ANTICIPATING ACTIVITY FROM MULTIMODAL SIGNALS

T. Rotondo¹, G. M. Farinella¹, D. Giacalone², S. M. Strano², V. Tomaselli², S. Battiato¹

¹Department of Mathematics and Computer Science, University of Catania, Italy ² STMicroelectronics, Catania, Italy

<u>tiziana.rotondo@unict.it, {gfarinella, battiato}@dmi.unict.it,</u> <u>{davide.giacalone, mauro.strano, valeria.tomaselli}@st.com</u>







Outline

- Problem definition
- ST Multimodal Dataset
- Proposed Pipeline and Results
- Conclusions and Future work

Problem Definition

 \circ Given the features vector x_t at time t as input, the goal is to predict the label of the next action, observing only data before the activity starts.



ST Multimodal Dataset -Modalities

Mobile phone

- Video
 - Camera Resolution: 720x1280
 - Frame rate: 29.95 fps
- Bluecoin:
 - Audio
 - Sampling rate: 32KHz
 - Acceleration, Gyroscope, Pressure, Temperature, Magnetic Field
 - Sampling rate: 52.63 Hz



ST Multimodal Dataset -Activities

• Activities: Desk, Reading, Sitting, Stairs, Standing, Typing, Walking.



ST Multimodal Dataset - # of sequences

Activity	# of sequences
Desk	2026
Reading	824
Typing	824
Walking	1353
Sitting	405
Stairs	420
Standing	405

"Desk" activity is more represented than "Stairs",
 "Sitting" or "Standing" activities because "Desk" is the past of "Typing", "Reading" and "Sitting" and the future of "Typing", "Reading" and "Standing".

• The collected dataset is balanced, i.e. the same number of sequences for each transition was acquired.

	\mathbf{Desk}	Reading	Typing	Walking	Stairs	Sitting	Standing
Desk		418	411				419
Reading	411		418				
Typing	418	411					
Walking					428	419	
Stairs				428			
Sitting	419						
Standing				419			

Cut- # of samples

- Time: 1,2160 sec
- Sensor: 64 samples

PAST FUTURE

WALKING STAIRS

- Video: 36 frames
- Audio: 32768 samples

The dataset contains 4874 transitions.

Pipeline



Baseline-SVM

	Classif	ication	Antici	pation
	Linear SVM	Rbf SVM	Linear SVM	Rbf SVM
Video	<u>66.70%</u>	<u>70.47%</u>	<u>61.35%</u>	<u>66.49%</u>
Audio	31.62%	35.49%	32.52%	34.26%
Sensors	46.27%	58.44%	43.89%	53.18%
Video and Sensors	<u>69.37%</u>	<u>72.89%</u>	<u>63.89%</u>	<u>70.57%</u>
Video and Audio	67.27%	70.72%	61.52%	67.58%
Audio and Sensors	46.29%	58.83%	43.03%	54.04%
Video Audio and Sensors	<u>69.55%</u>	<u>73.75%</u>	<u>64.06%</u>	<u>70.29%</u>

Baseline-K-NN

	Classification				Anticipation					
	K=1	K=3	K=5	K=7	K=9	K=1	K=3	K=5	K=7	K=9
Video	<u>60.02%</u>	<u>59.81%</u>	<u>59.51%</u>	<u>59.53%</u>	<u>59.61%</u>	<u>61.78%</u>	<u>61.31%</u>	<u>60.57%</u>	<u>60.45%</u>	<u>59.73%</u>
Audio	29.75%	30.92%	31.54%	32.30%	32.15%	30.94%	31.56%	32.44%	32.42%	32.79%
Sensors	52.44%	52.58%	54.16%	53.42%	53.81%	47.82%	48.28%	48.71%	49.30%	49.82%
Video and Sensors	<u>62.95%</u>	<u>62.66%</u>	<u>62.40%</u>	<u>61.72%</u>	<u>61.58%</u>	<u>64.18%</u>	<u>63.15%</u>	<u>63.09%</u>	<u>62.32%</u>	<u>61.87%</u>
Video and Audio	60.33%	60.00%	59.96%	59.76%	59.61%	62.05%	61.72%	60.72%	60.08%	60.29%
Audio and Sensors	54.49%	55.12%	55.57%	55.27%	54.63%	51.29%	52.11%	52.29%	52.19%	51.02%
Video Audio and Sensors	<u>63.24%</u>	<u>62.79%</u>	<u>62.42%</u>	<u>61.97%</u>	<u>61.76%</u>	<u>64.12%</u>	<u>63.61%</u>	<u>63.14%</u>	<u>62.27%</u>	<u>62.15%</u>

Triplet Network-SVM

Classif	ication	Anticipation				
Base	eline	Baseline Triplet			olet	
Linear Kernel	RBF	Linear Kernel RBF		Linear Kernel	RBF	
69.55%	73.75%	64.06%	70.29%	57.52%	63.32%	

Triplet Network-K-NN

	Classification	Predi	ction
К	Baseline	Baseline	Triplet
K=1	63.24%	64.12%	<u>64.65%</u>
K=3	62.79%	63.61%	<u>64.73%</u>
K=5	62.42%	63.14%	<u>64.14%</u>
K=7	61.97%	62.27%	<u>64.55%</u>
К=9	61.76%	62.15%	<u>64.18%</u>

Conclusions

 Our results suggest that multi-modality improves both classification and prediction.

 Considered activities can be anticipated with an accuracy close to the one obtained when the signals are fully observed (i.e., classification task).

 Future works could be devoted to collect bigger labelled multimodal datasets considering different environments and activities, as well as to model attention mechanisms among the different modalities

Thank you