

A VISION-FREE BASELINE FOR MULTIMODAL GRAMMAR INDUCTION (SUPPLEMENTARY MATERIALS)

Anonymous authors

Paper under double-blind review

1 VIDEO-ASSISTED PARSING

In Table 1 we present the unabridged version of comparisons for grammar induction with video and text, corresponding to Table 2 in the paper.

Table 1: Comparison with multi-modalities on three benchmark datasets. Note that our method, using features from OPT-2.7B, yields superior results despite not being regularized by multiple modalities. "W/ MD" denotes using multimodal data.

Method	W/ MD	DiDeMo					YouCook2		MSRVTT		
		NP	VP	PP	C-F1	S-F1	C-F1	S-F1	C-F1	S-F1	
LBranch	X	41.7	0.1	0.1	16.2	18.5	6.8	5.9	14.4	16.8	
RBranch	X	32.8	91.5	66.5	53.6	57.5	35.0	41.6	54.2	58.6	
Random	X	36.5 \pm 0.6	30.5 \pm 0.5	30.1 \pm 0.5	29.4 \pm 0.3	32.7 \pm 0.5	21.2 \pm 0.2	24.0 \pm 0.2	27.2 \pm 0.1	30.5 \pm 0.1	
C-PCFG	X	72.9 \pm 5.5	16.5 \pm 6.2	23.4 \pm 16.9	38.2 \pm 5.0	40.4 \pm 4.1	37.8 \pm 6.7	41.4 \pm 6.6	50.7 \pm 3.2	55.0 \pm 3.2	
VC-PCFG	Object	✓	70.5 \pm 15.3	25.7 \pm 15.9	36.5 \pm 24.6	42.6 \pm 10.4	44.0 \pm 10.4	39.9 \pm 8.7	44.9 \pm 8.3	52.2 \pm 1.2	56.0 \pm 1.6
	Action	✓	57.9 \pm 13.5	45.7 \pm 14.1	45.8 \pm 17.2	45.1 \pm 6.0	49.2 \pm 6.0	40.6 \pm 3.6	45.7 \pm 3.2	54.5 \pm 1.6	59.1 \pm 1.7
	R2PID	✓	61.2 \pm 8.5	38.1 \pm 5.4	62.1 \pm 4.1	48.1 \pm 4.4	50.7 \pm 4.2	39.4 \pm 8.1	44.4 \pm 8.3	54.0 \pm 2.5	58.0 \pm 2.3
	S3DG	✓	61.3 \pm 13.4	31.7 \pm 16.7	51.8 \pm 8.0	44.0 \pm 2.7	46.5 \pm 5.1	39.3 \pm 6.5	44.1 \pm 6.6	50.7 \pm 3.2	54.7 \pm 2.9
	Scene	✓	62.2 \pm 9.6	30.6 \pm 12.3	41.1 \pm 24.8	41.7 \pm 6.5	44.9 \pm 7.4	—	—	54.6 \pm 1.5	58.4 \pm 1.3
	Audio	✓	64.2 \pm 18.6	21.3 \pm 26.5	34.7 \pm 11.0	38.7 \pm 3.7	39.5 \pm 5.2	39.2 \pm 4.7	43.3 \pm 4.9	52.8 \pm 1.3	56.7 \pm 1.4
	OCR	✓	64.4 \pm 15.0	27.4 \pm 19.5	42.8 \pm 31.2	41.9 \pm 16.9	44.6 \pm 17.5	38.6 \pm 5.5	43.2 \pm 5.6	51.0 \pm 3.0	55.5 \pm 3.0
	Face	✓	60.8 \pm 16.0	31.5 \pm 17.0	52.8 \pm 9.8	43.9 \pm 4.5	46.3 \pm 5.5	—	—	50.5 \pm 2.6	54.5 \pm 2.6
	Speech	✓	61.8 \pm 12.8	26.6 \pm 17.6	43.8 \pm 34.5	40.9 \pm 16.0	43.1 \pm 16.1	—	—	51.7 \pm 2.6	56.2 \pm 2.5
	Concat	✓	68.6 \pm 8.6	24.9 \pm 19.9	39.7 \pm 19.5	42.2 \pm 12.3	43.2 \pm 14.2	42.3 \pm 5.7	47.0 \pm 5.6	49.8 \pm 4.1	54.2 \pm 4.0
MMC-PCFG	✓	67.9 \pm 9.8	52.3 \pm 9.0	63.5 \pm 8.6	55.0 \pm 3.7	58.9 \pm 3.4	44.7 \pm 5.2	48.9 \pm 5.7	56.0 \pm 1.4	60.0 \pm 1.2	
LC-PCFG	X	71.1 \pm 6.6	47.4 \pm 12.6	76.9 \pm 7.3	57.1 \pm 4.7	60.0 \pm 5.2	52.4 \pm 0.1	57.7 \pm 0.1	56.1 \pm 3.6	61.2 \pm 3.7	