

Real time hand gesture recognition in industry

Dumoulin William, Thiry Nicolas, Slama Rim, Bernard Michel
Haute Ecole d'ingénierie industrielle de l'Henallux Pierrard , Virton, Belgium

Abstract—With the 4th industrial revolution and the increased use of cobots in the industries comes many opportunities for new generation control panels. The work presented here aims to develop a deep learning model to recognise 10 different gestures that can be used to interact with a cobot. The data set fed to the model was completely created to be efficient in this context in particular. The videos were taken from a computer webcam and then processed to remove the noise created by the background by isolating the movement of the gray scale images. We proposed to extract the spatio-temporal features are extracted by the combination of 3D convolution and LSTM layers. We also proposed a real time method to recognize our gestures and we obtained for 8 out of 10 gestures, a recognition rate of more than 90%. Furthermore, an interface was created to test our method and to add new class of gestures to be recognized by our model.

I. INTRODUCTION

While working with image recognition by a computer, the main challenge that we can face is in the case of a real time application. The computational time has to be low enough to keep the program running without lags.

There are many applications for a gesture recognition program in different areas: To translate the sign language or to remove the need of a remote control for the television at home. Our program is developed for an application in industry. For instance, to control a cobot without a control panel.

In this paper, we proposed an interface that we used to acquire our dataset with specific gestures that can be used in industrial context. We also proposed to use a convolutional approach to extract the features and LSTM to classify gestures. Finally, we proposed a real time approach to recognize gestures using a friendly interface. The interface mainly shows a feed-back of the camera with the predictions from the model. There is also the possibility to create new videos and add new gestures.

II. STATE-OF-THE-ART

Many different techniques were developed in the past years. One of the first was an approach that separates the fingers from the palm [1]. The position of the hand on the static image was then deduced depending on the number of fingers found by the algorithm and their relative position. A much more developed method is to apply a skeletonization method [2], [3]. The idea is to reduce an object to lines that are only one pixel wide. The images are therefore easier to analyse by the model. This method is applied to the videos before the training and the tests. An algorithm to create the skeleton of a shape need a Voronoi tessellation [4] but it's a heavy algorithm that slows down the entire process. A solution would be to use the champfer distance transform [5] to reduce the computational cost.

N. Luthfil Hakim, T. K. Shih et al [6] created a model of gesture recognition for a television application. They used an RGB and a depth camera to get more valuable information. Their model is composed by multiple blocs of 3D convolution layers combined with dense and max pooling layers. The last part of their model is the Long Short-Term Memory layers [7] that finds features that are time related which are crucial when dealing with videos.

Another basic method is the haar-like features [8] that are used for face or eyes detections. This method is called a hand crafted feature extraction method. If the color changes are not sharp enough, this method will show no results. The LSTM layer is an advanced type of RNN [9]. Enhanced methods of Recurrent Neural Network that could be used for gesture recognition were developed : Differencial RNN [10], Hierarchical RNN [11], Bidirectionnal RNN [12].

Some articles [6], [13] use the kinect camera to get the depth information in addition to the 2D images.

Transformers [14], [15] are a type of layers that is designed to handle sequential data. Unlike LSTMs, they use a mechanism to give different weighs to the input data.

III. METHODS

A. Data set

In order to collect a dataset with gestures in industrial context in particular, we created our own dataset with gestures that can be used to interact with a cobot. To simplify the process, we made a python script to save the video in a folder and the label in a CSV file. The program is made to create a lot of data in the smallest amount of time. The user can switch between different labels by pressing a key and can start a new video with an other key. See details on Figure 1. There is a delay of one second before the video is taken to give the time to the user to place himself. The computers camera have a frame rate of 25 FPS so each video last 2.4 seconds and contains 61 frames. With a good rhythm, a new video can be made every 5 seconds.

A data set need recordings of different people otherwise the model would only recognise the gesture of the person that did all of the data set videos. The smallest data set need at least 12 different people and each one does each gesture twice. This first project has 10 different gestures which make a minimum of 240 videos for the training set. It is of course not enough so multiple people filmed more videos. The result is a training set of almost 1 500 videos.

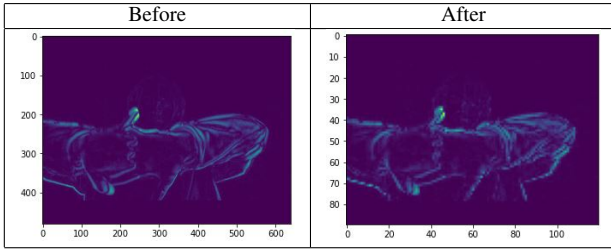
B. Data processing

The videos from the data set have 61 frames with a dimension of 640 pixels wide by 480 pixels height and 3



Fig. 1. New video program

TABLE I
NORMALIZATION



layers of depths for the red, green and blue channels. First step, conversion from RGB to gray image: Each pixel of the new gray image is calculated with the colors channels of source image with this formula: $Gray \leftarrow 0.299 * Red + 0.587 * Green + 0.114 * Blue$.

The second step is to isolate the movement in the images. To do that, each image is compared with the next one. The movement of the last frame can't be isolated because there are no more frame to do the comparison. That make 60 frames of movement. The pros of keeping only the movement is to remove the background. Of course, to make the training data, the background need to be still, otherwise it will be considered as part of the gesture. This technique still reduces the noise that can be produced by the background.

The model doesn't get 60 frames for each video because it would make an over complicated model. The third step is to isolate a sample of the frames. The choice was made to take only 10 fps, that make a sample of 24 frames.

The fourth step is to resize and normalize the videos. The normalization take all the pixels and map them from 0-255 to 0-1. After all the processing, the videos have 24 frames with a dimension of 120 pixels wide by 90 pixels height. All the pixels have a value between 0 and 1. See Table I for comparison.

IV. MODEL

A. Inspiration

Our first model is based on the work of N. L. Hakim, T. K. Shih, et al [6]. They have an RGB and a depth camera. The images from the two cameras are processed separately and combined after the LSTM units [16]. The image processing is made with multiple 3D convolution layers. The team had a data set with thousands of videos and got a result of 91%.

B. Model improvement and architecture

The results of our first architecture were not conclusive. The problem was probably the small data set and the large amount of parameters. We simplified their model to create our own. We trained 15 different models, each one had a small difference compare to the last one in order to understand how each parameter affects the model. Its accuracy raised from 26% to 83%. A clear improvement in the quality of the models can be seen by the increase of their accuracy : the 4th model showed an accuracy of 46% while the 12th showed an accuracy of 79%.

Our final model is composed of 3 blocs: The first two are an arrangement of Convolution 3d, dense, dropout, maxpooling 3d and batch normalization. The last one has two ConvLSTM2D followed by a long-short term memory and a dense layer to have as many outputs as we have classes.

TABLE II
EDGE DETECTION KERNEL

| | | |
|----|----|----|
| -1 | -1 | -1 |
| -1 | 8 | -1 |
| -1 | -1 | -1 |

1) *3D Convolution* [17], [18]: In the convolution layers, new channels are created with different kernels. A kernel is a squared matrix, whose size is usually 3x3 or 5x5. To analyse videos, it can be better to use 3D convolution layers, then the size are 3x3x3 or 5x5x5. The goal of this layer is to find the features of the images or videos. For example the Table II shows an edge detection kernel.

2) *Dense*: The dense layer is also called fully-connected layer. Each node of a layer is connected to all the nodes of the next layer. In our case a node is an entire image. Each channel of the last dimension is also an image. For example, the dimensions of the first dense layer are (?x24x90x120x256) which stands for 24 moments in time, 90 pixels height by 120 pixels wide and 256 images. The different images for one moment in time are created by the convolution layer.

The connection is defined by a weight with a value between minus one and one. In this architecture, the dense layer is placed after a convolution one. We can compare the dense layer to a filter. The meaningful channels created in the 3D convolution layer have a weight with a high value, the useless ones have a weight close to zero.

Of course this is a deduction of the system, we cannot prove it because the thousands of kernels and weights in the model does not make sense for a human.

3) *Dropout*: When this layer is placed after a dense layer, it removes a certain percentage of the weights by giving them new random values. The goal is to reduce the risk of overfitting.

4) *Max pooling 3D*: This layer reduces the size of the data set by keeping the max value of a section of an image (or multiple images when it's 3D). A 3D Max pooling layer with a matrix (3x3x3): The dimension go from (24,90,120,256) to (24,30,40,85) because the three last dimensions are divided by 3. This layer simplifies the data by keeping only the meaningful pixels.

5) *Batch Normalization*: The layer transforms inputs so that they are standardized. This means that they have a mean value of zero and a standard deviation of one.

6) *LSTM [7]*: Long Short-Term Memory layers are a type of recurrent neural network capable of finding time dependent features. This layer (and the ConvLSTM) is very important in this application because there are needed to differentiate "Swipe to the left" and "Swipe to the right" for example.

7) *ConvLSTM2D*: Like the classic LSTM, the convLSTM is capable of learning order dependence in sequence but this layer is specialized for the videos analyses because it's the combination of LSTM and convolution layer.

In this architecture, two ConvLSTM 2D are stacked. The first ConvLSTM layer provides an output for each input while the second layer gives only one output for a sequence of input.

8) *Reshape*: The layer doesn't learn anything in the training process. It only changes the dimension of the data to pass them to the LSTM layer which needs specific input.

C. Training

The videos are stored raw, they need to be processed when they are loaded to train the model. The data set of 1 500 videos is too heavy to be loaded in the RAM so the videos have to be load by package. The label of the videos are saved in a list which is randomly mixed. The videos can then be loaded in the order of the list. It's important that the model is trained on different gestures every time.

To train the model faster, the data processing and training are on two different threads. The process thread loads two packs of 40 videos and waits for the training part to be done with the previous pack. When the model is done training with a pack, the next pack is transferred to the training section and the thread loads the next one.

The 2 tasks are handled by the computer in parallel, that's what is called thread. It allows to reduce the overall time of computing. The videos are loaded by packs to avoid overloading the computer's memory. The packs have enough information to allow multiples training. However, if a model is trained multiple times on the same pack, the risk of overfitting increases.

To reduce this risk, the model should be trained on the entire data set once before being trained on videos that it already saw.

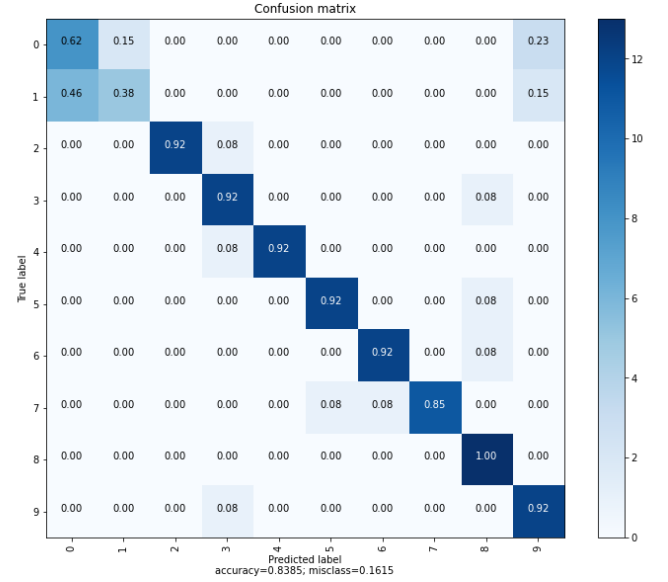


Fig. 2. Confusion matrix

The ideal process would be to train the model on the entire data set in a single pack. Since it is not possible due to memory size, the videos should be loaded by smaller packs. Each one of them should then train the model once and leave its place to the next pack. If multiple epochs are mandatory, the whole process needs to start again.

Multiple epochs can be necessary when the model has a bad accuracy.

Our algorithm doesn't use that process because each videos would be processed multiple times and the training would be unnecessarily long. Our model is trained on each pack three times in a row before moving to the next pack.

V. EXPERIMENTS RESULTS

A. Accuracy

The Figure 2 shows the confusion matrix of our best model. The numbers on the columns and lines of the confusion matrix correspond to this list:

TABLE III
LIST OF GESTURES

| | |
|---|--|
| 0 | Closing one hand |
| 1 | Opening one hand |
| 2 | Clap with both hands |
| 3 | The right hand swipes to the left |
| 4 | The left hand swipes to the right |
| 5 | Crossed fingers |
| 6 | The hands start closed in the middle of the image and they go in opposite directions while they open |
| 7 | The right arm starts horizontally and rotates around the elbow, it ends when the hand is pointing to the sky |
| 8 | The left arm starts horizontally and rotates around the elbow, it ends when the hand is pointing to the sky |
| 9 | Say "hello" with one hand |

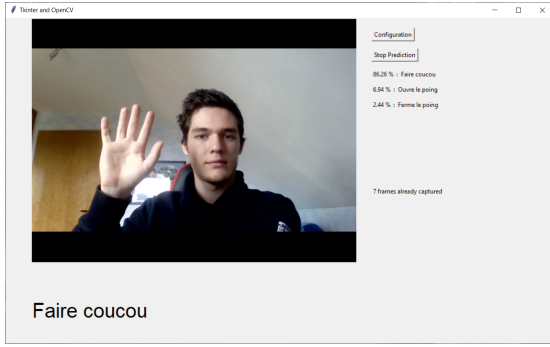


Fig. 3. First Window

The first two gestures are often confused with one another. This is probably due to the proximity of these gestures. We can see that closing a hand and opening a hand is done at the same places on the videos. Furthermore, the difference between closing and opening a hand is much more subtle than the difference between swipe left and swipe right.

For all the other gestures, the model is very efficient. On the live test we often got a precision of +90%.

VI. APPLICATION

A. Real time recognition

The application works with threads, there is one for the prediction and another one for the interface. The videos are recorded at 10 frames per seconds, which makes a pack of 24 frames, and then are sent to the prediction. 24 frames is the number of images that were given per video to train the model. The frames are processed directly when they are taken, the list of processed images is sent to the prediction thread. A CPU Intel i7-1065G7 takes 1.5s while a GPU Nvidia gtx1050 takes 0.4s to get the prediction done. More predictions could be done in parallel. For example send a video for prediction every 12 frames. The 12 other frames would be kept from the end of the previous video. However it could create some lag.

B. Interface

The interface was coded using the python libraries tkinter and openCV.

The main window, that can be seen at the Figure 3, shows the live feed from the camera. It also shows the prediction that have the highest score at the bottom of the screen. Furthermore the window shows the 3 first predictions for the same video with the percentage of recognition. If the first prediction given by the model is lower than 60%, the predictions are not displayed and a message is written on the screen explaining that the model did not recognise the gesture.

The last information that is shown is the number of frames that are already captured in the video. This allows the person in front of the camera to perform the full gesture on one video. As soon as the video is fully captured, another video is being captured.

The main window also shows two buttons. The first one freezes the prediction. In other words, the videos are no longer

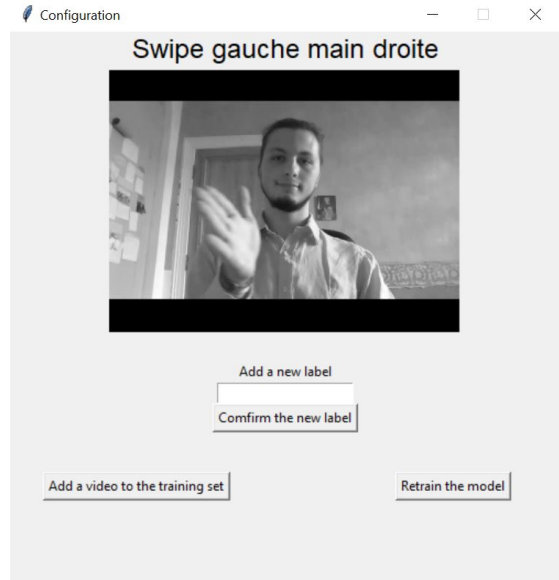


Fig. 4. Second window

captured and the prediction stays on the window. The second button opens a second window. While the second window is active, the first is completely frozen to avoid wasting the resources of the computer.

The second window, that can be seen at the Figure 4, displays an example of each gesture that the model is able to recognise with its name. It also allows the user to add a new gesture. The first step is to create a new label or title for the gesture and add it to the list through the text field on the window. If the label already exists on the list, a message is displayed on the screen.

On this window, there are two buttons. The first one opens program to record a video and add it to the training set. This program is explained at METHODS: A. Data set. The second button retrains the model. This operation takes a lot of time. Therefore, a safety needs to be installed to avoid miss-clicks.

VII. CONCLUSIONS AND FUTURE WORKS

A. Conclusions

The first method we considered was the skeletonization. The pros is a simpler model with less parameters because the skeleton is only one pixel wide with no colors information. The cons is that the process of skeletonize an image is very heavy and the gain of the simpler model is lost by the computational time of the skeletonization. Furthermore, the skeletonization is heavy in the training and the testing part while a big model takes time to train but is not a problem for the testing part.

To get rid of the noise made by the background, we isolate the movements in the videos. Then, we reduced the frame rate and the size of the images. The resulting data is fed to the model. The 1500 videos used to train the models were enough to get an overall accuracy of 84% and 90% for 8 gestures out of 10.

B. Future Works

While doing the tests of the complete project, we stumbled across ways to perform the gestures that our model had trouble recognising. These issues are the distance between the camera and the hands and speed at which the gesture is performed. The solution is data augmentation. It is the process to create data from already existing data. This method can then be used to increase the distance from the camera by resizing the images and adding black pixels on the sides to keep the same data size. The process is the same to solve the speed issue, it is the time that is resize instead of the space.

We can train models with different LSTM characteristics. After we have found the best solution, we can modify the characteristics of the data set fed to the model. These tests could also help us find if it is better to have a higher definition or a higher frame rate. Maybe there are both too computationally expensive and therefore not an interesting modification. When the characteristics of the new model are well defined, we can fix the limits of the model by increasing the size of the data set. The process is called data augmentation and is well described in the section "Limits of the model"

REFERENCES

- [1] Z.-h. Chen, J.-T. Kim, J. Liang, J. Zhang, and Y.-B. Yuan, "Real-time hand gesture recognition using finger segmentation," *The Scientific World Journal*, vol. 2014, 2014.
- [2] R. L. Ogniewicz and M. Ilg, "Voronoi skeletons: theory and applications," in *CVPR*, vol. 92, 1992, pp. 63–69.
- [3] B. Ionescu, D. Coquin, P. Lambert, and V. Buzuloiu, "Dynamic hand gesture recognition using the skeleton of the hand," *EURASIP Journal on Advances in Signal Processing*, vol. 2005, no. 13, pp. 1–9, 2005.
- [4] A. Beristain Iraola, *Skeleton based visual pattern recognition. Applications to tabletop interaction*. Servicio Editorial de la Universidad del País Vasco/Euskal Herriko ..., 2009.
- [5] M. A. Butt and P. Maragos, "Optimum design of chamfer distance transforms," *IEEE Transactions on Image Processing*, vol. 7, no. 10, pp. 1477–1484, 1998.
- [6] N. L. Hakim, T. K. Shih, S. P. Kasthuri Arachchi, W. Aditya, Y.-C. Chen, and C.-Y. Lin, "Dynamic hand gesture recognition using 3dcnn and lstm with fsm context-aware model," *Sensors*, vol. 19, no. 24, p. 5429, 2019.
- [7] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [8] Q. Chen, N. D. Georganas, and E. M. Petriu, "Real-time vision-based hand gesture recognition using haar-like features," in *2007 IEEE instrumentation & measurement technology conference IMTC 2007*. IEEE, 2007, pp. 1–6.
- [9] C. L. Giles, G. M. Kuhn, and R. J. Williams, "Dynamic recurrent neural networks: Theory and applications," *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 153–156, 1994.
- [10] V. Veeriah, N. Zhuang, and G.-J. Qi, "Differential recurrent neural networks for action recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4041–4049.
- [11] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1110–1118.
- [12] L. Pigou, A. Van Den Oord, S. Dieleman, M. Van Herreweghe, and J. Dambre, "Beyond temporal pooling: Recurrence and temporal convolutions for gesture recognition in video," *International Journal of Computer Vision*, vol. 126, no. 2, pp. 430–439, 2018.
- [13] Z. Ren, J. Meng, J. Yuan, and Z. Zhang, "Robust hand gesture recognition with kinect sensor," in *Proceedings of the 19th ACM international conference on Multimedia*, 2011, pp. 759–760.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [15] M. De Coster, M. Van Herreweghe, and J. Dambre, "Sign language recognition with transformer networks," in *12th International Conference on Language Resources and Evaluation*, 2020.
- [16] A. Ullah, J. Ahmad, K. Muhammad, M. Sajjad, and S. W. Baik, "Action recognition in video sequences using deep bi-directional lstm with cnn features," *IEEE access*, vol. 6, pp. 1155–1166, 2017.
- [17] S. Ji, W. Xu, M. Yang, and K. Yu, "3d convolutional neural networks for human action recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 1, pp. 221–231, 2012.
- [18] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4489–4497.