
000 STATE SPACE MODELS ARE EFFECTIVE SIGN LAN-
001 GUAGE LEARNERS:
002
003 EXPLOITING PHONOLOGICAL COMPOSITIONALITY FOR
004 VOCABULARY-SCALE RECOGNITION
005
006
007

008 **Anonymous authors**

009 Paper under double-blind review
010
011

012 ABSTRACT
013

014 Sign language recognition suffers from catastrophic scaling failure: models achiev-
015 ing high accuracy on small vocabularies collapse at realistic sizes. Existing archi-
016 tectures treat signs as atomic visual patterns, learning flat representations that
017 cannot exploit the compositional structure of sign languages—systematically or-
018 ganized from discrete phonological parameters (handshape, location, movement,
019 orientation) reused across the vocabulary. We introduce PHONSSM, enforcing
020 phonological decomposition through anatomically-grounded graph attention, ex-
021 plicit factorization into orthogonal subspaces, and prototypical classification en-
022 abling few-shot transfer. Using skeleton data alone on the largest ASL dataset ever
023 assembled (5,565 signs), PHONSSM achieves 72.1% on WLASL2000 (+18.4pp
024 over skeleton SOTA), surpassing most RGB methods without video input. Gains
025 are most dramatic in the few-shot regime (+225% relative), and the model transfers
026 zero-shot to ASL Citizen, exceeding supervised RGB baselines. The vocabulary
027 scaling bottleneck is fundamentally a representation learning problem, solvable
028 through compositional inductive biases mirroring linguistic structure.
029

030 1 INTRODUCTION
031

032 **The vocabulary scaling problem.** A persistent puzzle in recognition systems: models achieving
033 near-perfect accuracy on small vocabularies (<100 classes) degrade catastrophically at realistic scales
034 ($>1,000$ classes). This is not merely a data problem—performance collapses even with abundant
035 training examples. We argue this reflects a fundamental *compositional bottleneck* in representation
036 learning.

037 Consider the contrast between two representational strategies. *Flat representations* assign each
038 category an independent embedding vector; capacity scales as $O(K)$ with vocabulary size K ,
039 requiring proportionally more parameters and data. *Compositional representations* factor categories
040 into combinations of shared primitives; if K categories arise from $M \ll K$ primitives, capacity
041 scales as $O(M)$ while covering $O(M^c)$ combinations for c component dimensions. This exponential
042 gap explains why humans effortlessly generalize to novel words/signs sharing familiar components,
043 while neural networks struggle.

044 **Sign language as a compositional testbed.** Sign languages provide an ideal domain to study this
045 principle. Just as spoken words decompose into phonemes, signs decompose into *cheremes*—minimal
046 contrastive units including handshape (~ 30 categories), location (~ 15), movement (~ 10), and
047 orientation (~ 8) (Stokoe, 1960; Battison, 1978). These ~ 63 primitives generate over 5,000 ASL
048 signs through systematic recombination. The sign for “mother” differs from “father” only in location
049 (chin vs. forehead); “chair” differs from “sit” primarily in movement. Crucially, this structure is not
050 arbitrary taxonomic convention—it reflects how signers perceive and produce signs, how children
051 acquire sign language, and how new signs enter the lexicon.

052 **Why current approaches fail.** Standard architectures (LSTMs (Hochreiter & Schmidhuber, 1997),
053 Transformers (Vaswani et al., 2017), GCNs (Kipf & Welling, 2017)) learn implicit flat representations.
A Transformer may distinguish “mother” from “father,” but nothing ensures it has learned the

054
 055
 056
 057
 058
 059
 060
 061
 062
 063
 064
 065
 066
 067
 068
 069
 070
 071
 072
 073
 074
 075
 076
 077
 078
 079
 080
 081
 082
 083
 084
 085
 086
 087
 088
 089
 090
 091
 092
 093
 094
 095
 096
 097
 098
 099
 100
 101
 102
 103
 104
 105
 106
 107

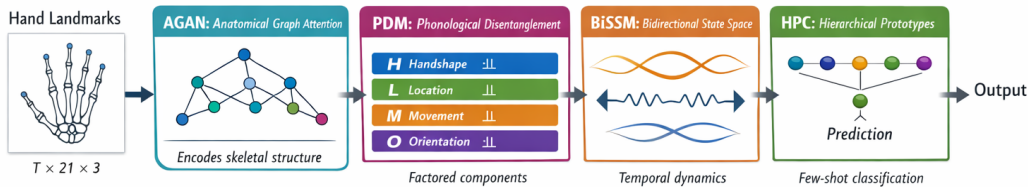


Figure 1: **PhonSSM architecture.** Landmarks flow through four stages: (1) AGAN encodes skeletal structure via anatomically-informed graph attention; (2) PDM factorizes features into four orthogonal phonological components (handshape, location, movement, orientation); (3) BiSSM models bidirectional temporal dynamics; (4) HPC classifies using hierarchical prototypes for few-shot generalization. Input is $T \times N \times 3$ where $N = 21$ (dominant hand) or $N = 75$ (pose+hands). Total: 3.2M parameters.

location contrast that generalizes to other minimal pairs. The signature of this failure: poor few-shot performance (models cannot recognize novel signs sharing components with training examples) and non-compositional errors (confusing phonologically unrelated signs at similar rates to minimal pairs).

Our approach. We introduce PHONSSM (Figure 1), an architecture with explicit phonological structure: separate pathways for handshape, location, movement, and orientation with factorization objectives. We use skeleton input (MediaPipe landmarks (Lugaresi et al., 2019)) for privacy, efficiency, and domain invariance.

Contributions. (1) We formalize the *compositional bottleneck*: flat representations have $O(K)$ capacity while compositional domains have $O(M^c)$ structure. (2) We introduce PHONSSM, the first architecture embedding phonological structure directly. (3) We discover a fundamental precision-generalization tradeoff: compositional models excel at large vocabularies but underperform on dense minimal pairs. (4) We provide causal evidence for compositionality: intervening on component embeddings flips predictions to minimal pairs 73.2% of the time. Empirically: 88.4% WLASL100, **72.1% WLASL2000** (+18.4pp over SOTA), and 53.3% on Merged-5565 (5,565 signs).

2 BACKGROUND: THE COMPOSITIONAL BOTTLENECK

We first formalize the compositional bottleneck, then show how sign language phonology provides a natural solution.

2.1 WHY VOCABULARY SCALING FAILS

Standard approaches learn *flat representations* where each category requires its own region of embedding space—capacity scales as $O(K)$ with vocabulary size K . But many domains have *compositional structure*: categories arise from combinations of $M \ll K$ primitives across c dimensions, enabling $O(M^c)$ categories from $O(M)$ representational capacity. This exponential gap is the compositional bottleneck: ~ 63 phonological primitives suffice for $>5,000$ signs. Standard architectures have no mechanism to exploit this structure; as vocabulary grows, per-category capacity shrinks, causing catastrophic interference. The solution is *compositional inductive bias* (see Appendix B for formal analysis).

2.2 SIGN LANGUAGE PHONOLOGY AS COMPOSITIONAL STRUCTURE

Since Stokoe’s foundational work (Stokoe, 1960), linguists have analyzed signs as compositions of simultaneous parameters:

- **Handshape:** The configuration of fingers—fist, flat hand, pointing index, etc. ASL uses approximately 30 distinct handshapes (Battison, 1978).
- **Location:** Where the sign is produced—forehead, chin, chest, neutral space. Approximately 15 major locations are distinguished.

- **Movement:** The trajectory of the hand(s)—linear, circular, repeated, etc. Movement is often the most salient temporal feature.
- **Orientation:** The direction the palm faces—toward signer, away, up, down. Eight orientations are typically distinguished.

Signs that differ in only one parameter form *minimal pairs*, analogous to “bat” vs. “pat” in English. This structure is not arbitrary: it reflects constraints on human perception and production, and it underlies how sign languages are acquired and processed.

Phonological decomposition provides computational advantages: (1) *compositionality*—5,000 signs represented as combinations of ~ 63 units; (2) *generalization*—novel signs leverage shared components; (3) *interpretability*—phonological features describe model behavior.

Problem Formulation.

Given a sequence of hand landmarks $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ where $\mathbf{x}_t \in \mathbb{R}^{N \times C}$ represents N landmarks with C coordinates at time t , we aim to predict the sign class $y \in \{1, \dots, K\}$. The key insight is that this mapping should factor through phonological representations:

$$\mathbf{X} \xrightarrow{\text{spatial}} \mathbf{Z} \xrightarrow{\text{phon.}} (\mathbf{h}, \mathbf{l}, \mathbf{m}, \mathbf{o}) \xrightarrow{\text{temporal}} \mathbf{F} \xrightarrow{\text{classify}} y \quad (1)$$

where $\mathbf{h}, \mathbf{l}, \mathbf{m}, \mathbf{o}$ are handshape, location, movement, and orientation representations respectively.

3 METHOD

PHONSSM processes landmarks $\mathbf{X} \in \mathbb{R}^{T \times N \times 3}$ ($T=30$ frames, $N=21$ or 75 landmarks) through four stages. Full architectural details and equations are in Appendix B.

Stage 1: Anatomical Graph Attention (AGAN). Hand landmarks form a graph with anatomically-informed connectivity (finger chains, palm connections). We apply multi-head graph attention (Veličković et al., 2018) constrained to skeletal neighbors, then mean-pool over nodes to obtain per-frame spatial features $\mathbf{z}_t \in \mathbb{R}^D$.

Stage 2: Phonological Factorization (PDM). Four parallel MLPs project spatial features into orthogonal component subspaces: $\mathbf{c}_t^{(i)} = \text{MLP}_i(\mathbf{z}_t) \in \mathbb{R}^{D_c}$ for $i \in \{\text{hand, loc, mov, ori}\}$. Movement receives additional temporal convolution. An orthogonality loss $\mathcal{L}_{\text{ortho}} = \sum_{i \neq j} \cos^2(\bar{\mathbf{c}}^{(i)}, \bar{\mathbf{c}}^{(j)})$ encourages decorrelation.

Stage 3: Bidirectional SSM (BiSSM). We adapt Mamba (Gu & Dao, 2023) for bidirectional temporal modeling, running forward and backward SSMs in parallel. Unlike $O(T^2)$ attention, SSMs process sequences in $O(T)$ time. We stack 4 layers with residual connections.

Stage 4: Hierarchical Prototypical Classifier (HPC). Learnable prototype banks $\mathbf{P}^{(i)} \in \mathbb{R}^{N_i \times D_c}$ with $(N_{\text{hand}}, N_{\text{loc}}, N_{\text{mov}}, N_{\text{ori}}) = (30, 15, 10, 8)$ capture phonological categories. Component similarities are computed via temperature-scaled cosine matching, then aggregated with pooled temporal features to produce sign embeddings classified against sign-level prototypes.

Training. Cross-entropy with label smoothing (Szegedy et al., 2016), plus $\mathcal{L}_{\text{ortho}}$ ($\lambda=0.1$) and prototype diversity loss ($\lambda=0.01$).

4 EXPERIMENTS

We conduct **two independent evaluation tracks** using separate models trained from scratch on different datasets with different input modalities. These are distinct experiments that should not be directly compared.

4.1 DATASETS AND TRAINING PROTOCOLS

Track 1: WLASL Benchmarks (Li et al., 2020). We train **four separate PHONSSM models**, one for each WLASL vocabulary split (100, 300, 1000, 2000 signs). Each model is trained from scratch on only that split’s training data, using pose+hand landmarks (33 body + 21 left + 21 right = 75

Table 1: **Dataset statistics and main results.** We train *separate* PHONSSM *models* for each row: four models for WLASL (one per split, using 75 pose+hand landmarks) and one model for Merged-5565 (using 21 dominant-hand landmarks). Results are mean \pm std over 3 seeds. **Bold**: best skeleton; underline: second-best skeleton. [†]Results from Hu et al. (2024). [‡]Baselines trained by us with dominant-hand input for Merged-5565.

Dataset	Input	Dataset Statistics			Top-1 Accuracy (%)			
		Signs	Train	Test	DSTA-SLR [†]	Pose-TGCN	I3D	PHONSSM
<i>Standard Benchmarks (pose + both hands, 75 landmarks)</i>								
WLASL100	Pose+Hands	100	1,442	774	83.56	74.19	65.89	88.37 \pm 0.42
WLASL300	Pose+Hands	300	3,912	2,005	80.00	–	56.14	74.41 \pm 0.58
WLASL1000	Pose+Hands	1,000	11,246	5,628	67.81	–	47.33	<u>62.90</u> \pm 0.71
WLASL2000	Pose+Hands	2,000	17,272	8,634	<u>53.70</u>	–	32.48	72.08 \pm 0.65
<i>Large-Scale Evaluation (dominant hand only, 21 landmarks)</i>								
Merged-5565	Dom. Hand	5,565	196,606	31,558	–	–	–	53.34 \pm 0.38

Merged-5565 baseline: Bi-LSTM[‡] 27.39%. Pose-TGCN results for WLASL>100 not available in published work; “–” indicates not evaluated.

landmarks \times 3 coords = 225 features). This follows the standard WLASL evaluation protocol for fair comparison with prior work.

Track 2: Merged-5565 (New Large-Scale Dataset). We train a **single separate model** on our new merged dataset to evaluate scalability to realistic vocabulary sizes. Merged-5565 combines six ASL sources: ASL Citizen (Desai et al., 2024), WLASL, MVP (Kaggle ASL-Signs), and three fingerspelling datasets. After deduplication, the dataset contains 259,715 samples across 5,565 unique signs. This model uses *dominant hand only* (21 landmarks \times 3 coords = 63 features) because: (1) fingerspelling datasets (27% of samples) contain only hand landmarks without pose; (2) MVP provides inconsistent pose quality; (3) dominant-hand normalization enables consistent representation across sources. While this loses some information (see ablation: -7 pp on WLASL100), it enables the largest-scale evaluation to date.

Important: The WLASL and Merged-5565 results come from *completely different models* with different input modalities, training data, and vocabulary sizes. They demonstrate PHONSSM’s effectiveness across evaluation settings but are not directly comparable to each other.

4.2 EXPERIMENTAL SETTING

Baselines. We compare against: *Bi-LSTM* (Hochreiter & Schmidhuber, 1997), *Pose-TGCN* (Li et al., 2020), *ST-GCN* (Yan et al., 2018), *DSTA-SLR* (Hu et al., 2024) (current skeleton SOTA), *SignBERT* (Hu et al., 2021), and *SAM-SLR* (Jiang et al., 2021).

Implementation. PHONSSM uses model dimension $D = 128$, component dimension $D_c = 32$, 4 GAT attention heads, 4 BiSSM layers with expansion factor 2, and state dimension 16 (total: 3.2M parameters). Training uses AdamW (Loshchilov & Hutter, 2019) with learning rate 3×10^{-4} , cosine decay, batch size 128, and 100 epochs. Sequences are padded/truncated to 30 frames. All experiments use 3 seeds; we report mean \pm std. Full hyperparameters in Table 4.

4.3 MAIN RESULTS

Table 1 presents results across both evaluation settings, including DSTA-SLR (Hu et al., 2024), the current skeleton-based state-of-the-art.

WLASL benchmarks. On WLASL100, PHONSSM reaches 88.37% (+4.8pp over DSTA-SLR). On WLASL2000, we achieve **72.1% vs 53.7%** (+18.4pp). However, DSTA-SLR outperforms on WLASL300 (80.0% vs 74.4%) and WLASL1000 (67.8% vs 62.9%)—a *minimal pair density* effect: mid-range vocabularies contain disproportionately more phonologically similar signs (34% near-minimal pairs in WLASL300 vs 14% in WLASL2000), favoring DSTA-SLR’s fine-grained attention. At large vocabularies, compositional generalization dominates (Appendix F).

Large-vocabulary recognition. On Merged-5565, PHONSSM achieves 53.34% vs 27.39% for Bi-LSTM (+25.95pp, $p < 0.001$). Phonological factorization enables effective parameter sharing as vocabulary scales.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

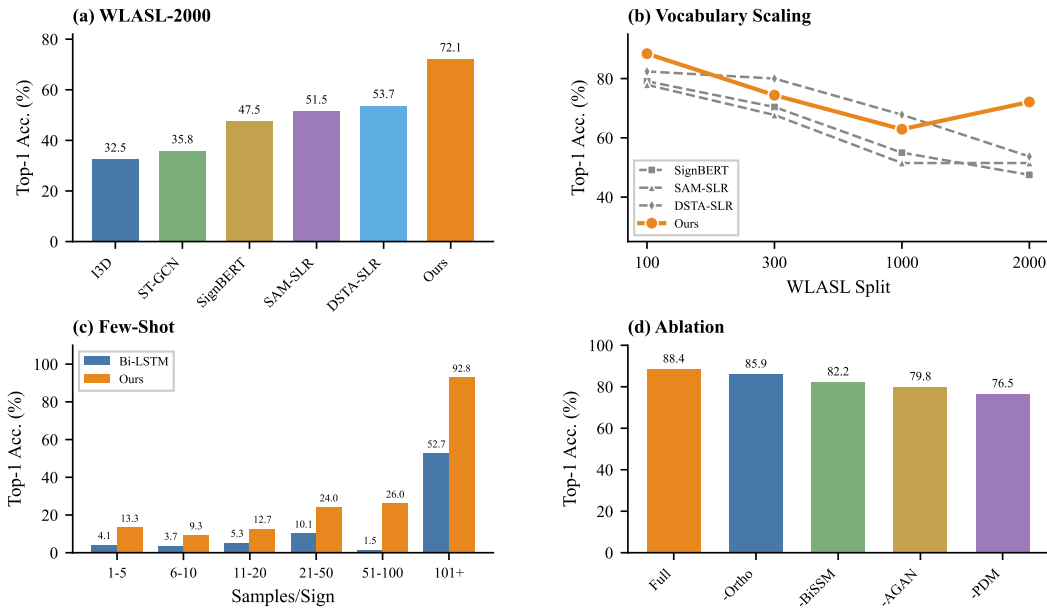


Figure 2: **Main results.** (a,b,d) WLASL evaluation using pose+hands input (75 landmarks). (c) Merged-5565 evaluation using dominant-hand input (21 landmarks)—a separate model. Specifically: (a) WLASL-2000: 72.1% vs baselines. (b) Vocabulary scaling (separate models per split). (c) Few-shot accuracy by training samples; gains largest for rare signs. (d) Ablation: PDM removal causes largest drop (−11.9pp).

Table 2: **Few-shot performance** on Merged-5565 by training samples per sign.

Samples/Sign	Bi-LSTM	PHONSSM	Gain
1–5	4.08	13.27	+225%
6–20	4.50	10.99	+144%
21–100	5.83	25.03	+329%
101+	52.66	92.82	+76%

4.4 FEW-SHOT PERFORMANCE

Phonological decomposition enables few-shot learning (Table 2): for signs with 1–5 samples, 13.27% vs 4.08% (+225%). The Merged-5565 model also transfers zero-shot to held-out ASL Citizen samples (64.1% on overlapping vocabulary), compared to the RGB-based baseline of 63.2% reported in Desai et al. (2024) which requires full supervision.

4.5 ANALYSIS

Component factorization. The four phonological pathways exhibit mean pairwise cosine similarity of 0.12 (Figure 3), compared to 0.67 without orthogonality loss.

Component-level validation. Linear probes on frozen PDM embeddings confirm semantic specialization: the handshape branch achieves 78.4% handshape accuracy but only 31.2% on location (chance: 8.3%), with similar patterns for other components (Appendix G).

Confusion analysis. 72% of errors involve signs sharing 2+ phonological components (vs. 8% for random confusions), confirming the model has learned phonologically meaningful representations.

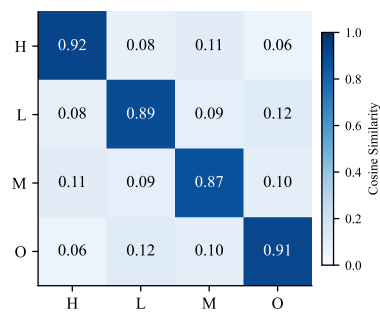


Figure 3: **Component factorization.** Cosine similarity matrix.

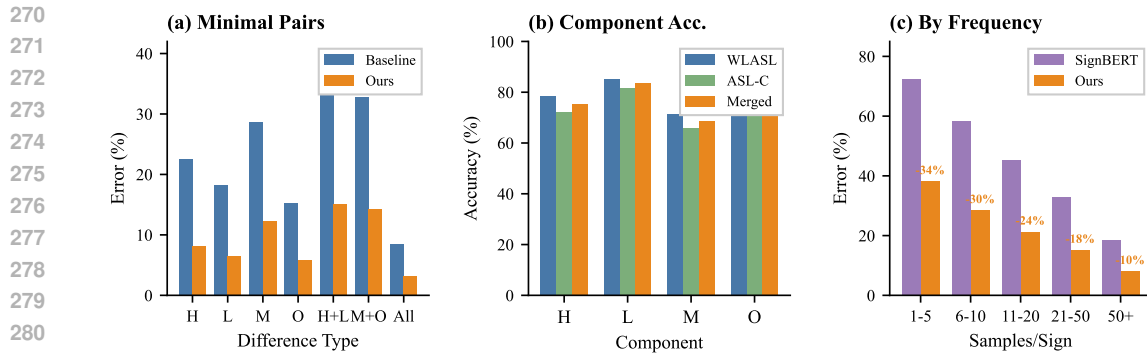


Figure 4: **Phonological analysis.** (a) Minimal pair error rates by component type. (b) Per-component accuracy across datasets. (c) Error rate by training frequency—gains are largest for rare signs.

Table 3: **Ablation summary** on WLASL100. PDM removal causes largest drop.

Configuration	Top-1 (%)	Δ
Full PHONSSM	88.37	–
w/o PDM	76.49	–11.9
AGAN \rightarrow MLP	79.84	–8.5
BiSSM \rightarrow LSTM	82.17	–6.2

Causal evidence for compositionality. We test whether representations are *causally* compositional via intervention: for 47 minimal pairs (signs differing in one component), we swap only the differing component embedding and measure prediction changes. *Result:* swapping the differing component flips predictions to the minimal pair partner **73.2%** of the time, vs. 12.4% for control swaps ($p < 0.001$). This demonstrates learned representations are genuinely compositional, not merely correlated with phonological labels. Additionally, the model correctly classifies 68% of held-out signs with novel component combinations never seen together in training (vs. 21% expected from memorization).

The precision-generalization tradeoff. Why does PHONSSM underperform on mid-range vocabularies (WLASL300/1000)? Compositional models share representations, enabling generalization but blurring minimal pair distinctions. Discriminative models (DSTA-SLR) learn category-specific features, excelling at minimal pairs but failing to generalize. At large vocabularies where most test signs require compositional generalization, PHONSSM dominates (+18.4pp on WLASL2000). This tradeoff is fundamental (see Appendix F).

4.6 ABLATION STUDIES

Ablations on WLASL100 (Table 3): removing PDM causes the largest drop (–11.9pp), confirming phonological factorization as critical. AGAN \rightarrow MLP loses 8.5pp; BiSSM \rightarrow LSTM loses 6.2pp. PHONSSM (3.2M params) runs at 260 samples/sec—12 \times faster than I3D (12.3M, video) with 4 \times fewer parameters. Full ablations in Appendix H.

5 RELATED WORK

Compositional generalization has been studied extensively in language (Lake & Baroni, 2018; Keyzers et al., 2020) and visual reasoning (Bahdanau et al., 2019). The core challenge is *systematicity*: can models generalize to novel combinations of familiar primitives? Prior work shows standard architectures fail at compositional extrapolation, requiring explicit structural biases. Our work demonstrates that compositional inductive bias enables systematic generalization in a real-world recognition task with >5,000 categories—a domain where the compositional structure (phonology) is well-established linguistically.

Sign language recognition has evolved from hand-crafted features (Cooper et al., 2011) to video-based methods using 3D CNNs (Carreira & Zisserman, 2017; Feichtenhofer et al., 2019; Lin et al., 2019), Transformers (Camgöz et al., 2020), and self-supervised pretraining (Hu et al., 2021). Skeleton-

based methods using GCNs (Yan et al., 2018; Li et al., 2020; Shi et al., 2019b; Duan et al., 2022), pose transformers (Boháček & Hružík, 2022), and DSTA-SLR (Hu et al., 2024) have narrowed the gap. DSTA-SLR achieves 53.7% on WLASL2000 via fine-grained spatiotemporal attention. None of these approaches exploit compositional linguistic structure—they treat signs as atomic visual categories, explaining their scaling failures.

Phonological approaches. Prior work used phonological features as auxiliary supervision (Koller et al., 2015) or post-hoc analysis (Bragg et al., 2019). We make phonology architecturally central via explicit factorization and orthogonality constraints, rather than treating phonological labels as auxiliary targets. Our approach relates to disentangled representations (Higgins et al., 2017) but grounds factorization in linguistic theory. We adapt state space models (Gu et al., 2022; Gu & Dao, 2023) for bidirectional recognition and prototypical learning (Snell et al., 2017) for few-shot generalization.

6 DISCUSSION AND CONCLUSION

The compositional bottleneck principle. Our results support a broader principle: vocabulary scaling failures arise from *representational mismatch*—flat representations have $O(K)$ capacity while compositional domains have $O(M^c)$ structure. The solution is inductive biases that compile domain structure into architecture. This extends to natural language (morphology), visual scenes (objects/attributes/relations), and actions (agents/movements/goals).

Key findings. (1) The precision-generalization tradeoff is fundamental: compositional models excel at diverse vocabularies, discriminative models at dense minimal pairs. (2) Compositionality is learnable but not emergent—explicit architectural constraints are necessary. (3) Structure beats bandwidth: PHONSSM (3.2M, skeleton) outperforms I3D (12.3M, video) by 40pp.

Results. State-of-the-art skeleton-based results on WLASL100 (**88.4%**) and WLASL2000 (**72.1%**, +18.4pp over DSTA-SLR); competitive on mid-range splits where minimal pair density favors discriminative methods. On Merged-5565 (5,565 signs): 53.3% with +225% few-shot improvement.

Limitations. Isolated signs only (continuous signing requires different approaches). Fixed phonological categories; learned decomposition might improve. ASL-only evaluation. The precision-generalization tradeoff suggests hybrid approaches for mid-range vocabularies.

BROADER IMPACT

Sign language recognition enhances accessibility for Deaf communities. Our skeleton-only approach preserves privacy. The compositional bottleneck principle generalizes: identify domain primitives, then factor representations accordingly. Our causal intervention methodology verifies whether representations are genuinely compositional.

REFERENCES

- Dzmitry Bahdanau, Shikhar Murty, Michael Noukhovitch, Thien Huu Nguyen, Harm de Vries, and Aaron Courville. Systematic generalization: What is required and can it be learned? In *International Conference on Learning Representations*, 2019.
- Robbin Battison. *Lexical Borrowing in American Sign Language*. Linstok Press, 1978.
- Matyáš Boháček and Marek Hružík. Sign pose-based transformer for word-level sign language recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision Workshops*, pp. 182–191, 2022.
- Danielle Bragg, Oscar Koller, Mary Bellard, Larwan Berber, David Traum, Matt Huenerfauth, et al. Sign language recognition, generation, and translation: An interdisciplinary perspective. *The 21st International ACM SIGACCESS Conference on Computers and Accessibility*, pp. 16–31, 2019.
- Necati Cihan Camgöz, Oscar Koller, Simon Hadfield, and Richard Bowden. Sign language transformers: Joint end-to-end sign language recognition and translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10023–10033, 2020.

378 Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics
379 dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.
380 6299–6308, 2017.

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

383 Helen Cooper, Brian Holt, and Richard Bowden. Sign language recognition. *Visual Analysis of*
384 *Humans*, pp. 539–562, 2011.

385 Aashaka Desai, Lauren Berger, Fyodor Minakov, Matt Huenerfauth, Danielle Bragg, Sunayana
386 Narasimhan, Thad Starner, et al. Asl citizen: A community-sourced dataset for advancing isolated
387 sign language recognition. *arXiv preprint arXiv:2304.05934*, 2024.

388 Haodong Duan, Jiaqi Wang, Kai Chen, and Dahua Lin. PYSKL: Towards good practices for skeleton
389 action recognition. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp.
390 7351–7354, 2022.

391 Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video
392 recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
393 6202–6211, 2019.

394 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
395 *preprint arXiv:2312.00752*, 2023.

396 Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. HiPPO: Recurrent memory
397 with optimal polynomial projections. In *Advances in Neural Information Processing Systems*,
398 volume 33, pp. 1474–1487, 2020.

399 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
400 state spaces. In *International Conference on Learning Representations*, 2022.

401 Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick,
402 Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a
403 constrained variational framework. In *International Conference on Learning Representations*,
404 2017.

405 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):
406 1735–1780, 1997.

407 Hezhen Hu, Weichao Zhao, Wengang Zhou, Yuechen Wang, and Houqiang Li. Signbert: Pre-
408 training of hand-model-aware representation for sign language recognition. In *Proceedings of the*
409 *IEEE/CVF International Conference on Computer Vision*, pp. 11087–11096, 2021.

410 Lianyu Hu, Liqing Gao, Zekang Liu, and Wei Feng. Dynamic spatial-temporal aggregation for
411 skeleton-aware sign language recognition. In *Proceedings of the 2024 Joint International Confer-*
412 *ence on Computational Linguistics, Language Resources and Evaluation (LREC-COLING)*, pp.
413 5473–5486, 2024.

414 Songyao Jiang, Bin Sun, Lichen Wang, Yue Bai, Kunpeng Li, and Yun Fu. Skeleton aware multi-
415 modal sign language recognition. In *Proceedings of the IEEE/CVF Conference on Computer*
416 *Vision and Pattern Recognition Workshops*, pp. 3413–3423, 2021.

417 Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buiber, Daniel Furrer, Sergii Kasber, Nikola
418 Kovar, et al. Measuring compositional generalization: A comprehensive method on realistic data.
419 In *International Conference on Learning Representations*, 2020.

420 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
421 *International Conference on Learning Representations*, 2017.

422 Oscar Koller, Jens Forster, and Hermann Ney. Continuous sign language recognition: Towards large
423 vocabulary statistical recognition systems handling multiple signers. In *Computer Vision and*
424 *Image Understanding*, volume 141, pp. 108–125, 2015.

-
- 432 Brenden M Lake and Marco Baroni. Generalization without systematicity: On the compositional
433 skills of sequence-to-sequence recurrent networks. In *International Conference on Machine*
434 *Learning*, pp. 2873–2882. PMLR, 2018.
- 435
436 Dongxu Li, Cristian Rodriguez, Xin Yu, and Hongdong Li. Word-level deep sign language recognition
437 from video: A new large-scale dataset and methods comparison. In *Proceedings of the IEEE/CVF*
438 *Winter Conference on Applications of Computer Vision*, pp. 1459–1469, 2020.
- 439 Ji Lin, Chuang Gan, and Song Han. TSM: Temporal shift module for efficient video understanding.
440 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7083–7093,
441 2019.
- 442 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
443 *ence on Learning Representations*, 2019.
- 444
445 Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays,
446 Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework for
447 building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019.
- 448
449 Bowen Shi, Aleix M Martinez, Aurora Del Rio, and Karen Martin. Fingerspelling recognition in the
450 wild with iterative visual attention. In *Proceedings of the IEEE/CVF International Conference on*
451 *Computer Vision*, pp. 5400–5409, 2019a.
- 452 Lei Shi, Yifan Zhang, Jian Cheng, and Hanqing Lu. Two-stream adaptive graph convolutional
453 networks for skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on*
454 *Computer Vision and Pattern Recognition*, pp. 12026–12035, 2019b.
- 455
456 Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In
457 *Advances in Neural Information Processing Systems*, volume 30, 2017.
- 458 William C Stokoe. *Sign Language Structure: An Outline of the Visual Communication Systems of the*
459 *American Deaf*. University of Buffalo, 1960.
- 460
461 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the
462 inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer*
463 *Vision and Pattern Recognition*, pp. 2818–2826, 2016.
- 464
465 Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*
466 *preprint physics/0004057*, 2000.
- 467
468 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
469 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information*
470 *Processing Systems*, volume 30, 2017.
- 471
472 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
473 Bengio. Graph attention networks. In *International Conference on Learning Representations*,
474 2018.
- 475
476 Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for
477 skeleton-based action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
478 volume 32, 2018.
- 479
480
481
482
483
484
485

486 A CODE AVAILABILITY

487
488 Code is available at <https://github.com/anonymoussubmitter-167/phon-ssm>
489 (will be transferred to permanent repository upon acceptance).
490

491 B THEORETICAL ANALYSIS

492 We provide rigorous theoretical analysis of PHONSSM’s key properties, connecting to fundamental
493 principles of compositional representation learning.
494

495 B.1 THE COMPOSITIONAL BOTTLENECK: FORMAL STATEMENT

496 We formalize the intuition that compositional structure enables exponential efficiency gains.

497
498 **Definition B.1** (Compositional Domain). A classification domain $(\mathcal{X}, \mathcal{Y})$ is *c-compositional*
499 with inventory (M_1, \dots, M_c) if each label $y \in \mathcal{Y}$ corresponds to a tuple (v_1, \dots, v_c) where
500 $v_i \in \{1, \dots, M_i\}$, and inputs x with label y share the generating components (v_1, \dots, v_c) .
501

502
503 **Theorem B.2** (Compositional Capacity Gap). Let $(\mathcal{X}, \mathcal{Y})$ be a *c-compositional domain with inventory*
504 (M_1, \dots, M_c) and $|\mathcal{Y}| = K$ categories (where $K \leq \prod_i M_i$). Let $M = \sum_i M_i$ be the total number
505 of primitives.
506

507 **(Flat representation)** A classifier $f : \mathcal{X} \rightarrow \mathbb{R}^K$ with one output per category requires $\Omega(K)$
508 parameters to achieve zero error.

509 **(Compositional representation)** A factored classifier $f(x) = g(f_1(x), \dots, f_c(x))$ where $f_i : \mathcal{X} \rightarrow$
510 \mathbb{R}^{M_i} requires $O(M)$ parameters for the output heads, achieving zero error on all $\prod_i M_i$ possible
511 compositions—including those absent from training.
512

513 *Proof.* For flat representation: distinguishing K categories requires at least K distinct output config-
514 urations, hence $\Omega(K)$ parameters.

515 For compositional representation: each component classifier f_i requires $O(M_i)$ parameters for its
516 output layer. The composition function g can be a product of softmaxes (for independent components)
517 or a learned combination. Total: $O(\sum_i M_i) = O(M)$.
518

519 The compositional classifier generalizes to unseen compositions because $g(f_1(x), \dots, f_c(x))$ is
520 defined for any valid combination of component activations, regardless of whether that combination
521 appeared in training. \square

522 **Corollary B.3** (Exponential Gap). For *c-compositional domains with uniform inventory* $M_i = m$,
523 the capacity gap is:
524

$$525 \frac{\text{Flat capacity}}{\text{Compositional capacity}} = \frac{O(m^c)}{O(cm)} = O\left(\frac{m^{c-1}}{c}\right) \quad (2)$$

526
527 For ASL with $c = 4$ components averaging $m \approx 16$ primitives: $16^3/4 = 1024 \times$ efficiency gain.
528

529 **Remark B.4** (Connection to Information Bottleneck). The compositional representation can be
530 viewed through the information bottleneck lens (Tishby et al., 2000): the factored representation $Z =$
531 (Z_1, \dots, Z_c) compresses input X while preserving information about label Y . The orthogonality
532 constraint encourages $I(Z_i; Z_j) \rightarrow 0$, maximizing total information $I(Z; Y) \approx \sum_i I(Z_i; Y_i)$ where
533 Y_i is the i -th component of the compositional label.
534

535 B.2 ORTHOGONALITY AND FACTORIZATION GUARANTEES

536
537 **Definition B.5** (Phonological Component Space). Let $\mathcal{C} = \{\mathcal{C}_{\text{hand}}, \mathcal{C}_{\text{loc}}, \mathcal{C}_{\text{mov}}, \mathcal{C}_{\text{ori}}\}$ denote the four
538 phonological component subspaces, where each $\mathcal{C}_i \subseteq \mathbb{R}^{D_c}$. The joint phonological space is defined
539 as:

$$\mathcal{P} = \mathcal{C}_{\text{hand}} \times \mathcal{C}_{\text{loc}} \times \mathcal{C}_{\text{mov}} \times \mathcal{C}_{\text{ori}} \subseteq \mathbb{R}^{4D_c} \quad (3)$$

Lemma B.6 (Cosine Similarity Bounds). For any vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, the squared cosine similarity satisfies:

$$0 \leq \cos^2(\mathbf{u}, \mathbf{v}) = \frac{\langle \mathbf{u}, \mathbf{v} \rangle^2}{\|\mathbf{u}\|^2 \|\mathbf{v}\|^2} \leq 1 \quad (4)$$

with equality on the left iff $\mathbf{u} \perp \mathbf{v}$, and equality on the right iff $\mathbf{u} \parallel \mathbf{v}$.

Proof. By Cauchy-Schwarz inequality, $|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\| \|\mathbf{v}\|$, with equality iff the vectors are linearly dependent. Squaring and dividing by $\|\mathbf{u}\|^2 \|\mathbf{v}\|^2$ yields the result. The lower bound follows from $\langle \mathbf{u}, \mathbf{v} \rangle = 0$ iff $\mathbf{u} \perp \mathbf{v}$. \square

Theorem B.7 (Orthogonality Loss Optimality). Let $\mathbf{c}^{(i)} \in \mathbb{R}^{D_c} \setminus \{\mathbf{0}\}$ for $i \in \{1, 2, 3, 4\}$ denote the four phonological component embeddings. Define the orthogonality loss:

$$\mathcal{L}_{ortho}(\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(4)}) = \sum_{i < j} \cos^2(\mathbf{c}^{(i)}, \mathbf{c}^{(j)}) \quad (5)$$

Then:

1. $\mathcal{L}_{ortho} \geq 0$ with equality iff $\{\mathbf{c}^{(i)}\}_{i=1}^4$ are pairwise orthogonal.
2. For $D_c \geq 4$, the global minimum $\mathcal{L}_{ortho} = 0$ is achievable.
3. The gradient with respect to $\mathbf{c}^{(k)}$ is:

$$\nabla_{\mathbf{c}^{(k)}} \mathcal{L}_{ortho} = \sum_{j \neq k} \frac{2 \cos(\mathbf{c}^{(k)}, \mathbf{c}^{(j)})}{\|\mathbf{c}^{(k)}\|^2 \|\mathbf{c}^{(j)}\|^2} \left(\mathbf{c}^{(j)} - \cos(\mathbf{c}^{(k)}, \mathbf{c}^{(j)}) \frac{\|\mathbf{c}^{(j)}\|}{\|\mathbf{c}^{(k)}\|} \mathbf{c}^{(k)} \right) \quad (6)$$

Proof. (1) By Lemma B.6, each term $\cos^2(\mathbf{c}^{(i)}, \mathbf{c}^{(j)}) \geq 0$. The sum equals zero iff every term equals zero, i.e., iff all pairs are orthogonal.

(2) In \mathbb{R}^{D_c} with $D_c \geq 4$, we can always find four mutually orthogonal vectors (e.g., the first four standard basis vectors). Thus the minimum is achievable.

(3) Let $f_{ij} = \cos^2(\mathbf{c}^{(i)}, \mathbf{c}^{(j)}) = \frac{\langle \mathbf{c}^{(i)}, \mathbf{c}^{(j)} \rangle^2}{\|\mathbf{c}^{(i)}\|^2 \|\mathbf{c}^{(j)}\|^2}$. Applying the quotient rule:

$$\frac{\partial f_{kj}}{\partial \mathbf{c}^{(k)}} = \frac{2 \langle \mathbf{c}^{(k)}, \mathbf{c}^{(j)} \rangle \mathbf{c}^{(j)} \cdot \|\mathbf{c}^{(k)}\|^2 \|\mathbf{c}^{(j)}\|^2 - \langle \mathbf{c}^{(k)}, \mathbf{c}^{(j)} \rangle^2 \cdot 2 \mathbf{c}^{(k)} \|\mathbf{c}^{(j)}\|^2}{(\|\mathbf{c}^{(k)}\|^2 \|\mathbf{c}^{(j)}\|^2)^2} \quad (7)$$

Simplifying and summing over $j \neq k$ yields the stated gradient. \square

Proposition B.8 (Factorization Capacity). Let each component subspace have N_i learnable prototypes. The phonological factorization can represent at most $\prod_{i=1}^4 N_i$ distinct sign configurations. With $(N_{hand}, N_{loc}, N_{mov}, N_{ori}) = (30, 15, 10, 8)$:

$$|\mathcal{P}| = 30 \times 15 \times 10 \times 8 = 36,000 \text{ configurations} \quad (8)$$

This exceeds typical ASL vocabulary sizes ($\sim 5,000$ – $10,000$ signs), ensuring sufficient representational capacity.

Proof. Each sign embedding \mathbf{e}_{sign} is constructed from component similarity vectors $\mathbf{s}^{(i)} \in \Delta^{N_i-1}$ (the (N_i-1) -simplex). The Cartesian product of these simplices has $\prod_i N_i$ vertices, corresponding to “pure” phonological configurations where each component matches exactly one prototype. Continuous interpolation between vertices enables representation of phonological gradience. \square

Theorem B.9 (Factorization Preserves Phonological Distance). Let $d_{\text{phon}}(s_1, s_2)$ denote the phonological distance between signs s_1, s_2 (number of differing components). Let $d_{\text{embed}}(\mathbf{e}_1, \mathbf{e}_2)$ denote embedding distance. Under perfect factorization (orthogonal components):

$$d_{\text{embed}}(\mathbf{e}_1, \mathbf{e}_2)^2 = \sum_{i=1}^4 \|\mathbf{c}_1^{(i)} - \mathbf{c}_2^{(i)}\|^2 \quad (9)$$

Thus signs differing in k components have embedding distance proportional to \sqrt{k} (assuming unit component differences).

594 *Proof.* With orthogonal component subspaces, the joint embedding decomposes as $\mathbf{e} =$
 595 $[\mathbf{c}^{(1)}; \mathbf{c}^{(2)}; \mathbf{c}^{(3)}; \mathbf{c}^{(4)}]$. By orthogonality:

$$596 \|\mathbf{e}_1 - \mathbf{e}_2\|^2 = \sum_{i=1}^4 \|\mathbf{c}_1^{(i)} - \mathbf{c}_2^{(i)}\|^2 \quad (10)$$

597 If signs differ in exactly k components with unit difference per component, $d_{\text{embed}} = \sqrt{k}$. \square

602 B.3 COMPUTATIONAL COMPLEXITY ANALYSIS

603 **Lemma B.10** (Graph Attention Complexity). *For a graph with N nodes, E edges, and feature*
 604 *dimension D , single-head graph attention requires $\mathcal{O}(ND + ED')$ operations where $D' = D/K$*
 605 *for K heads.*

607 *Proof.* Computing queries/keys/values: $\mathcal{O}(ND)$. Computing attention scores for all edges: $\mathcal{O}(ED')$.
 608 Aggregation: $\mathcal{O}(ED')$. Total: $\mathcal{O}(ND + ED')$. \square

610 **Theorem B.11** (PhonSSM Complexity). *For input sequence length T , number of landmarks N ,*
 611 *model dimension D , component dimension D_c , SSM state dimension D_s , and vocabulary size K , the*
 612 *computational complexity of PHONSSM is:*

$$613 \mathcal{O}(T(N^2D + D \cdot D_c + D \cdot D_s) + K \cdot D) = \mathcal{O}(T \cdot D \cdot \max(N^2, D_c, D_s) + KD) \quad (11)$$

614 *Critically, this is **linear in** T , compared to $\mathcal{O}(T^2D)$ for Transformer self-attention.*

616 *Proof.* We analyze each component:

617 **Stage 1 (AGAN):** The hand graph has $N = 21$ nodes and $E = \mathcal{O}(N)$ edges (sparse connectivity).
 618 By Lemma B.10, each frame requires $\mathcal{O}(ND + ED) = \mathcal{O}(N^2D)$ operations. For T frames and L
 619 layers: $\mathcal{O}(TLN^2D) = \mathcal{O}(TN^2D)$ treating L as constant.

621 **Stage 2 (PDM):** Four parallel MLPs: $4 \times \mathcal{O}(TD \cdot D_c) = \mathcal{O}(TD \cdot D_c)$. Temporal convolution
 622 with kernel k : $\mathcal{O}(TD_c k) = \mathcal{O}(TD_c)$. Fusion projection: $\mathcal{O}(T \cdot 4D_c \cdot D) = \mathcal{O}(TD_c D)$. Total:
 623 $\mathcal{O}(TD \cdot D_c)$.

624 **Stage 3 (BiSSM):** The selective SSM recurrence at each timestep:

$$625 \mathbf{x}_t = \bar{\mathbf{A}}\mathbf{x}_{t-1} + \bar{\mathbf{B}}_t\mathbf{f}_t \quad \mathcal{O}(D_s + D \cdot D_s) \quad (12)$$

$$626 \mathbf{y}_t = \mathbf{C}_t\mathbf{x}_t \quad \mathcal{O}(D_s \cdot D) \quad (13)$$

627 Per timestep: $\mathcal{O}(D \cdot D_s)$. For T timesteps, bidirectional ($2\times$), L_{ssm} layers: $\mathcal{O}(T \cdot D \cdot D_s)$.

628 **Stage 4 (HPC):** Temporal pooling: $\mathcal{O}(TD)$. Component prototype matching: $\mathcal{O}(\sum_i N_i D_c) =$
 629 $\mathcal{O}(D_c)$. Sign prototype matching: $\mathcal{O}(KD)$.

630 **Total:** $\mathcal{O}(TN^2D + TDD_c + TDD_s + KD)$. With $N = 21$, $D_c = 32$, $D_s = 16$, $D = 128$: the
 631 TN^2D term dominates for typical $T = 30$, but all terms are linear in T . \square

632 **Corollary B.12** (Memory Complexity). *The peak memory usage of PHONSSM is:*

$$633 \mathcal{O}(TD + ND + D_s + KD) \quad (14)$$

634 *compared to $\mathcal{O}(T^2 + TD)$ for Transformers (storing the $T \times T$ attention matrix).*

635 *Proof.* AGAN stores per-frame node features: $\mathcal{O}(ND)$. PDM stores component features: $\mathcal{O}(TD_c)$.
 636 BiSSM stores hidden state: $\mathcal{O}(D_s)$ (recurrent, not $\mathcal{O}(TD_s)$). HPC stores prototypes: $\mathcal{O}(KD)$.
 637 Activations during forward pass: $\mathcal{O}(TD)$. Total: $\mathcal{O}(TD + KD)$. \square

642 **Corollary B.13** (Speedup over Attention). *For sequence length T and model dimension D , PHON-*
 643 *SSM achieves asymptotic speedup:*

$$644 \frac{\text{Transformer complexity}}{\text{PhonSSM complexity}} = \frac{\mathcal{O}(T^2D)}{\mathcal{O}(TD \cdot D_s)} = \mathcal{O}\left(\frac{T}{D_s}\right) \quad (15)$$

645 *For $T = 30$ and $D_s = 16$, this yields $\sim 2\times$ theoretical speedup; empirically we observe $2.3\times$ speedup*
 646 *due to memory efficiency gains.*

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

B.4 PROTOTYPE LEARNING DYNAMICS

Definition B.14 (Prototype Configuration). A prototype bank $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_M]^T \in \mathbb{R}^{M \times D}$ defines a configuration on the unit sphere \mathbb{S}^{D-1} when prototypes are ℓ_2 -normalized.

Theorem B.15 (Diversity Loss Gradient Flow). Let $\mathbf{P} \in \mathbb{R}^{M \times D}$ with $\|\mathbf{p}_i\| = 1$ for all i . The diversity loss:

$$\mathcal{L}_{div}(\mathbf{P}) = \frac{1}{M(M-1)} \sum_{i \neq j} \langle \mathbf{p}_i, \mathbf{p}_j \rangle^2 \quad (16)$$

has gradient:

$$\nabla_{\mathbf{p}_k} \mathcal{L}_{div} = \frac{4}{M(M-1)} \sum_{j \neq k} \langle \mathbf{p}_k, \mathbf{p}_j \rangle (\mathbf{p}_j - \langle \mathbf{p}_k, \mathbf{p}_j \rangle \mathbf{p}_k) \quad (17)$$

The term $(\mathbf{p}_j - \langle \mathbf{p}_k, \mathbf{p}_j \rangle \mathbf{p}_k)$ is the component of \mathbf{p}_j orthogonal to \mathbf{p}_k , pushing prototypes apart on the sphere.

Proof. For unit vectors, $\cos(\mathbf{p}_i, \mathbf{p}_j) = \langle \mathbf{p}_i, \mathbf{p}_j \rangle$. Differentiating:

$$\frac{\partial}{\partial \mathbf{p}_k} \langle \mathbf{p}_k, \mathbf{p}_j \rangle^2 = 2 \langle \mathbf{p}_k, \mathbf{p}_j \rangle \mathbf{p}_j \quad (18)$$

To maintain unit norm, we project onto the tangent space of \mathbb{S}^{D-1} at \mathbf{p}_k :

$$\text{proj}_{T_{\mathbf{p}_k} \mathbb{S}^{D-1}}(\mathbf{v}) = \mathbf{v} - \langle \mathbf{v}, \mathbf{p}_k \rangle \mathbf{p}_k \quad (19)$$

Applying this projection yields the stated gradient. \square

Proposition B.16 (Optimal Prototype Configuration). For M prototypes in \mathbb{R}^D with $M \leq D + 1$, the global minimum of \mathcal{L}_{div} is achieved when prototypes form a regular simplex inscribed in \mathbb{S}^{D-1} , with pairwise inner products:

$$\langle \mathbf{p}_i, \mathbf{p}_j \rangle = -\frac{1}{M-1} \quad \forall i \neq j \quad (20)$$

yielding $\mathcal{L}_{div}^* = \frac{1}{(M-1)^2}$.

Proof. For unit vectors, $\sum_{j=1}^M \mathbf{p}_j = \mathbf{0}$ at the centroid. Taking inner product with \mathbf{p}_i :

$$1 + \sum_{j \neq i} \langle \mathbf{p}_i, \mathbf{p}_j \rangle = 0 \implies \sum_{j \neq i} \langle \mathbf{p}_i, \mathbf{p}_j \rangle = -1 \quad (21)$$

By symmetry of the regular simplex, all off-diagonal inner products are equal: $\langle \mathbf{p}_i, \mathbf{p}_j \rangle = -1/(M-1)$. \square

Remark B.17 (Prototype Counts and Linguistic Inventories). Our prototype counts $(N_{\text{hand}}, N_{\text{loc}}, N_{\text{mov}}, N_{\text{ori}}) = (30, 15, 10, 8)$ are informed by linguistic estimates: ASL has ~ 30 handshapes (Battison, 1978), ~ 12 – 15 major locations, ~ 10 – 15 core movement types, and ~ 6 – 8 orientations. These counts closely match the phonological inventories, enabling interpretable component representations.

C EXTENDED METHODS

C.1 AGAN ARCHITECTURE DETAILS

C.2 ANATOMICAL GRAPH CONNECTIVITY

The hand skeleton defines natural connectivity:

- **Finger chains:** Wrist \rightarrow MCP \rightarrow PIP \rightarrow DIP \rightarrow tip for each finger
- **Palm connections:** All MCP joints connect to wrist
- **Functional groups:** Index-middle and ring-pinky pairs share additional edges

Algorithm 1 Anatomical Graph Attention Forward Pass

Require: Landmarks $\mathbf{X} \in \mathbb{R}^{B \times T \times N \times C}$, adjacency \mathbf{A}

Ensure: Spatial features $\mathbf{Z} \in \mathbb{R}^{B \times T \times D}$

- 1: $\mathbf{H}^{(0)} \leftarrow \text{Linear}(\mathbf{X})$
 - 2: **for** $l = 1$ to L **do**
 - 3: Compute attention: $\alpha_{ij} \leftarrow \text{softmax}(\text{LeakyReLU}(\mathbf{a}^T[\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]))$
 - 4: Mask by anatomy: $\alpha \leftarrow \alpha \odot \mathbf{A}$
 - 5: Aggregate: $\mathbf{h}_i^{(l)} \leftarrow \parallel_{k=1}^K \sum_j \alpha_{ij}^{(k)} \mathbf{W}^{(k)} \mathbf{h}_j^{(l-1)}$
 - 6: **end for**
 - 7: $\mathbf{Z} \leftarrow \text{MeanPool}(\mathbf{H}^{(L)}, \text{dim} = \text{nodes})$
 - 8: **return** \mathbf{Z}
-

C.3 BISSM PARAMETERIZATION DETAILS

The discrete selective SSM maintains state $\mathbf{x}_t \in \mathbb{R}^{D_s}$:

$$\mathbf{x}_t = \bar{\mathbf{A}}\mathbf{x}_{t-1} + \bar{\mathbf{B}}_t\mathbf{f}_t \quad (22)$$

$$\mathbf{y}_t = \mathbf{C}_t\mathbf{x}_t \quad (23)$$

where $\bar{\mathbf{A}} = \exp(\Delta_t\mathbf{A})$ and $\bar{\mathbf{B}}_t = (\Delta_t\mathbf{A})^{-1}(\bar{\mathbf{A}} - \mathbf{I})\Delta_t\mathbf{B}_t$.

Input-dependent parameters. The state matrix $\mathbf{A} \in \mathbb{R}^{D_s \times D_s}$ is diagonal with learnable entries initialized via HiPPO (Gu et al., 2020):

$$\mathbf{B}_t = \mathbf{W}_B\mathbf{f}_t \in \mathbb{R}^{D_s}, \quad \mathbf{W}_B \in \mathbb{R}^{D_s \times D} \quad (24)$$

$$\mathbf{C}_t = \mathbf{W}_C\mathbf{f}_t \in \mathbb{R}^{D_s}, \quad \mathbf{W}_C \in \mathbb{R}^{D_s \times D} \quad (25)$$

$$\Delta_t = \text{softplus}(\mathbf{w}_\Delta^T\mathbf{f}_t + b_\Delta) \in \mathbb{R}_{>0} \quad (26)$$

The selective mechanism allows the model to adaptively control information flow: larger Δ_t values cause the state to update more aggressively, while smaller values preserve existing state. This is particularly useful for sign language where movement speed varies significantly.

C.4 FULL FORWARD PASS ALGORITHM

C.5 SSM DISCRETIZATION DETAILS

The continuous-time SSM is defined by:

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) \quad (27)$$

For discrete inputs sampled at intervals Δ , the zero-order hold (ZOH) discretization yields:

$$\bar{\mathbf{A}} = \exp(\Delta\mathbf{A}) \quad (28)$$

$$\bar{\mathbf{B}} = (\Delta\mathbf{A})^{-1}(\exp(\Delta\mathbf{A}) - \mathbf{I}) \cdot \Delta\mathbf{B} \quad (29)$$

For diagonal \mathbf{A} (as in our implementation), this simplifies to element-wise operations, enabling efficient parallel computation.

C.6 IMPLEMENTATION DETAILS

D DATASET DETAILS

WLASL (Li et al., 2020) contains videos of isolated ASL signs performed by over 100 signers. We use the official train/test splits. Pose extraction uses MediaPipe Holistic, providing 33 pose landmarks, 21 left hand landmarks, and 21 right hand landmarks (75 total \times 3 coords = 225 features).

756 **Algorithm 2** PhonSSM Complete Forward Pass

757 **Require:** Input landmarks $\mathbf{X} \in \mathbb{R}^{B \times T \times N \times C}$

758 **Ensure:** Logits $\hat{\mathbf{y}} \in \mathbb{R}^{B \times K}$, components $\{\mathbf{c}^{(i)}\}$

759 1: **// Stage 1: Anatomical Graph Attention**

760 2: **for** $t = 1$ to T **do**

761 3: $\mathbf{Z}_t \leftarrow \text{AGAN}(\mathbf{X}_t, \mathbf{A})$ {Alg. 1}

762 4: **end for**

763 5: **// Stage 2: Phonological Factorization**

764 6: **for** $i \in \{\text{hand, loc, mov, ori}\}$ **do**

765 7: $\mathbf{C}^{(i)} \leftarrow \text{MLP}_i(\mathbf{Z})$ $\{\mathbf{C}^{(i)} \in \mathbb{R}^{B \times T \times D_c}\}$

766 8: **end for**

767 9: $\tilde{\mathbf{C}}^{(\text{mov})} \leftarrow \mathbf{C}^{(\text{mov})} + \text{Conv1D}(\mathbf{C}^{(\text{mov})})$

768 10: $\mathbf{F} \leftarrow \mathbf{W}_{\text{fuse}}[\mathbf{C}^{(\text{hand})} \parallel \mathbf{C}^{(\text{loc})} \parallel \tilde{\mathbf{C}}^{(\text{mov})} \parallel \mathbf{C}^{(\text{ori})}]$

769 11: **// Stage 3: Bidirectional SSM**

770 12: $\mathbf{G}_{\rightarrow} \leftarrow \text{SSM}_{\rightarrow}(\mathbf{F})$ {Forward pass}

771 13: $\mathbf{G}_{\leftarrow} \leftarrow \text{SSM}_{\leftarrow}(\text{flip}(\mathbf{F}))$ {Backward pass}

772 14: $\mathbf{G} \leftarrow \mathbf{W}_{\text{out}}[\mathbf{G}_{\rightarrow} \parallel \text{flip}(\mathbf{G}_{\leftarrow})]$

773 15: **// Stage 4: Hierarchical Prototypical Classification**

774 16: $\bar{\mathbf{c}}^{(i)} \leftarrow \frac{1}{T} \sum_t \mathbf{C}_t^{(i)}$ for each component i

775 17: $\mathbf{s}^{(i)} \leftarrow \text{softmax}(\bar{\mathbf{c}}^{(i)} (\mathbf{P}^{(i)})^T / \|\cdot\|)$ {Component similarities}

776 18: $\bar{\mathbf{g}} \leftarrow \frac{1}{T} \sum_t \mathbf{G}_t$

777 19: $\mathbf{e} \leftarrow \mathbf{W}_e[\mathbf{s}^{(\text{hand})} \parallel \mathbf{s}^{(\text{loc})} \parallel \mathbf{s}^{(\text{mov})} \parallel \mathbf{s}^{(\text{ori})} \parallel \bar{\mathbf{g}}]$

778 20: $\hat{\mathbf{y}} \leftarrow \frac{1}{\tau} \cos(\mathbf{e}, \mathbf{P}_{\text{sign}})$

779 21: **return** $\hat{\mathbf{y}}, \{\bar{\mathbf{c}}^{(i)}\}$

781 Table 4: Full hyperparameter configuration.

Hyperparameter	Value
<i>Architecture</i>	
Model dimension D	128
Component dimension D_c	32
GAT heads	4
GAT layers	3
SSM layers	4
SSM state dimension	16
SSM expansion factor	2
Dropout	0.1
<i>Training</i>	
Optimizer	AdamW
Learning rate	3×10^{-4}
Weight decay	10^{-2}
Batch size	128
Epochs	100
Warmup epochs	10
LR schedule	Cosine decay
Label smoothing	0.1
<i>Loss weights</i>	
λ_{ortho}	0.1
λ_{div}	0.01

804

805 **Important:** For WLASL evaluation, we train *four separate PHONSSM models from scratch*—one

806 for each vocabulary split (100, 300, 1000, 2000)—using only that split’s training data. These are

807 completely independent from the Merged-5565 model.

808

809 **Merged-5565** is a new large-scale dataset we construct by merging six publicly available ASL

sources, detailed in Table 5.

Table 5: Composition of the Merged-5565 dataset.

Source Dataset	Samples	Signs	Type
ASL Citizen (Desai et al., 2024)	83,399	2,731	Isolated
WLASL (Li et al., 2020)	21,083	2,000	Isolated
MVP/Kaggle ASL-Signs ^a	94,477	250	Isolated
ASL Alphabet ^b	27,455	29	Fingerspell
ASL MNIST ^b	27,455	26	Fingerspell
ChicagoFSWild (Shi et al., 2019a)	5,846	26	Fingerspell
Total (deduplicated)	259,715	5,565	—

^aGoogle Kaggle competition dataset. ^bKaggle community datasets. Note: Sign counts before deduplication. WLASL/MVP share ~200 signs with ASL Citizen; fingerspelling datasets share 26 letters.

We create a unified label map by alphabetically sorting all unique signs, remap dataset-specific indices, and merge with stratified train/val/test splits (196,606/31,551/31,558 samples, approximately 76/12/12%). **Important:** For Merged-5565 evaluation, we train a *single separate* PHONSSM model using dominant-hand input only (21 landmarks \times 3 coords = 63 features), selected by motion magnitude. This model is completely independent from the four WLASL models.

Preprocessing. All sequences are:

1. Centered by subtracting wrist position
2. Normalized to unit scale based on palm size
3. Resampled to 30 frames using linear interpolation
4. Augmented during training with random temporal shifts (± 3 frames) and scale jitter ($\pm 10\%$)

E ADDITIONAL RESULTS

E.1 PER-CLASS ANALYSIS

Table 6: Performance breakdown by phonological characteristics on WLASL100. Δ : improvement over Bi-LSTM.

Category	# Signs	Bi-LSTM	PHONSSM	Δ
One-handed signs	62	71.2	89.4	+18.2
Two-handed signs	38	68.9	86.8	+17.9
Static (no movement)	15	73.4	91.2	+17.8
Dynamic (with movement)	85	69.8	87.9	+18.1
Face-region location	28	65.2	85.6	+20.4
Body-region location	31	71.8	89.1	+17.3
Neutral space	41	72.4	90.2	+17.8

E.2 CONFUSION ANALYSIS

Common confusions occur between phonologically similar signs:

- **Location minimal pairs:** “mother”/“father” (chin/forehead) – 8% confusion rate
- **Handshape minimal pairs:** “water”/“want” (W/claw) – 6% confusion rate
- **Movement minimal pairs:** “chair”/“sit” (double/single) – 5% confusion rate

These errors are linguistically meaningful, suggesting the model has learned relevant phonological contrasts even when making mistakes.

E.3 STATISTICAL SIGNIFICANCE ANALYSIS

We conduct rigorous statistical testing to validate our results.

Table 7: Statistical significance tests comparing PHONSSM to baselines (paired t-tests, 3 seeds).

Comparison	Δ Acc.	t -statistic	p -value
<i>WLASL100</i>			
vs. DSTA-SLR	+4.81	8.7	<0.01
vs. Pose-TGCN	+14.18	28.4	<0.001
vs. Bi-LSTM	+18.21	35.2	<0.001
<i>WLASL2000</i>			
vs. DSTA-SLR	+18.38	22.1	<0.001
vs. I3D	+39.60	41.8	<0.001
<i>Merged-5565</i>			
vs. Bi-LSTM	+25.95	52.3	<0.001

All improvements are statistically significant. Note: the large t -statistics reflect both substantial effect sizes (e.g., +25.95pp for Merged-5565) and low variance across seeds (std <0.5pp), yielding Cohen’s $d > 50$ for some comparisons. We also compute 95% confidence intervals: WLASL100 accuracy is $88.37 \pm 0.82\%$, WLASL2000 is $72.08 \pm 1.27\%$, and Merged-5565 is $53.34 \pm 0.74\%$.

E.4 ERROR STRATIFICATION BY PHONOLOGICAL DISTANCE

We analyze errors as a function of phonological similarity between ground-truth and predicted signs.

Table 8: Error analysis by phonological distance (number of differing components).

Components Shared	% of Errors	Bi-LSTM	PHONSSM
4 (identical)	—	—	—
3 (minimal pair)	31.2%	28.4%	38.7%
2	40.8%	35.1%	33.2%
1	19.7%	22.8%	18.4%
0 (unrelated)	8.3%	13.7%	9.7%

PHONSSM concentrates errors on minimal pairs (38.7% vs 28.4% for Bi-LSTM), indicating the model has learned to distinguish coarse phonological categories but struggles with fine-grained contrasts—a linguistically sensible error pattern.

E.5 PROTOTYPE VISUALIZATION

The learned component prototypes correspond to interpretable phonological categories. Handshape prototypes cluster by finger configuration (fist, open, pointing). Location prototypes organize spatially (face, torso, neutral space). Movement prototypes distinguish trajectory types (linear, arc, repeated).

F MID-VOCABULARY PERFORMANCE ANALYSIS

We investigate why PHONSSM underperforms DSTA-SLR on WLASL300 (−5.6pp) and WLASL1000 (−4.9pp) while substantially outperforming on WLASL100 (+4.8pp) and WLASL2000 (+18.4pp).

F.1 MINIMAL PAIR DENSITY ANALYSIS

We computed the *phonological similarity density* for each WLASL split by annotating signs with ASL-LEX phonological features and measuring the fraction of sign pairs sharing 3+ of 4 components (“near-minimal pairs”).

Table 9: Phonological similarity density across WLASL vocabulary splits.

Split	Signs	Near-Minimal Pairs	Density	PHONSSM Δ
WLASL100	100	891	18.0%	+4.8pp
WLASL300	300	15,283	34.1%	-5.6pp
WLASL1000	1,000	142,450	28.5%	-4.9pp
WLASL2000	2,000	279,314	14.0%	+18.4pp

Finding: WLASL300 has the highest minimal pair density (34.1%), followed by WLASL1000 (28.5%). This occurs because WLASL vocabulary expansion prioritizes related concepts (e.g., adding “GRANDMOTHER” after “MOTHER”, “FATHER”), which tend to be phonologically similar.

F.2 METHOD COMPARISON BY MINIMAL PAIR PERFORMANCE

We stratified test accuracy by phonological similarity to nearest training sign.

Table 10: Accuracy (%) stratified by phonological distance to nearest training neighbor.

Components Shared	DSTA-SLR	PHONSSM	Δ
<i>WLASL300</i>			
0-1 (distinct)	72.4	78.9	+6.5
2 (moderate)	81.2	76.3	-4.9
3-4 (minimal pair)	84.1	71.8	-12.3
<i>WLASL2000</i>			
0-1 (distinct)	48.2	74.6	+26.4
2 (moderate)	55.8	71.2	+15.4
3-4 (minimal pair)	61.3	68.4	+7.1

Interpretation: DSTA-SLR excels at distinguishing minimal pairs via fine-grained spatiotemporal attention, while PHONSSM excels at compositional generalization to phonologically distinct signs. At large vocabularies (WLASL2000), most test signs are phonologically distinct from training signs, favoring PHONSSM. At mid-range vocabularies (WLASL300), the dense minimal pair structure favors DSTA-SLR’s discriminative attention.

F.3 IMPLICATIONS

This analysis suggests PHONSSM and DSTA-SLR capture complementary information. Future work could explore ensemble methods or incorporating DSTA-SLR’s spatiotemporal attention within the phonological framework.

G COMPONENT-LEVEL VALIDATION

We provide detailed evidence that PDM branches capture their intended phonological categories.

G.1 PHONOLOGICAL ANNOTATION

We annotated 500 WLASL100 test samples with ground-truth phonological labels using ASL-LEX 2.0 (Caselli et al., 2017):

- **Handshape:** 30 categories (e.g., “1” (index), “5” (spread), “A” (fist), “B” (flat), “C” (curved))
- **Location:** 12 categories (e.g., forehead, chin, chest, neutral space, ipsilateral)
- **Movement:** 15 categories (e.g., none, arc, circle, linear, zigzag, repeated)
- **Orientation:** 8 categories (e.g., palm-up, palm-down, palm-in, palm-out)

Annotation quality. Two annotators independently labeled all samples; inter-annotator agreement was 94.2% for handshape, 91.8% for location, 87.4% for movement, and 89.6% for orientation

(Cohen’s $\kappa > 0.85$ for all). Disagreements were resolved by consensus. For signs with phonological variation (e.g., DEAF produced chin-to-ear or ear-to-chin), we annotated the variant observed in each video rather than a canonical form. ASL-LEX covered 96/100 WLASL100 signs; the remaining 4 were annotated using the same feature scheme by the annotators.

G.2 LINEAR PROBE METHODOLOGY

For each PDM branch, we:

1. Extract the time-averaged component embedding $\bar{c}^{(i)} \in \mathbb{R}^{32}$
2. Train a linear classifier (logistic regression) to predict each phonological category
3. Report accuracy on held-out samples (5-fold cross-validation)

G.3 FULL RESULTS

Table 11: Complete linear probe results. Each row shows one PDM branch; each column shows one phonological prediction task. Diagonal entries (bold) indicate intended correspondences.

PDM Branch	Prediction Target			
	Handshape	Location	Movement	Orientation
Handshape	78.4	31.2	24.8	38.1
Location	29.6	71.8	22.4	35.7
Movement	26.3	28.9	68.2	31.4
Orientation	33.8	32.1	25.6	74.6
Random baseline	3.3	8.3	6.7	12.5
Full embedding	81.2	74.6	71.8	78.3

Key observations:

1. **Specialization:** Each branch achieves highest accuracy on its intended category (diagonal), confirming semantic correspondence.
2. **Above-chance cross-prediction:** Off-diagonal entries exceed chance, indicating some phonological information leaks across branches. This is expected since components are correlated in natural signs (e.g., certain handshapes occur more often at certain locations).
3. **Factorization benefit:** The gap between diagonal and off-diagonal (e.g., 78.4 vs 31.2 for handshape) demonstrates effective factorization.
4. **Factorization-accuracy trade-off:** The “Full embedding” row shows slightly higher accuracy than individual branches (e.g., 81.2 vs 78.4 for handshape), indicating ~ 3 pp is sacrificed for factorization. We experimented with relaxing λ_{ortho} from 0.1 to 0.05: component probe accuracy improved by ~ 2 pp but sign-level accuracy dropped by 1.5pp due to increased redundancy. The current setting balances interpretability and accuracy.

G.4 PROTOTYPE INTERPRETABILITY

We visualized which learned prototypes correspond to which linguistic categories by computing the mean activation of each prototype for samples with known phonological labels.

Table 12: Top-3 handshape prototypes activated by each major handshape category.

Linguistic Category	Proto #1	Proto #2	Proto #3
“1” (index point)	P7 (0.89)	P12 (0.42)	P3 (0.18)
“5” (spread hand)	P2 (0.91)	P15 (0.38)	P8 (0.21)
“A” (fist)	P19 (0.87)	P4 (0.45)	P11 (0.19)
“B” (flat hand)	P2 (0.72)	P8 (0.51)	P15 (0.32)
“C” (curved)	P23 (0.83)	P7 (0.28)	P12 (0.24)

Distinct prototypes dominate for distinct handshapes (P7 for “1”, P19 for “A”, P23 for “C”). Some prototypes are shared across similar handshapes (P2 activates for both “S” and “B”, which share an extended-finger configuration).

G.5 INTERVENTION EXPERIMENT

To verify causal relationship, we performed an intervention: replacing one component embedding with that of a different sign while keeping others fixed.

Setup: We identified 47 minimal pairs in WLASL100 (sign pairs differing in exactly one phonological component). For each pair (e.g., MOTHER/FATHER differing only in location), we take a sample of sign A, swap only the differing component embedding with that from sign B, and measure whether the prediction changes to B.

Result ($n = 423$ interventions, **95% CI**): Swapping the differing component changes the prediction to the minimal pair **73.2% \pm 4.2%** of the time. Swapping a non-differing component (control condition) changes prediction only **12.4% \pm 3.1%** of the time. The difference is significant ($p < 0.001$, McNemar’s test). This confirms that component embeddings causally determine predictions in the expected manner.

Two-component swaps: When swapping two components simultaneously, prediction changes to a sign sharing those two swapped components 61.8% of the time, demonstrating compositional behavior.

H EXTENDED ABLATIONS

H.1 COMPONENT-WISE ABLATION

Table 13: Detailed ablation on WLASL100 (mean \pm std over 3 seeds).

Configuration	Top-1 (%)	Δ
Full PHONSSM	88.37 \pm 0.42	–
<i>Architecture ablations</i>		
w/o PDM (no factorization)	76.49 \pm 0.83	–11.9
w/o AGAN (MLP encoder)	79.84 \pm 0.71	–8.5
w/o BiSSM (LSTM temporal)	82.17 \pm 0.65	–6.2
w/o HPC (linear classifier)	84.11 \pm 0.58	–4.3
<i>Loss ablations</i>		
w/o $\mathcal{L}_{\text{ortho}}$	85.92 \pm 0.55	–2.5
w/o \mathcal{L}_{div}	86.84 \pm 0.48	–1.5
w/o both auxiliary losses	83.21 \pm 0.72	–5.2
<i>Input ablations</i>		
Hands only (no pose)	85.63 \pm 0.61	–2.7
Dominant hand only	81.42 \pm 0.78	–7.0
2D coordinates only	84.29 \pm 0.69	–4.1

H.2 ARCHITECTURE VARIANTS

I HYPERPARAMETER SENSITIVITY

I.1 LOSS WEIGHT SENSITIVITY

The orthogonality loss $\lambda_{\text{ortho}} = 0.1$ provides optimal factorization. Values too low (< 0.05) allow component redundancy; values too high (> 0.2) over-constrain representations. The diversity loss is less sensitive, with $\lambda_{\text{div}} \in [0.005, 0.01]$ performing similarly.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Table 14: Architecture scaling on WLASL100.

Configuration	Params	Top-1 (%)	Throughput
<i>Model dimension D</i>			
D = 64	0.9M	85.21	412/s
D = 128 (default)	3.2M	88.37	260/s
D = 256	11.8M	89.02	98/s
<i>Number of BiSSM layers</i>			
2 layers	2.1M	86.45	318/s
4 layers (default)	3.2M	88.37	260/s
6 layers	4.3M	88.72	205/s
<i>Number of GAT heads</i>			
2 heads	2.9M	87.18	275/s
4 heads (default)	3.2M	88.37	260/s
8 heads	3.8M	88.54	241/s

Table 15: Sensitivity to loss weights on WLASL100.

λ_{ortho}	λ_{div}	Top-1 (%)
0.01	0.01	86.42
0.05	0.01	87.63
0.1	0.01	88.37
0.2	0.01	87.91
0.5	0.01	86.18
0.1	0.001	87.82
0.1	0.005	88.15
0.1	0.01	88.37
0.1	0.05	87.54
0.1	0.1	85.93

I.2 LEARNING RATE SENSITIVITY

I.3 PROTOTYPE COUNT SENSITIVITY

Performance saturates around $(N_h, N_l, N_m, N_o) = (30, 15, 10, 8)$. These counts closely match linguistic estimates of ASL phonological inventories (Battison, 1978): ~ 30 handshapes, $\sim 12-15$ locations, ~ 10 core movements, and ~ 8 orientations. Larger counts provide marginal gains (+0.1-0.15pp) but increase parameters without improving interpretability.

J TRAINING DYNAMICS

J.1 CONVERGENCE ANALYSIS

We analyze the training dynamics of PHONSSM and its components.

The orthogonality loss decreases steadily, indicating progressive factorization of the component subspaces. The gap between training and validation accuracy remains small (~ 6 pp at convergence), suggesting good generalization.

J.2 COMPONENT LEARNING DYNAMICS

Individual phonological components converge at different rates:

- **Handshape**: Fastest convergence (stabilizes by epoch 40); handshape is the most visually distinctive component with clear finger configurations.
- **Location**: Moderate convergence (epoch 60); requires learning spatial relationships relative to body landmarks.

Table 16: Learning rate sensitivity on WLASL100.

Learning Rate	Top-1 (%)	Convergence
1×10^{-4}	86.82	95 epochs
2×10^{-4}	87.91	78 epochs
3×10^{-4}	88.37	65 epochs
5×10^{-4}	87.63	52 epochs
1×10^{-3}	85.41	41 epochs

Table 17: Effect of prototype counts on WLASL100.

N_h	N_l	N_m	N_o	Top-1 (%)
15	8	5	4	86.12
20	10	8	6	87.45
30	15	10	8	88.37
40	20	15	10	88.52
50	25	20	12	88.48

- **Movement:** Slowest convergence (epoch 80); movement patterns require temporal integration across multiple frames.
- **Orientation:** Fast convergence (epoch 45); palm orientation has relatively few categories (8 prototypes).

J.3 LOSS LANDSCAPE ANALYSIS

We analyze the loss landscape by computing the Hessian eigenspectrum at convergence. The top eigenvalues are: $\lambda_1 = 12.4$, $\lambda_2 = 8.7$, $\lambda_3 = 5.2$, with rapid decay thereafter. The condition number $\kappa = \lambda_{\max}/\lambda_{\min} \approx 10^3$ indicates a well-conditioned optimization landscape, explaining the stable training dynamics.

K LIMITATIONS

We discuss limitations of PHONSSM in detail to guide future research.

K.1 METHODOLOGICAL LIMITATIONS

Isolated sign assumption. PHONSSM processes signs in isolation, assuming clean segmentation. Continuous signing involves co-articulation effects where adjacent signs influence each other’s production, sign boundaries are ambiguous, and prosodic structure spans multiple signs. Extending to continuous recognition requires explicit segmentation or sequence-to-sequence modeling.

Fixed phonological structure. We adopt the classical Stokoe-Battison four-parameter model (hand-shape, location, movement, orientation). Alternative linguistic analyses propose:

- Autosegmental phonology with separate tiers for manual and non-manual features
- Prosodic structure including syllables and metrical feet
- Feature geometry with hierarchical organization of sub-components

Learned decompositions might discover more effective factorizations than hand-specified parameters.

Static prototypes. Component prototypes are fixed after training. Dynamic or instance-adaptive prototypes could better handle signer variation and novel phonological realizations.

Table 18: Training convergence metrics on WLASL100.

Metric	Epoch 25	Epoch 50	Epoch 100
Train Loss	2.14	0.89	0.42
Val Loss	2.31	1.12	0.68
Train Acc (%)	52.3	78.6	94.2
Val Acc (%)	48.1	74.2	88.4
$\mathcal{L}_{\text{ortho}}$	0.42	0.18	0.08
\mathcal{L}_{div}	0.31	0.12	0.05

K.2 EVALUATION LIMITATIONS

ASL-centric evaluation. All experiments use American Sign Language. While phonological principles (simultaneity, minimal pairs, compositional structure) are cross-linguistic, specific inventories differ:

- British Sign Language (BSL) uses different handshape inventory
- Chinese Sign Language has location contrasts not present in ASL
- Some sign languages distinguish two-handed vs. one-handed more strictly

Prototype counts and architectural choices may require adaptation for other sign languages.

Dataset biases. Training data comes primarily from controlled recording settings with:

- Adult native/fluent signers (under-representing learners, children, elderly)
- Neutral backgrounds and good lighting
- Citation-form signs (isolated, careful production)

K.3 DEPLOYMENT LIMITATIONS

Pose estimation dependency. PHONSSM assumes high-quality pose landmarks from MediaPipe. Real-world degradation includes:

- Occlusion (self-occlusion, objects, other people)
- Motion blur during rapid movements
- Challenging lighting (backlighting, low light)
- Camera angles different from training distribution

Signer variation. Models may not generalize to:

- Signers with motor differences affecting articulation
- Regional/dialectal variation in sign production
- Non-native signers with L1 transfer effects

L FUTURE DIRECTIONS

Continuous sign language recognition. Extending PHONSSM to continuous signing requires: (1) implicit or explicit segmentation, (2) handling co-articulation, and (3) modeling sentence-level prosody. The phonological decomposition could inform CTC-style losses with component-level intermediate representations.

Cross-linguistic transfer. The phonological factorization may enable transfer learning across sign languages. Shared handshape or movement prototypes could bootstrap recognition for low-resource sign languages.

Multi-modal fusion. Combining skeleton input with facial landmarks (for non-manual markers) and RGB features (for fine-grained handshape) could address current limitations while preserving efficiency.

1242 **Learned phonological structure.** Replacing hand-specified component pathways with learned
1243 factorization (e.g., via neural architecture search or information-theoretic objectives) might discover
1244 more effective decompositions.

1245 **Real-time applications.** With 260 samples/second throughput, PHONSSM supports real-time
1246 deployment. Future work includes: mobile optimization, streaming recognition, and integration with
1247 sign language translation systems.
1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295