

- Task: Proposing a generalized model for multi-label action recognition in humans and animals. for action recognition. It leverages visual and textual cues to eliminate the actor-specific information required in prior art.
- **Results:** MSQNet outperforms the prior art on five single and multi-label datasets by up to **50%**.

- **Different actors** (*e.g.* humans and animals) required **separate action recognition models**.
- Those models had **customized designs** to accommodate **specific pose information** of the actors.
- Incorporating multiple actors in a single model is challenging since animals exhibit different shapes, sizes, and appearances.

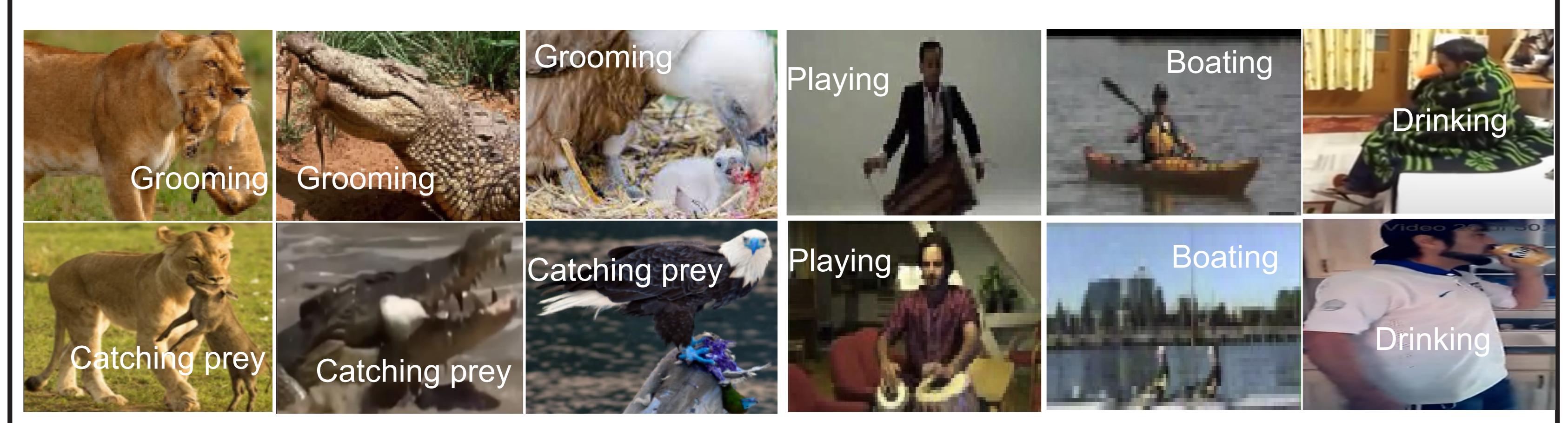
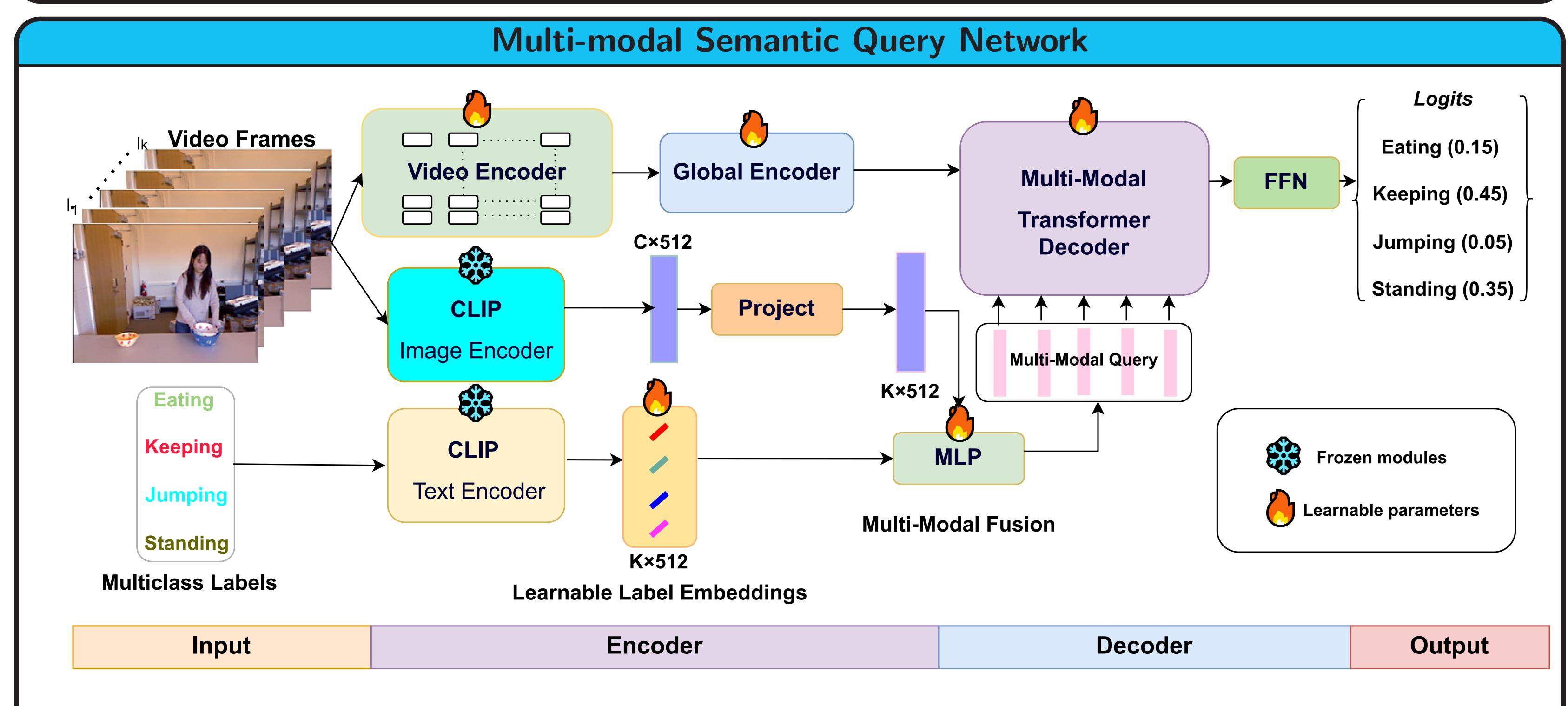


Illustration of large action variation across different actors (e.g., animals and humans). Such differences often motivate the development of actor-specific action recognition models, such as using actor-specific pose estimation.



Our proposed Multi-modal Semantic Query Network (MSQNet) for multi-modal multi-label action recognition has three main components: • a **spatio-temporal video encoder**, that extracts the spatio-temporal features from an input video, • a vision-language query encoder, that merges the visual and textual information, • a multi-modal decoder, transforms the video encoding to make multi-label classification with a feed-forward network (FFN). We evaluate MSQNet on three multi-label (Charades, Hockey, and Animal Kingdom) and two single-label (Thumos 14 and HMDB51) action recognition datasets.

Actor-agnostic Multi-label Action Recognition with Multi-modal Query Sauradip Nag^{*1,2,5} Joaquin M Prada^{1,4} Xiatian Zhu^{1,2,3} Anjan Dutta^{*1,2,3,4} Anindya Mondal* 1,2,3

⁵iFlyTek-Surrey Joint Research Center on Al

*Authors have equal contributions

TL:DR

• Action: Formulating a novel Multi-modal Semantic Query Network (MSQNet) model in a transformer-based object detection framework (e.g., DETR)

Motivation

• Previous models primarily focused on single-label action classification, whereas multiple actions usually occur in a single video in the real world.

		Quar	itita	tive	Per	form	nan	ce	
	Chara	des						Thu	mos
Method	Backbo	ne Pretrain	e Pretrain MMQ		Method		Backb	Backbone	
Zhang et al. CVPR '21	Nonlocal	-101 -	No	44.20	Zhao et al. ICCV '17		C3I	D	
Fan et al. ICCV '21	SlowFa	st K600	No	43.90	Xu et al. TPAMI '19		C3I	D	
Wang et al. ArXiv '21	ViT-E	3 -	No	44.30	Lin et al. ICCV '19		C3I	D	
MSQNet	ViT-E	3 K400	No	43.99	MSQNet		C3I	D	
MSQNet	TS	K400	No	44.11	MSQNet		TS	5	
MSQNet	TS	K400	Yes	47.57	MSQNet		TS	5	
	Animal Ki	ingdom						H	ockey
Method	Backbo	ne Pretrain	MMQ	mAP	Method		Backb	one	
Ng et al. CVPR '22	X3D	-	No	25.25	Carbonneau. ETS'17		-		
Ng et al. CVPR '22	I3D	-	No	16.48	Zhang et al. CVPR '21		1 CSN-	152	
MSQNet	I3D	-	No	55.59	MSQNet		130)	
MSQNet	TS	K400	No	71.63	MSQNet		T5	5	
MSQNet	TS	K400	Yes	73.10		MSQNe	t	TS	5
				HMDE	351				
	Method Backbor				ne Pi	retrain	MM	Accurac	:y
	Wu et al. CVPR		'23	ViT	WI	T-400M	Yes	84.31	
	Kalfaoglu et al. CVPR '23			R(2+1)	D IO	G65M	No	85.10	
	Wang et al. CVPR		R '23	ViT	K40	0/K600	No	88.10	
		MSQNet		TS		K400	Yes	93.25	
	Split	Method			ımos 14	Charad	les H	MDB51	
	-	Wang et al. ArXiv '21		-		21.1		58.70	
	Reported	Bain et al. ArXiv '22			-	25.8		-	
_		Wu et al. CVPR '23 Vanilla MSQNet			-	- 15.87	,	61.40 45.66	7
	50% Seen	Vanilla MSQNet + Text Ini			49.37 53.76	18.30		45.00 51.22	Zero
	JU/0 Jeen	MSQNet + Text III			53.98	30.9 1		59.24	
I —		Vanilla MSQNet			52.02	17.43		48.37	
і 	75% Seen	Vanilla MSQNet				18.62 59			
I		MSQN	7	75.33		35.59 69			

Horse Eating + Walking Walking Eating Standing

