



Jniversity



MARLIN: Masked Autoencoder for facial video Representation LearnINg Zhixi Cai¹, Shreya Ghosh², Kalin Stefanov¹, Abhinav Dhall^{3,1}, Jianfei Cai¹, Hamid Rezatofighi¹, Reza Haffari¹, Munawar Hayat¹

Contributions

We introduce **MARLIN**, a universal facial encoder that learns robust representations from non-annotated web-crawled facial videos in a self-supervised manner.

2. We propose *Fasking*, a facial region-guided tube masking strategy that reconstructs facial regions from densely masked areas. This approach captures both local and global aspects in facial videos, aiding in the acquisition of generic and transferable features.

3. Through thorough analysis, we demonstrate that MARLIN learns rich, consistent, and versatile facial representations, performing well across various tasks such as *Facial* Attribute Recognition, Lip Synchronization, and even in few shot settings.





¹Monash University, ²Curtin University, ³Indian Institute of Technology Ropar



i I	Method				LSE-D↓	LSE-C↑	FID↓
	Speech2V	id [41]			14.230	1.587	12.320
I	LipGAN [42]			10.330	3.199	4.861
I	Wav2Lip [Wav2Lip [57]		7.521	6.406	4.887	
I I	AttnWav2	Lip [<mark>74</mark>]			7.339	6.530	_
I I	Wav2Lip -	Wav2Lip + ViT [28]		8.996	2.807	13.352	
I I	Wav2Lip -	+ ViT + V	VideoMA	E [71]	7.316	5.096	4.097
1	Way2L in			-	7 1 2 7	5 500	2 450
	wav2Lip-				7.127	5.528	3.452
	Few- Data→	Shot	Ada MOSEI [apta 7]	FF++ [62]	5.528	3.452 HQ [85]
	Few- Data→ Task→	Shot Emo.	Ada Ada AOSEI [7-Sen.	apta 7] 2-Sen.	FF++ [62] DeepFake	CelebV- Appr.	3.452 HQ [85] Act.
	Few-, Data→ Task→ Anno.%	+ v11 + M Shot Shot N Emo. Acc.↑	Ada Ada AOSEI [7-Sen. Acc.↑	apta 7] 2-Sen. Acc.↑	FF++ [62] DeepFake AUC↑	5.528 CelebV- Appr. AUC↑	3.452 HQ [85] Act. AUC↑
	Few- Data→ Task→ Anno.%	+ vni + r Shot Emo. Acc.↑ 80.60	IARLIN Ada IOSEI [7-Sen. Acc.↑ 34.63	apta 7] 2-Sen. Acc.↑ 73.70	7.127 tion FF++ [62] DeepFake AUC↑ 0.9305	5.328 CelebV- Appr. AUC↑ 0.9373	3.452 HQ [85] Act. AUC↑ 0.9278
	Few- Task \rightarrow Anno.%	+ VII + N Shot Emo. Acc.↑ 80.60 80.59	Adda Adda Adda Adda 7-Sen. Acc.↑ 34.63 33.73	apta 7] 2-Sen. Acc.↑ 73.70 73.33	7.127 tion FF++ [62] DeepFake AUC↑ 0.9305 0.8681	5.328 CelebV- Appr. AUC↑ 0.9373 0.9273	3.452 HQ [85] Act. AUC↑ 0.9278 0.9270
	Few- Data→ Task→ Anno.% 100% 10%	+ VII + N Shot Emo. Acc.↑ 80.60 80.59 79.89	Adda Adda Adda Acc.↑ 34.63 33.73 33.56	7] 2-Sen. Acc.↑ 73.70 73.33 72.26	7.127 tion FF++ [62] DeepFake AUC↑ 0.9305 0.8681 0.7459	5.328 CelebV- Appr. AUC↑ 0.9373 0.9273 0.8996	3.452 HQ [85] Act. AUC↑ 0.9278 0.9270 0.9201



