

On the Components That Enable Robust Generalization in HAR Models

Otávio Napoli and Edson Borin *

Institute of Computing - State University of Campinas
Av. Albert Einstein, 1251 - Cidade Universitária, Campinas - SP, 13083-889 - Brazil

Abstract. Generalizing human activity recognition (HAR) models across datasets remains challenging due to variations in sensors, environments, and user behavior. Domain Generalization (DG) methods attempt to address these shifts through objective-level modifications, architecture-level augmentations, and model pretraining strategies, but prior HAR studies often evaluate these components in isolation using suboptimal baselines. We systematically assess the contribution of each DG component across multiple HAR architectures, from CNNs to Transformers, using the DAGHAR benchmark. Our results show that pretraining the model with a self-supervised learning technique provides the most substantial and consistent gains in cross-dataset generalization, while architecture-level augmentations offer complementary improvements, and objective-level methods alone yield limited benefits across architectures. This suggests that DG studies should treat model pretraining as a standard baseline rather than an optional enhancement.

1 Introduction

Human activity recognition (HAR) using smartphone sensors has advanced with deep learning, yet most models are still evaluated within a single dataset using random splits or cross-validation, assuming that train and test samples share the same distribution (*i.i.d.* assumption). In real-world scenarios, however, models frequently encounter distribution shifts arising from changes in sensor placement, device type, user behavior, and environmental conditions [1]. Figure 1 illustrates these shifts with *t*-SNE projections from the DAGHAR benchmark [2], which contains six datasets from different domains (collected with different devices, body positions, and protocols). The plots highlight cross-position and cross-device differences, as well as the more severe cross-dataset shift combining both factors and additional variability in collection protocols and demography.

Domain Generalization (DG) methods aim to improve robustness to such shifts by training on multiple source domains and evaluating on unseen targets [3]. DG techniques can be grouped into three families: (a) objective-level methods, which modify the optimization objective, while treating the model as

*This project was supported by the Ministry of Science, Technology, and Innovation of Brazil, with resources granted by the Federal Law 8.248 of October 23, 1991, under the PPI-Softex. The project was coordinated by Softex and published as Intelligent agents for mobile platforms based on Cognitive Architecture technology [01245.003479/2024-10]. The authors also thank CNPq (315399/2023-6) and Fapesp (2013/08293-7) for their financial support, and Discovery Laboratory for their computational resources.

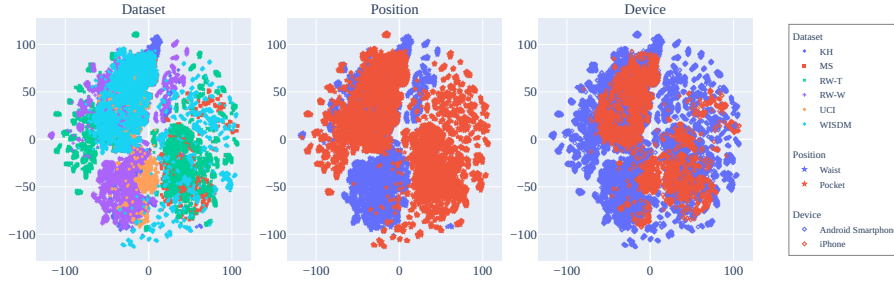


Fig. 1: t -SNE projections of DAGHAR samples in the frequency domain. The same projection is shown across panels, colored by dataset, body position, and device type to highlight different sources of domain shift. Clear cluster separation indicates strong distribution differences.

a black box; (b) architecture modifiers, which insert lightweight modules into models to perturb intermediate representations; and (c) initialization strategies, typically based on Self-Supervised Learning (SSL) pretraining, which provide more robust starting points. Although complementary, these components are often evaluated in isolation in prior HAR work [1, 4], leading to conclusions drawn from suboptimal baselines [5].

In this work, we systematically assess the contribution of each DG component to cross-dataset generalization across multiple HAR architectures, from CNNs to Transformers. Using the DAGHAR benchmark under a leave-one-dataset-out setup, we show that SSL-based initialization is one of the most important component in DG performance and should be treated as a standard component of evaluation. Feature-level augmentations provide complementary improvements, whereas objective-level methods alone offer limited benefits.

2 Domain Generalization Problem

DG aims to learn predictors that remain reliable under distribution shifts. Let \mathcal{X} and \mathcal{Y} denote the input and output spaces, and let a domain be a dataset $S = \{(x_i, y_i)\}_{i=1}^n$ drawn from a joint distribution $P_{\mathcal{X}\mathcal{Y}}$. In DG, the training set S_{train} contains M source domains $\{S^i\}_{i=1}^M$, each drawn from a distinct distribution $P_{\mathcal{X}\mathcal{Y}}^i$. The goal is to learn a predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ that generalizes to an unseen target domain $S_{\text{test}} \sim P_{\mathcal{X}\mathcal{Y}}^u$, whose distribution is unavailable during training ¹.

A model (predictor) consists of a feature extractor f_θ and a classifier g_ϕ , combined as $h(x) = g_\phi(f_\theta(x))$. Training minimizes a prediction loss \mathcal{L} . The standard baseline, Empirical Risk Minimization (ERM), optimizes this loss over all source domains. Although ERM does not explicitly enforce domain invariance, it remains a central baseline in DG studies [3], as training with data from

¹The label space \mathcal{Y} is assumed to be shared across domains.

diverse domains can implicitly encourage representations that generalize to unseen distributions, depending on data diversity and model capacity.

2.1 Domain Generalization Methods

DG methods aim to extract invariances shared across domains while suppressing domain-specific variation.²

Domain-Invariant Feature Exploration (DIFEX) [4] combines frequency-based guidance with cross-domain alignment. A teacher network processes the phase-only FFT of x to obtain invariant features z_T . The student representation $z = f_\theta(x)$ is split into (z_1, z_2) , where z_1 is guided toward z_T and z_2 is aligned across domains using CORAL [6]. The objective is

$$\mathcal{L}_{\text{DIFEX}} = \mathcal{L}_{\text{cls}} + \alpha \|z_1 - z_T\|_2^2 + \beta (-\|z_1 - z_2\|_2^2) + \lambda \mathcal{L}_{\text{CORAL}}(z_2), \quad (1)$$

encouraging frequency invariance, complementary representations, and second-order alignment.

MixStyle [7] perturbs intermediate features by mixing channel-wise statistics. For a feature map f with per-channel mean μ and std σ , and another sample (μ', σ') , a mixing coefficient λ sampled from a Beta distribution produces mixed statistics $\tilde{\mu} = \lambda\mu + (1-\lambda)\mu'$ and $\tilde{\sigma} = \lambda\sigma + (1-\lambda)\sigma'$, yielding the stylized feature

$$f_{\text{mix}} = \tilde{\sigma} \frac{f - \mu}{\sigma} + \tilde{\mu}.$$

During training, the non-trainable MixStyle layer is applied with probability p to features from selected layers, exposing the model to a broader range of domain styles and improving its robustness.

Time-Frequency Consistency (TF-C) [8] is an SSL method originally designed for representation learning, but also evaluated under DG conditions. TF-C enforces consistency between temporal and spectral views of the same input. Given $x \in \mathbb{R}^{C \times T}$ (C sensor axes and T time steps), augmentations generate $x_t = \mathcal{T}_{\text{time}}(x)$ and $x_f = \mathcal{T}_{\text{freq}}(\text{FFT}(x))$, where \mathcal{T} is a set of stochastic transformations. Two encoders G_t and G_f produce latent features $h_t = G_t(x_t)$ and $h_f = G_f(x_f)$, which are projected to embeddings $z_t = R_t(h_t)$ and $z_f = R_f(h_f)$.

TF-C combines three contrastive objectives: a time-view loss \mathcal{L}_t , a frequency-view loss \mathcal{L}_f , and a consistency loss aligning the two views,

$$\mathcal{L}_{\text{TF-C}} = \lambda(\mathcal{L}_t + \mathcal{L}_f) + (1 - \lambda)\mathcal{L}_{\text{consistency}}(z_t, z_f), \quad (2)$$

typically implemented with NT-Xent loss. The λ parameter balances view-specific and cross-view alignment. By aligning temporal and spectral structure, TF-C yields encoders that capture domain-stable information and substantially improve downstream HAR performance after fine-tuning.

²We assume the target distribution shares some structure with at least one source domain so that transferable features exist.

2.2 Limitations of Prior HAR DG Studies

Although DIFEX, MixStyle, and TF-C target complementary aspects of DG (objective-level, architecture-level, and initialization strategies, respectively), prior HAR studies typically evaluate them in isolation [1, 4], leading to incomplete conclusions and weak baselines [5]. To address this, we systematically assess each DG component, individually and in combination, across multiple HAR architectures using the DAGHAR benchmark, clarifying their contributions to domain generalization across HAR datasets.

3 Methodology

We follow a standardized evaluation protocol inspired by DomainBed [3], ensuring consistent batching, training, and model selection across all configurations.

Datasets Experiments use the DAGHAR benchmark [2], which unifies six public HAR datasets under a common format and split structure. We adopt a leave-one-dataset-out (LODO) strategy: in each run, one dataset is held out as the target domain, while the remaining five serve as sources. The train and validation splits of all source datasets form the training set, and their test splits are used for model selection (our validation set); the target dataset remains unseen until evaluation. We report performance across all target splits rather than only the test split as in prior work [2].

Training and Optimization Models are trained with domain-balanced mini-batches of 60 samples (12 per source domain). Training lasts 5000 steps using the ADAM optimizer with learning rates 10^{-3} and 10^{-4} , reporting the best result. The model with the highest validation accuracy is selected for evaluation. Each configuration is run with three random seeds. Unlike prior work [1, 2], which uses early stopping or domain truncation, we follow DomainBed’s sampling and model-selection procedures.

Models and DG Techniques We evaluate four representative HAR backbones [1, 2]: CNN-PFF, ResNet-SE-5, TS2Vec, and IMUTransformer, covering convolutional and transformer-based architectures. Each backbone is combined with DG methods spanning objective-level strategies (ERM, DIFEX), architecture modifiers (MixStyle), and initialization strategies (Random or TF-C). For DIFEX, we set $\alpha = 1$, $\beta = 1$, and $\lambda = 1$, following the authors’ reference implementation. For MixStyle, we apply the module after the first three convolutional layers (CNNs) or transformer blocks (Transformers) using a mixing probability of $p = 0.5$. For TF-C, each backbone is pretrained for 10000 steps on the combined source datasets, and the weights of the time-view encoder G_t is used to initialize the backbone of the target model. Pretraining excludes the held-out dataset, requiring six independent runs to maintain the LODO protocol.

| | CNN-PFF | | | | | IMU-Transformer | | | | | ResNet-SE-5 | | | | | TS2Vec | | | | | | | | | | | | |
|------------------|---------|------|------|------|------|-----------------|------|------|------|------|-------------|------|-------|------|------|--------|------|------|------|-------|------|------|------|------|------|------|-------|------|
| | KH | MS | RW-T | RW-W | UCI | WISDM | mean | KH | MS | RW-T | RW-W | UCI | WISDM | mean | KH | MS | RW-T | RW-W | UCI | WISDM | mean | KH | MS | RW-T | RW-W | UCI | WISDM | mean |
| DIFEX + MixStyle | 59.7 | 67.9 | 66.7 | 70.1 | 68.1 | 57.1 | 64.9 | 55.7 | 58.8 | 50.6 | 58.8 | 50.0 | 57.4 | 55.2 | 60.5 | 65.6 | 63.3 | 67.2 | 61.5 | 55.9 | 62.3 | 60.5 | 67.9 | 59.7 | 71.1 | 61.6 | 52.7 | 62.2 |
| Random Init. | +1.7 | +2.1 | +1.5 | +0.2 | +1.2 | +3.2 | +1.6 | +1.4 | +1.5 | +2.4 | +1.1 | +2.3 | +4.1 | +2.1 | +1.0 | +4.5 | +1.6 | +1.8 | +4.7 | +2.5 | +2.7 | +0.6 | +5.3 | +3.6 | +1.2 | +0.9 | +3.6 | +2.5 |
| DIFEX | 62.6 | 76.9 | 66.6 | 70.2 | 67.2 | 59.1 | 67.1 | 59.4 | 64.8 | 54.7 | 62.8 | 48.2 | 62.3 | 58.7 | 60.7 | 73.9 | 64.1 | 71.0 | 66.4 | 62.3 | 66.4 | 61.3 | 68.2 | 68.6 | 71.1 | 62.5 | 55.5 | 64.5 |
| TF-C Init. | +3.2 | +0.7 | +1.7 | +1.7 | +1.0 | +1.4 | +1.6 | +2.2 | +0.8 | +1.0 | +1.0 | +1.2 | +1.6 | +1.3 | +1.1 | +2.5 | +1.0 | +0.5 | +1.5 | +1.6 | +1.4 | +1.2 | +4.2 | +2.4 | +1.0 | +3.5 | +1.6 | +2.3 |
| DIFEX | 61.9 | 72.2 | 67.2 | 68.9 | 65.0 | 58.4 | 65.6 | 58.8 | 56.1 | 50.3 | 57.3 | 52.2 | 56.3 | 55.2 | 64.5 | 69.1 | 63.0 | 71.3 | 65.0 | 55.0 | 64.7 | 63.8 | 62.3 | 56.4 | 68.9 | 63.6 | 50.7 | 61.0 |
| Random Init. | +2.1 | +0.6 | +1.4 | +2.7 | +1.5 | +2.1 | +1.7 | +1.8 | +0.6 | +1.9 | +3.2 | +4.7 | +2.9 | +2.5 | +0.9 | +3.7 | +0.1 | +1.6 | +1.7 | +3.0 | +1.8 | +1.4 | +2.7 | +1.8 | +1.1 | +3.7 | +2.6 | +2.2 |
| DIFEX | 64.9 | 74.5 | 64.8 | 67.2 | 70.7 | 60.7 | 67.1 | 55.8 | 66.5 | 55.2 | 61.9 | 46.6 | 57.7 | 57.3 | 64.3 | 72.4 | 65.9 | 68.9 | 66.4 | 58.9 | 66.1 | 62.9 | 70.8 | 68.2 | 71.1 | 62.7 | 58.3 | 63.7 |
| TF-C Init. | +1.7 | +1.9 | +1.7 | +1.9 | +1.1 | +1.2 | +1.6 | +2.5 | +0.6 | +1.4 | +0.6 | +2.0 | +1.1 | +1.4 | +2.7 | +2.9 | +0.6 | +0.4 | +0.8 | +1.7 | +1.5 | +2.0 | +2.2 | +1.4 | +1.0 | +1.1 | +3.9 | +2.0 |
| ERM + MixStyle | 60.7 | 70.6 | 65.8 | 71.1 | 64.5 | 57.4 | 65.0 | 60.3 | 60.2 | 53.2 | 63.7 | 59.2 | 56.4 | 58.8 | 67.4 | 70.6 | 66.3 | 71.6 | 70.4 | 56.3 | 67.1 | 69.2 | 70.9 | 65.5 | 74.5 | 68.8 | 61.1 | 68.3 |
| Random Init. | +0.7 | +1.1 | +0.6 | +1.9 | +4.3 | +1.2 | +1.1 | +0.7 | +1.7 | +2.4 | +2.2 | +5.0 | +2.6 | +2.4 | +3.2 | +1.8 | +1.8 | +1.1 | +1.7 | +1.0 | +1.8 | +2.0 | +1.5 | +1.0 | +4.6 | +4.2 | +1.3 | +1.9 |
| ERM + MixStyle | 62.7 | 77.3 | 67.9 | 72.0 | 63.3 | 60.9 | 67.3 | 63.7 | 64.4 | 57.5 | 64.7 | 53.0 | 60.7 | 60.7 | 65.3 | 75.3 | 67.1 | 71.6 | 69.7 | 59.4 | 68.1 | 66.9 | 73.9 | 70.4 | 73.4 | 70.9 | 60.2 | 69.3 |
| TF-C Init. | +1.8 | +0.3 | +1.2 | +0.2 | +3.0 | +1.1 | +1.3 | +6.5 | +1.5 | +2.2 | +0.2 | +1.7 | +1.0 | +2.2 | +1.2 | +0.7 | +0.8 | +2.0 | +0.6 | +0.7 | +1.0 | +0.5 | +4.2 | +0.4 | +1.4 | +1.0 | +1.9 | +1.5 |
| ERM | 66.5 | 70.6 | 67.2 | 71.1 | 71.1 | 61.9 | 68.0 | 62.2 | 55.2 | 53.0 | 60.1 | 58.1 | 55.6 | 57.4 | 68.6 | 73.3 | 65.5 | 71.8 | 69.6 | 57.7 | 67.7 | 75.6 | 70.8 | 65.2 | 73.3 | 72.2 | 60.2 | 69.6 |
| Random Init. | +0.7 | +2.1 | +1.0 | +1.0 | +4.4 | +3.1 | +2.1 | +1.9 | +1.5 | +4.2 | +1.2 | +1.9 | +0.8 | +1.9 | +1.5 | +2.5 | +2.1 | +3.0 | +1.7 | +1.2 | +2.0 | +1.5 | +2.1 | +0.9 | +2.4 | +3.1 | +2.4 | +2.1 |
| ERM | 67.6 | 76.0 | 67.5 | 70.9 | 64.0 | 63.0 | 68.2 | 67.6 | 64.2 | 57.1 | 65.6 | 51.2 | 59.0 | 60.8 | 68.1 | 76.3 | 68.3 | 72.8 | 72.3 | 60.0 | 69.6 | 66.2 | 72.4 | 67.9 | 72.6 | 71.9 | 60.0 | 68.5 |
| TF-C Init. | +1.9 | +1.5 | +0.8 | +0.5 | +1.9 | +0.5 | +1.3 | +5.4 | +0.6 | +2.7 | +2.6 | +3.2 | +3.5 | +3.0 | +1.1 | +1.2 | +0.9 | +0.5 | +2.8 | +1.6 | +1.4 | +3.8 | +0.7 | +0.4 | +0.7 | +2.7 | +1.1 | +1.5 |

Fig. 2: Average accuracy (%) for each backbone (subplots), DG technique, and initialization strategy (rows) across all target datasets (columns) in the DAGHAR benchmark under a LODO protocol. Best results for each model-dataset pair are highlighted with a red dashed outline.

4 Experimental Results

Figure 2 summarizes the results across all backbones, DG techniques, and target domains. Some observations emerge: (i) TF-C pretraining provides the largest and most consistent gains, improving mean accuracy by almost 2% across all backbones. (ii) MixStyle yields some pontual improvements on CNN-PFF and TS2Vec but has limited overall effect, with a mean decrease of about 0.39%. (iii) DIFEX does not consistently outperform simple ERM, even when combined with TF-C or MixStyle, and on average decreases performance by roughly 3.14%.

These observations were further validated through Wilcoxon paired tests, comparing all combinations of DG techniques independent of backbone and dataset. Figure 3 shows these results, in which arrows denote statistically significant differences ($p < 0.05$), pointing from the higher to the lower-performing combination. The analysis confirms that TF-C pretraining yields significant improvements, that MixStyle is an optional complementary component and that DIFEX shows no significant advantage under any condition.

These findings suggest that robust SSL initialization is a primary component of DG in HAR and should be treated as a standard evaluation component rather than an optional enhancement, contrary to previous assumptions [1, 4].

5 Conclusion

This work presents a systematic evaluation of DG techniques for smartphone-based HAR across multiple backbones and datasets in the DAGHAR benchmark. Using a leave-one-dataset-out setup, we show that SSL initialization consistently provides the largest gains in cross-dataset generalization, outperforming objective-level methods such as DIFEX. Feature-level augmentations like MixStyle offer modest complementary improvements but cannot substitute for strong initialization. Thus, robust SSL initialization is a key component for

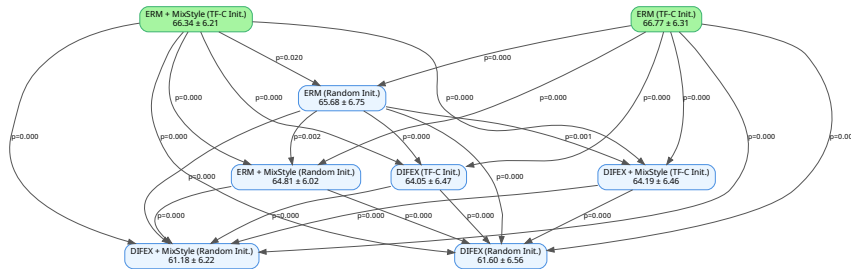


Fig. 3: Wilcoxon paired comparisons between DG technique combinations (nodes), independent of model and dataset. Arrows indicate significant differences ($p < 0.05$), pointing from higher- to lower-performing combinations. Green nodes indicate combinations with no incoming arrows (best ones).

DG in HAR and should be treated as a standard evaluation element rather than an optional enhancement [1, 4].

Although our work evaluates a selected set of DG methods, it spans a wider range of architectures than prior studies (which is typically limited to a single model or architecture type) and shows that these conclusions hold for both convolutional and transformer-based backbones. This foundation can guide future research on more effective DG strategies and support more reliable HAR systems in real-world applications such as healthcare and fitness monitoring.

References

- [1] O. Napoli and E. Borin. On domain generalization for human activity recognition with mix-based methods. In *Proceedings of 33th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pages 491–496, 2025.
- [2] O. Napoli, D. Duarte, P. Alves, D. H. P. Soto, H. E. de Oliveira, A. Rocha, L. Boccatto, and E. Borin. A benchmark for domain adaptation and generalization in smartphone-based human activity recognition. *Scientific Data*, 11(1):1192, 2024.
- [3] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International Conference on Learning Representations*, 2021.
- [4] Wang Lu, Jindong Wang, Haoliang Li, Yiqiang Chen, and Xing Xie. Domain-invariant feature exploration for domain generalization. *Transactions on Machine Learning Research*, 2022, 2022.
- [5] P. Teterwak, K. Saito, T. Tsiligkaridis, K. Saenko, and B. A Plummer. Erm++: An improved baseline for domain generalization. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 8525–8535. IEEE, 2025.
- [6] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *European conference on computer vision*, pages 443–450. Springer, 2016.
- [7] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang. Domain generalization with mixstyle. *arXiv preprint arXiv:2104.02008*, 2021.
- [8] X. Zhang, Z. Zhao, T. Tsiligkaridis, and M. Zitnik. Self-supervised contrastive pre-training for time series via time-frequency consistency. *Advances in neural information processing systems*, 35:3988–4003, 2022.