

1 Supplementary material

2 This supplementary material presents the baseline results for *fear* detection based on the intelligence engine proposed in [1] by
 3 using the physiological and speech data publicly available in WEMAC. In this baseline, 88 volunteers are employed and 12 of
 4 them are discarded due to the high unbalance of the fear labels' distribution, following a similar criterion as in [1]. Note that
 5 the fear detection results in [1] comprise the use of only 42 volunteers of the WEMAC database.

6 Detailed Results of Fear Classification

7 Based on the same mono-modal and data fusion architectures presented in [1], the same time arrangements used for the
 8 alignment of the physiological and speech signals (Bindi 1.0, Bindi 2.0a, Bindi 2.0b) are employed in this case. Figure 2
 9 represents the F1-score (1) results by considering the 88 volunteers in WEMAC for the mono-modal fear detection systems
 10 (only based on physiological or speech data) in addition to the fusion strategies for the merging of both modalities. Both
 11 performance metrics are given to get more interpretable results from the slight imbalance between the binary *fear* labels.

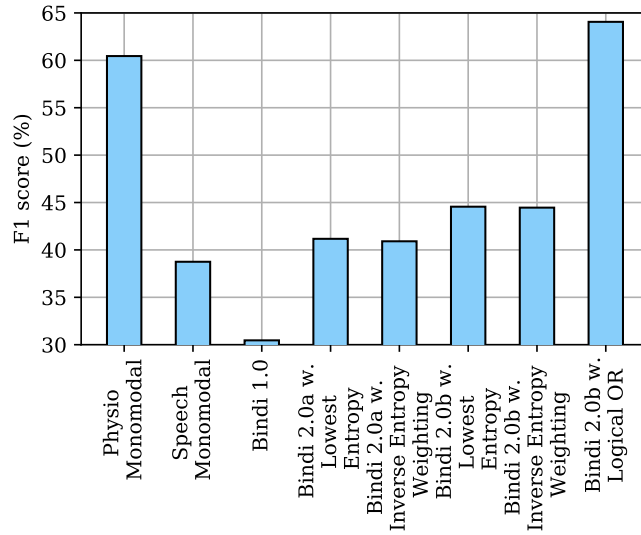


Figure 1

Figure 2. F1 score average performance analysis predicting over the 87 testing groups for the different architecture configurations. From left to right, the configurations are: physiological monomodal subsystem, the speech monomodal subsystem, Bindi 1.0, Bindi 2.0a with lowest entropy data fusion, Bindi 2.0a with inverse entropy weighting data fusion, Bindi 2.0b with lowest entropy data fusion, Bindi 2.0b with inverse entropy weighting data fusion, and Bindi 2.0b with logical OR data fusion. Note that Bindi 2.0a was not combined with logical OR data fusion because it is equivalent to Bindi 1.0.

		Physiological Monomodal	Speech Monomodal	BINDI 1.0	Bindi 2.0a Lowest Entropy	Bindi 2.0a Inverse Entropy Weighting	Bindi 2.0b Lowest Entropy	Bindi 2.0b Inverse Entropy Weighting	Bindi 2.0b Logical OR
F1-score	Mean	60.45	38.75	30.46	41.17	40.91	44.56	44.46	64.06
	Std	17.70	29.26	30.06	29.47	29.60	29.95	30.01	14.94

Table 1. Average performance analysis predicting over the 88 testing groups. Mean and standard deviations (Std).

12 Compared to the results in [1], we now double the amount of user data, adding 45 volunteers in these experiments. We use
 13 a *LASO* (Leave hAlf Subject Out) approach in which we train 88 models, one per each group or volunteers, using as fine-tuning
 14 data half of the data belonging to each user, and using as blind testing data the other half. Note that this is done with the
 15 intention of developing a general *fear* detection model personalized to each particular user.

16 As for the performances, these results are very similar to the ones achieved in [1], even slightly lower in some cases.
 17 This leads us to think that the addition of more data by doubling the number of users is still far from achieving higher and

18 more reliable rates for the detection of *fear*, leaving the door open for the research community to test new fusion methods,
19 personalisation strategies for each user, study the correlation of temporal alignment of the two data modalities available, and
20 test other suitable methodologies for this particular dataset.

21 **References**

22 **1.** Miranda, J. A. *et al.* Bindi: Affective internet of things to combat gender-based violence. *IEEE Internet Things J.* (2022).