

ÉCOLE NORMALE SUPÉRIEURE
DÉPARTEMENT D'INFORMATIQUE
RAPPORT DE STAGE DE L3

Multimodal Learning

A case study for Gesture and Audio-Visual Speech Recognition

HSIEH YU-GUAN

Sous la direction de

AMÉLIE CORDIER¹ & MATHIEU LEFORT²

¹Hoomano, Équipe R&D

²LIRIS, Équipe SMA

14 Juin 2017 – 11 Août 2017

1 Introduction

1.1 Cadre de travail, outils et environnement logiciel

Mon stage s'inscrit dans le cadre du laboratoire commun BEHAVIORS.AI¹ qui est un projet de recherche collaboré entre la société Hoomano et LIRIS. Hoomano est une entreprise créée en 2014 spécialisée dans le développement de logiciel pour les robots sociaux (ex: Nao et Pepper). Ils commercialisent un moteur d'interaction qui vise à offrir une interaction homme-robot la plus naturelle possible. Le LIRIS (Laboratoire en InfoRmatique, Image et Système d'Information) est une unité mixte de recherche en informatique située à Lyon qui compte 14 équipes de recherche et 330 membres. Le laboratoire commun est particulièrement lié à l'équipe SMA (Systèmes Multi-Agents) qui met l'accent sur l'intelligence artificielle, les systèmes multi-agents, et les systèmes cognitifs.

Pendant mon stage, j'ai passé les deux premières semaines au laboratoire et pour le reste du temps j'étais à Hoomano car j'entraînais les modèles sur le serveur de l'entreprise. Tous les codes ont été réalisés en python et précisément à l'aide de la bibliothèque open source TensorFlow² développée par Google. Cela incluait la construction des différents architectures et réseaux de neurones, l'entraînement et l'évaluation. TensorBoard était un outil assez pratique pour visualiser l'apprentissage, et je pouvais même facilement visualiser les données et les représentations apprises grâce au plugin qui est proposé. Toute la suite du rapport sera rédigée en anglais à cause de la difficulté de traduire tous les termes techniques en français.

1.2 Scientific overview

Historically, expert systems have been widely used for the design of artificial intelligence. Agents rely heavily on handcrafted ontologies and symbolic rules to make decisions and act in the world. Although impressive results have been obtained with this line of research, it requires a lot of engineering efforts and human knowledge. Furthermore, engineering in this way all the behaviors that are needed to display general human level intelligence is out of reach.

In addition, we face also the symbol grounding problem [11]. The symbols manipulated by the AI systems have no meaning to them. To solve these problems, the embodied paradigm [6] first argues an agent must be able to perceive and act from its own perspective, so it should have a body, e.g. a robot. Developmental robotics [51], which draws inspiration from developmental psychology [36], tries to implement child development in robots. It puts emphasis on the importance of endowing the robot with the capacity of learning new knowledge from its own sensorimotor experience.

In [41], Smith and Gasser discerned six fundamental principles for the development of embodied intelligence: multimodality, incremental development, physical interaction with the environment, exploration, social guidance and symbolic language acquisition. I was particularly interested in the multimodal point in my internship. The everyday concepts that a human acquires is intrinsically multimodal. For example, the word "cat" can be quickly associated with the visual appearance of a cat, its vocalization, and the tactile feeling that we have while petting a cat. The existence of neurons that receive early information from different modalities has also been proven [28].

In machine learning, recently, deep learning has achieved a great success in various domains, such as image recognition [18], text generation [10] or machine translation [46]. It has also been more and more often applied to multimodal inputs [2, 33]. Its capacity of learning hierarchical representations of data in a fully unsupervised way [40, 50] is particularly interesting in the domain of robotics.

In my internship, I studied mainly the application of multimodal learning using deep networks in the fields of multimodal gesture recognition and audio-visual speech recognition (AVSR). I focused especially on the problem of shared representation learning and knowledge transfer. I wanted to show that the availability of multiple modalities could enable the model to learn a better representation for

¹ <http://behaviors.ai/>

² <https://www.tensorflow.org/>

each single modality and that one can leverage information from one modality to be used for another when there is an imbalance of amount of data among different sources.

The report is organized as follows. In section 2 I briefly review related work on gesture recognition, AVSR and other fields of multimodal learning. In section 3, I present several basic deep neural network architectures playing an important role in my work. In sections 4 and 5, I describe respectively the datasets and common experimental setups that were used in my internship. Experimental details and results are given in sections 6 and 7. Finally, conclusions and perspectives are presented in section 8.

The source code of all of the work described here can be found on my github: https://github.com/cybermeow/internship_2017.

2 Related Work

2.1 Gesture recognition

Gesture recognition has been studied for a while within the fields of computer vision and pattern recognition [26, 44]. Recently, deep learning algorithms have led to significant advances in this domain [1]. The use of Convolutional Neural Networks (CNNs, see 3.1) is the most common [30]. 3d CNNs are used to deal with image sequences in [27]. Architectures that take multimodal inputs have also grown in popularity. In [31] Neverova *et al.* copes with color, depth, skeleton and audio information using CNNs and Recurrent Neural Networks (RNNs). Some other works use 3d CNNs and deep belief networks (DBNs) while facing similar multimodal inputs [32, 37, 53].

2.2 Audio-visual speech recognition (AVSR)

AVSR is probably one of the earliest examples of multimodal research. Most early works were based on various hidden Markov model (HMM) extensions [38], but the use of neural networks were also explored [5, 54]. In the past few years, AVSR has drawn attention from the deep learning community [16, 33, 34]. In [33] the authors use a deep network to learn a joint representation of visual and auditory input. They show that better features for one modality can be learned if multiple modalities are present at feature learning time. A transfer deep learning framework applied in AVSR is proposed in [29] and this forms the base of my study in 7.2. In [7] and [12] several examples of how deep architectures can be used to deal with audio data are given.

2.3 Multimodal learning in other fields

Multimodal learning also finds its application in many other domains. For example, there has been a surge of interest in the exploitation of multimodal information in the fields of image annotation [52], zero-shot learning [9, 42], and automatic image caption generation [15]. They first train visual and language models separately and as a next step they try to learn either a common embedding [52], a mapping [9, 42] or an alignment [15] between the two models.

More generally, in [2] Baltrusaitis *et al.* offer an overview of recent advances in multimodal machine learning. They identify five core technical challenges surrounding this research field: shared representation learning, translation from one modality to another, multimodal alignment that aims to find correspondences between sub-components of instances from two or more modalities, multimodal fusion that integrates information from multiple modalities with the goal of predicting an outcome measure and co-learning which aids the modeling of a modality by exploiting knowledge from another modality.

Concerning the use of multimodal learning in robotics, in [8] Droniou *et al.* propose a network that is able to learn both a symbolic representation of data and a fine discrimination between two similar stimuli in an unsupervised way. The authors apply their method to visual, auditory and proprioceptive data of the humanoid robot iCub. A multimodal embedding that is able to combine information coming from vision, language and motion trajectories is defined in [45] and endows the robot with the capacity of manipulating an unseen object.

3 Presentation of Some Deep Network Architectures

3.1 Convolutional Neural Networks

CNNs are an early family of deep learning architectures inspired from the human vision system [19]. Generally we have convolutional layers alternating with pooling (subsampling) layers, but fully connected layers can also be introduced (see Figure 5). CNNs have been shown to achieve state-of-the-art performance in image processing tasks such as image classification [18] and object detection [20]. However, they can be equally applied in other fields like speech recognition [7].

For a convolutional layer with a mono-channel input x , the latent representation of the k -th feature map is given by

$$h^k = \sigma(x * W^k + b^k)$$

where W^k is the k -th kernel, b^k the bias is broadcasted to the whole feature map, $*$ denotes the 2d convolution and σ is an activation function. Currently, ReLu (Rectified Linear Units) is probably the activation function that is the most commonly used for a CNN [18]. It is defined by

$$\sigma(x) = x^+ = \max(0, x),$$

when x is a single scalar and otherwise the above function is applied to each element of the input. The convolution kernel is extended to the full depth of input when there are multiple input channels. On the contrary, no parameters are required for defining a pooling layer. It just take some $d \times d$ region (supposing that the kernel has the same height and width) and output a single value, which for example, is the maximum in that region when using max-pooling.

There are yet two other important factors that are used to define these operations: *stride* and *padding*.³ Concerning this, two different zero paddings are defined in TensorFlow, ‘SAME’ and ‘VALID’. For ‘SAME’ padding, $h_o = \lceil h_i / s_h \rceil$, and for ‘VALID’ padding, $h_o = \lceil (h_i - h_f + 1) / s_h \rceil$, where h_i , h_o , h_f are respectively the height of the input feature map, the output feature map and the filter (kernel) and s_h is the stride of the height dimension. Similar formulas also exist for the width.⁴

Suppose all the parameters of a CNN model have been learned, the model is then ready to be used for a specific task. One just needs to run the model forward with the operations described above. Therefore, the purpose of the learning phase is to find the appropriate parameters that can help us solve the task. Mathematically, this is achieved by optimizing the weights of the network to minimize some loss function (e.g. cross entropy or L2 distance, see section 5). The optimization is usually done by a SGD-based algorithm. SGD is the abbreviation of stochastic gradient descent. At each training step, we compute the gradients of the loss with respect to all weights in the network and then update the weights according to these gradients. Using the chain rule, gradients are computed layer by layer (starting from the output of the network) and the error is thus *back-propagated* to the whole network.⁵

In [14], 3d CNNs are proposed to be used for video inputs. They’re like traditional CNNs except for the fact that we replace 2d feature maps by 3d feature maps and 2d kernels by 3d kernels.

3.2 Auto-Encoder

Auto-Encoders are networks that are trained to minimize the reconstruction error by back-propagating it from the output layer to hidden layers. In the simplest model with one hidden layer, an auto-encoder takes an input $\mathbf{x} \in \mathbb{R}^d$ and maps it to the latent representation $\mathbf{h} \in \mathbb{R}^{d'}$ given by $\mathbf{h} = \sigma(W\mathbf{x} + \mathbf{b})$ where W is a weight matrix, \mathbf{b} is a bias vector and σ is an activation function. Then the network

³ Due to the space limit, more details of the CNN architecture, including the role of these two hyperparameters, can be found here: <http://cs231n.github.io/convolutional-networks/>.

⁴ https://www.tensorflow.org/api_guides/python/nn#Convolution

⁵ <http://andrew.gibiansky.com/blog/machine-learning/convolutional-neural-networks/>

tries to reconstruct the input by a reverse mapping $\mathbf{x}' = \sigma(W'\mathbf{h} + \mathbf{b}')$. The loss function can be for instance $\|x' - x\|$ where $\|\cdot\|$ is some distance.

To prevent the auto-encoder from learning the identity function as a trivial solution, several regularization techniques have been proposed. The bottleneck approach forces dimensionality reduction by having fewer neurons in hidden layers than in the input layer. For example, in the above case, we must have $d' < d$. Sparse auto-encoders impose sparsity on hidden units [21]. Denoising auto-encoders, which played an important role in my work, try to reconstruct the clean input from its corrupted version [3, 49]. Binomial noise (switching pixels on or off) adding to input or hidden layers were used in my case.

Intuitively, auto-encoders are useful for data reconstruction. Nevertheless, the true interest lies in fact in its capacity to learn a representation (encoding) for a set of data in a purely unsupervised fashion [50]. Recently, auto-encoders have also been more and more used as a generative model [4].

3.3 Convolutional Auto-Encoder

Fully connected auto-encoders ignore the 2d image structure. This can cause problems when dealing with real-world size inputs, and introduce redundancy in the parameters. Convolutional Auto-Encoders [22, 48] (CAEs) are intuitively similar to architectures described in 3.2. However, convolutional and transposed convolutional layers are used instead. Transposed convolutional layers are also called deconvolutional layers [55]. We perform in fact also a convolution operation but with zero paddings around the image and sometimes around each pixel to upsample the image and to inverse the previous convolution (imagine that if the value of one pixel comes from 9 pixels of the previous layer, when doing a transposed convolution this pixel contributes to the same 9 pixels for reconstruction). This github directory https://github.com/vdumoulin/conv_arithmetic is quite useful for the understanding of the concept.

Pooling and unpooling layers are also sometimes considered in a CAE architecture [48]. More details on this aspect, which was not used in my work, can be found in in [35, 48].

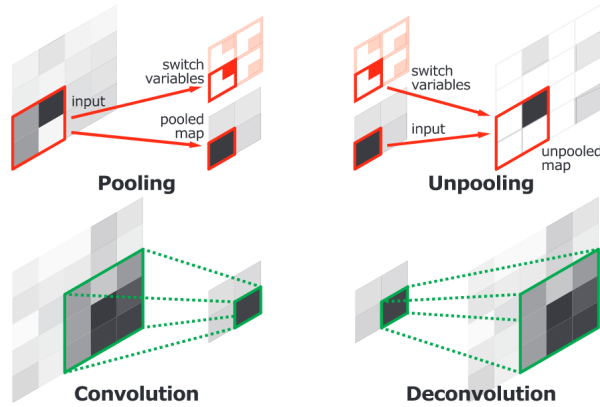


Figure 1: Illustration of deconvolution and unpooling layers [35]

4 Datasets and Preprocessing

Many datasets were explored during my internship. In this report, I will only detail the three that I used the most. Two of them are for gesture recognition: Creative Senz3D [24, 25] and ASL Finger Spelling [39], and one is for AVSR: AVLetters [23].

4.1 Creative Senz3D

The dataset contains gestures coming from 4 different people, each performing 11 different static gestures repeated 30 times each, for a total of 1320 samples. For each sample, color, depth and confidence frames are available. I only used the color and depth frames of this dataset since my

architectures took at most two modalities in input. The original size of each image is 480×640 and they were resized to 299×299 pixels before being fed to the network since this is also the input size of the Inception model [47]. No other preprocessing were done. For both color and depth images I used the three color channels (even though a priori only one channel was needed for depth maps).

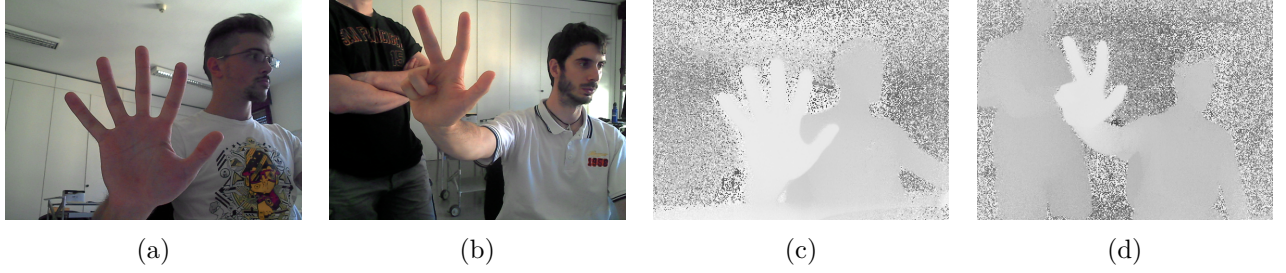


Figure 2: Example images in the Creative Senz3D dataset.

a, b) Color images.

c, d) Corresponding depth images. (c) and (d) are respectively the depth maps of (a) and (b).

All of the images are of size 480×640 and contain the the entire upper body of the subject.

4.2 ASL Finger Spelling

The dataset is composed of more than 60000 images in each modality (RGB and depth images are provided). Five subjects are asked to perform the 24 static signs in the American Sign Language (ASL) alphabet (excluding j and z which involve motion) a certain number of times, captured with similar lighting and background.

Images of this dataset are of variable sizes. The data preprocessing included resizing each image to 83×83 pixels, converting them to grayscale and Z-normalization (normalizing to zero mean and unit of variance).

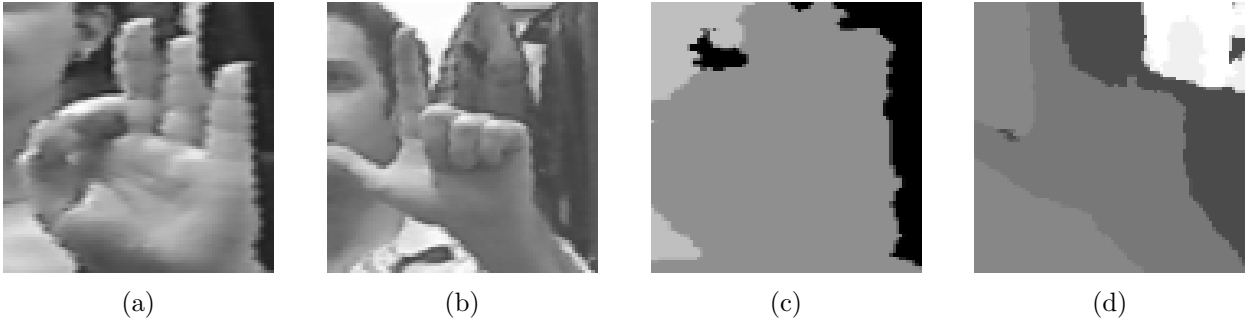


Figure 3: Example images in the ASL Finger Spelling dataset (after preprocessing).

a, b) Grayscale intensity images.

c, d) The corresponding depth maps. (c) and (d) are respectively the depth images of (a) and (b).

Images of this dataset have variable sizes, and they're all resized to 83×83 before being fed to the network. Generally only the hand region is contained in the image.

4.3 AVLetters

The dataset comprises video and audio recordings of 10 speakers uttering the letters A to Z, three times each. We count therefore 780 samples in total. For video data, image sequences of pre-extracted lip regions are provided. Each single image is of size 60×80 . For audio data, only the mel-frequency cepstrum coefficients (MFCCs⁶) are given, and each audio frame is represented by 26 MFCCs. The lack of raw audio data is a strong constraint on what we're able to do on this dataset.

⁶ MFCCs are a feature widely used in automatic speech and speaker recognition. Their job is to extract spectral envelope in order to discriminate the phoneme being pronounced.

Since all utterances do not have the same time duration, I used Fourier resampling to force every video input to be of length 12 and every audio input to be of length 24. These numbers were chosen according to [29]. Video frames were Z-normalized. Several data augmentation techniques were also considered, including random brightness adjusting, random contrast adjusting and random cropping (but at least 80% of the original image was kept). The data was augmented on the fly.



Figure 4: Example visual input for the AVletters dataset.

Pre-extracted lip regions of 60×80 pixels are provided. Each image sequence is resampled to be of length twelve in order to give an input of fixed size to the network. In the above figure, time goes from left to right then top to bottom.

5 Experimental Setup

Each classifier used in my work was trained with the cross entropy cost function and each auto-encoder was trained with the L2 distance between the input and output vector as the loss. The cross entropy measures the similarity between two probability distributions p and q over the same underlying set of events. When p and q are discrete, which is always true in my case, it's defined by

$$H(p, q) = - \sum_x p(x) \log q(x).$$

For the sake of preventing overfitting, L2 regularization [3] was applied to all the weights of network with a regularization coefficient λ . The Adam algorithm [17], which is a variant of the standard SGD, was then introduced for minimizing the loss function. An exponential decay was further used for the stepsize α of this algorithm with an initial stepsize α_0 varying from 0.01 to 0.0001 depending on experiment. The decaying rate γ was generally close to 0.8 and the decay took place every 100 training steps. That is, the stepsize was multiplied by γ every 100 steps.

Inputs of the network were normally fed as mini-batches of size 24 since rather satisfying results could be obtained with this number and when larger batches were used training became slower without a significant improvement of the performance. Batch normalization [13] were introduced after every convolutional and transposed convolutional layer. Therefore, the real operations used to compute neural activations were more complicated than what are described above. The presented settings and hyperparameter values were found with a manual heuristic search as I focused on principles and not on better performance achievement.

Here are some more details of the network architectures: ReLu activations were added to all the hidden layers of all the architectures [18] while no activation function was used for the output layers. For classification model dropout [43] was always applied to the second to last layer during training. The output of the classification layer was mapped to the probabilities that one data example belongs to each class by the softmax function. Therefore, if the output vector is $\mathbf{z} = (z_1, \dots, z_n)$, for $j \in \{1, \dots, n\}$, the probability that the input has the j^{th} label is

$$\text{softmax}(\mathbf{z})_j = \frac{\exp(z_j)}{\sum_{i=1}^n \exp(z_i)}.$$

For classification experiments, except for the one described in 7.2, the dataset was always separated

into training and test set. The classifier was first trained on training data and then tested on test data once the training phrase was finished. Unless stated otherwise, the prediction accuracies mentioned hereinafter are always evaluated on the test set.

Moreover, we can consider two different possible settings:

Subject Dependent. In a subject dependent setting, data samples are separated randomly into training set and test set. Therefore, during the testing phase, the classifier does not need to deal with data from an individual that it has never seen before.

Subject Independent. On the contrary, the classifier faces individuals never seen before during testing in a subject independent setting. In my case, I always separated one subject from the others to form the test set while the training set consisted of the rest of the data.

6 Experiments and Results: Unimodal Cases

6.1 Classification

For every dataset, I began with training a classifier on it in a totally supervised manner. This gave me an insight into its data quality, the preprocessing effectiveness and ensured that further experiments could be conducted. CNN was then one of the most suitable architecture for this purpose.

First, it's worth mentioning the problem of overfitting. It was observed for almost all the classification experiments that I carried out, and it was particularly severe for CNNs with many layers. In fact, CNN architectures could usually classify perfectly the training data, but experienced a drop of performance when evaluating on test set.

It is well-known that by reducing the number of hidden layers and increasing the weight regularization coefficient λ we may be able to cure this problem. For example, when λ was augmented from 0.0004 to 0.1, the classification accuracy for the audio input of the AVLetters dataset increased by about 10 points (from 63.22% to 77.84%) while using the architecture detailed on the left side of [Figure 10](#) with the hyperparameters given in [Table 2](#). By using fewer hidden layers in the 3d CNN architecture overfitting was also alleviated when dealing with the video input of this same dataset. However, these techniques did not always work and more often I got a poorer performance during training without an improvement of performance for test.

Below I'll briefly describe the final results of this part dataset by dataset.

6.1.1 Creative Senz3d

The tested architectures were perceptron (in all of this report by perceptron I denote a single-layer perceptron), pre-trained InceptionV4 model [\[47\]](#) and several hand-coded CNNs⁷. The lack of data quantity, variety, and the fact that the head is also contained in the image seems to increase the classification difficulty. In conclusion, I found that it was more interesting to use a subject dependent setting for this dataset. In this case, for RGB images, all of the classifiers were able to have a classification accuracy that was closed to 100%. For depth images, the classification accuracy was between 60% and 70% using a perceptron and near 90% for other architectures.

6.1.2 ASL Finger Spelling

The large number of data contained in this dataset and the relatively simple image content (single hand instead of the entire upper body) makes the classification task much easier. By using the CNN architecture shown in [Figure 5](#), I could achieve a classification accuracy of respectively 79.73% and 64.46% for intensity and depth images (see [Table 1](#)) in a subject independent setting (four subjects for training and one subject for testing). Fewer layers in the CNN architecture may also allow us to achieve the same performance. More experiments need to be carried out to find the most suitable hyperparameters.

⁷ For those who are interested, please refer to my github directory for more details.

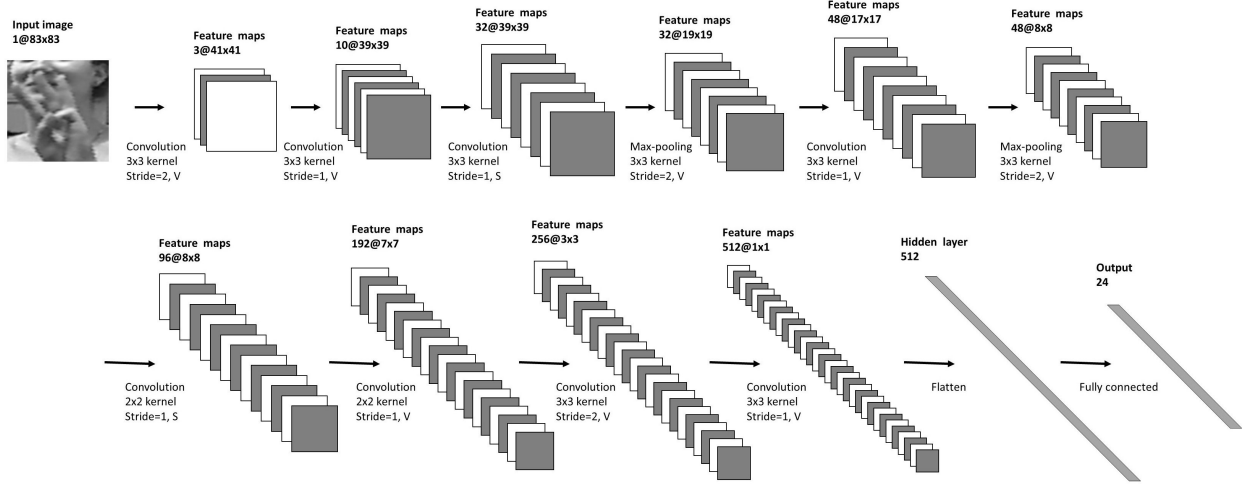


Figure 5: CNN architecture used for the Finger Spelling dataset.

The input of the network is a one-channel image of size 83×83 . It contains ten hidden layers. S stands for ‘SAME’ padding and V stands for ‘VALID’ padding (see 3.1).

6.1.3 AVLetters

One can refer to Figure 10 for the main CNN architectures that were used in this dataset. Notice that 3d CNNs were employed to deal with video inputs. Considering the small number of available data, a speaker dependent setting was used. With some carefully chosen hyperparameter values and network architectures (see Figure 10 and Table 2), the prediction accuracy was of 77.84% for audio and 54.32% for video.

6.2 Convolutional auto-encoder

Several distinct CAE hyperparameters were tested during my internship. Here I present the one with five hidden layers; therefore it contains three convolutional layers and three transposed convolutional layers as illustrated in Figure 6. I used end-to-end training for learning my architectures.

The proposed architecture was then trained on the two gesture recognition datasets. First of all, I was interested in the denoising capacity of the auto-encoder. An example is given in Figure 7. The auto-encoder is effectively able to reconstruct the clean image in a way, even though the result is blurred and sometimes distorted.

What’s more important, can meaningful high-level features be learned in this way? In fact, when autonomous robots interact with the environment, they need to extract high-level knowledge from raw perception on its own. This should be done in an unsupervised way since no labeled data is provided. That is why I am particularly interested in the problem of representation learning here. The output of the middle layer was taken as a new representation of the input data. I wanted to show that it could better represent the data than the raw input.

By doing principal component analysis, activation values of different layers could be projected into three dimensions for visualization. However, to quantify the learned features, I further trained two perceptrons respectively on top of the raw input and the high-level representation and compared their performance.

As suggested in 6.1, I used a subject dependent setting for Creative Senz3d while a subject independent setting was employed for ASL Finger Spelling. As already mentioned earlier, a perceptron built on raw data could classify perfectly the input while dealing with color images of Creative Senz3d. In other cases, the use of the learned data representation always led to an improvement of 10-20 points for the prediction accuracy (refer to Table 1 for results on ASL Finger Spelling). Useful high-level data representation can thus be learned in a fully unsupervised manner. However, if the whole architecture used for classification here (three layers of convolution and one layer of classification) was trained in

a totally supervised manner, better performance could be achieved (see the column CAE architecture of Table 1).

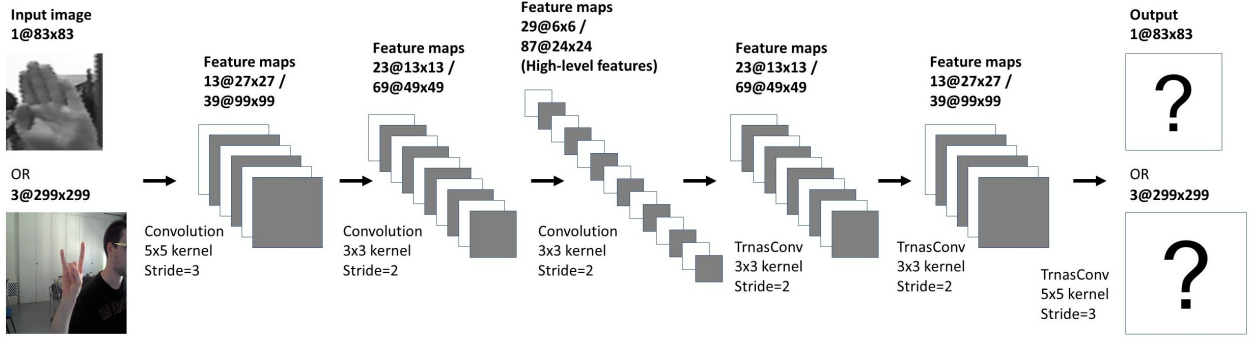


Figure 6: Convolutional auto-encoder architecture with three convolutional layers and three transposed convolutional layers.

Activation values of the middle layer are taken as high-level features of the input image. Inputs of the network can be of different sizes. We only use ‘VALID’ paddings here. For the results that are presented in this work, the training hyperparameters are $\alpha_0 = 0.01$, $\gamma = 0.8$ and $\lambda = 0.0004$ (see the first paragraph of section 5).

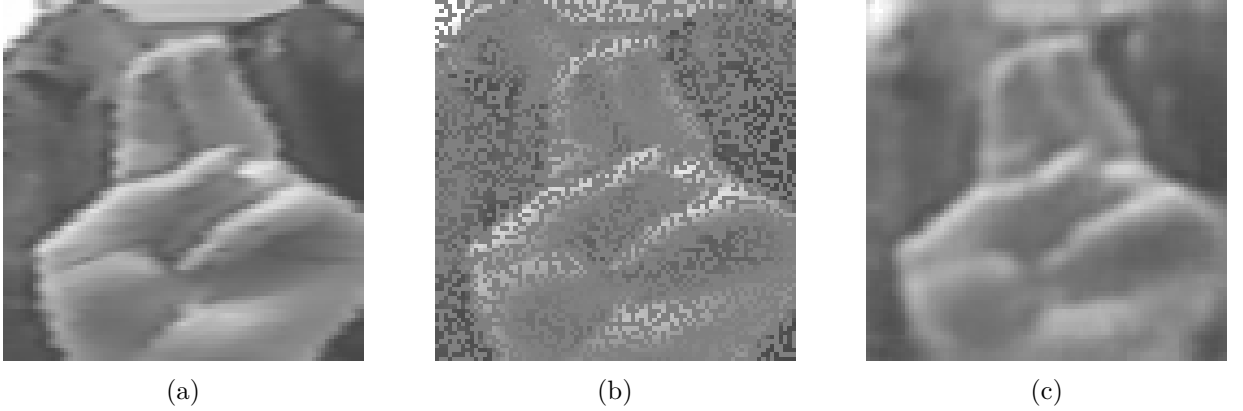


Figure 7: Image restoration using convolutional auto-encoder.

- a) Clean Image.
- b) Noisy image [input].
- c) Restored image [output].

7 Experiments and Results: Multimodal Cases

7.1 Shared representation learning

A first fundamental challenge in multimodal learning is the problem of representing and summarizing data from several modalities. This is also a key problem in developmental robotics as robots normally possess multiple perceptual channels. So how can we relate information from multiple sources? Mainly inspired from [8, 33], I generalized the CAE architecture introduced in 6.2 to a multimodal setup. As presented in Figure 8, two CAEs of different modalities shared a common middle layer by doing a weighted sum of corresponding values. I first pre-trained the first two layers of each modality using a unimodal CAE. In a second stage, I trained the rest of the network to reconstruct the two modalities of the input. This is basically the same model as the one proposed in [33] except for a different method that is applied for unimodal representation learning.

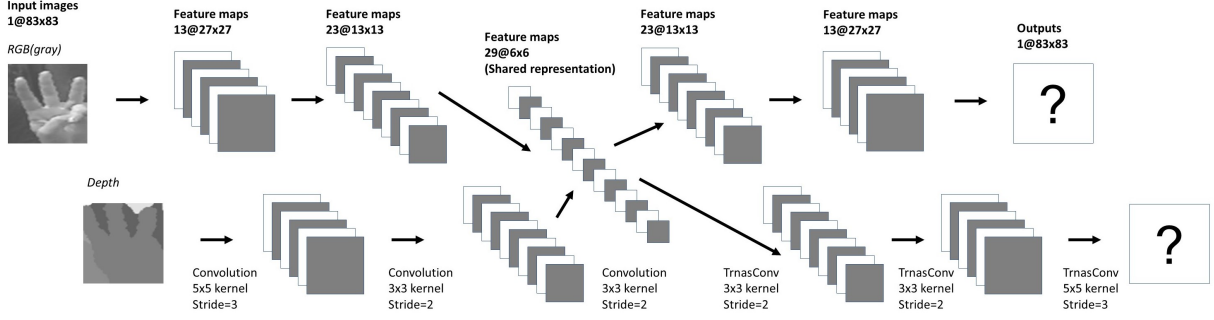


Figure 8: The bimodal convolutional auto-encoder model that is used to learn shared multimodal representation.

We simply take the CAE architecture that is introduced [Figure 6](#) for each modality but force them to have a shared middle layer by adding the corresponding activation values. We then try to reconstruct the two images separately through two disjoint paths. During training, I used $\alpha_0 = 0.01$, $\gamma = 0.8$ and $\lambda = 0.0004$ (see the first paragraph of [section 5](#)) for the presented results.

To prevent the network from finding representations such that different hidden units are tuned for different modalities separately, random dropouts were added to inputs. I first drew a random number p between 0 and 1. Each pixel of the color modality was then kept with probability p while for the depth modality they were kept with probability $1 - p$. In particular, sometimes one modality could be totally absent whereas the whole clean image was given for the other modality. In this way I gauranteed also that I always had enough information in input for the reconstruction if the two modalities were used jointly by the network. In addition, proper scaling were considered to keep the expected sum of the activations at each layer to be the same.

Ideally, we expect that the network is able to capture correlations across different modalities and learn a better single modality representation thanks to the presence of the other modality during feature learning. To verify this hypothesis, I built a perceptron on the top of the middle layer of the architecture and trained it in a supervised way while only one modality was given in input. For example, if I wanted to train a classifier for color images, zeros were fed as depth inputs.

The results are shown in [Table 1](#). We can observe an improvement of the classification accuracy for a classifier that is based on the learned shared representation compared to the one that is built on the features that are learned by an unimodal CAE. This reveals the utility of using multiple modalities for representation learning. However, more tests and statistics are necessary for validating this preliminary result.

How about exploiting the information from the two modalities in a totally supervised way? I took the CNN architecture of [Figure 5](#) for color and depth images separately until the seventh layer where a fusion was carried out. After training the network, the prediction accuracy was always at about 80% and no improvement was obtained compared with a CNN trained only on color inputs. The little or no improvement may be due to the similarity between the two modalities. Supposing that color images contain already all the necessary information for the task at hand, then depth maps would not bring any supplementary information that benefits the specific purpose I defined. Again, still more tests, and perhaps with different data, are needed to verify this hypothesis.

There is yet another interesting result that's worth mentioning. In an experiment, I trained a CAE (the one presented in [Figure 6](#)) for noisy depth inputs but instead of comparing the outputs with clean depth images they were compared with clean color images. In other words, I tried to construct the color modality from noisy depth maps. When I did the classification based on the learned representation, the performance was comparable with the perceptron that was built on the CAE features that were learned by a normal depth CAE. This result is particularly useful when one or several modalities are noisy whereas clean information is available for the others.

Table 1: Classification performance on the ASL Finger Spelling dataset

Raw) Perceptron that reads raw input data.

CAE features) Perceptron stacked on the middle layer of the CAE. (Figure 6 and 6.2).

Shared) Perceptron that exploits the shared representation learned by a bimodal CAE (Figure 8 and 7.1).

CAE architecture) Perceptron stacked on the middle layer of the CAE but train the whole network in a supervised way as a CNN. (Figure 6 and 6.2).

CNN) The CNN architecture in Figure 5 (Figure 5 and 6.1.2).

The used hyperparameters are $\alpha_0 = 0.005$, $\gamma = 0.8$ and $\lambda = 0.0004$ (see the first paragraph of section 5). We notice that we have exactly the same network architecture for the middle three experimental setups and only the training process differs one from another.

		Raw	CAE features	Shared	CAE architecture	CNN
Intensity	train	69.47 %	78.87 %	85.85 %	91.29 %	99.69 %
	test	32.64 %	50.24 %	53.38 %	65.44 %	79.73 %
Depth	train	63.64 %	79.61 %	81.83 %	88.80 %	97.24 %
	test	29.93 %	41.64 %	42.85 %	55.62 %	64.46 %

We may also want to ask if this architecture, when a partial input with only one modality is given, is able to infer the values of the missing modality. Several examples can be seen in Figure 9. Visually speaking, the reconstruction result seems better when both modalities are available in input even though they're both very noisy.

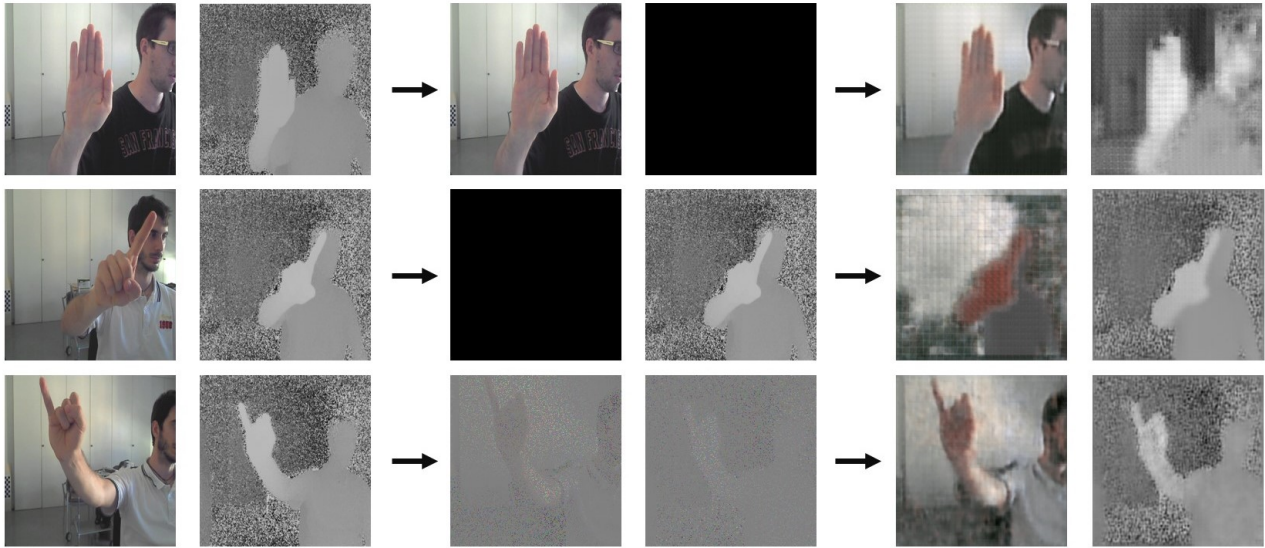


Figure 9: Restore color and depth images from incomplete input information.

Top) Only the color image is given.

Middle) Only the depth image is given.

Bottom) Both modalities are given but with little information (10% of pixels).

From left to right: original images \rightarrow images with dropout \rightarrow reconstructed images.

7.2 Transfer learning

In real world, we face frequently an imbalance in the amount of data among different modalities, and acquiring a parallel multimodal dataset is extremely resource consuming. The same defect equally exists for an autonomous robot because multimodal information received by its different sensors are not always related. Consequently, I'll discuss the knowledge transfer problem between different modalities in the second part of this section. The goal is to exploit information from one modality to improve the learning of another modality. This work is very similar to what is done in [29], but at the same

time it can also be viewed as a form of zero-shot learning [9, 42] since as we will see, not every label used for testing is directly presented during training.

I studied particularly the transfer between speech and lip-reading video data using the AVLetters dataset. First the dataset was separated randomly into training and test set. They are respectively noted Tr and Te . Tr contained 600 data samples while Te comprised the left 180 instances. For X an arbitrary subset of data, X^a and X^v denote respectively the audio and video data in X .

Next, to simulate the imbalance of data quantity between different modalities in the real world, I further splitted Tr , Te into $Tr_{A \sim T}$, $Tr_{U \sim Z}$, $Te_{A \sim T}$ and $Te_{U \sim Z}$ according to the label of each sample. For instance, a lip-reading video of a speaker saying E is in either $Tr_{A \sim T}^v$ or $Te_{A \sim T}^v$. During the first training stage, the video model was only accessible to a partial label space. The network (see Figure 10 and Table 2 for architecture and hyperparameters) was trained on $Tr_{A \sim T}^v$ which contained 466 samples. On the other hand, the audio model was trained on the whole label space Tr^a .

After that, I wanted to leverage speech data to fine-tune the network trained for video recognition. For a data sample $x^a \in Tr_{U \sim Z}^a$, I first computed its audio representation $h^a(x^a)$ that was taken from some hidden layer of audio model. A transfer function \mathcal{T} was used to approximate the video representation

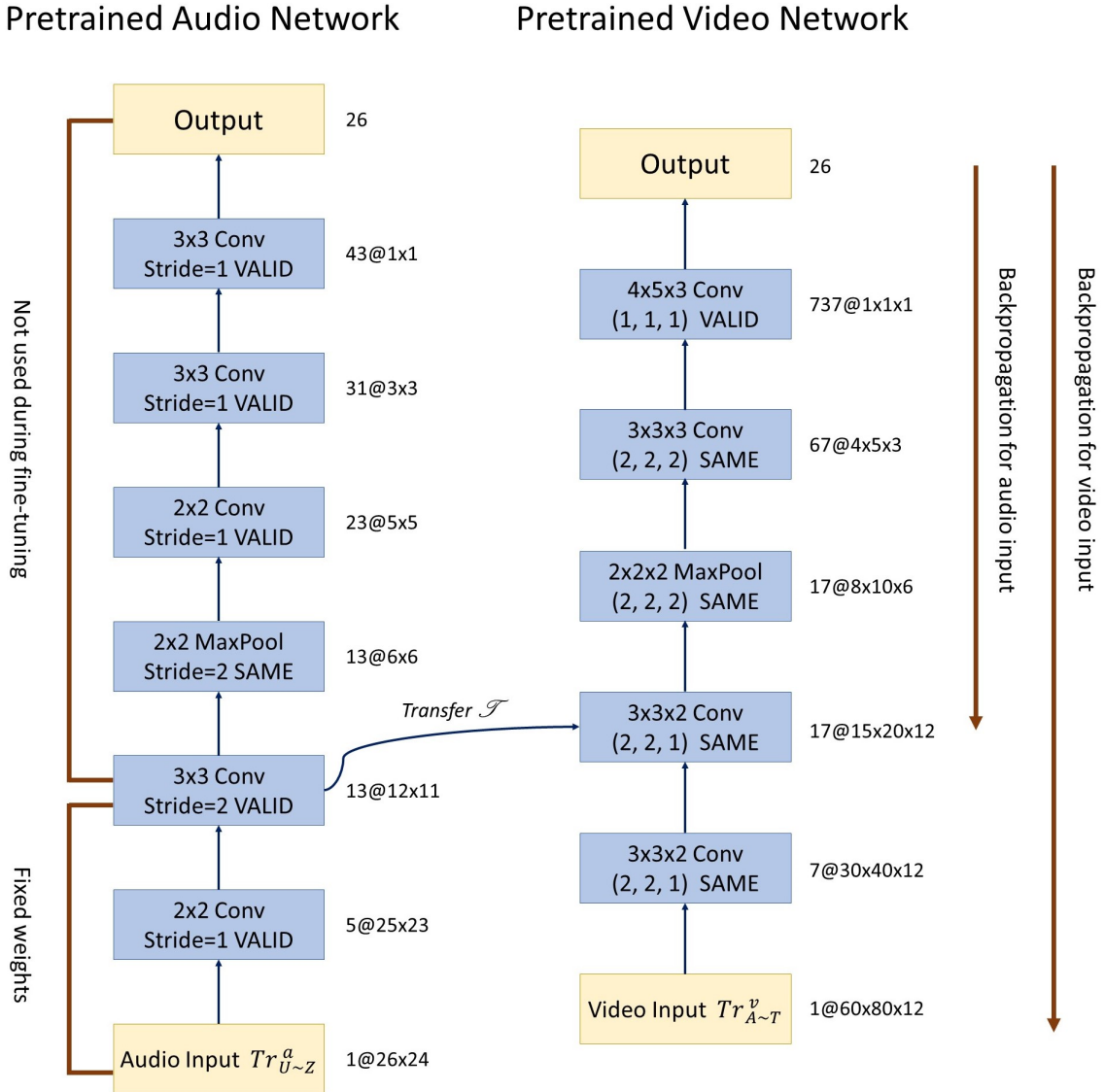


Figure 10: Illustration of the transfer learning approach applied in audio and lip-reading speech recognition tasks.

We first learn two separated models for audio and visual inputs (in my case two CNNs) and try to fine-tune the video network with transferred audio data.

of this data sample. That is, we want $\mathcal{T}(h^a(x^a)) \approx h^v(x^v)$. Finally, I fine-tuned the subnetwork situated after the hidden video layer that was taken as representation via a standard backpropagation algorithm. I chose audio as the source modality for two principal reasons. First, generally the quantity of audio data surpasses that of video data in the field of AVSR. Second, the lip-reading task appears to be more difficult than that of audio speech recognition.

The detailed experiment is presented in Figure 10. For both audio and video network I choose the second hidden layer for transfer. For the sake of simplicity, an instance-based approach (KNN-based mapping) was employed to define \mathcal{T} . I obtained mapping of a new audio input x^a by first finding the K -closest audio samples in $Tr_{A \sim T}^a$ in the representation space, and then returning the average values of the corresponding video samples (also in the representation space).

Nonetheless, if only audio inputs were presented during fine-tuning, it might cause a bias of the fine-tuned network towards the six new classes. To avoid this problem, data from $Tr_{A \sim T}^v$ were also fed as input from time to time to fine-tune the whole network. Concretely, with probability p_a audio data was presented and otherwise video input was given. At the end, the new network was tested on different parts of video data (e.g. $Tr_{A \sim T}^v$, $Te_{U \sim Z}^v$) for evaluation.

Many experiments with various hyperparameter values were conducted and some results are shown in Table 2. We see that despite the fact that the overall performance is not certainly improved by knowledge transfer, the fine-tuned network is now able to classify samples whose labels range from U to Z with an accuracy of 10-20% (**Exp1**, 21.64% for $Tr_{U \sim Z}^v$ and 15.22% for $Te_{U \sim Z}^v$). It shows the utility of the transfer learning approach which enables us to classify labels that are not at all available for the target modality. During fine-tuning, the magnitude of learning rate is equally important. As we can see, if the learning rate is too high, it can deteriorate the pre-trained network and the performance may drop drastically (**Exp2**). On the other hand, as predicted before, if only audio data labeled from U to Z are provided for fine-tuning, the classification performance drops for the first twenty labels even though a higher accuracy can be achieved for the last six labels (**Exp3**, 45.52% for $Tr_{U \sim Z}^v$ and 41.34% for $Te_{U \sim Z}^v$).

Table 2: Some results of the audio-visual transfer experiment.

The audio model was pre-trained with $\lambda = 0.1$, $\alpha_0 = 0.005$ and $\gamma = 0.8$ while the video model was pre-trained with $\lambda = 0.0004$, $\alpha_0 = 0.002$ and $\gamma = 0.96$. The transfer learning experiment was carried out under the setting $\lambda = 0.0004$, $\gamma = 0.8$, $K = 12$ and was trained for 160 steps. For **Exp1**, **Exp2** and **Exp3**, I used respectively $\alpha_0 = 0.001, 0.005, 0.001$ and $p_a = 0.85, 0.85, 1$.

For the meaning of different variables see the first paragraph of section 5 and also the text of 7.2. Notice that since $p_a = 1$ for **Exp3** no video data is given in input during fine-tuning.

	Tr^v	$Tr_{A \sim T}^v$	$Tr_{U \sim Z}^v$	Te^v	$Te_{A \sim T}^v$	$Te_{U \sim Z}^v$
No transfer	77.67 %	100 %	0 %	40.56 %	54.48 %	0 %
Exp1	81.17 %	98.28 %	21.64 %	39.44 %	47.76 %	15.22 %
Exp2	40.83 %	51.07 %	5.22 %	23.89 %	30.60 %	4.35 %
Exp3	19.67 %	12.23 %	45.52 %	12.22 %	2.24 %	41.34 %

8 Conclusion and Perspectives

In this report, I studied two different multimodal learning architectures and evaluated them on three different datasets. Several preliminary results were obtained and my work ensured the possibility and the possible benefit of taking into multiple modalities into account while dealing with real world problems. To go further in this direction, my work can be improved in several ways. First, datasets with bigger size and higher variety should give more reliable results. Second, a more exhaustive search of hyperparameters may lead to better outcomes. Finally, we should also pay attention on the modalities that we deal with and if the adopted method is really appropriate.

Ideally, we want to apply these techniques to real data coming from human-robot interactions. In this way, we may be able to improve the perception of the robot thanks to the presence of multiple sensors. Learning the shared representation in a totally unsupervised way can be useful for the development of an autonomous robot and absence of constant relation between different modalities can be compensated by the knowledge transfer approach.

References

- [1] M. Asadi-Aghbolaghi, A. Clapés, M. Bellantonio, H. J. Escalante, V. Ponce-López, X. Baró, I. Guyon, S. Kasaei, and S. Escalera. A Survey on Deep Learning Based Approaches for Action and Gesture Recognition in Image Sequences. In *Automatic Face & Gesture Recognition (FG 2017) 2017 12th IEEE International Conference*, 2017.
- [2] T. Baltrusaitis, C. Ahuja and L-P. Morency. Multi-modal Machine Learning: A Survey and Taxonomy. In *arXiv preprint arXiv: 1705.09406*, 2017.
- [3] Y. Bengio. Practical recommendations for gradient-based training of deep architectures. In *Neural networks: Tricks of the trade*, pages 437–478, 2012.
- [4] Y. Bengio, L. Yao, G. Alain and P. Vincent. Generalized denoising auto-encoders as generative models. In *Advances in Neural Information Processing Systems*, pages 899–907, 2013.
- [5] C. Bregler and Y. Konig. “Eigenlips” for robust speech recognition. In *ICASSP*, 1994.
- [6] A. Clark. Being there: Putting brain, body, and world together again. In *MIT press*, 1997.
- [7] L. Deng, G. Hinton and B. Kingsbury. New types of deep neural network learning for speech recognition and related applications: An overview. In *ICASSP*, 2013.
- [8] A. Droniou, S. Ivaldi, and O. Sigaud. A deep unsupervised network for multimodal perception, representation and classification. In *Robotics and Autonomous System*, 2014.
- [9] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A deep visual-semantic embedding model. In *NIPS*, 2013.
- [10] A. Graves. Generating sequences with recurrent neural networks. In *arXiv preprint arXiv:1308.0850*, 2013.
- [11] S. Harnad. The symbol grounding problem. In *Physica D*, vol. 42, pages 335–346, 1990.
- [12] G. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. Sainath and B. Kingsbury. Deep neural networks for acoustic modeling in speech recognition. In *IEEE Signal Processing Mag* Vol. 29, 2012.
- [13] S. Ioffe and C. Szegedy. Batch normalization, accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015.
- [14] S. Ji, W. Xu, M. Tang and K. Yu. 3D Convolutional Neural Networks for Human Action Recognition. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, 2013.
- [15] A. Karpathy, and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- [16] A. K. Katsaggelos, S. Bahaadini and R. Molina. Audiovisual fusion: Challenges and new approaches. In *Proceedings of the IEEE*, vol. 103, 2015.
- [17] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *arXiv preprint arXiv:1412.6980*, 2014.
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, 2012.
- [19] Y. LeCun, L. Bottou, Y. Bengio and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, vol. 86, pages 2278–2324, 1998.
- [20] Y. LeCun, K. Kavukcuoglu and C. Farabet. Convolutional networks and applications in vision. In *Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium*, pages 253–256, 2010.
- [21] A. Makhzani and B. Frey, K-Sparse autoencoders. In *International Conference on Learning Representations (ICLR)*, 2014.
- [22] J. Masci, U. Meier, D. C. San, and J. Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *Artificial Neural Networks and Machine Learning-ICANN*, 2011.
- [23] I. Matthews, T. F. Cootes, J. A. Bangham and S. Cox. Extraction of visual features for lipreading. In *PAMI*, 2002.
- [24] A. Memo, L. Minto and P. Zanuttigh. Exploiting Silhouette Descriptors and Synthetic Data for Hand Gesture Recognition. In *STAG: Smart Tools & Apps for Graphics*, 2015.
- [25] A. Memo and P. Zanuttigh. Head-mounted gesture controlled interface for human-computer interaction. In *Multimedia Tools and Applications*, 2017.
- [26] S. Mitra and T. Acharya. Gesture recognition: A survey. In *IEEE Systems, Man, and Cybernetics*, 37:311–324, 2007.
- [27] P. Molchanov, S. Gupta, K. Kim, and J. Kautz. Hand gesture recognition with 3D convolutional neural networks. In *CVPRW*, 2015.

- [28] S. Molholm, W. Ritter, M. M. Murray, D. C. Javitt, C. E. Schroeder, and J. J. Foxe. Multisensory auditory–visual interactions during early sensory processing in humans: a high-density electrical mapping study. In *Cognitive brain research*, 14:115–128, 2002.
- [29] S. Moon, S. Kim and H. Wang. Multimodal transfer deep learning with applications in audio-visual recognition. In *NIPS MMML Workshop*, 2015.
- [30] J. Nagi, F. Ducatelle, G. Di Caro, D. Ciresan, U. Meier, A. Giusti, F. Nagi, J. Schmidhuber, and L. Gambardella. Max-pooling convolutional neural networks for vision-based hand gesture recognition. In *Proceedings of the IEEE International Conference on Signal and Image Processing Applications*, 2011.
- [31] N. Neverova, C. Wolf, G. Paci, G. Somnavilla, G. W. Taylor, and F. Nebout. A multi-scale approach to gesture detection and recognition. In *Computer Vision Workshops (ICCVW)*, 2013.
- [32] N. Neverova, C. Wolf, G. W. Taylor, and F. Nebout. Multiscale deep learning for gesture detection and localization. In *ECCVW*, 2014.
- [33] J. Ngiam, A. Khosla, M. Kim, J. Nam, and A. Y. Ng. Multimodal deep learning. In: In *International Conference on Machine Learning*. 28. Bellevue, Washington, USA, 2011.
- [34] K. Noda, Y. Yamaguchi, K. Nakadai, H. G. Okuno and T. Ogata. Audio-visual speech recognition using deep learning. In *Applied Intelligence*, vol. 42, no. 4, pp. 722–737, 2014.
- [35] H. Noh, S. Hong and B. Han. Learning deconvolution network for semantic segmentation. In 2015 IEEE International Conference on Computer Vision (ICCV), pages 1520–1528, 2015.
- [36] J. Piaget. *The origins of intelligence in children*. WW Norton & Co, 1952.
- [37] L. Pigou, S. Dieleman, P. J. Kindermans, and B. Schrauwen. Sign language recognition using convolutional neural networks. In *Workshop at the European Conference on Computer Vision*, pages 572–578, 2014.
- [38] G. Potamianos, C. Neti, J. Luetttin and I. Matthews. Audio-visual automatic speech recognition: An overview. In *Issues in Visual and Audio-Visual Speech Processing*, MIT Press, 2004.
- [39] N. Pugeault, and R. Bowden. Spelling It Out: Real-Time ASL Fingerspelling Recognition. In *Proceedings of the 1st IEEE Workshop on Consumer Depth Cameras for Computer Vision*, 2011.
- [40] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *arXiv preprint arXiv:1511.06434*, 2015.
- [41] L. Smith and M. Gasser. The development of embodied cognition: Six lessons from babies. In *Artificial life*, 11:13–29, 2005.
- [42] R. Socher, M. Ganjoo, H. Sridhar, O. Bastani, C. D. Manning, and A. Y. Ng. Zero-shot learning through cross-modal transfer. In *International Conference on Learning Representations (ICLR)*, Scottsdale, Arizona, USA, 2013.
- [43] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. In *Journal of Machine Learning Research*, vol. 15, 2014.
- [44] T. Starner, A. Pentland, and J. Weaver. Real-time american sign language recognition using desk and wearable computer based video. In *PAMI*, 20(12):1371–1375, 1998.
- [45] J. Sung, I. Lenz, and A. Saxena. Deep multimodal embedding: Manipulating novel objects with point-clouds, language and trajectories. In *Robotics and Automation (ICAR), 2017 IEEE International Conference*, pages 2794–2801, 2017.
- [46] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- [47] C. Szegedy, S. Ioffe, V. Vanhoucke. Inception-v4, inception-ResNet and the impact of residual connections on learning. In *AAAI*, pages 4278–4284, 2017.
- [48] V. Turchenko, E. Chalmers and A. Luczak. A deep convolutional auto-encoder with pooling – unpooling layers in Caffe. In *arXiv preprint arXiv:1701.04949*, 2017.
- [49] P. Vincent, H. Larochelle, Bengio and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine Learning (ICML)*, pages 1096–1103; 2008.
- [50] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio and P.-A. Manzagol. Stacked denoising autoencodes: Learning useful representations in a deep network with a local denoising criterion? In *The Journal of Machine Learning Research*, vol. 11, pages 3371–3408, 2010.
- [51] J. Weng, J. McClelland, A. Pentland, O. Sporns, I. Stockman, M. Sur, and E. Thelen. Autonomous mental development by robots and animals. In *Science*, vol. 291, no. 5504, pages 599–600, 2001.
- [52] J. Weston, S. Bengio, and N. Usunier. Large scale image annotation: learning to rank with joint word-image embeddings. In *Machine Learning*, 81(1):21–35, 2010.
- [53] D. Wu, L. Pigou, P.-J. Kindermans, N. D.-H. Le, L. Shao, J. Dambre, and J.-M. Odobez. Deep Dynamic Neural Networks for Multimodal Gesture Segmentation and Recognition. In *Pattern Analysis and Machine Intelligence IEEE Transactions*, vol. 38, 2016.
- [54] B. P. Yuhas, M. H. Goldstein, and T. J. Sejnowski. Integration of acoustic and visual speech signals using neural networks. In *IEEE Comm. Magazine*, pp. 65–71, 1989.
- [55] M. D. Zeiler, G. W. Taylor and R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In *Computer Vision (ICCV), 2011 IEEE International Conference*, pages 2018–2025, 2011.