EGO-ONLY: EGOCENTRIC ACTION DETECTION WITHOUT EXOCENTRIC TRANSFERRING

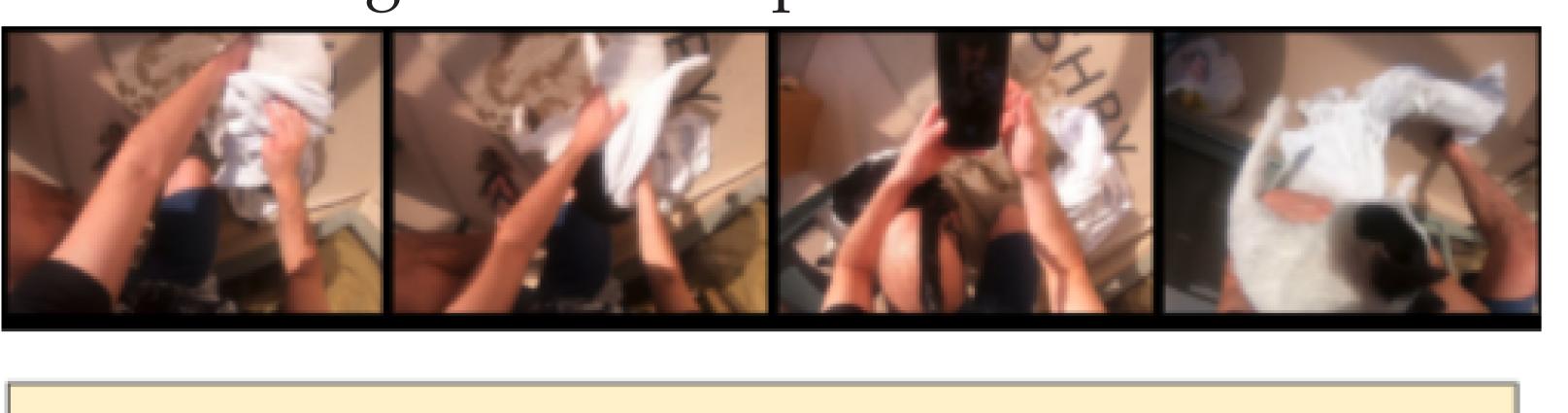
Huiyu Wang¹ Mitesh Kumar Singh¹ Lorenzo Torresani¹

¹Meta AI



GOAL

• Egocentric action detection and recognition. a.k.a. egocentric temporal action localization.

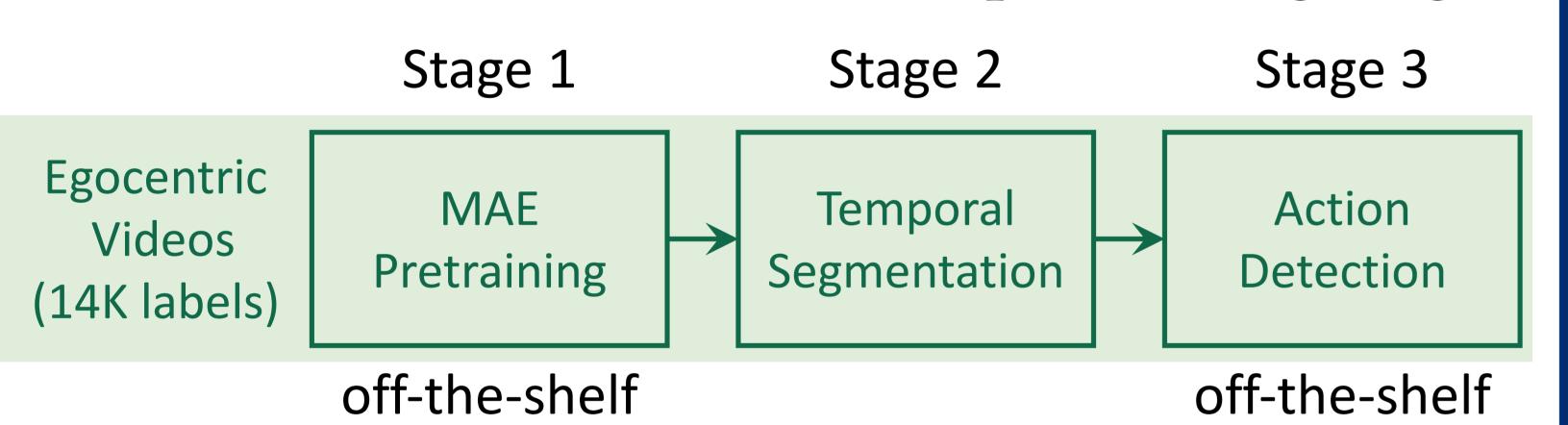


Interact with pet

Use phone

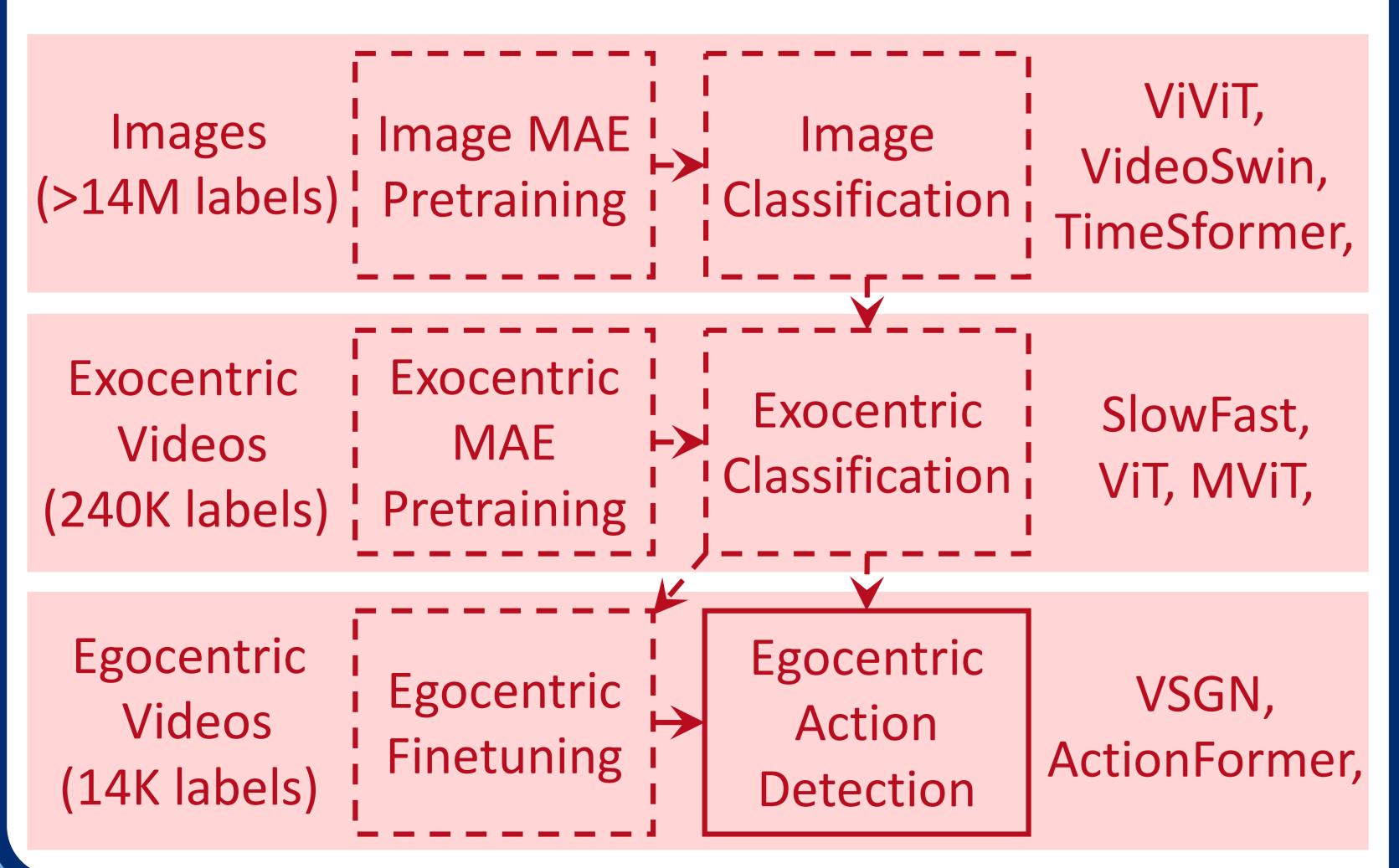
METHOD

• Our Ego-Only approach learns from ego videos only and simplifies the existing pipeline that uses extra data, extra labels, and extra pretraining stages.



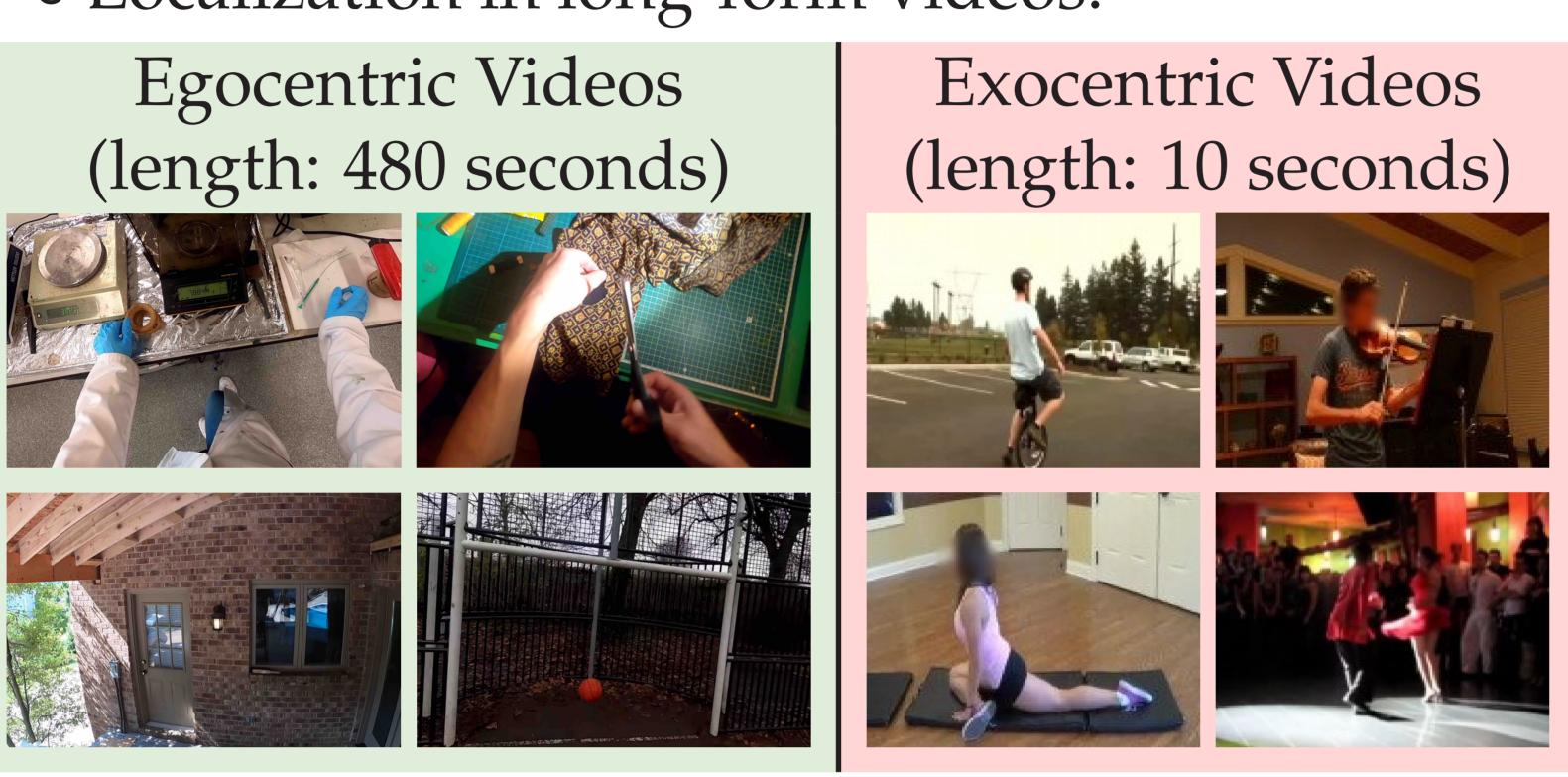
PREVIOUS METHODS

• **Previous methods** depend heavily on transferring from exo to ego, despite the content and appearance gap separating the two domains.



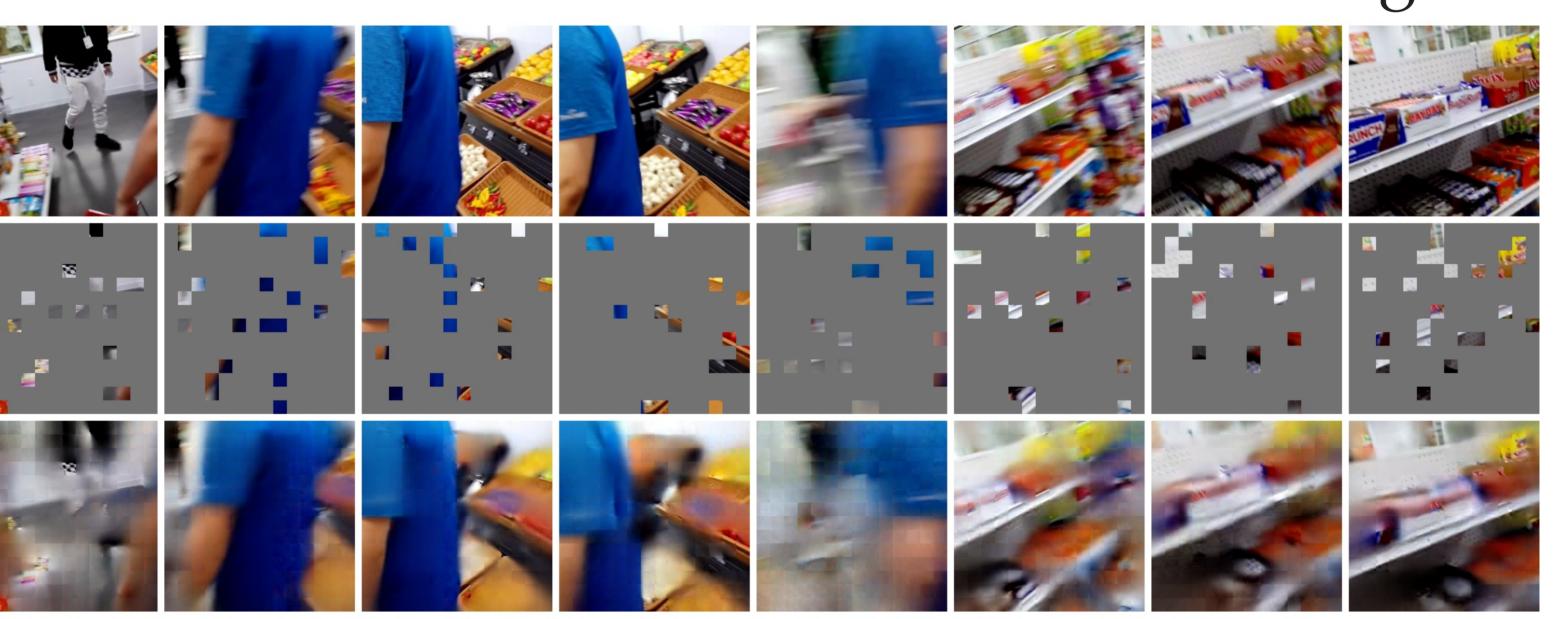
EGO-EXO DOMAIN GAP

- Egocentric domain specific challenges:
- o No actor visible but only hand object interaction.
- o Fine-grained actions, long-tail class distribution.
- o Domain shift between ego and exo datasets.
- o Localization in long-form videos.



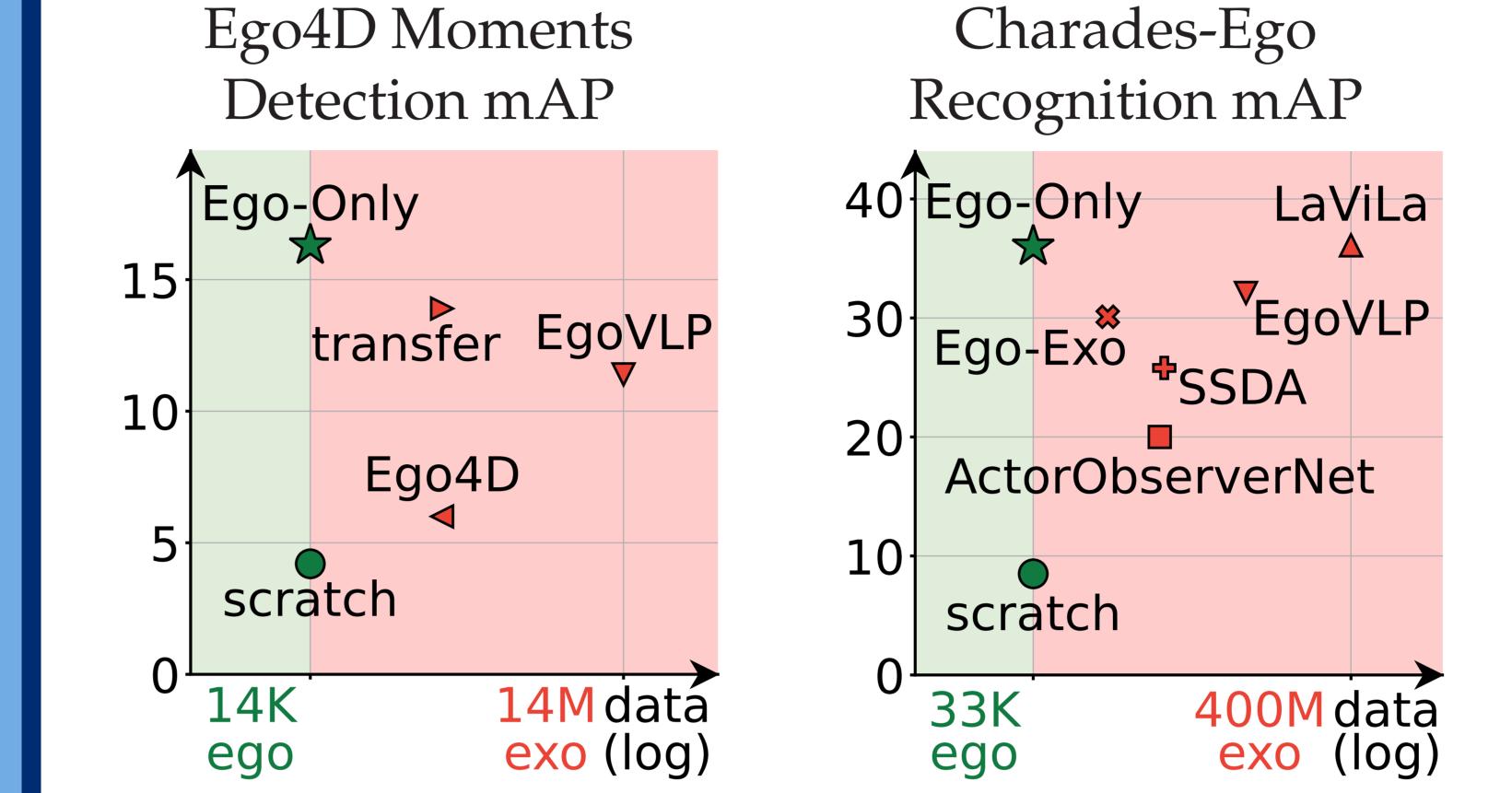
STAGE 1 MAE PRETRAINING

Visualization of MAE reconstruction on Ego4D.



MAIN RESULTS

• Ego-Only uses orders of magnitude fewer labels, simplifies the pipeline, and improves the results.



• EPIC-Kitchens-100 Action Detection mAP.

method	backbone	exo data	verb	noun	# labels
ActionFormer	SlowFast	K400	23.5	21.9	307K
Ego-Only	ViT-B		28.1	26.5	67K

• EPIC-Kitchens-100 Recognition top-1 accuracy.

method	variant	verb	noun
ViViT	ViViT-L/16x2, IN-21k+K400	66.4	56.8
MoViNet	MoViNet-A6, 120 frames	72.2	57.3
MeMViT	, , , , , , , , , , , , , , , , , , , ,	71.4	60.3
LaViLa	TSF-L, WebImageText+Ego4D	72.0	62.9
Ego-Only	ViT-L, 32 frames, 3.2 sec	73.3	59.4

ABLATION STUDY

• Varying the pretraining stage with Video MAE.

method	self-sup. MAE	sup. exo	sup. ego	Ego4D mAP	# labels seen
exo-sup	_	K400	Ego4D	13.9	254K (18×)
ours	Ego4D	-	Ego4D	16.3	14K (1×)
scratch	_	_	Ego4D	4.2	14K (1×)
exo-MAE	K400	_	Ego4D	13.4	14K (1×)
exo-FT	K400	K400	Ego4D	16.2	254K (18×)

Varying the finetuning stage with segmentation.

method	backbone	self-sup. MAE	sup. exo	sup.	Ego4D mAP
exo-sup	SlowFast	_	K400	_	13.2
exo-MAE	ViT-B	K400	_	_	6.7
ego-MAE	ViT-B	Ego4D	_	_	7.8
exo-FT	ViT-B	K400	K400	_	13.5
ours	ViT-B	Ego4D	-	Ego4D	16.3

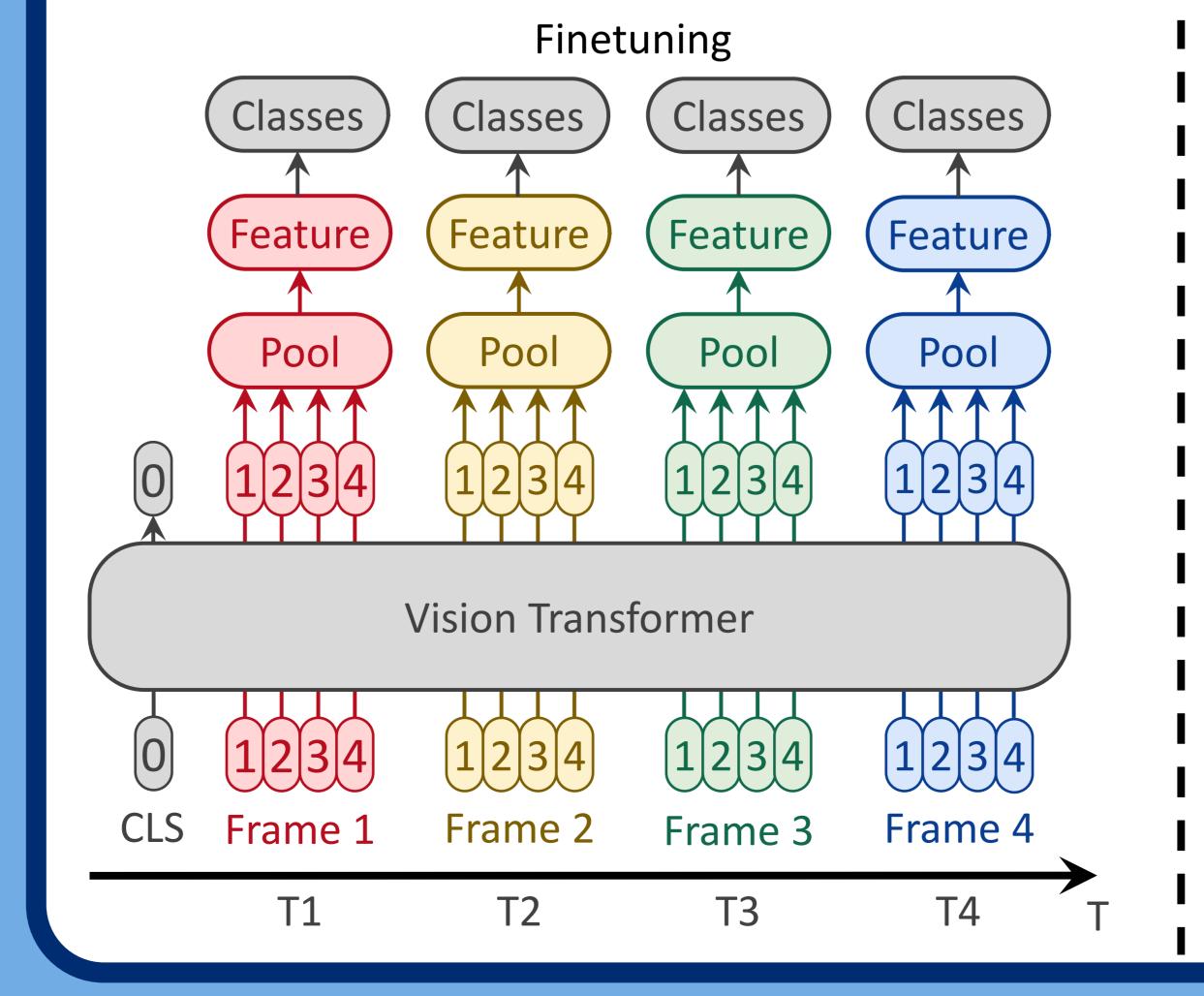
• Joint ego-exo pretraining on Ego4D and Kinetics.

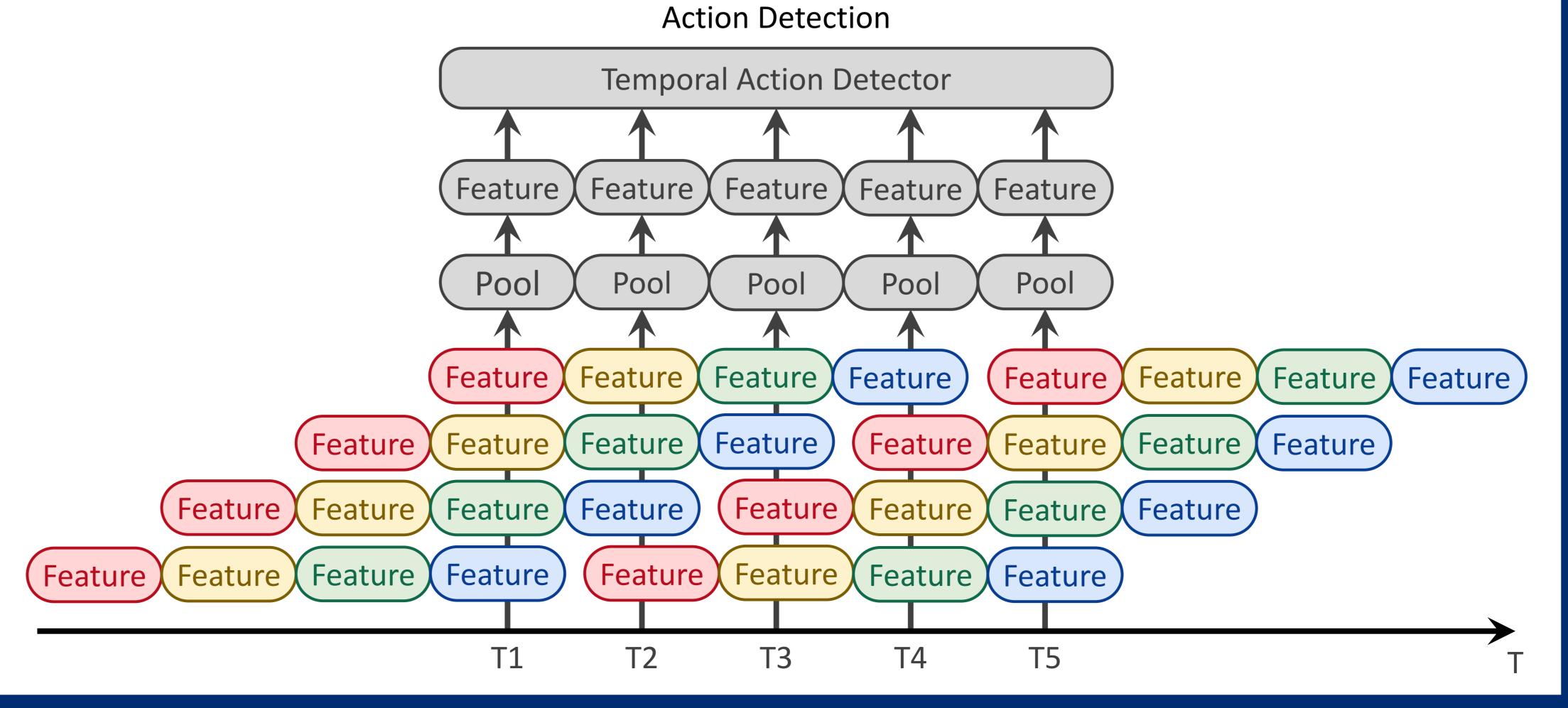
method	self-sup. MAE	sup. exo	sup. ego	Ego4D mAP
exo-MAE	K400	_	Ego4D	13.4
joint-MAE	K400 & Ego4D	_	Ego4D	16.0
ours	Ego4D	_	Ego4D	16.3

- Joint ego-exo finetuning on KEEC.
- o Kinetics-600, Ego4D, EPIC-Kitchens-100, COIN.

method	self-sup. MAE	sup.	sup.	verb mAP	noun mAP	# labels seen
joint-FT	EPIC		KEEC			
ours	EPIC	-	EPIC	29.0	28.1	67K

STAGE 2 TEMPORAL SEGMENTATION & STAGE 3 ACTION DETECTION





RUNTIME

mothod	FLOPs (G)			training time (hours) MAE exo ego total			
method	MAE	exo	ego	MAE	exo	ego	total
exo-sup	_	598	598	_	200.9	10.0	210.9
exo-FT	81	598	598	100.5	50.2	10.0	160.7
ours	81	-	598	100.5	-	10.0	110.5