

Self-Supervised Multimodal Sensor Fusion for Human Activity Recognition



Robert Piechocki¹, Kevin Chetty², N.D. Lane³, Muhammad J. Bocus¹

¹ School of Computer Science, Electrical and Electronic Engineering, and Engineering Maths, University of Bristol

² Department of Security and Crime Science, University College London

³ Department of Computer Science & Technology, University of Cambridge



1. Introduction

- Sensors such as Wi-Fi are particularly promising for in-home healthcare applications.
- Most passive Human Activity Recognition (HAR) deep learning-based systems are unimodal, i.e., they use information from only one type of sensor.
- Propose to use multiple synchronised sensors, or views, to improve the performance of a passive HAR system.
- Learn shared representation across different data types in a fully self-supervised manner, thus avoiding the need to label a huge amount of data, which is time-consuming and expensive.

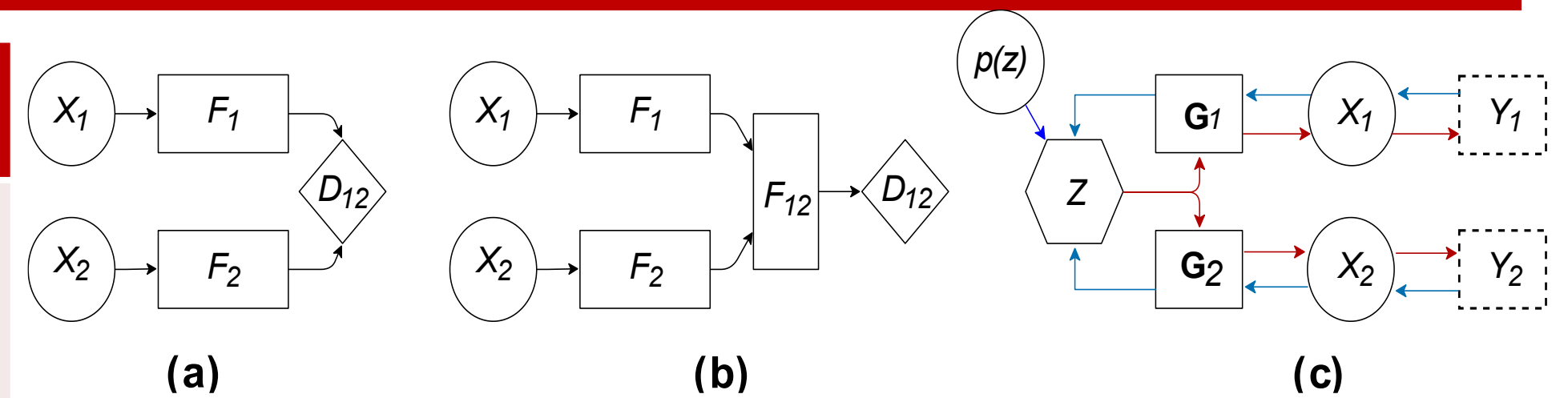


Figure 1. Multimodal Sensor Fusion: (a) Decision fusion, (b) Feature fusion, (c) Our technique: fusion in the latent representation with optional compressed sensing measurements; F features, $p(z)$ prior model, G generators, X complete data, Y subsampled data.

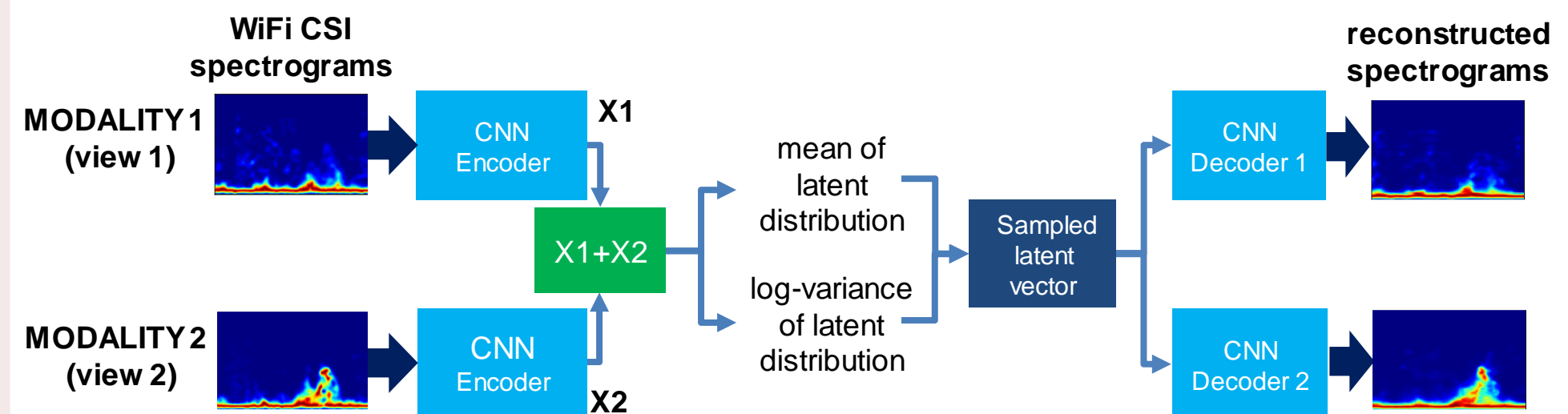


Figure 2. Multimodal Variational Autoencoder (M-VAE)

2. Results - Multimodal Sensor Fusion in the Latent Representation Space

Algorithm 1 Multimodal Sensor Fusion in the Latent Representation Space (SFLR)

```

1: Training data:  $\mathcal{D}_T \equiv \{X_{1:M}^{(1:T)}\}$ , Test data  $\mathcal{D}_T \equiv \{X_{1:M}^{(1:T)}\}$ , Samplers  $\{X_{1:M}\}$ 
2: Stage 1: Train M-VAE using  $\mathcal{D}_T$  Full data in first stage
3: Output:  $p(z)$ , Encoders  $\{\phi_{1:M}\}$ , Decoders  $\{\psi_{1:M}\}$ 
4: Stage 2: Fusion
5:  $y_{1:M}^{(i)} \sim \mathcal{D}_T$  subsampled data
6: Sample the initial point  $z^0 \sim p(z)$  Sample from prior
7: while not converged do
8:    $z \leftarrow z - \eta_0 \nabla_z (\|z\|^2) - \eta_1 \nabla_z (\|y_1^{(i)} - \chi_1(\psi_1(z))\|^2)$  One or several
9:    $z \leftarrow z - \eta_0 \nabla_z (\|z\|^2) - \eta_2 \nabla_z (\|y_2^{(i)} - \chi_2(\psi_2(z))\|^2)$  SGD steps are
10:    $\vdots$  taken for each
11:    $z \leftarrow z - \eta_0 \nabla_z (\|z\|^2) - \eta_M \nabla_z (\|y_M^{(i)} - \chi_M(\psi_M(z))\|^2)$  modality in turn.
12: end while
13:  $\hat{z}_{MAP} \leftarrow z$ 
14: Downstream tasks:  $\hat{x}_m = \psi_m(\hat{z}_{MAP})$ , classification tasks  $K$ -NN( $\hat{z}_{MAP}$ )
    
```

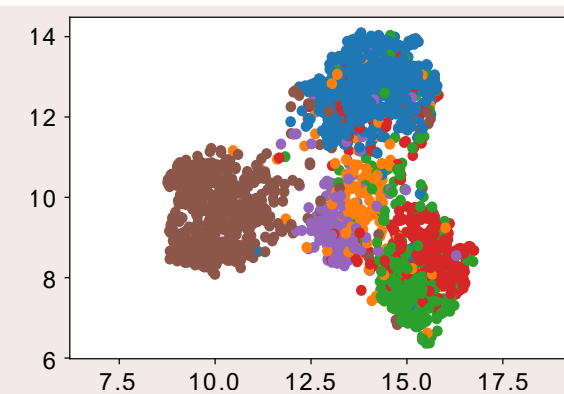


Figure 3. Trained Latent Space (6 clusters = 6 activities)

Table 1. Few-shot learning classification results (F1 macro) for HAR

	1 example per class	5 examples per class	10 examples per class
2-channel CNN	0.4273	0.5709	0.6185
1-channel CNN (Modality 1)	0.3491	0.4513	0.5045
1-channel CNN (Modality 2)	0.4466	0.6000	0.6057
Probability fusion (product rule)	0.4404	0.5847	0.6419
Dual-branch CNN	0.5082	0.5688	0.5759
SFLR (ours)	0.6527	0.7182	0.7375

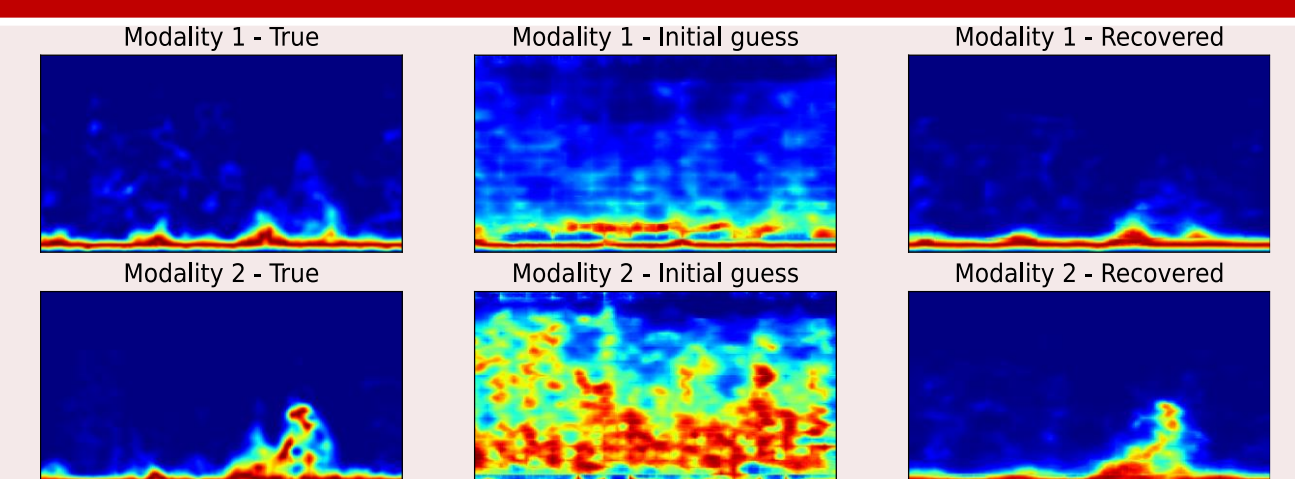


Figure 4. Recovery with compressed sensing measurements as low as 784 out of 50,176 (1.56%).

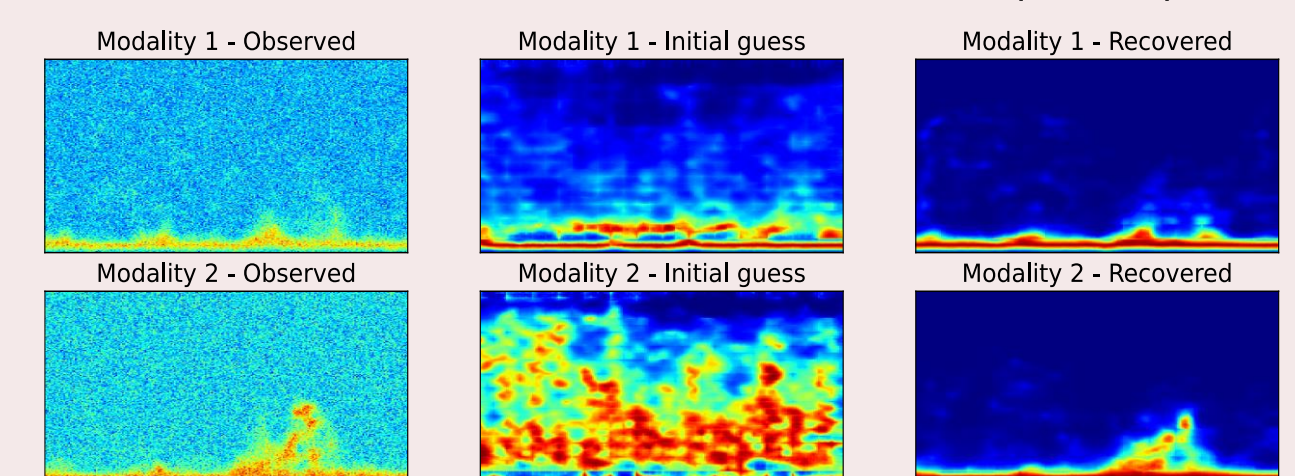
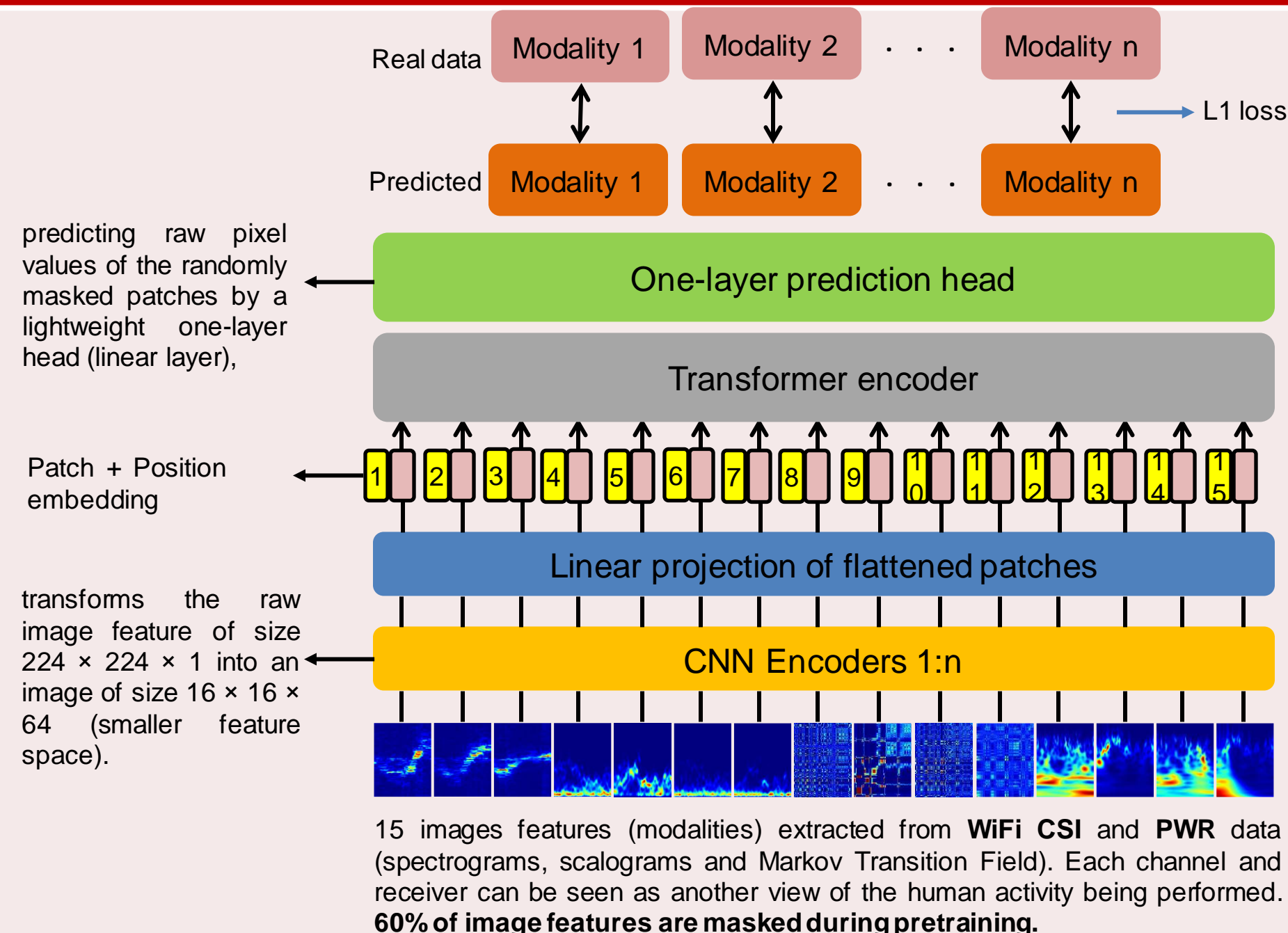


Figure 5. Recovery under additive Gaussian noise (std dev. = 0.4).

3. Results - Multimodal Fusion Transformer for Passive HAR



1. mask 60% of image features and pre-train model to predict the masked image features.
2. pretrain model with both PWR and CSI spectrograms, using all different views and image-based features
3. fine-tune the pretrained model in a supervised way using few labelled examples.

	F1-score for different amount of labelled samples in training set				
	1 sample	5%	10%	20%	Full
Fusion Transformer with SSL	56.3%	84.5%	89.7%	91.2%	95.9%
ResNet-34 (CNN)	32.6%	56.9%	62.7%	73.8%	94.9%

4. Conclusions

- Multimodal sensor fusion brings improved performance for downstream tasks such as Human Activity Recognition (HAR), which serves a vital role in the E-Health paradigm.
- Pretraining models such as Multimodal Variational Autoencoder (MVAE) and multimodal Vision Transformer (ViT) in a self-supervised fashion outperforms non-pretrained models under few-shot learning (i.e. under the condition that few labelled samples are available).

References

- [1] Bocus, M.J., Li, W., Vishwakarma, S. et al. OPERAnet, a multimodal activity recognition dataset acquired from radio frequency and vision-based sensors. *Sci Data* 9, 474 (2022).
- [2] Piechocki, R. J., Wang, X. and Bocus, M. J., "Multimodal sensor fusion in the latent representation space", 2022, <https://arxiv.org/abs/2208.02183>.
- [3] Koupai, A. K., Bocus, M. J., Santos-Rodríguez, R., Piechocki, R. J. and McConville, R., "Self-Supervised Multimodal Fusion Transformer for Passive Activity Recognition", 2022, <https://arxiv.org/abs/2209.03765>.



Engineering and Physical Sciences Research Council

