

Signing as Input for a Dictionary Query: Matching Signs Based on Joint Positions of the Dominant Hand

Manolis Fragkiadakis, Victoria Nyst, Peter van der Putten

Leiden University

Nonnensteeg 1-3 2311VJ, P.N. van Eyckhof 3 2311BV, Niels Bohrweg 1 2333CA

m.fragkiadakis@hum.leidenuniv.nl, v.a.s.nyst@hum.leidenuniv.nl, p.w.h.van.der.putten@liacs.leidenuniv.nl

Abstract

This study presents a new method to search sign language lexica, using a full sign as input for a query. Thus, a dictionary user can look up information about a sign by signing the sign to a webcam. The recorded sign is then compared to potential matching signs in the lexicon. As such, it provides a new way of searching sign language dictionaries to complement existing methods based on (spoken language) glosses or phonological features, like handshape or location. The “find the sign” method analyzes the recorded sign using OpenPose to extract the body and finger joint positions. To compare the recorded sign with the signs in the database, the variation in trajectories of the dominant hand and of the fingers is quantified and compared, using Dynamic Time Warping (DTW). The method was tested with ten people with various degrees of sign language proficiency. Each subject viewed a set of 20 out of 100 total signs from the newly compiled Ghanaian Sign Language lexicon and was asked to replicate the signs. The results show that our method can predict the matching sign with 87% and 74% accuracy at the Top-10 and Top-5 ranking level respectively by using only the trajectory of the dominant hand. Additionally, more proficient signers obtain 90% accuracy at the Top-10 ranking. The methodology has the potential to be used also as a variation measurement tool to quantify the difference in signing between different signers or sign languages in general.

Keywords: sign language lexica, search functionality, variation measurement

1. Introduction

In most sign language dictionaries, users can search a sign through a written gloss, a unique identifier that by definition refers to a sign. In some cases, the lexica offer the possibility to specify formal parameters of the target sign, for instance, its handshape and location (Figure 1). The Flemish Sign Language (VGT) dictionary (Van Herreweghe et al., 2004), the Swedish Sign Language (Institutionen for Lingvistik, 2009) and the Danish Sign Language (Center for Tegnsprog, 2008) are some examples of such dictionaries. After the input, the user is offered a set of signs that match the selected properties which can be then viewed individually.

Although sign search functionality on the basis of a sign parameter value is a useful attribute of sign language lexica, dictionary compilers still have to link these values to the videos. Also, as Zwitterlood discusses, the users of such dictionaries must “abstract away from the sign as a whole” if they want to use the parameter search functionality (Zwitterlood, 2010). Even then, only signs that match the query 100% are returned, and there is no concept of an ordered set of results that match to some degree. A thorough overview of sign language lexica and their features can be found in Zwitterlood’s review (2010).

In this paper we describe our “find the sign” methodology that allows inputting a full video-recorded sign to search for entries in a dictionary. This method requires no training of any kind of model such as the ones used for sign language recognition tasks. In its core, it is a comparison method to quantify the difference in the movement between signs. As a result, it can be used for any sign language. By utilizing a pre-trained pose estimation framework we extract the body and hand joint positions from users using their webcam. Subsequently, by employing Dynamic Time Warping we find the closest matching signs from a compiled lexicon.

To date, this methodology has only been applied to sign language classification tasks (Jangyodsuk et al., 2014; Schneider et al., 2019; Ten Holt et al., 2007) and not as a mode to complement sign search possibly solving the problem of ordering retrieval previously discussed. Additionally, we have developed a visualization tool to allow researchers to view the rendered paths of the dominant hand to further explore the overall difference in signing movements.

The paper is structured as follows: in Section 2 we give an overview of methods that utilize Dynamic Time Warping in the gestural and sign language domain. In Section 3 we describe our methodology regarding the extraction of the body joint coordinates as well as the experimental setup, analysis, and visualization tool. In Section 4 we present the results of our experiments. We discuss them in Section 5 and conclude and motivate future research in Section 6.



Figure 1: Traditional search functionality as seen in the online Danish Sign Language dictionary (Center for Tegnsprog, 2008).

2. Related Work

Dynamic Time Warping (DTW) is a dynamic programming based time series comparison algorithm to produce a distance metric between two inputs. It has been widely used in the speech recognition domain since the early 1970's (Abdulla et al., 2003; Axelrod and Maison, 2004; Myers et al., 1980). While the original algorithm can be computationally expensive, different variations have been developed over the years to reduce the overall complexity, with most notably the works of Itakura-Parallelgram (Itakura, 1975), Ratanamahatana-Koegh-Band (Ratanamahatana and Keogh, 2004) and Sakoe-Chiba-Band (Sakoe and Chiba, 2013).

As a technique, it has been long-established in the gesture and sign language recognition domain as well (Ahmed et al., 2016; Jambhale and Khaparde, 2014; Jangyodsuk et al., 2014). Due to the fact that it is a distance metric it requires no training and it is a perfect choice for applications where limited training samples are available.

Ten Holt and her colleagues presented an algorithm for Dynamic Time Warping (DTW) on multi-dimensional time series (MDDTW) to perform classification on 121 gestures recorded with two cameras in stereo position (Ten Holt et al., 2007). In Jangyodsuk et al. (2014) the authors investigated the use of DTW and Histogram of Oriented Gradient (HOG) to compare a query sign with those in a database of ASL signs using Kinect data. Their results showed an accuracy of 82% in a Top-10 ranking level.

Recent developments in the field of machine and deep learning have lead to advances in sign language and gesture recognition. However, these approaches pose restrictions to their overall applicability as they require large amount of data and computational power in order to be trained. Furthermore, proposed methods for sign language classification have been based on special sensor hardware, such as Microsoft's Kinect presenting additional challenges in their duplicability as well difficulty in their technical set-up. Our proposed method does not require the use of depth data to extract the pose key-points as this is being held by the pre-trained pose estimation framework OpenPose. This makes our approach suitable for any kind of sign language lexicon. Most recently, Schneider et al. (2019) used Dynamic Time Warping in conjunction with One-Nearest-Neighbor algorithm and OpenPose to perform classification on six gestures. Their results suggested an accuracy of 77.4%. A major advantage of their methodology is the necessity for very little training data. However, a considerable drawback of their study is that they have only tested a small amount of gestures. As a result, such as pipeline shows a major deterioration of the overall accuracy when an additional gesture is added into the classification task.

Our study repurposes the work of Schneider et al. by:

- considering signs instead of gestures as inputs in DTW
- extending significantly the number of signs used in the experiment
- adding the finger joints extracted by OpenPose as additional data

- testing whether signing proficiency influences the accuracy of the method

3. Methodology

In this section we describe the pose estimation framework (i.e. OpenPose) as well as the apparatus and materials used in this study.

3.1. Pose Estimation

OpenPose is a real-time, open source for academic purposes library for multi-person 2D pose estimation (Cao et al., 2017). It can detect body, foot, hand and facial key-points. It is a bottom-up approach meaning that it does not recognize first where a person is in an image and then extract the body joints but from the detection of the various key-points predicts the overall pose. In general, it exceeds in performance similar 2D body pose estimation libraries like Mask R-CNN (He et al., 2017) and Alpha-Pose (Li et al., 2018). Its major advantage lies in its high accuracy regardless of the number of people in an image or video.

OpenPose is able to run on different operating systems and hardware architectures while providing all the necessary tools for acquisition, visualization and output file generation. Its output consists of multiple json formatted files containing the pixel x, y coordinates of the body, hand and face joints. In this study only the body and hand predictions were used as the face joints were irrelevant for our purposes.

3.2. Preprocessing

The output of OpenPose consists of x,y pixel coordinates. As the people in each frame can potentially be in different locations, it is important to normalize their keypoints. Rotational invariance is omitted in this study as most people are expected to be in an upright position in front of the web camera. The normalization is done in two steps. Firstly, all the key points are translated in such way so that the neck key point shifts to the origo at (0,0). To accomplish the shift, the neck key points coordinates are subtracted from all other key points. Secondly, the key points are scaled in such way so that the distance between the left and the right shoulder key point becomes 1. This is achieved by dividing all key points' coordinates by the distance between the left and right shoulder key point. The scale normalization method is based on previous studies by Celebi et al. (2013), Schneider et al. (2019) and Östling et al. (2018).

One additional step added to the pipeline is the horizontal flip of the videos when a participant was left-handed. This step is achieved by measuring the average velocity of each hand. In cases where the left hand's velocity is greater than the respective of the right hand, a horizontal flip is applied. Such a process allows an independent handedness feature of the overall methodology.

3.3. Participants

Ten people were asked to participate in the research. Four of them have no experience with sign language whatsoever while the rest are experienced signers. Additionally, they were all informed about the general purpose of the research and gave their consent to participate. This study was approved by the Faculty ethics committee.

Identifier	Sign Gloss
s1	ABOUT
s2	BED
s3	BOOK
s4	CAPTAIN
s5	DREAM
s6	EAT
s7	ELEPHANT
s8	HISTORY
s9	HOTEL
s10	IF
s11	LAPTOP
s12	LATER
s13	LUNCH
s14	MEET
s15	MIND
s16	NEAR
s17	NOSE
s18	OPEN
s19	TALK
s20	TRUE

Table 1: List of signs shown to the participants of our experiment

3.4. Data

Each participant viewed only once a selection of 20 signs from the newly compiled Ghanaian Sign Language lexicon (HANDS!Lab, 2020). While the overall lexicon has more than 1300 signs we selected randomly 100 of them to be used in our experiments due to time limitations. The order was randomized for each participant to avoid potential biases. A full list can be seen in Table 1. Each video had a 1000 by 580 pixel resolution at 30 frames per second and lasted approximately 5 (± 2) seconds. Recordings were made with a Macbook Pro’s webcam at 1280 by 720 pixel resolution and 30 frames per second.

We employ the soft DTW method by Cuturi and Blondel (2017) deployed by the `tslearn` python package (Tavenard et al., 2017) to perform DTW on the normalized trajectories of the dominant hand. Their work takes advantage of a smoothed formulation of DTW that computes the soft-minimum of all alignment costs. In a pilot test we observed that soft DTW performed better compared to other DTW variants, and was thus used in the rest of the experiment. Furthermore, a DTW variant created by Sakoe and Chiba (2013) used by the same python module was utilized to measure the distance of the trajectories of all finger coordinates.

Most signs in our lexicon are one-handed where the left hand is inert either by being “absent” or passively fixed at a location. In the two-handed signs, the left hand mostly copies the movement of the right hand. As a result, we employed DTW only on the dominant hand features as the left hand would either be less informative or equally informative.

Finally, the limited resolution of the output from OpenPose had an undesired effect producing sudden spikes in the signal. This attribute has been previously acknowledged by

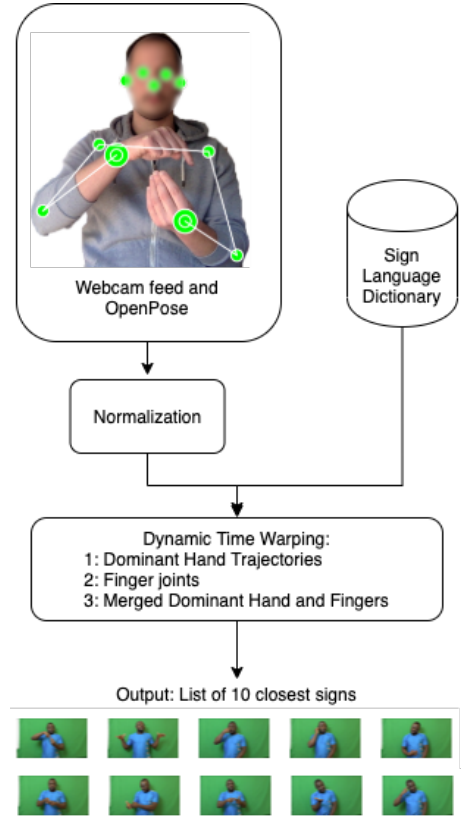


Figure 2: Overview of the overall pipeline of our methodology.

Schneider et al. (2019) and was present here too. The videos in the lexicon were blurry when the hand was moving fast making OpenPose to mispredict the proper joint locations between consecutive frames. As such, failed to create a smooth path. To compensate for this behavior we included two additional steps. Firstly, all the dominant hand’s wrist x,y coordinates that had a confidence level lower than 0.3 were deleted. Additionally, we used a median filter with radius $r = 3$ for smoothing the remaining signal. Moreover, we noticed that due to the good lighting conditions in the GSL lexicon there was a mismatch on the body joint’s coordinates predicted by OpenPose. The lighting conditions of the videos captured with the participants were of poor quality making it hard for the DTW algorithm to operate properly. To solve that problem we decided to include in the lexicon the data from a random participant every time we tested the methodology. This step seems to add the necessary noise in the database that is nevertheless similar to the noise in the participants’ data. As a result, the data of each participant’s sign was compared with 120 signs in our database (100 from the GSL lexicon and 20 from another random participant). The overall pipeline can be seen in Figure 2.

3.5. Visualization

To further explore the outputs of OpenPose and how they are rendered in our methodology, we have created an interactive visualization tool. Developed with the python module “bokeh” (Bokeh Development Team, 2014), the user is able

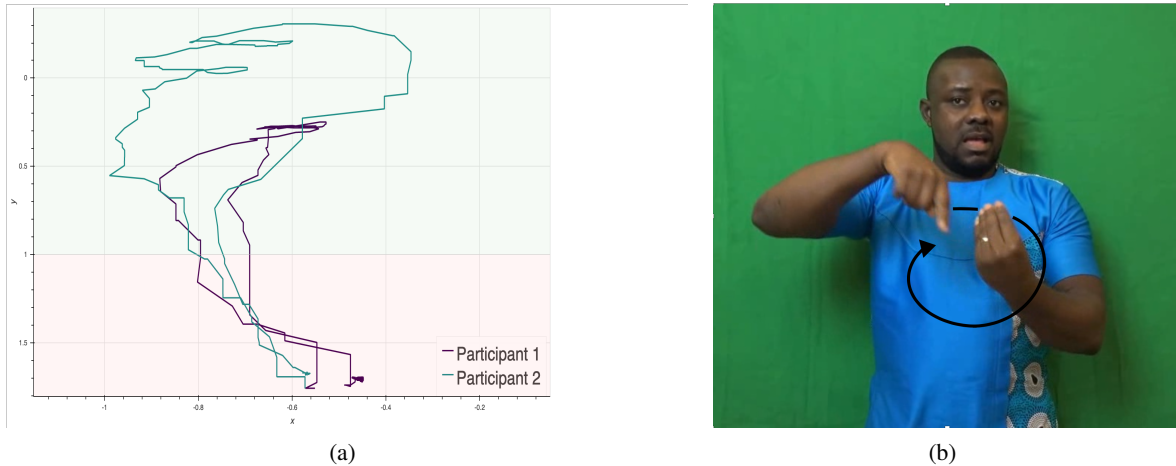


Figure 3: Visualization of the dominant hand trajectories between two participants (a) for the sign ABOUT (b).

to view the extracted dominant hand trajectories from the participants as a whole or individually. As all participants started and ended each sign in the same position, we have color coded as red the preparation and retraction phase and as green the stroke of each sign.

While the motivation behind the creation of this tool was to solely verify the output of openpose and the normalization part in our method, its potential reaches beyond the scope of this study. Such a tool, in combination with the DTW output, can potentially be used as a metric to quantify the variation in the movement and location of signers and sign languages in general. An example of the trajectories of two participants for the same sign can be seen in Figure 3a. It is evident that one participant produced the sign in a larger space with more distinctive movements. Moreover, it can be deduced that the location parameter is different as Participant 2 made the sign at a higher plane (almost in front of the face) while Participant 1 in front of the torso.

4. Results

Table 2 presents the overall accuracy of our methodology. Top-k refers to the number of signs a user must look up before finding a correct match. Accuracy indicates whether the target sign is present in the Top-k retrieved signs and is averaged across all participants and signs.

It is evident that the highest accuracy is apparent at a Top-10 rank level at 87%. Furthermore, Top-5 rank shows an adequate accuracy at 74%. Contrary to expectations, using DTW in the joints of fingers extracted by OpenPose did not yield significant results with a highest accuracy at the Top-10 rank at approximately 52%. Merged DTW distances from the dominant hand trajectories and the finger joints also did not generate compelling results.

If only the experienced signers' data is considered then the accuracy at the Top-10 rank raises at 90% and the Top-5 at 78% (Table 2 row 4). On the other hand, the accuracy on the non-experienced signers drops at 82% and 0.67% at the Top-10 and 5 rank respectively (Table 2 row 5). Moreover, DTW on the finger's trajectories shows a significant drop at the Top-10 rank between the experienced and non-experienced signers of approximately 22% (Table 2 column 7).

The most striking observation to emerge from the analysis was that four out of 20 signs were consistently recognized with almost 100% accuracy at the Top-1 level rank. These signs were: CAPTAIN, DREAM, ELEPHANT and OPEN. Such behavior is justified as these signs have large, distinctive movements and locations that are hard to misinterpreted by the DTW.

5. Discussion

In this study we have investigated the use of OpenPose and Dynamic Time Warping as a ranking pipeline to retrieve matching signs from a sign language dictionary. Our results demonstrated that such a task can be achieved with an adequate accuracy rate.

This is in good agreement with the results obtained by Jangyodsuk et al. (2014). Although the accuracy rate does not match the one from Schneider et al. (2019) we have tested a larger vocabulary and lexicon. Additionally, we are not aiming at classifying each sign but rather create a suggestion ranking system. As such, our results suggest that approximately 9 out of 10 times the matched sign will be present in the first 10 retrieved signs.

Moreover, the results have further strengthened our hypothesis that signing proficiency is an influencing factor for

Condition	Dominant hand trajectory			Fingers' trajectories			Merged trajectories		
	Top 1	Top 5	Top 10	Top 1	Top 5	Top 10	Top 1	Top 5	Top 10
Accuracy of all participants	0,27	0,74	0,87	0,23	0,40	0,52	0,29	0,55	0,71
Accuracy of experienced signers	0,29	0,78	0,90	0,30	0,47	0,61	0,39	0,64	0,79
Accuracy of non-experienced signers	0,25	0,67	0,82	0,12	0,30	0,39	0,14	0,42	0,59

Table 2: Sign retrieval accuracy. Top k refers to number of best matches.

classification efforts. Although our sample size was limited there was a significant drop in the accuracy rates between the experienced and non-experienced signers. The former, produced well structured signs matching more appropriately the ones from the lexicon, which made DTW perform in a more excellent matter.

Our research failed to account for the low values of accuracy on the finger joints. This was probably as a result of the low performance of OpenPose in accurately predicting the finger joints due to low lighting conditions in the videos. It was often the case that joint predictions would disappear between frames or mis-predicted in wrong locations. Thus, caution must be exercised when OpenPose is being used for such trivial tasks.

6. Conclusion

To sum up, we have obtained satisfactory results demonstrating the use of OpenPose and Dynamic Time Warping for a new, sign-based search functionality in reduced sign language dictionaries. We showed that our “find the sign” methodology can be used as a suggestion tool for sign retrieval in a small lexicon by using only the trajectory of the dominant hand. Additionally, our research has highlighted the importance of considering the level of signing proficiency when it comes to classification tasks. The significance of this study lies on the fact that the methodology in question can be easily used in any kind of sign language lexicon, irrespective of its quality and language. Additionally, no prior training of any kind of model is required. As such, this approach, in combination with the developed visualization module, has the potential to be used also as a metric tool to quantify the variation between signers and overall languages.

Furthmore, a number of things is left for future work; first and foremost, to investigate how extracted finger joints can be utilized more efficiently in the overall pipeline. Moreover, different variants of the original DTW algorithms need to be tested. Finally, we intend to evaluate the use of other pose estimation frameworks, such as PoseNet, to further enhance the web and mobile user-friendliness of the method used.

7. Acknowledgements

We would like to thank all the people who participated in the study, without whose help this work would have never been possible.

8. Bibliographical References

- Abdulla, W., Chow, D., and Sin, G. (2003). Cross-words reference template for dtw-based speech recognition systems. In *TENCON 2003. Conference on Convergent Technologies for Asia-Pacific Region*, volume 4, pages 1576–1579, Bangalore, India. Allied Publishers Pvt. Ltd.
- Ahmed, W., Chanda, K., and Mitra, S. (2016). Vision based hand gesture recognition using dynamic time warping for indian sign language. In *2016 International Conference on Information Science (ICIS)*, pages 120–125, Kochi, India. IEEE.
- Axelrod, S. and Maison, B. (2004). Combination of hidden Markov models with dynamic time warping for speech recognition. In *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages I–173–6, Montreal, Que., Canada. IEEE.
- Bokeh Development Team. (2014). Bokeh: Python library for interactive visualization. <http://bokeh.pydata.org>.
- Cao, Z., Simon, T., Wei, S.-E., and Sheikh, Y. (2017). Realtime multi-person 2d pose estimation using part affinity fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1302–1310, Honolulu, HI. IEEE.
- Celebi, S., Aydin, A. S., Talha, T. T., and Tarik, A. (2013). Gesture recognition using skeleton data with weighted dynamic time warping. In *Proceedings of the International Conference on Computer Vision Theory and Applications*, pages 620–625, Barcelona, Spain. SciTePress - Science and Technology Publications.
- Center for Tegnsprog. (2008). Ordbog over Dansk Tegnsprog. <http://www.tegnsprog.dk/>.
- Cuturi, M. and Blondel, M. (2017). Soft-dtw: a differentiable loss function for time-series. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 894–903. JMLR. org.
- HANDS!Lab. (2020). Ghanaian Sign Language. https://play.google.com/store/apps/details?id=com.ljsharp.gsldictionary&hl=es_US.
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. (2017). Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969.
- Institutionen for Lingvistik. (2009). Svenskt teckensprakslexikon. <https://teckensprakslexikon.su.se>.
- Itakura, F. (1975). Minimum prediction residual principle applied to speech recognition. *IEEE Transactions on acoustics, speech, and signal processing*, 23(1):67–72.
- Jambhale, S. S. and Khaparde, A. (2014). Gesture recognition using DTW & piecewise DTW. In *2014 International Conference on Electronics and Communication Systems (ICECS)*, pages 1–5, Coimbatore. IEEE.
- Jangyodsuk, P., Conly, C., and Athitsos, V. (2014). Sign language recognition using dynamic time warping and hand shape distance based on histogram of oriented gradient features. In *Proceedings of the 7th International Conference on Pervasive Technologies Related to Assistive Environments - PETRA '14*, pages 1–6, Rhodes, Greece. ACM Press.
- Li, J., Wang, C., Zhu, H., Mao, Y., Fang, H.-S., and Lu, C. (2018). Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. *arXiv preprint arXiv:1812.00324*.
- Myers, C., Rabiner, L., and Rosenberg, A. (1980). Performance tradeoffs in dynamic time warping algorithms for isolated word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(6):623–635.
- Ratanamahatana, C. A. and Keogh, E. (2004). Making time-series classification more accurate using learned constraints. In *Proceedings of the 2004 SIAM Interna-*

- tional Conference on Data Mining*, pages 11–22. Society for Industrial and Applied Mathematics.
- Sakoe, H. and Chiba, S. (2013). Gesture recognition using skeleton data with weighted dynamic time warping. In *Proceedings of the International Conference on Computer Vision Theory and Applications*, pages 620–625, Barcelona, Spain. SciTePress - Science and Technology Publications.
- Schneider, P., Memmesheimer, R., Kramer, I., and Paulus, D. (2019). Gesture recognition in rgb videos using human body keypoints and dynamic time warping. *arXiv:1906.12171 [cs]*. arXiv: 1906.12171.
- Tavenard, R., Faouzi, J., Vandewiele, G., Divo, F., Androz, G., Holtz, C., Payne, M., Yurchak, R., Russwurm, M., Kolar, K., and Woods, E. (2017). tslearn: A machine learning toolkit dedicated to time-series data. <https://github.com/rtavenar/tslearn>.
- Ten Holt, G. A., Reinders, M. J., and Hendriks, E. (2007). Multi-dimensional dynamic time warping for gesture recognition. In *Thirteenth annual conference of the Advanced School for Computing and Imaging*, volume 300, page 1.
- Van Herreweghe, M., Slembrouck, S., and Vermeerbergen, M. (2004). Digitaal Vlaamse Gebarentaal-Nederlands/Nederlands-Vlaamse Gebarentaal woordenboek. <https://woordenboek.vlaamsegebarentaal.be>.
- Zwitsersloot, I. (2010). Sign language lexicography in the early 21st century and a recently published dictionary of sign language of the netherlands. *International Journal of Lexicography*, 23(4):443–476.
- Östling, R., Börstell, C., and Courtaux, S. (2018). Visual iconicity across sign languages: Large-scale automated video analysis of iconic articulators and locations. *Frontiers in Psychology*, 9:725.