

Sign Language ASL Recognition using Bag of Visual Words and Support Vector Machines

Samuel Ortega Cuadra^a and Laia Alexandra Sjöberg Cerezo^b

^a1669776

^b1667894

Abstract—This project focuses on developing a system for recognizing American Sign Language (ASL) alphabet and numbers. The proposed approach combines image preprocessing, feature extraction, and prediction models to accurately classify hand gestures. The dataset consists of images of hands representing ASL letters and numbers, which are processed by isolating the hand region, converting it in grayscale and applying a mask for feature extraction. Dense SIFT and a Bag of Words model are used to create feature representations, and a Support Vector Machine (SVM) classifier is trained and optimized through grid search for maximum accuracy. The finished system can recognize hands by retrieving frames from video input, and translating ASL gesture sequences into coherent text. This work lays a foundation for future developments in sign language recognition and improving accessibility for the deaf community by providing a tool for real-time translation of sign language to text.

Contents

1	Introduction	1
2	Finding The Dataset	1
2.1	Dataset Preprocessing	1
2.2	Masking the Hand	1
3	Feature extraction	2
3.1	Keypoints	2
3.2	The Codebook	2
4	Predicting using the codebook	2
4.1	From Descriptors to Histograms	2
4.2	SVM	3
5	Video Processing	3
6	Hand Detection	3
7	Final Results	4
8	Future Improvements	4
8.1	Category Coherence	4
8.2	Transition Frames	4
8.3	Model Improvement	4

1. Introduction

This document presents the development of a system for recognizing American Sign Language (ASL) alphabet and numbers. It outlines the steps taken to achieve this goal, from the **dataset collection** and preprocessing to the implementation of a **Bag of Visual Words** to help building a **SVM model** for classification. The system is then extended to **process videos**, extract frames, identify ASL symbols, and **reconstruct** the words being spelled.

2. Finding The Dataset

When looking for a good ASL dataset, we found the [American Sign Language Dataset](#) dataset on Kaggle. This dataset contains 2515 images of 400x400 pixels, each showing a hand making an ASL letter or number. The dataset is divided into 36 folders, one for each letter of the alphabet and one for each of the numbers from 0 to 9. The images are in JPEG format and have a black background with a white hand. The dataset is well-organized and easy to use, making it a good choice for our project.

2.1. Dataset Preprocessing

The preprocessing of the dataset involves the following steps:

- Dataset Organization:** The dataset is organized into folders, one for each letter of the alphabet and one for each number from 0 to 9. Each folder contains 65 to 70 images of hands making the corresponding ASL symbol. A *categories* dictionary was created to map each label to its corresponding folder name. The filenames of all images were collected and combined with their respective category labels to construct a DataFrame.
- Shuffling the Dataset:** To avoid any bias in the training and testing data, the dataset was shuffled. This ensures that the data is randomly distributed and that the model will not be trained on any specific order of images.
- Image Loading:** Each image was loaded from its respective path and converted into a NumPy array. Providing a new variable with a structured representation of pixels called *pixel_data*.
- Grayscale conversion:** The images were converted to grayscale to reduce the complexity of the model and improve training speed.

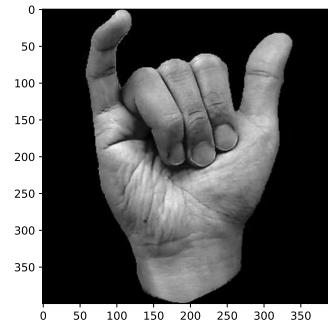


Figure 1. Grayscale original image

2.2. Masking the Hand

In order to find the best features we tried to mask the hand in the image using different methods. The optimal resulting method was surprisingly the easiest. We took a look at the first 20 pixels of the image before converting it to grayscale and we found the most frequent color. Subsequently, we applied a standard deviation and removed all of the background with the same colors, which will be really important when extracting the keypoints in the next chapter.

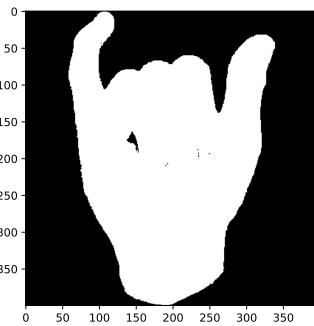


Figure 2. Mask of the image

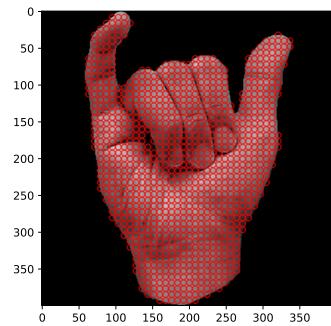


Figure 4. Mask applied Dense SIFT

Information

We tried other methods such as brightness threshold, gaussian blur and several others from the computer vision library. They all ended up being less efficient than the one we finally used.

3.2. The Codebook

Once we have all the descriptors from our images, we can create a codebook. The idea behind this process is to create a "dictionary of visual words" that will represent the most common features in our dataset. This will help us find the meaning behind the keypoints we have detected.

In order to achieve this we implemented the KMeans algorithm. This will group all the descriptors found into 200 clusters which will define our codebook. This will later be used to create the histograms that will represent our images in the prediction step.

Information

The first number of clusters used was 36, which is the number of classes in the dataset. This was a mistake, as the codebook was not able to differentiate between the different classes. We then tried with increasing the number of clusters, which gave much better results.

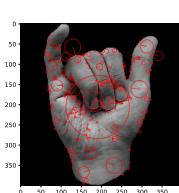
45 3. Feature extraction

46 3.1. Keypoints

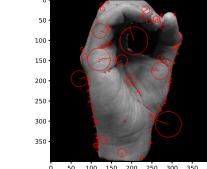
47 In the implementation of Bag of Visual Words there are two main
48 steps to follow once you have your images: Detecting the keypoints
49 and extracting the descriptors.

51 We used the SIFT algorithm to detect the keypoints in the image.
52 This algorithm is one of the most implemented in the field of
53 computer vision because of how robust it is. We first tried to use
54 `detectAndCompute()` (see 3a) but this created a huge problem,
55 when taking the background into consideration. This process is
56 non-positional, it does not matter if it finds a hole in the middle of the
57 image, (see 3b), it will be treated the same as a space in the left.
58

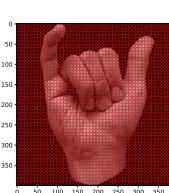
59 After this first approach we tried to implement **Dense SIFT** (3c),
60 which improved the results considerably (increasing the final accuracy
61 by almost 15%) but in the end the game changer was the mask
62 we applied in the previous chapter. This allowed us to remove the
63 background and focus on the hand, which allowed us to detect only
64 the keypoints we were most interested in.



(a) Simple detect and compute

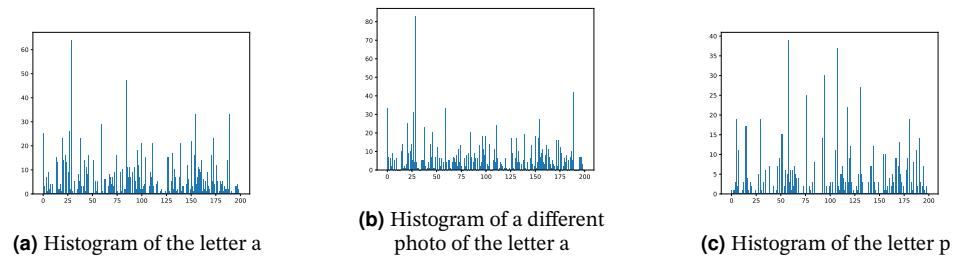


(b) Background mistaken keypoints



(c) Dense SIFT

Figure 3. Several keypoints detection methods considered

**Figure 5.** Comparison of histograms of different images**4.2. SVM**

There are lots of different models that can be used to predict the categories of these histograms, the one we ended up choosing is SVM (Support Vector Machines). After doing a split between the train and the test we tried a starting model with a linear kernel. After this first approach, we tried to improve the model by using a grid search to find the best hyperparameters. After looking at the classification reports, this are the results:

	Precision	Recall	F1-score
Simple model	0.96	0.96	0.96
Grid search	0.97	0.97	0.97

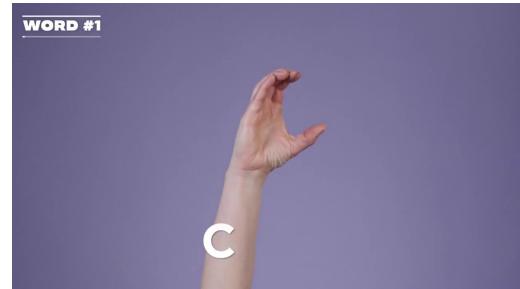
Once we have the complete process, we compressed all the operations in one unique function `predict_image`

```

1 def predict_image(img_path, normalizer, km,
2     best_model):
3     # Load the image from a path
4     img = Image.open(img_path)
5     numpydata = asarray(img)
6
7     # Get the mask of the image
8     mask = get_mask(numpydata)
9
10    # Convert to grayscale
11    if len(numpydata.shape)==3:
12        gray_data = numpydata.mean(axis=2)
13        gray_data = gray_data.astype(np.uint8)
14    else:
15        gray_data = numpydata.astype(np.uint8)
16
17    # Get dense keypoints
18    dense_kp = generate_dense_keypoints(
19        gray_data, mask)
20
21    # Compute descriptors
22    _, des = sift.compute(gray_data, dense_kp,
23        mask)
24
25
26    # Create the histogram of the image
27    im_features = np.zeros(k, dtype=np.float32)
28
29    if des is not None:
30        for feature in des:
31            feature = feature.reshape(1, -1)
32            idx = kmeans.predict(feature)
33            im_features[idx] += 1
34
35    # Normalize the histogram
36    im_features = normalizer.transform(
37        im_features.reshape(1,-1))
38
39    # Predict the category using the svm model
40    y_pred = best_model.predict(im_features)
41
42    return categories[y_pred[0]]
```

Code 1. Predict Image Function.**5. Video Processing**

Now we have a model that can correctly predict the category of an image. We can use this model to predict the category of each frame of a video¹. To do this we will use the `cv2.VideoCapture()` function together with a ratio variable to take 1 of every few frames. This will allow us to process the video faster and will help us to avoid the transition frames between signs.

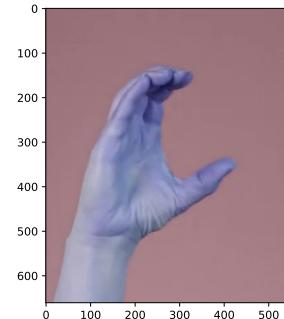
**Figure 6.** Frame of the spelling video

Once we have extracted some frames from the video, we need to process it into images we can predict.

6. Hand Detection

As it seems, there is a lot of noise in the frames extracted (text, coloured background, etc.). In order to have a better accuracy when predicting the ASL symbols, we first need to detect the hand in the frame.

To do this we used the hands library from MediaPipe in order to detect the hand in the frame.

**Figure 7.** Hand detected in the frame

Once we have the cropped image we apply the same mask as before

¹The video used corresponds to [ASL Fingerspelling Exercise](#) from [Learn How to Sign](#)

120 to remove the background and focus on the hand, in order to predict
121 the result with better accuracy.

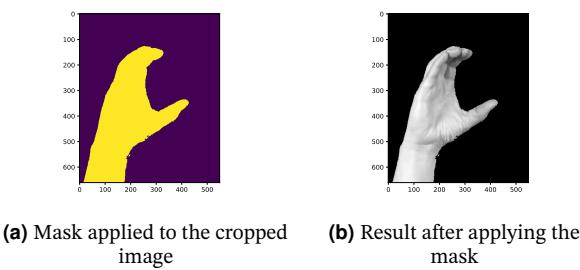


Figure 8. Applying the mask to the cropped image

122 Once we have removed the background, just as way to standardize
123 our photos, the images will be resized into 400x400px, the size of the
124 dataset images.

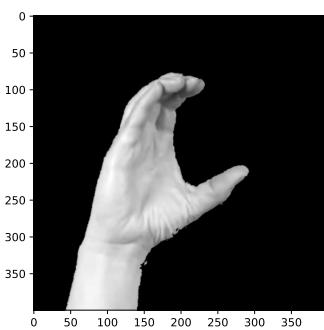


Figure 9. Final preprocessed hand

7. Final Results

126 After doing a prediction of all the images the results presented some
127 problems. One of the main ones was the transition and "no sign"
128 frames that appeared between the different signs. This derived in
129 a lot of noise in the final prediction. In order to adjust this, one
130 solution was to adjust the ratio of frames, or take a majority voting of
131 several sections of the prediction.

133 However, a second problem appeared. Some signs, such as the
134 "A" and the "T" or the "I" and the "J" are really similar which created
135 confusion in the prediction on the model.

136 The predictor is able to detect and classify particularly selected
137 frames, nevertheless, the solution to the noise and confusion prob-

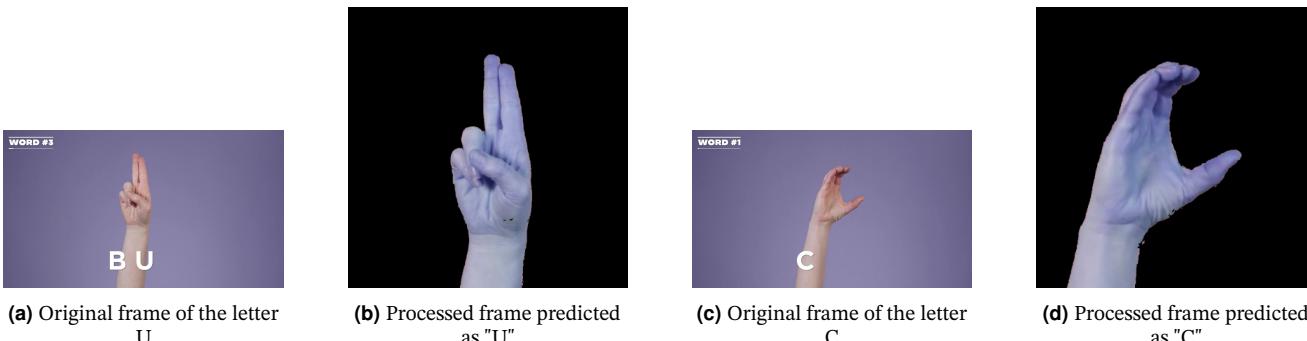


Figure 10. Image predictions of the video

139 lems is still under development. Here are some of the results of the
140 prediction (see Figure 10).

8. Future Improvements

8.1. Category Coherence

142 One of the main problems we have found when predicting is the
143 different ways that exist to sign the same letter. This is a problem that
144 can be solved by adding more variability to the images that train the
145 model. We can appreciate this problem if we compare the letter "G"
146 in the video and the dataset.

8.2. Transition Frames

148 The problem with the transition frames has many solutions, but none
149 of them has ended working for the purpose of the project. A future
150 upgrade to the frame extraction step will be selecting only the frames
151 that present a really small variance with the previous one. The idea
152 behind this solution is to have into account that when signing, the
153 hand will not move a lot between the frames corresponding to the
154 same letter.

8.3. Model Improvement

156 The model can be improve by fine tuning the model even more or
157 trying different classifiers. The solution we have in mind is imple-
158 menting a Neural Network to predict the sign, but that solution will
159 be applied in future approaches to this problem.