

WLASL-LEX: a Dataset for Recognising Phonological Properties in American Sign Language

Anonymous ACL submission

Abstract

Signed Language Processing (SLP) concerns the automated processing of signed languages, the main means of communication of Deaf and hearing impaired individuals. SLP features many different tasks, ranging from sign recognition to translation and production of signed speech, but has been overlooked by the NLP community thus far. In this paper, we bring to attention the task of modelling the phonology of sign languages. We leverage existing resources to construct a large-scale dataset of American Sign Language signs annotated with six different phonological properties. We then conduct an extensive empirical study to investigate whether data-driven end-to-end and feature-based approaches can be optimised to automatically recognise these properties. We find that, despite the inherent challenges of the task, graph-based neural networks that operate over skeleton features extracted from raw videos are able to succeed at the task to a varying degree. Most importantly, we show that this performance pertains even on signs unobserved during training.

1 Introduction

Around 200 languages in the world are signed rather than spoken, featuring their own vocabulary and grammatical structures. For example the American Sign Language (ASL) is not a mere translation of English into signs and is unrelated to the British Sign Language (BSL). This introduces many novel challenges to their automated processing. Research on Sign Language Processing (SLP) encompasses tasks such as sign language detection, i.e. recognising if and which signed language is performed (Moryossef et al., 2020) and sign language recognition (SLR) (Koller, 2020), i.e. the identification of signs either in isolation or in continuous speech. Other tasks concern the translation from signed to spoken (or written) (Camgoz et al., 2018) language or the production of signs from text

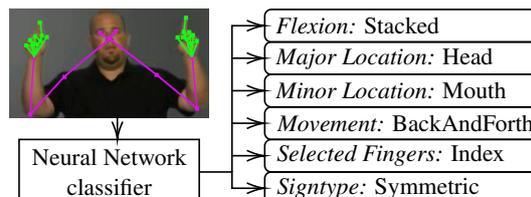


Figure 1: We annotate ASL sign videos with their corresponding phonological information and skeleton features of the speakers, and train neural networks to recognise the former from the latter.

(Rastgoo et al., 2021). With the recent success of deep learning-based approaches in computer vision (CV), as well as advancements in —from the CV perspective—related tasks of action and gesture recognition (Asadi-Aghbolaghi et al., 2017), SLP is gaining more attention in the CV community (Zheng et al., 2017).

Some recent approaches to various SLP tasks rely on *phonological* features, perhaps due to the complexity of the tasks (Tornay, 2021; Metaxas et al., 2018; Gebre et al., 2013; Tavella et al., 2021). Surprisingly, however, little work has been carried out on explicitly modelling the phonology of signed languages. This presents a timely opportunity to investigate signed languages from a linguist’s perspective (Yin et al., 2021). In the context of signed languages, phonology typically distinguishes between manual features, such as usage, position and movement of hands and fingers, and non-manual features, such as facial expression. Sign language phonology is a matured field with well-developed theoretical frameworks (Liddell and Johnson, 1989; Fenlon et al., 2017; Sandler, 2012). These phonological features, or *phonemes*, are drawn from a fixed inventory of possible configurations which is typically much smaller than the vocabulary of signed languages (Borg and Camilleri, 2020). For example, there is only a limited number of fingers that can be used to perform a sign due to anatomical constraints. Hence, different signs share phonolog-

ical properties and well performing classifiers can be used to predict those properties for signs unseen during training. This potentially holds even across different languages, because, while different languages may dictate different combinations of phonemes, there are also significant overlaps (Tornay et al., 2020).

Finally, these phonological properties have a strong discriminatory power when determining signs. For example, in ASL-Lex (Caselli et al., 2017), a lexicon which also captures phonology information, the authors report that more than 50% of its 994 described signs have a unique combination of only six phonological properties and more than 80% of the signs share their combination with at most two other signs. By relying on additional (i.e., phonological) information from resources such as ASL-Lex, many signs can be determined from (predicted) phonological properties alone, without encountering them in training data. This is a capability that current data-driven approaches to SLR lack by design (Koller, 2020). Thus, in combination, mature approaches to phonology recognition can facilitate the development of sign language resources. This is an important task for both documenting low-resource sign languages as well as rapid developing of large-scale datasets, to fully harness data-driven CV approaches.

To spur research in this direction, we extend the preliminary work by Tavella et al. (2021) and introduce the task of Phonological Property Recognition (PPR). More specifically, this paper contributes (i) WLASLlex2001, a large-scale, automatically constructed PPR dataset, (ii) an analysis of the dataset quality, and (iii) an empirical study of the performance of different deep-learning based baselines thereon.

2 Methodology

We address PPR as a classification problem based on features extracted from videos of people speaking SL. Albeit manual annotation approaches are generally adopted, an automated approach would be less time and resource consuming, allowing researchers to limit their efforts to data validation. To extract such features, we take advantage of pretrained deep models from the computer vision community (Rong et al., 2021; Wang et al., 2019). Finally, we train several deep models to classify them as phonological classes.

Dataset construction: As previously men-

tioned, ASL-Lex (Caselli et al., 2017) contains phonological features of American Sign Language, such as where the sign is executed, the movement performed by the hand or the number of hands involved. The latter properties were coded by 3 ASL-versed people. In our work, we are interested in recognising phonological classes from videos of people speaking ASL. Consequently, we aim to construct a dataset suitable for supervised learning, containing videos labelled with 6 phonological properties. We choose: (i) *flexion*, aperture of the selected fingers of the dominant hand at sign onset, (ii) *major location*, general location of the dominant hand at sign onset, (iii) *minor location*, specific location of the dominant hand at sign onset, (iv) *movement*, path movement of the first morpheme in the sign, (v) *selected fingers*, fingers that are moving or foregrounded in the first morpheme of the sign, and (vi) *sign type*, symmetry of the hands according to Battison (1978). A detailed description of all the properties is provided in the appendix. We selected these manual properties as they have a strong discriminatory power to predict signs based on their configuration (Caselli et al., 2017). One of the limitations of ASL-Lex is the small number of examples and its limited variety: its first iteration (ASL-Lex 1.0) contains less than 1000 videos, all signed by the same person. While sufficient for educational purposes, these videos are of limited suitability for developing robust classifiers that can capture the diversity of ASL speakers (Yin et al., 2021). To this end, we source videos from WLASL (Li et al., 2020) (Word Level-ASL), one of the largest available SL datasets, featuring more than 2000 glosses demonstrated by over 100 people, for a total of more than 20000 videos. Each sign is performed by at least 3 different signers, which implies greater variability compared to having one gloss performed by only one user. By cross referencing ASL-Lex and WLASL2000 based on corresponding glosses, we can increase the number of samples available to train our models. Finally, to leverage state of the art SLR architectures that operate over structured input, we enrich each raw video with its extracted keypoints that represent the joints of the speaker. To do so, we use two pretrained models, FrankMocap (Rong et al., 2021) and HRNet (Wang et al., 2019). While these tracking algorithms follow different paradigms, the former extracting 3D coordinates based on a predicted human model and the latter predicting keypoints as

coordinates from videos directly, they produce similar outputs. An important distinction is that while FrankMocap estimates the 3D keypoints, HRNet outputs 2D keypoints with associated prediction confidence scores. We use these different models to explore whether different tracking algorithms affect the recognition of phonological classes. We select a subset of features of the upper body, namely: nose, eyes, shoulders, elbows, wrists, thumbs and first/last knuckles of the fingers. These manual features were determined to be the most informative while performing sign language recognition (Jiang et al., 2021b).

Our final dataset, WLASL-Lex2001 (WLASL2000 + ASL-Lex 1.0), is composed of 10017 videos corresponding to 800 glosses, 3D skeletons (x , y , z from FrankMocap and x , y and $score$ from HRNet) labelled with their phonological properties. A characteristic of this dataset is that it follows a long tailed distribution. Due to the nature of language, some phonological properties are more common than others, which means that some classes are more represented than others. On the one hand, the training setup for our models should take this factor into account, but on the other hand, the advantage of training over phonological classes instead of glosses is that different glosses can share phonological classes.

Models: To estimate the complexity of the dataset, we use the majority-class baseline and the Multi-Layer Perceptron (MLP) as a basic deep model. We further use Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) as models capable of capturing the temporal component of videos. As state-of-the-art SLP architectures that have been used to perform SLR, we use the I3D 3D Convolutional Neural Network (Carreira and Zisserman, 2017; Li et al., 2020) able to learn from raw videos, and the Spatio-Temporal Graph Convolutional Network (STGCN) (Jiang et al., 2021b) that captures both spatial and temporal components from the extracted keypoints.

Experimental Setup: We generate one dataset and train different models for each phonological property. While this might not be the optimal way, as opposed to a multiclass multilabel approach, it is the best one in order to understand which features can and cannot be singularly learned, making the error analysis much easier. From now on, when we cite the *dataset*, we refer to an instance of the WLASL-Lex 2001 dataset, whose labels are the

values of a single phonological class. We make this distinction because we split the dataset into train, validation and test sets (with a 70 : 15 : 15 ratio) using a stratified strategy based on the selected phonological class (*Phoneme*). By doing so, we make sure that all the different splits contain all possible values for a phonological class. Because our dataset features multiple videos per gloss, glosses in the test set appear in the training set as well. Thus, to investigate how well the models can predict properties on unseen glosses, we also produce label-stratified splits on gloss-level (*Gloss*), such that videos of glosses in the validation and test set do not appear in training data and vice versa.

The I3D is pre-trained on Kinetics-400 (Carreira and Zisserman, 2017) and fine-tuned on our datasets. The other models are trained from scratch using keypoints as input. We fix the length of all input to 150 frames, longer sequences are truncated while shorter sequences are looped to reach the fixed length. We select the best performing model based on performance on the validation set and for the final test set performance we train the models on both train and validation set. For more details on model selection, consult the appendix. We measure both accuracy, to investigate how well models perform in general, and class-balanced accuracy to take into account how well they are able to model different classes of the phonological properties.

3 Results and discussion

The upper half of Table 1 presents the results for the six datasets split in a stratified fashion, not taking into account the corresponding glosses. The poor performance of the simple MLP architecture suggests that the tasks are in fact challenging and do not exhibit easily exploitable regularities. Due to its simplicity, for some properties it is barely able to reach the baseline (34% vs. 35% and 44% vs. 50% for *movement* and *flexion* respectively). In particular, MLP classifying based on FrankMocap (MLP_F) output is often the worst performing combination. Conversely, STGCN using HRNet output ($STGCN_H$) outperforms other models on all six tasks. In some cases, for example when predicting *movement* or *flexion*, it is the only model which significantly surpasses the majority class baseline. This superior performance is expected, as specifically this combination of the STGCN operating over HRNet-extracted keypoints has been shown to be the largest contributor to the SLR performance

| | FLEXION | | MAJLOCATION | | MINLOCATION | | MOVEMENT | | FINGERS | | SIGNTYPE | | |
|----------------|--------------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|-------------------|-------------|
| | A | \bar{A} | A | \bar{A} | A | \bar{A} | A | \bar{A} | A | \bar{A} | A | \bar{A} | |
| <i>Phoneme</i> | Baseline | 50.3 | 11.1 | 34.4 | 20.0 | 33.9 | 3.1 | 35.5 | 16.7 | 48.2 | 11.1 | 39.3 | 20 |
| | MLP _H | 50.1 ± 2.5 | 11.1 | 70.3 ± 2.3 | 64.0 | 51.6 ± 2.5 | 28.2 | 34.3 ± 2.4 | 18.7 | 59.4 ± 2.5 | 25.0 | 73.9 ± 2.2 | 52.6 |
| | MLP _F | 50.3 ± 2.5 | 11.1 | 57.8 ± 2.5 | 46.8 | 34.3 ± 2.4 | 9.1 | 34.3 ± 2.4 | 18.7 | 43.4 ± 2.5 | 12.9 | 67.0 ± 2.4 | 42.8 |
| | RNN _H | 49.0 ± 2.5 | 30.0 | 75.8 ± 2.2 | 72.4 | 64.3 ± 2.4 | 46.0 | 35.1 ± 2.4 | 29.5 | 71.0 ± 2.3 | 46.5 | 78.7 ± 2.1 | 58.8 |
| | RNN _F | 50.3 ± 2.5 | 11.1 | 64.6 ± 2.4 | 54.2 | 30.3 ± 2.3 | 4.0 | 35.4 ± 2.4 | 18.1 | 46.5 ± 2.5 | 12.4 | 70.9 ± 2.3 | 46.8 |
| | STGCN _H | 62.3 ± 2.4 | 45.0 | 83.2 ± 1.9 | 78.6 | 74.5 ± 2.2 | 63.5 | 63.6 ± 2.4 | 58.2 | 73.8 ± 2.2 | 56.0 | 84.5 ± 1.8 | 69.6 |
| | STGCN _F | 43.4 ± 2.5 | 20.8 | 70.5 ± 2.3 | 62.1 | 53.0 ± 2.5 | 40.0 | 45.7 ± 2.5 | 37.8 | 63.1 ± 2.4 | 32.8 | 73.0 ± 2.2 | 53.1 |
| | 3DCNN | 46.5 ± 2.5 | 13.2 | 64.3 ± 2.4 | 55.2 | 42.3 ± 2.5 | 18.6 | 32.0 ± 2.4 | 20.8 | 47.5 ± 2.5 | 14.5 | 69.5 ± 2.3 | 44.8 |
| <i>Gloss</i> | Baseline | 53.1 | 11.1 | 35.7 | 20.0 | 42.0 | 5.0 | 35.2 | 16.7 | 47.4 | 12.5 | 38.3 | 20.0 |
| | MLP _H | 44.6 ± 2.5 | 15.5 | 68.1 ± 2.3 | 56.6 | 47.3 ± 2.5 | 19.7 | 28.4 ± 2.2 | 19.8 | 56.2 ± 2.5 | 22.9 | 75.3 ± 2.2 | 50.7 |
| | MLP _F | 50.3 ± 2.5 | 11.1 | 56.6 ± 2.5 | 42.9 | 38.3 ± 2.4 | 10.7 | 37.1 ± 2.4 | 21.7 | 39.3 ± 2.5 | 12.5 | 68.4 ± 2.4 | 41.2 |
| | RNN _H | 49.0 ± 2.5 | 30.0 | 72.8 ± 2.2 | 67.3 | 49.3 ± 2.5 | 26.3 | 32.2 ± 2.3 | 24.9 | 60.7 ± 2.5 | 32.5 | 75.4 ± 2.2 | 53.5 |
| | RNN _F | 50.3 ± 2.5 | 11.1 | 64.1 ± 2.4 | 52.6 | 44.4 ± 2.4 | 17.8 | 36.7 ± 2.4 | 20.1 | 27.3 ± 2.3 | 12.7 | 72.0 ± 2.3 | 46.9 |
| | STGCN _H | 49.1 ± 2.5 | 21.6 | 77.3 ± 2.1 | 70.0 | 55.1 ± 2.4 | 32.7 | 52.5 ± 2.5 | 46.5 | 65.7 ± 2.4 | 34.4 | 76.6 ± 2.1 | 54.4 |
| | STGCN _F | 39.0 ± 2.5 | 14.4 | 66.7 ± 2.3 | 60.1 | 45.1 ± 2.4 | 21.1 | 43.1 ± 2.5 | 34.9 | 60.0 ± 2.5 | 29.2 | 71.3 ± 2.3 | 47.5 |
| | 3DCNN | 46.0 ± 2.5 | 12.8 | 64.9 ± 2.4 | 52.0 | 10.8 ± 1.5 | 13.6 | 32.0 ± 2.3 | 19.3 | 45.9 ± 2.5 | 14.7 | 71.6 ± 2.3 | 46.3 |

Table 1: Accuracy (A.) and per-class averaged accuracy (\bar{A}) of various models on the test tests of the six tasks. For accuracy, we report the error margin as a confidence interval at $\alpha = 0.05$ using asymptotic normal approximation. We omit error margins for balanced accuracy as the low number of classes results in a small sample size.

on the WLASL2000 dataset (Jiang et al., 2021a). Models that operate over structured input often outperform the 3D CNN, demonstrating the utility of additional information provided by the skeleton features. The results also suggest that models using the HRNet skeleton output outperform those who use FrankMocap, possibly due to confidence scores produced by HRNet and associated with the coordinates. This difference in performance suggests to conduct a more rigorous study to investigate the impact of different feature extraction methods as a possible future research direction.

The lower half of Table 1 shows the evaluation results on unseen glosses (*Gloss*). The performance of all tasks and all models deteriorates, suggesting that their success is partly derived from exploiting the similarities of videos that appear in training and test data and refer to the same gloss. However, the best model, STGCN_H, performs comparably to the *Phoneme*-split, with a drop of less than 10 accuracy points for five of the six tasks.

Often, automatically constructed datasets such as ours, have a performance ceiling, for example due to incorrectly assigned ground truth labels or low quality of input data (Chen et al., 2016). To investigate the former, we measure the agreement on videos that all models misclassify using Fleiss’ κ . Intuitively, if all models agree on a label different than the ground truth, the ground truth label might be wrong. We find that averaged across the six tasks, the agreement is negligible: 0.09 ± 0.06 and 0.11 ± 0.09 for *Phoneme* and *Gloss* split, respectively. Similarly, for the latter, if all models consistently fail to assign any correct label for a

given video (e.g. all models err on a video appearing in the test sets of *movement* and *flexion*), this can hint at low quality of the input, exacerbating processing it correctly. We find that this is not the case with WLASL-LEX2001, as videos appearing in test sets of different tasks tend to have a low mutual misclassification rate: 1% and 0.7% of videos appearing in test sets of two and three tasks were misclassified by all models for all associated tasks for the *Phoneme* split. For the *Gloss* split the numbers are 3 and 0% for two and three tasks, respectively. Together, these observations suggest that the models presented in this paper are unlikely to reach the performance ceiling on WLASL-Lex2001 and more advanced approaches could obtain even higher accuracy scores.

4 Conclusion

In this paper, we discuss the task of Phonological Property Recognition (PPR). We automatically construct a dataset for the task featuring six phonological properties and analyse it extensively. We find that there is potential for improvement over our presented data-driven baseline approaches. Researchers pursuing this direction can focus on developing better-performing models, for example by relying on jointly learning all properties, as labels for different properties can be mutually dependent.

Another possible avenue is to investigate the feasibility of using PRR to perform *tokenisation* of continuous sign language speech, by decomposing it into multiple phonemes, which is identified as one of the big challenges of SLP (Yin et al., 2021).

340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395

References

Maryam Asadi-Aghbolaghi, Albert Clapes, Marco Belantoni, Hugo Jair Escalante, Victor Ponce-Lopez, Xavier Baro, Isabelle Guyon, Shohreh Kasaei, and Sergio Escalera. 2017. [A Survey on Deep Learning Based Approaches for Action and Gesture Recognition in Image Sequences](#). *Proceedings - 12th IEEE International Conference on Automatic Face and Gesture Recognition, FG 2017 - 1st International Workshop on Adaptive Shot Learning for Gesture Understanding and Production, ASLAGUP 2017, Biometrics in the Wild, Bwild 2017, Heteroge*, pages 476–483.

Robbin Battison. 1978. Lexical borrowing in american sign language.

Mark Borg and Kenneth P. Camilleri. 2020. [Phonologically-Meaningful Subunits for Deep Learning-Based Sign Language Recognition](#). *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12536 LNCS:199–217.

Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. 2018. [Neural Sign Language Translation](#).

João Carreira and Andrew Zisserman. 2017. [Quo vadis, action recognition? a new model and the kinetics dataset](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733.

Naomi K. Caselli, Zed Sevcikova Sehyr, Ariel M. Cohen-Goldberg, and Karen Emmorey. 2017. [Asllex: A lexical database of american sign language](#). *Behavior Research Methods*, 49(2):784–801.

Danqi Chen, Jason Bolton, and Christopher D. Manning. 2016. [A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 4, pages 2358–2367, Stroudsburg, PA, USA. Association for Computational Linguistics.

J Fenlon, Kearsy A Cormier, and Diane Brentari. 2017. [Sign language phonology](#). In *Routledge Handbook of Phonological Theory*.

Binyam Gebrekidan Gebre, Peter Wittenburg, and Tom Heskes. 2013. [Automatic sign language identification](#). *2013 IEEE International Conference on Image Processing, ICIP 2013 - Proceedings*, pages 2626–2630.

Songyao Jiang, Bin Sun, Lichen Wang, Yue Bai, Kunpeng Li, and Yun Fu. 2021a. [Sign Language Recognition via Skeleton-Aware Multi-Model Ensemble](#).

Songyao Jiang, Bin Sun, Lichen Wang, Yue Bai, Kunpeng Li, and Yun Fu. 2021b. [Skeleton Aware Multimodal Sign Language Recognition](#). *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 3408–3418.

Oscar Koller. 2020. [Quantitative Survey of the State of the Art in Sign Language Recognition](#). 396
397

Dongxu Li, Cristian Rodriguez, Xin Yu, and Hongdong Li. 2020. [Word-level deep sign language recognition from video: A new large-scale dataset and methods comparison](#). In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1459–1469. 398
399
400
401
402

Scott K. Liddell and Robert E. Johnson. 1989. [American Sign Language: The Phonological Base](#). *Sign Language Studies*, 1064(1):195–277. 403
404
405

B.W. Matthews. 1975. [Comparison of the predicted and observed secondary structure of t4 phage lysozyme](#). *Biochimica et Biophysica Acta (BBA) - Protein Structure*, 405(2):442–451. 406
407
408
409

Dimitris Metaxas, Mark Dilsizian, and Carol Neidle. 2018. [Scalable ASL sign recognition using model-based machine learning and linguistically annotated corpora](#). 410
411
412
413

Amit Moryossef, Ioannis Tsochantaridis, Roei Aharoni, Sarah Ebling, and Srini Narayanan. 2020. [Real-Time Sign Language Detection Using Human Pose Estimation](#). *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 12536 LNCS:237–248. 414
415
416
417
418
419
420

Razieh Rastgoo, Kouros Kiani, and Sergio Escalera. 2021. [Sign Language Recognition: A Deep Survey](#). *Expert Systems with Applications*, 164:113794. 421
422
423

Yu Rong, Takaaki Shiratori, and Hanbyul Joo. 2021. [Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration](#). In *IEEE International Conference on Computer Vision Workshops*. 424
425
426
427
428

Wendy Sandler. 2012. [The Phonological Organization of Sign Languages](#). *Language and Linguistics Compass*, 6(3):162–182. 429
430
431

Federico Tavella, Aphrodite Galata, and Angelo Cangelosi. 2021. [Phonology recognition in american sign language](#). 432
433
434

Sandrine Tornay. 2021. [Explainable Phonology-based Approach for Sign Language Recognition and Assessment](#). Ph.D. thesis, Lausanne, EPFL. 435
436
437

Sandrine Tornay, Marzieh Razavi, and Mathew Magimai.-Doss. 2020. [Towards Multilingual Sign Language Recognition](#). In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6304–6308. IEEE. 438
439
440
441
442
443

Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. 2019. [Deep high-resolution representation learning for visual recognition](#). *TPAMI*. 444
445
446
447
448

- 449 Kayo Yin, Amit Moryossef, Julie Hochgesang, Yoav
450 Goldberg, and Malihe Alikhani. 2021. [Including](#)
451 [Signed Languages in Natural Language Processing](#).
452 pages 7347–7360.
- 453 Lihong Zheng, Bin Liang, and Ailian Jiang. 2017. [Re-](#)
454 [cent Advances of Deep Learning for Sign Language](#)
455 [Recognition](#). *DICTA 2017 - 2017 International Con-*
456 *ference on Digital Image Computing: Techniques*
457 *and Applications*, 2017-Decem:1–7.

A Hyperparameters optimization

Table 2 contains all the hyperparameters explored during our experiment over each different model. The best model is the one that maximises the Matthew’s correlation coefficient

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

with TP, TN, FP, FN being true/false positive/negative. For the STGCN we use hyperparameters chosen by Jiang et al. (2021a), because initial experiments on our data showed a difference of at most 2% accuracy, which is within the uncertainty estimate. To find the optimal hyperparameters for the other models, we perform Bayesian optimisation over a pre-defined set We maximise Matthew’s correlation coefficient (MCC) (Matthews, 1975) on the validation sets of all six tasks. We choose MCC as it provides a good trade-off between overall and class-level accuracy which is necessary due to the unbalance inherently present in our dataset.

| Model | Parameters |
|--------|----------------------|
| MLP | number of layers |
| | hidden dimension |
| | dropout |
| | learning rate |
| | scheduler step size |
| RNN | gamma |
| | number of RNN layers |
| | RNN hidden dimension |
| | RNN dropout |
| STGCN | learning rate |
| | number of groups |
| | block size, |
| | window size |
| | scheduler step size |
| 3D CNN | dropout |
| | learning rate |
| | gamma |
| | scheduler step size |
| | window size |

Table 2: Set of explored hyperparameters for each different model

B Seed dependency

Table 3 illustrates the performance on the test set for each model with respect to chance as measured by training 5 models from different random seeds. The performance difference is negligible suggesting that model training is largely stable with regard to chance.

| Model | Accuracy |
|--------|------------------|
| MLP | 74.39 ± 0.35 |
| RNN | 79.12 ± 0.46 |
| STGCN | 84.12 ± 0.29 |
| 3D CNN | 69.23 ± 0.93 |

Table 3: Mean and standard deviation of accuracy of all architectures trained with the HRNet output, measured on the SIGNTYPE test set and averaged over 5 different random seeds. Results for the 3D CNN are obtained from the validation set.

C Phonological classes description

Tables 4 to 9 describe in detail the meaning of values for all the phonological classes according to ASL-Lex (Caselli et al., 2017).

The cardinality is calculated on WLASL-Lex, which is why some classes that are in ASL-Lex are not represented (i.e., cardinality equal to 0).

| Value | Definition | Cardinality |
|--------------|-----------------------------------|--------------------|
| imrp | index, middle, ring, pinky finger | 4824 |
| imr | index, middle, ring finger | 95 |
| mrp | middle, ring, pinky finger | 28 |
| im | index, middle finger | 1296 |
| ip | index, pinky finger | 51 |
| mr | middle, ring finger | 0 |
| mp | middle, pinky finger | 0 |
| rp | ring, pinky finger | 0 |
| i | index finger | 2547 |
| m | middle finger | 259 |
| r | ring finger | 0 |
| p | pinky | 407 |
| thumb | thumb | 510 |

Table 4: Values and relative definitions for selected fingers

| Value | Definition | Cardinality |
|--------------|--|--------------------|
| Head | Sign is produced on or near the head | 3137 |
| Arm | Sign is produced on or near the arm | 219 |
| Body | Sign is produced on or near the trunk | 1019 |
| Hand | Sign is produced on or near the non-dominant hand | 2194 |
| Neutral | Sign is not produced in another location on the body | 3448 |
| Other | Sign is produced in another unspecified location on the body | 0 |

Table 5: Values and relative definitions for major location

| Value | Definition | Cardinality |
|--------------|--|--------------------|
| 1 | Fully open: no joints of selected fingers are flexed | 5037 |
| 2 | Bent (closed): non-base joints are flexed | 693 |
| 3 | Flat-open: base joints flexed less than 90 degrees | 909 |
| 4 | Flat-closed: base joints flexed equal to or more that 90 degrees | 507 |
| 5 | Curved open: base and non-base joints flexed without contact | 1130 |
| 6 | Curved closed: base and non-base joints flexed with contact | 642 |
| 7 | Fully closed: base and non-base joints fully flexed | 795 |
| Stacked | Stacked: Flexion of selected fingers differs | 123 |
| Crossed | Crossed | 181 |

Table 6: Values and relative definitions for flexion

| Value | Definition | Cardinality |
|--------------|---|--------------------|
| HeadTop | Sign is produced on top of the head | 20 |
| Forehead | Sign is produced at the forehead | 246 |
| Eye | Sign is produced near the eye | 616 |
| CheekNose | Sign is produced on the cheek or nose | 511 |
| UpperLip | Sign is produced on the upper lip | 53 |
| Mouth | Sign is produced on the mouth | 431 |
| Chin | Sign is produced on the chin | 717 |
| UnderChin | Sign is produced under the chin | 74 |
| UpperArm | Sign is produced on the upper arm | 39 |
| ElbowFront | Sign is produced in the crook of the elbow | 0 |
| ElbowBack | Sign is produced on the outside of the elbow | 13 |
| ForearmBack | Sign is produced on the outside of the forearm | 32 |
| ForearmFront | Sign is produced on the inside of the forearm | 10 |
| ForearmUlnar | Sign is produced on the ulnar side of the forearm | 56 |
| WristBack | Sign is produced on the back of the wrist | 23 |
| WristFront | Sign is produced on the front of the wrist | 0 |
| Neck | Sign is produced on the neck | 68 |
| Shoulder | Sign is produced on the shoulder | 101 |
| Clavicle | Sign is produced on the clavicle | 419 |
| TorsoTop | Sign is produced in the upper third of the torso | 0 |
| TorsoMid | Sign is produced in the middle third of the torso | 0 |
| TorsoBottom | Sign is produced in the bottom third of the torso | 19 |
| Waist | Sign is produced at the waist | 34 |
| Hips | Sign is produced on the hips | 59 |
| Palm | Sign is produced on the palm of the non-dominant hand | 925 |
| FingerFront | Sign is produced on the front of the fingers of the non-dominant hand | 99 |
| PalmBack | Sign is produced on the back of the palm of the non-dominant hand | 218 |
| FingerBack | Sign is produced on the back of the fingers of the non-dominant hand | 186 |
| FingerRadial | Sign is produced on the radial side of the non-dominant hand | 410 |
| FingerUlnar | Sign is produced on the ulnar side of the non-dominant hand | 40 |
| FingerTip | Sign is produced on the tip of the fingers of the non-dominant hand | 158 |
| Heel | Sign is produced on the heel of the non-dominant hand | 88 |
| Other | Sign is produced in an unspecified location on the body | 707 |
| Neutral | Sign is not produced on or near the body | 3390 |

Table 7: Values and relative definitions for minor location

| Value | Definition | Cardinality |
|-------------------------------------|--|--------------------|
| One Handed | Sign only recruits one hand | 3939 |
| Symmetrical Or Alternating | Sign recruits both hands Phonological specifications for both hands are identical Movement of both hands is either symmetrical or alternating | 3358 |
| Asymmetrical Same Handshape | Sign recruits both hands Only the dominant hand moves The location and orientation of the hands may differ, but the other specifications of handshape are the same Non-Dominant hand must be an unmarked handshape (B A S I C O 5) | 938 |
| Asymmetrical Different Handshape | Sign recruits both hands Only the dominant hand moves The location and orientation of the hands may differ, and the other specifications of handshape are not the same Non-Dominant hand must be an unmarked handshape (B A S I C O 5) | 1639 |
| Other | Sign violates Battison's Symmetry and Dominance Conditions | 143 |

Table 8: Values and relative definitions for sign type

| Value | Definition | Cardinality |
|--------------|---|--------------------|
| Straight | Straight movement of the dominant hand through xyz space | 1938 |
| Curved | Single arc movement of the dominant hand through xyz space Hands may or may not make contact with multiple locations | 1255 |
| BackAndForth | Sequence of more than one straight or curved movements | 3549 |
| Circular | Circular movement of the dominant hand through space Rotation alone does not constitute a circular movement | 1129 |
| None | Entire sign (or first free morpheme) does not have a path movement | 1748 |
| Other | Sign has another unspecified path movement | 398 |

Table 9: Values and relative definitions for movement