# Multi-Surrogate-Teacher Assistance for Representation Alignment in Fingerprint-based Indoor Localization
# Supplementary Material

Son Minh Nguyen[1], Linh Duy Tran[2], Duc Viet Le[1], Paul J.M Havinga[1]

[1]Department of Computer Science, University of Twente
[2]Viettel AI, Viettel Group

{m.s.nguyen, v.d.le, p.j.m.havinga}@utwente.nl
linhtd15@viettel.com.vn

## A. Hyperparemeter Selection

For the angular margin $\alpha$, we searched within a range from 0 degrees (0 radians) to 30 degrees (0.52 radians) on UIJIndoorLoc, with a step size of 0.1 radians, to determine the optimal angular margin. Our experimentation revealed that an angular margin of 0.2 radians consistently yielded stable results. Higher values of the angular margin $\alpha$ were not recommended, as they led to significant increases in model loss during training, making convergence extremely difficult.

The weighting factors mentioned are also determined through a grid search on UIJIndoorLoc and subsequently applied to other datasets. This process is expedited by leveraging prior knowledge that the primary task of indoor localization (denoted as $J_{MAE}$) should be accorded greater emphasis, necessitating a larger range and higher amplitude for the weighting factor $\lambda_1$. Conversely, the high sensitivity of $J_{FI}$ to the learning ability necessitates a much narrower range for the weighting factor $\lambda_4$. To ensure balanced impacts, these weighting factors are normalized together.

The rationale behind conducting grid-search hyperparameter selection exclusively on the UJIIndoorLoc dataset before applying the chosen hyperparameters to other datasets lies in UJIIndoorLoc's notable generalizability. This dataset has a collection period spanning months and encompasses expansive campus coverage across three buildings, totaling 110,000 square meters, effectively capturing the dynamic nature of real-world environments. During the grid search process, we set specific search ranges for each hyperparameter, such as [1-5] for $\lambda_1$ with a step size of 0.2, and [0.1-1] for $\lambda_{2,3,4}$ with a step size of 0.1. Through this rigorous exploration of over 7 days, we identified that the set of values [3,0.5,0.5,0.5] yielded the best performance during testing on the UJIIndoorLoc dataset.

## B. Incompatibly with other knowledge-transfer in RSS-Fingerprint-based Indoor Localization

Compounded by the nature of radio propagation and multipath effects, the distinctive characteristics of RSS datasets, including variabilities in building structure, occupancy levels, and the arrangement and number of input WiFi anchors, give rise to uncompromising discrepancies in both appearance (input size) and content (locations). Unlike other data types such as images or text that easily achieve a common input size with minimal content alteration using standard interpolation techniques, RSS fingerprints cannot be resized as their arrangements of the disparate number of anchors are just unknown. Even if the input sizes were reluctantly synchronized, the content representing specific locations would undergo significant alterations, leading to deviations from the true data distribution. Consequently, traditional domain adaptation approaches [1,6], such as meta-learning and adversarial learning, face limitations in their applicability to such datasets.

Meta-learning approaches optimize a common meta-learner for the target task through the learning abilities of its versions trained on sub-tasks. However, implementing this approach to achieve a unified meta-learner for different RSS fingerprint datasets presents challenges. These datasets often vary significantly in the number of anchors, with differences of hundreds observed between datasets. Additionally, standard resizing techniques cannot be applied due to the unique characteristics of RSS fingerprints.

Adversarial domain adaptation offers a potential solution to address differences in input size by employing separate feature generators for different datasets. However, this approach requires substantial modifications to existing architectures to ensure the delivery of homogeneous-sized features to domain discriminators for evaluation. Moreover,

blindly learning domain-invariant features solely through artificial coarse-grained domain labels, especially in an adversarial learning framework, is inadequate for capturing fine-grained information particularly essential for precise localization. Additionally, this method is susceptible to instability, including notorious issues of model collapse, where discriminators fail to keep track of distribution changes in generated data.

## C. Impacts of Target Relevance

The robustness of the framework is additionally evaluated on the target side where the target distribution is changed with the proportion of training data. Specifically, experiments are first carried out using only $1\%$ and $10\%$ of the training data for UJIIndoorLoc, and then expanded to $10\%$ for the other datasets for general examination. This random partitioning is designed to simulate real-world scenarios for which all the models are subjected to the same conditions, and repeated for ten rounds to achieve statistical results. As demonstrated in Table 1, the framework exhibits its tolerance to target constrictions and consistently empowers state-of-the-art models to achieve strong performance.

## D. Specific steps to the final $J_{MI}$ in Eq.4

Cross-Mutual Information Maximization Constraint $J_{MI}$ in Eq.4 can be represented by JS Divergence $D_{JS}$ in retrospect as follows.

$$D_{JS} = \mathbb{E}_{Z_{G_i}^E, Z_S^E \sim p\left(Z_{G_i}^E, Z_S^E\right)} \left(\log \frac{p_{Z_{G_i}^E, Z_S^E}}{m\left(Z_{G_i}^E, Z_S^E\right)}\right)$$
$$+ \mathbb{E}_{Z_{G_i}^E \sim p\left(Z_{G_i}^E\right), Z_S^E \sim p\left(Z_S^E\right)} \left(\log \frac{p_{Z_{G_i}^E} p_{Z_S^E}}{m\left(Z_{G_i}^E, Z_S^E\right)}\right)$$

where $m\left(Z_{G_i}^E, Z_S^E\right) = \frac{1}{2} p\left(Z_{G_i}^E\right) p\left(Z_S^E\right) + \frac{1}{2} p\left(Z_{G_i}^E, Z_S^E\right)$

$$\Rightarrow 2 D_{JS} =$$

$$\mathbb{E}_{Z_{G_i}^E \sim p\left(Z_S^E\right)} \left[ \begin{matrix} \mathbb{E}_{Z_{G_i}^E \sim p\left(Z_{G_i}^E | Z_S^E\right)} \left[ \begin{matrix} \log \frac{p(Z_{G_i}^E | Z_S^E)}{p(Z_{G_i}^E)} \\ -\log \left(1 + \frac{p(Z_{G_i}^E | Z_S^E)}{p(Z_{G_i}^E)}\right) \end{matrix} \right] \\ + \mathbb{E}_{Z_{G_i}^E \sim p\left(Z_{G_i}^E\right)} \left[ -\log \left(1 + \frac{p(Z_{G_i}^E | Z_S^E)}{p(Z_{G_i}^E)}\right) \right] \end{matrix} \right]$$

(1)

In addition, we make use of the mutual information estimator $\Psi_\theta$ [3] to estimate the logarithm ratio between $p\left(Z_{G_i}^E | Z_S^E\right)$ and $p\left(Z_{G_i}^E\right)$, represented by $\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right) = \log \frac{p(Z_{G_i}^E | Z_S^E)}{p(Z_{G_i}^E)}$. In the reverse direction, it can be interpreted as $\frac{p(Z_{G_i}^E | Z_S^E)}{p(Z_{G_i}^E)} = e^{\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right)}$. Put it all together in Eq.1, the JS Divergence is elaborated further, which is exactly the Mutual Information constraint $J_{MI}$ presented in this work:

$$D_{JS} = \mathbb{E}_{p\left(Z_{G_i}^E, Z_S^E\right)} \left[\log \frac{e^{\Psi_\theta\left(z_{G_i}^E, z_S^E\right)}}{1 + e^{\Psi_\theta\left(z_{G_i}^E, z_S^E\right)}}\right]$$
$$- \mathbb{E}_{p\left(Z_{G_i}^E\right) p\left(Z_S^E\right)} \left[\log \left(1 + e^{\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right)}\right)\right]$$
$$= \mathbb{E}_{p\left(Z_{G_i}^E, Z_S^E\right)} \left[\log \frac{1}{e^{-\Psi_\theta\left(z_{G_i}^E, z_S^E\right)} + 1}\right]$$
$$- \mathbb{E}_{p\left(Z_{G_i}^E\right) p\left(Z_S^E\right)} \left[\log \left(1 + e^{\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right)}\right)\right]$$
$$= \mathbb{E}_{p\left(Z_{G_i}^E, Z_S^E\right)} \left[-\log \left(1 + e^{-\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right)}\right)\right]$$
$$- \mathbb{E}_{p\left(Z_{G_i}^E\right) p\left(Z_S^E\right)} \left[\log \left(1 + e^{\Psi_\theta\left(Z_{G_i}^E, Z_S^E\right)}\right)\right]$$

(2)

## E. Specific steps expanded in Eq.9

We elaborate on the expression of the transmitting matrix $T_j$, which is mentioned in Eq.8 and is reduced to a simple form in Eq.9 as follows:

$$T_i \triangleq \left[\left(\widehat{Z}_{j-1}\right)^\Gamma \left(\widehat{Z}_j\right)\right]^\Gamma \left[\left(\widehat{Z}_{j-1}\right)^\Gamma \left(\widehat{Z}_j\right)\right]$$
$$= \left(W_j \widehat{Z}_{j-1}\right)^\Gamma \left(\widehat{Z}_{j-1}\right) \left(\widehat{Z}_{j-1}\right)^\Gamma \left(W_i \widehat{Z}_{j-1}\right)$$
$$= \left(W_j \widehat{Z}_{j-1}\right)^\Gamma \left(W_j \widehat{Z}_{j-1}\right)$$
$$= \widehat{Z}_{j-1}^\Gamma \left(W_j^\Gamma W_j\right) \widehat{Z}_{j-1}$$

(3)

## F. Power Iteration Algorithm

The proposed framework estimates the spectral norm of the blocks in neural networks using the Power Iteration Algorithm. This method is chosen for its lightweight computation and continuous differentiability. Here is the pseudocode for the algorithm:

---

**Algorithm 1:** Numerical Estimation of Spectral Norm with Tensorflow Pseudocode, called *top_eigenvalue*

---

**Input:** Transmitting Matrix $T_i$, power iteration $n$
**Output:** The Largest Eigenvalue $\sigma(\cdot)$

1   v = tf.random.normal([$T_i$.shape[0], $T_i$.shape[1]])
2   **for** $i = 0 \rightarrow n$ **do**
3     m = tf.matmul($T_i$, $v$)
4     $\mu$ = tf.sqrt( tf.reduce_sum(tf.square($m$), axis = 1))
5     $v = m/\mu$
6   $v\_norm$ = tf.sqrt(tf.reduce_sum(tf.square($v$), axis = 1))
7   $\sigma(T_i)$ = tf.sqrt($\mu/v\_norm$)
8   **return** $\sigma(T_i)$

---

Table 1. The impact of the target relevance to the source datasets on Expert Distilling phase

| Method | UJIIndoorLoc | | UTS | Tampere |
| | 1% (199/19937) | 10% (1993/19937) | 10% (910/9108) | 10% (69/697 ) |
| | MAE (m)↓ | MAE (m)↓ | MAE (m)↓ | MAE (m) ↓ |
|---|---|---|---|---|
| DNN [2] | 160.56±0.24 | 18.28±0.41 | 10.65±5.52 | 25.04±14.99 |
| DNN+++ | **36.36±2.32** | **17.42±0.45** | **8.18±0.68** | **24.28±10.87** |
| CNNLoc [8] | 22.14±1.09 | 14.07±0.27 | 7.34±0.16 | 15.99±9.58 |
| CNNLoc+++ | **21.93±1.17** | **13.73±0.36** | **6.96±0.34** | **12.78±1.09** |
| BayesCNN [7] | 26.84±2.21 | 16.25±0.47 | 8.53±0.55 | 15.14±1.25 |
| BayesCNN+++ | **25.3±1.93** | **15.68±0.72** | **8.05±0.46** | **15.10±0.9** |
| bAaT [4, 5] | 19.56±1.09 | 13.06±0.26 | 6.79±0.22 | 13.30±1.03 |
| bAaT++ (wo $J_{FI}$) | **18.18±0.86** | **12.88±0.18** | **6.58±0.18** | **12.33±0.89** |

## G. Pseudo code of the training pipeline

The framework executes the alignment through two primary phases. The first phase called *Expert Training*, involves modeling the representations established by the specialized networks on their respective source datasets using surrogate teacher networks. In the second phase, *Expert Distilling*, these modeled representations are collectively distilled into essential knowledge for alignment with representations learned on the target dataset.

**Algorithm 2:** Expert Training Phase

**Input:** Surrogate teachers
$$\{G_i\}_{i=1}^N = \{F_i^G, E_i^G\}_{i=1}^N,$$
Critics $\{C_i\}_{i=1}^N$, $c\_step$, Source dataset $D_{i=1}^N$,
Specialized models $\{S_i\}_{i=1}^N = \{F_i^S, E_i^S, R_i^S\}_{i=1}^N$,
Epoch $E$, gradient weight $\alpha$, loss weights $\{\beta_i\}_{i=1}^3$

**Output:** Pre-trained surrogate teachers $\{G_i\}_{i=1}^N$

**1** Initialize $T\_list$
**2 for** $i = 1 \to N$ **do**
**3**    **for** $epoch\ e = 1 \to E$ **do**
**4**      Load $D_i$
**5**      Initialize $G_i, C_i, S_i$
**6**      **for** $batch\ b \in D_i$ **do**
        /* First, training Critics for c steps */
**7**        **for** $k$ to $c\_step$ **do**
**8**          Initialize $Noise$
**9**          $Z_{S_i}^F, Z_{S_i}^E, \hat{Y}_{S_i} = S_i(b)$
**10**         $Z_{G_i}^F, Z_{G_i}^E = G_i(Noise)$
**11**         $r\_logits = C_i(Z_{S_i}^E)$
**12**         $f\_logits = C_i(Z_{G_i}^E)$
**13**         $gp = grad\_penalty(C_i, b, Z_{S_i}^E, Z_{G_i}^E)$
**14**         $L_c = critic\_loss(r\_logits, f\_logits)$
**15**         $L_{tc} = L_c + \alpha * gp$
**16**         $C_i = update(L_{tc}, C_i)$
        /* Then, training Generators and Specialized model */
**17**        Initialize $Noise$
**18**        $Z_{S_i}^F, Z_{S_i}^E, \hat{Y}_{S_i} = S_i(b)$
**19**        $Z_{G_i}^F, Z_{G_i}^E = G_i(Noise)$
**20**        $\hat{Y}_{G_i} = R_i^S(Z_{G_i}^E)$
**21**        $L_{S\_MAE} = J_{MAE}(\hat{Y}_{S_i}, Y_{S_i})$
**22**        $L_{G\_MAE} = J_{MAE}(\hat{Y}_{G_i}, Y_{S_i})$
**23**        $L_{Sim} = J_{Sim}(Z_{S_i}^E, Z_{G_i}^E)$
**24**        $L_{tg} = \beta_1 * L_{S\_MAE} + \beta_2 * L_{G\_MAE} + \beta_3 * L_{Sim}$
**25**        $G_i = update(L_{tg}, G_i)$
**26**        $S_i = update(L_{S\_MAE}, S_i)$
**27**    $T\_list.append(G_i)$
**28 return** $T\_list$

---

**Algorithm 3:** Expert Distilling Phase

**Input:** Surrogate teachers
$$\{G_i\}_{i=1}^N = \{F_i^G, E_i^G\}_{i=1}^N,$$
Target dataset $D_t$
Mutual Information Estimator $\Psi_\theta$,
Specialized models $S = \{F^S, E^S, R^S\}$,
Epoch $E$, loss weights $\lambda_{i=1}^4$

**Output:** Specialized models $S = \{F^S, E^S, R^S\}$

**1** Load $D_t$
**2** Initialize $\Psi_\theta, S$
**3 for** $e = 1 \to E$ **do**
**4**    **for** $b \in D_t$ **do**
**5**      Initialize $Noise, L_t^{MI}, L_t^{Sim}, L_t^{FI}$
**6**      $Z_S^F, Z_S^E, \hat{Y}_S = S(b)$
**7**      $L_{S\_MAE} = J_{MAE}(\hat{Y}_{S_i}, Y_{S_i})$
      /* Computing Functional Information in the specialized model **S** for comparison with other surrogate teachors */
**8**      $TM_S = transmitting\_matrix(Z_S^F, Z_S^E)$
**9**      $FI_S = top\_eigenvalue(TM_S)$
**10**      **for** $i = 1 \to N$ **do**
**11**        Load $G_i$
**12**        $Z_{G_i}^F, Z_{G_i}^E = G_i(Noise)$
        /* Computing Mutual Information Constraint */
**13**        $product\_examples = concat(Z_{S_i}^E[1:], Z_{S_i}^E[0])$
**14**        $joint\_stat = \Psi_\theta(Z_{S_i}^E, Z_{G_i}^E)$
**15**        $product\_stat = \Psi_\theta(product\_examples, Z_{G_i}^E)$
**16**        $L_i^{MI} = J_{MI}(joint\_stat, product\_stat)$
**17**        $L_t^{MI} += L_i^{MI}$
        /* Computing Angular Similarity Constraint */
**18**        $L_i^{Sim} = J_{Sim}(Z_{S_i}^E, Z_{G_i}^E)$
**19**        $L_t^{Sim} += L_i^{Sim}$
        /* Computing Functional Information Constraint */
**20**        $TM_{G_i} = transmitting\_matrix(Z_{G_i}^F, Z_{G_i}^E)$
**21**        $FI_{G_i} = top\_eigenvalue(TM_{G_i})$
**22**        $L_i^{FI} = J_{FI}(FI_S, FI_{G_i})$
**23**        $L_t^{FI} += L_i^{FI}$
**24**      $L_t = J_{overall}(\lambda_{i=1}^4, L_{S\_MAE}, L_t^{Sim}, L_t^{MI}, L_t^{FI})$
**25**      $S = update(L_t, S)$
**26 return** $S = \{F^S, E^S, R^S\}$

# References

[1] Abolfazl Farahani, Sahar Voghoei, Khaled Rasheed, and Hamid R Arabnia. A brief review of domain adaptation. *Advances in data science and information engineering: proceedings from ICDATA 2020 and IKE 2020*, pages 877–894, 2021. 1

[2] Gibrán Félix, Mario Siller, and Ernesto Navarro Alvarez. A fingerprinting indoor localization algorithm based deep learning. In *2016 eighth international conference on ubiquitous and future networks (ICUFN)*, pages 1006–1011. IEEE, 2016. 3

[3] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018. 2

[4] Son Minh Nguyen, Duc Viet Le, and Paul JM Havinga. Learning the world from its words: Anchor-agnostic transformers for fingerprint-based indoor localization. In *2023 IEEE International Conference on Pervasive Computing and Communications (PerCom)*, pages 150–159. IEEE, 2023. 3

[5] Son Minh Nguyen, Duc Viet Le, and Paul JM Havinga. Seeing the world from its words: All-embracing transformers for fingerprint-based indoor localization. *Pervasive and Mobile Computing*, 100:101912, 2024. 3

[6] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009. 1

[7] Shreya Sinha and Duc V Le. Completely automated cnn architecture design based on vgg blocks for fingerprinting localisation. In *2021 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pages 1–8. IEEE, 2021. 3

[8] Xudong Song, Xiaochen Fan, Chaocan Xiang, Qianwen Ye, Leyu Liu, Zumin Wang, Xiangjian He, Ning Yang, and Gengfa Fang. A novel convolutional neural network based indoor localization framework with wifi fingerprinting. *IEEE Access*, 7:110698–110709, 2019. 3