🕐 NTT

Unsupervised Co-Segmentation for Athlete Movements and Live Commentaries Using Crossmodal Temporal Proximity

<u>Yasunori Ohishi</u>, Yuki Tanaka, and Kunio Kashino NTT Corporation, Japan



🕐 NTT







Model 🕐 NTT Guided attention scheme to efficiently detect and utilize temporal co-occurrences of audio and video information Existing approaches (Baseline) Visual feature Spatial and Video Dot product temporal (ECO) pooling 32 video frames Time axis (10-second video) Temporal information is • averaged or discarded. Audio feature Similarity Audio Tempora network pooling DAVEnet Mel spectrogram (10-second audio) Time axis

Dataset		🕐 NTT		
170 hours of NHK broadcast of grand sume tournaments	Winning techniques	Training	Validation	
grand sumo tournaments	Frontal push out	365	10	
 1,218 matches of nine frequent 	Frontal force out	362	10	
winning techniques	Slap down	141	10	
	Thrust down	77	10	
 IU-second video clips and their 	Over arm throw	45	10	
raw audio waveforms centered	Frontal thrust out	42	10	
around labeled times as audio-	Frontal crush out	34	10	
visual pairs	Rear push out	34	10	
	Frontal push down	28	10	
		1,128	90	
10-second video 10-second audio			5	

Audio-visual retrieval recall scores when the correct result was defined as the clips with the same winning techniques as the query

σ_g	Audio to Video			Video to Audio		
	R@1	R@3	R@5	R@1	R@3	R@5
0.001	.289	.600	.739	.294	.611	.717
0.01	.348	.656	.770	.304	.604	.785
0.1	.304	.648	.763	.307	.581	.733
1	.289	.600	.711	.211	.511	.622
10	.211	.461	.611	.144	.389	.561
100	.122	.389	.511	.056	.211	.411
Baseline	.256	.422	.589	.233	.511	.633



