

Learning Signer-Invariant Representations with Adversarial Training

Pedro M. Ferreira^{a,b}, Diogo Pernes^{a,c}, Ana Rebelo^{a,d}, and Jaime S. Cardoso^{a,b}

^aCentre for Telecommunications and Multimedia, INESC TEC, 4200-465 Porto, Portugal

^bFaculdade de Engenharia da Universidade do Porto, 4200-465 Porto, Portugal

^cFaculdade de Ciências da Universidade do Porto, 4169-007 Porto, Portugal

^dUniversidade Portucalense, 4200-072 Porto, Portugal

ABSTRACT

Sign Language Recognition (SLR) has become an appealing topic in modern societies because such technology can ideally be used to bridge the gap between deaf and hearing people. Although important steps have been made towards the development of real-world SLR systems, signer-independent SLR is still one of the bottleneck problems of this research field. In this regard, we propose a deep neural network along with an adversarial training objective, specifically designed to address the signer-independent problem. Concretely speaking, the proposed model consists of an *encoder*, mapping from input images to latent representations, and two classifiers operating on these underlying representations: (i) the *sign-classifier*, for predicting the class/sign labels, and (ii) the *signer-classifier*, for predicting their signer identities. During the learning stage, the *encoder* is simultaneously trained to help the *sign-classifier* as much as possible while trying to fool the *signer-classifier*. This adversarial training procedure allows learning signer-invariant latent representations that are in fact highly discriminative for sign recognition. Experimental results demonstrate the effectiveness of the proposed model and its capability of dealing with the large inter-signer variations.

Keywords: Sign Language Recognition, Gesture Recognition, Adversarial Neural Networks, Adversarial Training, Deep Learning

1. INTRODUCTION

Sign language is a nonverbal form of communication especially used by hearing impaired people within deaf communities worldwide. It combines articulated hand gestures along with facial expressions to convey meaning. Contrary to the popular belief, sign language is not universal and, just like spoken languages, it has its own lexicon, syntax and grammar. This is why most of hearing people are unfamiliar with sign language, which obviously creates a serious communication barrier between deaf communities and the hearing majority.

As a key technology to help bridging the gap between deaf and hearing people, Sign Language Recognition (SLR) has become one of the most active research topics in the human-computer interaction field. Its main purpose is to automatically translate the signs, from video or images, into the corresponding text or speech. Although recent SLR methods have demonstrated remarkable performances in signer-dependent scenarios, i.e. when training and test data come from the same signers, their recognition rates typically decrease significantly when the signer is new to the system. This performance drop is the result of the large inter-signer variability in the manual signing process of sign languages (see Figure 1). However, a practical SLR system must operate in a signer-independent scenario, which means that the signer of the probe must not be seen during the training routine of the models. Therefore, signer-independent SLR has become one of the bottleneck problems for the development of a real-world and practical SLR system.

Borrowing from recent works on adversarial neural networks [3, 5] and domain transfer [4], we introduce a deep neural network along with a novel adversarial training objective to specially tackle the signer-independent SLR problem. The underlying idea is to preserve as much information as possible about the signs, while discarding the signer-specific information that is implicitly present in the manual signing process. For this purpose, the proposed deep model is composed by an *encoder* network, which maps from the input images to latent representations,



Figure 1. The inter-signer variability: it is possible to observe not only phonological variations (i.e., different handshapes, palm orientations, and sign locations) but also a large physical variability (i.e., different hand sizes) when six signers are performing the same sign.

as well as two discriminative classifiers operating on top of these underlying representations, namely the *sign-classifier* network and the *signer-classifier* network. While the *sign-classifier* is trained to predict the sign labels, the *signer-classifier* is trained to discriminate their signer identities. In addition, the parameters of the *encoder* network are optimized to minimize the loss of the *sign-classifier* while trying to fool the *signer-classifier* network. This adversarial and competitive training scheme encourages the learned representations to be signer-invariant and highly discriminative for the sign classification task. To further constrain the latent representations to be signer-invariant, we introduce an additional training objective that operates on the hidden representations of the *encoder* network in order to enforce the latent distributions of different signers to be as similar as possible.

Although this adversarial training framework is similar to those initially introduced by Ganin *et al* [4], in the context of domain adaptation, and then by Feutry *et al* [3] to learn anonymized representations, our main contributions on top of these works are two-fold:

- The application of the adversarial training concept to the signer-independent SLR problem;
- A novel adversarial training objective that differs from the ones of Ganin *et al* [4] and Feutry *et al* [3] in two ways. First, our training objective is minimum if and only if the adversarial classifier, which in our case corresponds to the *signer-classifier*, produces a uniform distribution over the signer identities, meaning that our model is completely invariant to the signer identity of the training data. Second, we introduce an additional term to the adversarial training objective that further discourages the learned representations of retaining any signer-specific information, by explicitly imposing similarity in the latent distributions of different signers.

The remainder of the paper is organized as follows. Section 2 presents the related work. The proposed model along with its adversarial training scheme are fully described in Section 3. Experimental results and conclusions are reported in Sections 4 and 5, respectively.

2. RELATED WORK

According to the amount of data required from the test signers, previous signer-independent SLR works can be roughly classified into two main groups, namely (i) signer adaptation approaches, where a previously trained model is adapted to a new test signer by using a small amount of signer specific data, and (ii) truly signer-independent methods, in which a generic model robust for new test signers is built without using data of those test signers.

Greatly inspired by speaker adaptation methods from the speech recognition research, Von Agris *et al* [21] proposed the combination of the eigenvoice (EV) approach [11] with maximum likelihood linear regression (MLLR) and maximum a posteriori (MAP) estimation to adapt trained Hidden Markov Models (HMMs) to new signers. More recently, Kim *et al* [8] investigated the potential of different deep neural network adaptation strategies for the signer-independence problem. Yin *et al* [25] proposed an interesting weakly-supervised signer adaptation approach, in which the adaptation data from the new signer has not to be labeled. Specifically, a generic metric is first learnt from the available labeled data of several different signers and, then, adapted to the new signer by considering clustering and manifold constraints along with the collected unlabeled data. Although signer adaptation is a reasonable approach, in practice, collecting enough training data from each new signer to retrain and adapt the model may not be feasible. In this regard, several works focused on the development of truly-signer independent models that do not require any data from the new signers [1, 7, 10, 18, 22, 24, 26]. Most

of them involved a huge feature engineering effort in order to build normalized hand-crafted feature descriptors robust to the large inter-signer variations. A major weakness across all the aforementioned methods is related to the fact that representation and metric learning is not jointly performed. Motivated by the inherent difficulty of designing reliable hand-crafted features to the large inter-signer variability, recent SLR systems are mostly based on deep neural networks [9, 12, 15, 16, 23]. It is well-known that deep neural networks are remarkably good in figuring out reliable high-level feature representations from the data. However, in previous deep SLR methodologies, the learned representations are not explicitly constrained to be signer-invariant. Therefore, there is nothing to prevent the learned representations of different signers and the same class of being far apart in the representation space and, hence, signer invariance is not ensured.

This paper presents a novel adversarial training objective, based on representation learning and deep neural networks, specifically designed to address the signer-independent SLR problem. Different from the aforementioned methodologies, our model jointly learns the representation and the classifier from the data, while explicitly imposing signer invariance in the high-level representations for a robust and truly signer-invariant sign recognition.

3. PROPOSED METHOD

The ultimate goal of our model is to learn signer-invariant latent representations that preserve the relevant part of the information about the signs while discarding the signer-specific traits that may hamper the sign classification task. To accomplish this purpose, we introduce a deep neural network along with an adversarial training scheme that is able to learn feature representations that combine both sign discriminativeness and signer-invariance.

More specifically, let $\mathbb{X} = \{\mathbf{X}_i, y_i, s_i\}_{i=1}^N$ denote a labeled dataset of N samples, where \mathbf{X}_i represents the i -th colour image, and y_i and s_i denote the corresponding class (sign) label and signer identity, respectively. To induce the model to learn signer-invariant representations, the proposed model comprises three distinct sub-networks:

- an *encoder* network, which aims at learning an encoding function $h(\mathbf{X}; \theta_h)$, parameterized by θ_h , that maps from an input image \mathbf{X} to a latent representation \mathbf{h} ;
- a *sign-classifier* network, which operates on top of this underlying latent representation \mathbf{h} to learn our task-specific function $f(\mathbf{h}; \theta_f)$, parameterized by θ_f , that maps from \mathbf{h} to the predicted probabilities $p(y|\mathbf{h}; \theta_f)$ of each sign class.
- a *signer-classifier* network, with the purpose of learning a signer-specific function $g(\mathbf{h}; \theta_g)$, parameterized by θ_g , that maps the same hidden representation \mathbf{h} to the predicted probabilities $p(s|\mathbf{h}; \theta_g)$ of each signer identity.

During the learning stage, the parameters of both classifiers are optimized in order to minimize their errors on their specific tasks on the training set. In addition, the parameters of the *encoder* network are optimized in order to minimize the loss of the *sign-classifier* network while forcing the *signer-classifier* of being a random guessing predictor. In the course of this adversarial training procedure, the learned latent representations \mathbf{h} are encouraged to be signer-invariant and highly discriminative for sign classification. To further discourage the latent representations of retaining any signer-specific traits, we introduce an additional training objective that enforces the latent distributions of different signers to be as similar as possible. The result is a truly signer-independent model robust to new test signers.

3.1 Architecture

As illustrated in Figure 2, the architecture of the proposed model is composed by three main sub-networks or blocks, i.e. an *encoder*, a *sign-classifier* and a *signer-classifier*.

The *encoder* network attempts to learn a mapping from an input image \mathbf{X} to a latent representation \mathbf{h} . It simply consists of a sequence of L_e pairs of consecutive 3×3 convolutional layers with Rectified Linear Units (ReLU) as non-linearities. For downsampling, the last convolutional layer of each pair has a stride of 2. On

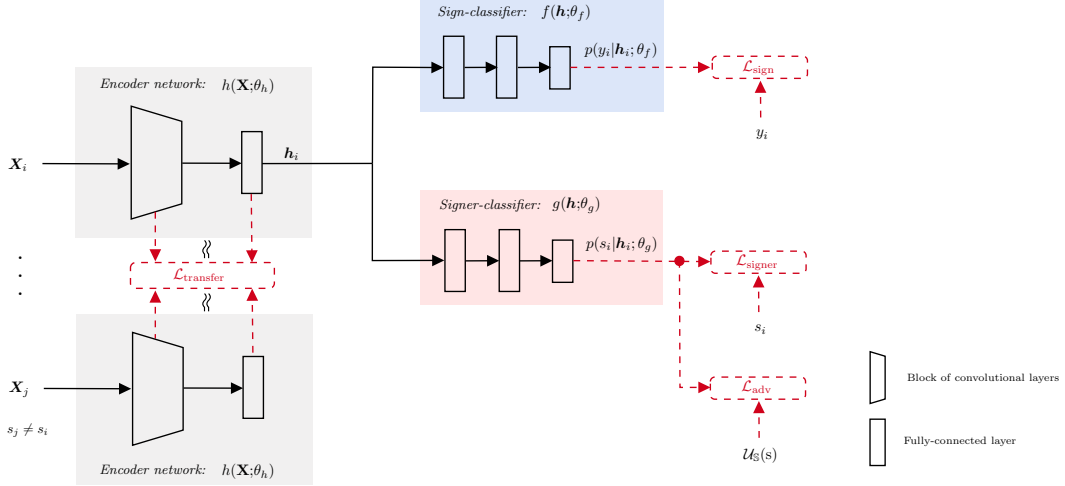


Figure 2. The architecture of the proposed signer-invariant neural network. It comprises three main sub-networks or blocks, i.e. an *encoder*, a *sign-classifier* and a *signer-classifier*.

top of that, there is a fully-connected layer, also with a ReLU, representing the desired signer-invariant latent representations \mathbf{h} .

Taking the latent representations \mathbf{h} as input, the *sign-classifier* block is composed by a sequence of L_s fully-connected layers, with ReLUs as the non-linear functions, for predicting the sign class $\hat{y} = \arg \max f(\mathbf{h}; \theta_f)$. Therefore, the last fully-connected layer has a softmax activation function which outputs the probabilities for each sign class.

The *signer-classifier* network has exactly the same topology as the *sign-classifier* net. However, it maps the latent representations \mathbf{h} to the predicted signer identity $\hat{s} = \arg \max g(\mathbf{h}; \theta_g)$. Therefore, the number of nodes of the output layer is defined accordingly to the number of signers in the training set.

3.2 Adversarial training

By definition, signer-invariant representations discard all signer-specific information and, as such, no function (i.e., classifier) exists that maps such representations into the correct signer identity. This naturally leads to an adversarial problem, in which: (i) a *signer-classifier* network $g(\cdot; \theta_g)$ receives latent representations $\mathbf{h} = h(\mathbf{X}; \theta_h)$ from an *encoder* network $h(\cdot; \theta_h)$ and tries to predict the signer identity s corresponding to image \mathbf{X} and (ii) the *encoder* network tries to fool the *signer-classifier* network while still providing good representations for the *sign-classifier* network $f(\cdot; \theta_f)$, which in turn receives the same representations \mathbf{h} and aims to predict the sign label y corresponding to image \mathbf{X} .

Therefore, the *signer-classifier* network shall be trained to minimize the negative log-likelihood of correct signer predictions:

$$\min_{\theta_g} \mathcal{L}_{\text{signer}}(\theta_h, \theta_g) = -\frac{1}{N} \sum_{i=1}^N \log p(s_i | h(\mathbf{X}_i; \theta_h); \theta_g) \quad (1)$$

In the perspective of the *encoder*, the predictions of the *sign-classifier* should be as accurate as possible and the predictions of the *signer-classifier* should be kept close to uniform, meaning that this latter classifier is not capable of doing better than random guessing the signer identity. Formally, this may be translated into the following constrained objective:

$$\min_{\theta_h, \theta_f} \mathcal{L}_{\text{sign}}(\theta_h, \theta_f) = -\frac{1}{N} \sum_{i=1}^N \log p(y_i | h(\mathbf{X}_i; \theta_h); \theta_f), \quad (2)$$

$$\text{subject to } \frac{1}{N} \sum_{i=1}^N D_{\text{KL}}(\mathcal{U}_{\mathbb{S}}(s) || p(s | h(\mathbf{X}_i; \theta_h); \theta_g)) \leq \epsilon, \quad (3)$$

where D_{KL} is the Kullback-Leibler (KL) divergence and $\mathcal{U}_{\mathbb{S}}(s)$ denotes the discrete uniform distribution on the random variable s , defined over the set of identities \mathbb{S} in the training set. Here, $\epsilon \geq 0$ determines how far from uniform the *signer-classifier* predictions are allowed to be (as measured by the KL divergence). The choice of the uniform distribution implies the underlying assumption that the training set is balanced relatively to the number of examples per signer (which should be true for most practical datasets). When this is not the case, the empirical distribution of signer identities in the training set may be used instead.

The constraint inequality (3) may be rewritten as:

$$\mathcal{L}_{\text{adv}}(\theta_h, \theta_g) = -\frac{1}{N|\mathbb{S}|} \sum_{i=1}^N \sum_{s \in \mathbb{S}} \log p(s|h(\mathbf{X}_i; \theta_h); \theta_g) \leq \epsilon + \log |\mathbb{S}|, \quad (4)$$

and the constrained optimization problem may be equivalently formulated as:

$$\min_{\theta_h, \theta_f} \mathcal{L}(\theta_h, \theta_f, \theta_g) = \mathcal{L}_{\text{sign}}(\theta_h, \theta_f) + \lambda \mathcal{L}_{\text{adv}}(\theta_h, \theta_g), \quad (5)$$

where $\lambda \geq 0$ depends on ϵ and \mathcal{L}_{adv} plays the role of an adversarial loss with respect to the signer classification loss $\mathcal{L}_{\text{signer}}$.

This objective and the structure of our model are similar to those used in [4], in the context of domain adaptation, and in [3], to learn anonymized representations for privacy purposes. However, the former uses the negative signer classification loss as the adversarial term (i.e., $\mathcal{L}_{\text{adv}} \leftarrow -\mathcal{L}_{\text{signer}}$), which is not lower bounded, leading to high gradients and difficult optimization. The latter addresses this problem by replacing this term with the absolute difference between the adversarial loss as defined in equation (4) and the signer classification loss (i.e., $\mathcal{L}_{\text{adv}} \leftarrow |\mathcal{L}_{\text{adv}} - \mathcal{L}_{\text{signer}}|$). This option has a nice information theoretic interpretation as being an empirical upper-bound for the mutual information between the distribution of signer identities and the distribution of latent representations. Nonetheless, there exist infinitely many (non-uniform) distributions for which this loss vanishes. Our choice, besides being clearly lower bounded by the entropy of the uniform distribution, $\log |\mathbb{S}|$, is minimum if and only if $p(s|h(\mathbf{X}_i; \theta_h); \theta_g) \equiv \mathcal{U}_{\mathbb{S}}(s)$, $\forall i$, meaning that the *signer-classifier* block is completely agnostic relatively to the signer identity of the training data.

3.3 Signer-transfer training objective

To further encourage the latent representations \mathbf{h} to be signer-invariant, we introduce an additional term in objective (5), the so-called signer-transfer loss $\mathcal{L}_{\text{transfer}}$. The core idea of $\mathcal{L}_{\text{transfer}}$ is to enforce the latent distributions of different signers to be as similar as possible. In practise, this is achieved by minimizing the difference between the hidden representations of different signers, at each layer of the *encoder* network. To measure the signer’s distribution difference at the m -th layer, $m = 1, \dots, M$, we compute a distance $\mathcal{D}^{(m)}$ between the hidden representations $h^{(m)}(\cdot; \theta_h)$ of two signers s and t at the output of that layer, such that:

$$\mathcal{D}^{(m)}(s, t; \theta_h) = \left\| \frac{1}{N_s} \sum_{i: s_i=s} h^{(m)}(\mathbf{X}_i; \theta_h) - \frac{1}{N_t} \sum_{j: s_j=t} h^{(m)}(\mathbf{X}_j; \theta_h) \right\|_2^2, \quad (6)$$

where $\|\cdot\|_2$ is the ℓ^2 -norm, and N_s and N_t denote the number of training examples of signers s and t , respectively. Accordingly, the signer-transfer loss at the m -th layer is the sum of the pairwise distances between all signers, i.e.:

$$\mathcal{L}_{\text{transfer}}^{(m)}(\theta_h) = \sum_{s \in \mathbb{S}} \sum_{\substack{t \in \mathbb{S}, \\ t \neq s}} \mathcal{D}^{(m)}(s, t; \theta_h) \quad (7)$$

The overall signer-transfer loss $\mathcal{L}_{\text{transfer}}$ is then a weighted sum of the losses computed at each layer of the *encoder* network, such that:

$$\mathcal{L}_{\text{transfer}}(\theta_h) = \sum_{m=1}^M \beta^{(m)} \mathcal{L}_{\text{transfer}}^{(m)}(\theta_h), \quad (8)$$

where $\beta^{(m)} \geq 0$ is a hyperparameter that controls the relative importance of the loss obtained at the m -th layer. By combining (5) and (8), the *encoder* and *sign-classifier* networks are trained to minimize the following loss function:

$$\min_{\theta_h, \theta_f} \mathcal{L}(\theta_h, \theta_f, \theta_g) = \mathcal{L}_{\text{sign}}(\theta_h, \theta_f) + \lambda \mathcal{L}_{\text{adv}}(\theta_h, \theta_g) + \gamma \mathcal{L}_{\text{transfer}}(\theta_h), \quad (9)$$

where $\gamma \geq 0$ is the weight that controls the relative importance of the signer-transfer term.

Summing up, the adversarial training procedure is organized by alternatively either training both the *encoder* and the *sign-classifier* in order to minimize objective (9) or training the *signer-classifier* in order to minimize objective (1).

4. EXPERIMENTAL EVALUATION

The experimental evaluation of the proposed model was performed using two publicly available SLR databases: the Jochen-Triesch database [19], and the Microsoft Kinect and Leap Motion American sign language (MKLM) database [13, 14]. Jochen-Triesch [19] is a dataset of 10 hand signs performed by 24 signers against three different types of backgrounds: uniform light, uniform dark and complex. Experiments on Jochen-Triesch were conducted using the standard evaluation protocol of this dataset [6], in which 8 signers are used for the training and the remaining 16 signers are used for the test. MKLM [13, 14] contains a total of 10 signs, each one repeated 10 times by 14 different signers. In this dataset, the performance of the models is assessed using 5 random splits, created with signer-independence, yielding at each split a training set of 10 signers, a validation set of 2 signers and a test set of 2 signers.

4.1 Implementation details

In order to extract the manual signs from the noisy background of the images, the automatic hand detection algorithm [2] is used as a pre-processing step. The images are then cropped, resized to the average sign size of the training set, and normalized to be in the range $[-1, 1]$. Throughout this section, the proposed model is compared with state-of-the-art methods for each dataset. Nevertheless, to further attest the robustness of the proposed model, two different baselines are also implemented:

- (Baseline 1) A CNN trained from scratch with ℓ^2 regularization. For a fair comparison, the architecture of the baseline CNN corresponds to the architecture of the *encoder* network followed by the *sign-classifier* network of the proposed model.
- (Baseline 2) A CNN with the baseline 1 topology, but trained with the triplet loss [17].

Here, the triplet loss concept is explored in order to impose signer-independence in the representation space and, hence, build up a more robust baseline. The underlying idea is to minimize the distance between an *anchor* and a *positive* latent representation, \mathbf{h}_{y_i, s_i} and \mathbf{h}_{y_p, s_p} , respectively; while maximizing the distance between the *anchor* \mathbf{h}_{y_i, s_i} and a *negative* representation \mathbf{h}_{y_n, s_n} . It is important to note that while *anchor* and *positive* latent representations have to be from the same sign class, their signer identity may or not change. On the other hand, *anchor* and *negative* representations are from different sign classes, whereas their signer identity may also change. In order to train baseline 2 in an end-to-end fashion for sign classification, the overall loss function to be minimized is a trade-off between the triplet loss $\mathcal{L}_{\text{triplet}}$ and a classification loss $\mathcal{L}_{\text{sign}}$, such that:

$$\mathcal{L} = \mathcal{L}_{\text{sign}} + \frac{\rho}{N} \sum_{i=1}^N \left[\|\mathbf{h}_{y_i, s_i} - \mathbf{h}_{y_p, s_p}\|_2^2 - \|\mathbf{h}_{y_i, s_i} - \mathbf{h}_{y_n, s_n}\|_2^2 + \alpha \right], \quad (10)$$

where $\mathcal{L}_{\text{sign}}$ corresponds to the categorical cross-entropy as defined in equation (2). The second term denotes the $\mathcal{L}_{\text{triplet}}$, where $y_p = y_i$ and $y_n \neq y_i$, and $\rho \geq 0$ is a hyperparameter controlling its the relative importance. The margin enforced between *positive* and *negative* pairs was fixed as $\alpha = 1$. In addition, following [17], an *online* triplet generation strategy, by selecting the hardest *positive/negative* samples within every mini-batch, was adopted.

Table 1. Hyperparameters sets.

Hyperparameters	Acronym	Set
Learning rate	-	$\{1e^{-04}, 1e^{-03}\}$
ℓ^2 -norm coefficient	-	$\{1e^{-05}, 1e^{-04}\}$
$\mathcal{L}_{\text{triplet}}$ weight	ρ	$\{0.1, 0.5, 1, 5, 10\}$
\mathcal{L}_{adv} weight	λ	$\{0.1, 0.5, 0.8, 1, 3\}$
$\mathcal{L}_{\text{transfer}}$ weight	γ	$\{1.5e^{-04}, 2e^{-04}, 4e^{-04}, 1e^{-03}\}$

Table 2. Experimental results on: (a) Jochen-Triesch dataset, and (b) MKLM dataset.

3*Method	Classification accuracy (%)			2*Method	Classification accuracy (%)		
	Background				average (std)	min	max
	Uniform	Complex	Both				
Just <i>et al</i> [6]	92.79	81.25	87.92	Marin <i>et al</i> [13]	89.71 (-)	-	-
Kelly <i>et al</i> [7]	91.80	-	-	Ferreira <i>et al</i> [2]	93.17 (-)	-	-
Dahmani <i>et al</i> [1]	93.10	-	-	CNN (Baseline 1)	89.90 (8.81)	73.00	98.00
CNN (Baseline 1)	97.50	74.38	89.79	CNN with Triplet loss (Baseline 2)	91.40 (3.93)	86.50	96.50
CNN with Triplet loss (Baseline 2)	98.13	75.63	90.63	Proposed method	94.80 (3.53)	89.50	100.00
Proposed method	98.75	91.25	96.25				

All deep models were implemented in PyTorch and trained with the Adam optimization algorithm using a batch size of 32 samples. For reproducibility purposes, the source code as well as the weights of the trained models are publicly available online*. The hyperparameters that are common to all the implemented models (i.e., learning rate and ℓ^2 regularization weight) as well as some hyperparameters that are specific to the proposed model (i.e., λ and γ) and to the implemented baseline 2 (i.e., ρ) were optimized by means of a grid search approach and cross-validation on the training set (see Table 1 for more details). The signer-transfer penalty $\mathcal{L}_{\text{transfer}}$ is applied to the last two layers of the *encoder* network with a relative weight of 1. Regarding the model’s architecture, the number of consecutive convolutional layers pairs L_e was set to 3, which results in a total of 6 convolutional layers. The number of filters starts as 32 which is then doubled after each convolutional pair. The dense layer on top of the *encoder* network has 128 neurons. The number of dense layers of both classifiers L_s was set to 3, and the number of nodes of each hidden layer was set as 128.

4.2 Results and discussion

Experiments on the Jochen-Triesch and MKLM databases are summarized in Tables 2(a) and 2(b), respectively. The results on the Jochen-Triesch database are presented in terms of average classification accuracy in the overall test set as well as against each specific background type (i.e., uniform and complex). For the MKLM database, Table 2(b) depicts the average classification accuracy computed across all the 5 test splits, as well as the minimum and maximum accuracy value achieved by each method.

The most interesting observation is the superior performance of the proposed model. Specifically, the proposed model provides the best overall classification accuracy on both SLR databases, clearly outperforming both implemented baselines and all the previous state-of-the-art models. In complex scenarios, as reported in Table 2(a), the proposed model surpasses all the other methods by a large margin (i.e., 91.25% against 81.25%, 74.38% and 75.63%). In addition, by analyzing the standard deviation as well as the minimum and maximum accuracy values, it possible to observe that the proposed model is the method with the lowest variability, yielding consistently high accuracy rates across all test splits of the MKLM dataset (see Table 2(b)). These results attest the robustness of the proposed model and its capability of better dealing with the large inter-signer variability that exists in the manual signing process of sign languages. Interestingly, the obtained results also reveal that the implemented baselines are in fact fairly strong models, both of them outperforming most of the state-of-the-art methods on both datasets.

Table 3 illustrates the effect of each proposed training scheme by itself. For this purpose, the proposed model was trained either (i) with just the adversarial procedure, without the signer-transfer $\mathcal{L}_{\text{transfer}}$ loss, or (ii) with just the $\mathcal{L}_{\text{transfer}}$ penalty on the *encoder* network without adversarial training. The results clearly demonstrate the complementary effect between the two training procedures, as their combination provides the best overall classification accuracy. Interestingly, each training scheme outperforms on its own both baselines and state-of-the-art methods.

*<https://github.com/pmmf/SI-SLR>

Table 3. The effect of each training procedure in the proposed model. The results in the last column are replicated from Tables 2(a) and 2(b) as they include both training procedures.

2*Dataset	Classification accuracy (%)		
	Only adversarial training	Only $\mathcal{L}_{\text{transfer}}$ penalty	Both
Jochen-Triesch	95.21	94.38	96.25
MKLM	94.00	94.10	94.80

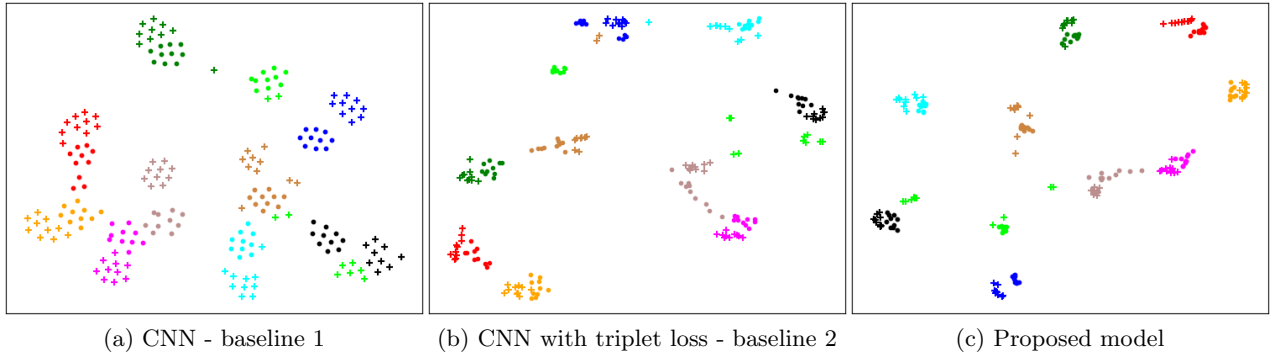


Figure 3. Two-dimensional projection of the latent representation space using the t-distributed stochastic neighbor embedding (t-SNE) [20]. Markers \bullet and $+$ represent 2 different test signers, while the different colors denote the 10 sign classes.

4.3 Latent space visualization

To further demonstrate the effectiveness of the proposed model in promoting signer-invariant latent representation spaces, we have performed a visual inspection of the latent representations through the t-distributed stochastic neighbor embedding (t-SNE) [20] (see Figure 3). These plots clearly demonstrate the better capability of the proposed model of imposing signer-independence in the latent representations. The proposed model yields a latent representation space in which representations of the same signer and different classes are close to each other and well mixed, while it keeps latent representations of different classes far apart. By analyzing the t-SNE plot of baseline 1, it is possible to observe that the latent representations of different signers and the same class tend to be far apart in the latent space. In addition, there is some overlapping between clusters of different classes. Although baseline 2 (CNN with the triplet loss) promoted slightly improvements over the standard baseline CNN, the proposed model achieved by far the best signer-invariance and class separability.

5. CONCLUSION

This paper presents a novel adversarial training objective, based on representation learning and deep neural networks, specifically designed to tackle the signer-independent SLR problem. The underlying idea is to learn signer-invariant latent representations that preserve as much information as possible about the signs, while discarding the signer-specific traits that are irrelevant for sign recognition. For this purpose, we introduce an adversarial training procedure for simultaneously training an *encoder* and a *sign-classifier* over the target sign variables, while preventing the latent representations of the *encoder* to be predictive of the signer identities. To further discourage the underlying representations of retaining any signer-specific information, we propose an additional training objective that enforces the latent distributions of different signers to be as similar as possible. Experimental results demonstrate the effectiveness of the proposed model in several SLR databases.

References

- [1] Djamila Dahmani and Slimane Larabi. User-independent system for sign language finger spelling recognition. *Journal of Visual Communication and Image Representation*, 25(5):1240 – 1250, 2014.
- [2] Pedro M. Ferreira, Jaime S. Cardoso, and Ana Rebelo. On the role of multimodal learning in the recognition of sign language. *Multimedia Tools and Applications*, Sep 2018.
- [3] Clément Feutry, Pablo Piantanida, Yoshua Bengio, and Pierre Duhamel. Learning anonymized representations with adversarial neural networks. *arXiv preprint arXiv:1802.09386*, 2018.

- [4] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1180–1189, Lille, France, 07–09 Jul 2015. PMLR.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- [6] A. Just, Y. Rodriguez, and S. Marcel. Hand posture classification and recognition using the modified census transform. In *7th International Conference on Automatic Face and Gesture Recognition (FGR06)*, pages 351–356, April 2006.
- [7] Daniel Kelly, John McDonald, and Charles Markham. A person independent system for recognition of hand postures used in sign language. *Pattern Recognition Letters*, 31(11):1359 – 1368, 2010.
- [8] Taehwan Kim, Weiran Wang, Hao Tang, and Karen Livescu. Signer-independent fingerspelling recognition with deep neural network adaptation. *CoRR*, abs/1602.04278, 2016.
- [9] O. Koller, H. Ney, and R. Bowden. Deep hand: How to train a cnn on 1 million hand images when your data is continuous and weakly labelled. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3793–3802, June 2016.
- [10] W.W. Kong and Surendra Ranganath. Towards subject independent continuous sign language recognition: A segment and merge approach. *Pattern Recognition*, 47(3):1294 – 1308, 2014. Handwriting Recognition and other PR Applications.
- [11] R. Kuhn, J. . Junqua, P. Nguyen, and N. Niedzielski. Rapid speaker adaptation in eigenvoice space. *IEEE Transactions on Speech and Audio Processing*, 8(6):695–707, Nov 2000.
- [12] Pradeep Kumar, Himaanshu Gauba, Partha Pratim Roy, and Debi Prosad Dogra. A multimodal framework for sensor based sign language recognition. *Neurocomputing*, 259:21 – 38, 2017. Multimodal Media Data Understanding and Analytics.
- [13] G. Marin, F. Dominio, and P. Zanuttigh. Hand gesture recognition with leap motion and kinect devices. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 1565–1569, Oct 2014.
- [14] Giulio Marin, Fabio Dominio, and Pietro Zanuttigh. Hand gesture recognition with jointly calibrated leap motion and depth sensor. *Multimedia Tools and Applications*, 75(22):14991–15015, Nov 2016.
- [15] N. Neverova, C. Wolf, G. Taylor, and F. Nebout. Moddrop: Adaptive multi-modal gesture recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(8):1692–1706, Aug 2016.
- [16] Lionel Pigou, Sander Dieleman, Pieter-Jan Kindermans, and Benjamin Schrauwen. Sign language recognition using convolutional neural networks. In Lourdes Agapito, Michael M. Bronstein, and Carsten Rother, editors, *Computer Vision - ECCV 2014 Workshops*, pages 572–578, Cham, 2015. Springer International Publishing.
- [17] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. *CoRR*, abs/1503.03832, 2015.
- [18] T. Shanableh and K. Assaleh. User-independent recognition of arabic sign language for facilitating communication with the deaf community. *Digital Signal Processing*, 21(4):535 – 542, 2011.
- [19] Jochen Triesch and Christoph von der Malsburg. A system for person-independent hand posture recognition against complex backgrounds. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(12):1449–1453, December 2001.

- [20] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [21] U. von Agris, C. Blomer, and K. Kraiss. Rapid signer adaptation for continuous sign language recognition using a combined approach of eigenvoices, mllr, and map. In *2008 19th International Conference on Pattern Recognition*, pages 1–4, Dec 2008.
- [22] Ulrich von Agris, Jörg Zieren, Ulrich Canzler, Britta Bauer, and Karl-Friedrich Kraiss. Recent developments in visual sign language recognition. *Universal Access in the Information Society*, 6(4):323–362, Feb 2008.
- [23] D. Wu, L. Pigou, P. Kindermans, N. D. Le, L. Shao, J. Dambre, and J. Odobez. Deep dynamic neural networks for multimodal gesture segmentation and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(8):1583–1597, Aug 2016.
- [24] Fang Yin, Xiujuan Chai, and Xilin Chen. Iterative reference driven metric learning for signer independent isolated sign language recognition. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 434–450, Cham, 2016. Springer International Publishing.
- [25] Fang Yin, Xiujuan Chai, Yu Zhou, and Xilin Chen. Weakly supervised metric learning towards signer adaptation for sign language recognition. In Mark W. Jones Xianghua Xie and Gary K. L. Tam, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 35.1–35.12. BMVA Press, September 2015.
- [26] Jörg Zieren and Karl-Friedrich Kraiss. Robust person-independent visual sign language recognition. In Jorge S. Marques, Nicolás Pérez de la Blanca, and Pedro Pina, editors, *Pattern Recognition and Image Analysis*, pages 520–528, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.

AUTHORS’ BACKGROUND

Your Name	Title*	Research Field	Personal website
Pedro M. Ferreira	Phd candidate	Computer Vision (e.g., gesture and emotion recognition), Deep Learning	https://www.inesctec.pt/en/people/pedro-martins-ferreira
Diogo Pernes	Phd candidate	Machine Learning Fundamental Topics, with some application to Computer Vision	https://www.inesctec.pt/en/people/diogo-pernes-cunha
Ana Rebelo	assistant professor	Computer Vision, Image Processing, Biometrics, Document Analysis	http://www.inescporto.pt/~arebelo/
Jaime S. Cardoso	associate professor	Computer Vision, Machine Learning, Medical Decision Support Systems	http://www.inescporto.pt/~jsc/

*This form helps us to understand your paper better, **the form itself will not be published.**

*Title can be chosen from: master student, Phd candidate, assistant professor, lecture, senior lecture, associate professor, full professor