmarks while being **much faster**

Model	LVU								Breakfast	COIN
	Content			Metadata				Avg		
	Relation	Speak	Scene	Director	Genre	Writer	Year			
FACT [21]	-	-	-	-	-	-	-	-	86.1	-
Obj. Transformer [43]	53.1	39.4	56.9	52.1	54.6	34.5	39.1	47.1	-	-
VIS4mer [15]	57.1	40.8	67.4	62.6	54.7	48.8	44.8	53.7	88.2	88.4
TranS4mer [16]	59.5	39.2	70.9	63.9	55.9	46.9	45.5	54.5	90.3	89.2
S5 [41]	67.1	42.1	73.5	67.3	<u>65.4</u>	51.3	48.0	59.2	90.7	90.8
Movies2Scenes [5]	<b>71.2</b>	42.2	68.2	70.9	57.8	55.9	<u>53.7</u>	60.0	-	-
MA-LMM [14]	58.2	<u>44.8</u>	<u>80.3</u>	<u>74.6</u>	61.0	<u>70.4</u>	51.9	<u>63.0</u>	<u>93.0</u>	<u>93.2</u>
HERMES (Ours)	<u>67.6</u>	<b>47.5</b>	<b>90.0</b>	<b>82.6</b>	<b>69.5</b>	<b>77.2</b>	<b>57.7</b>	<b>70.3</b>	<del>95.2</del>	<b>93.5</b>

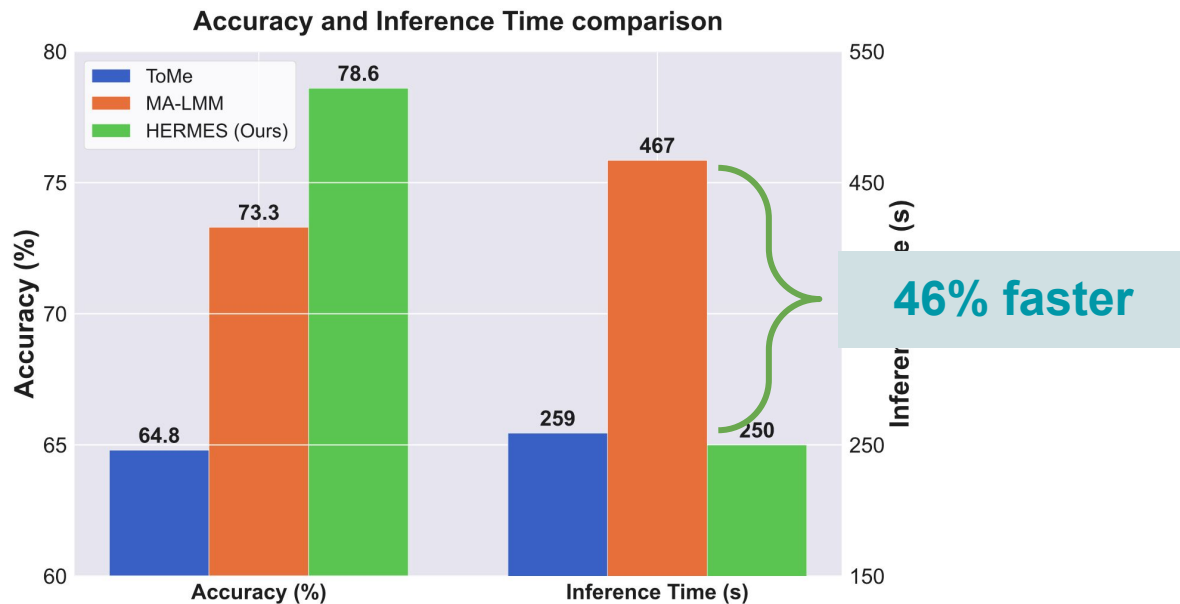
**+7.3%**  
**Accurate**

Model	Acc.	Time	Mem. (GB)
LongVA (7B)	54.11	1	42.5
+ ECO	54.19	0.700 (-30%)	<b>22.9</b>
+ SeTR	<b>54.56</b>	0.726 (-27%)	32.7

**46% less**  
**GPU Usage**



**HERMES** achieves **SOTA** results in several benchmarks while being **much faster**





**HERMES is open-source**, feel free to use it yourself



**Code**



**Project Page**