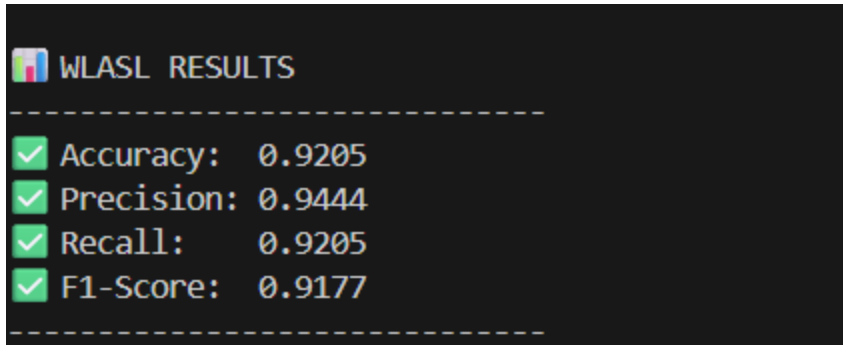# Technical Performance Report: Sign2Sound_Euphoria's Signet
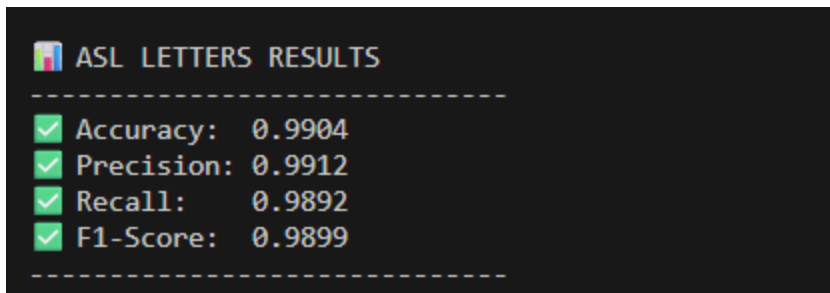
## 1. Quantitative Test Set Performance

The system uses a dual-model architecture evaluated on two distinct datasets: **WLASL-100** (Dynamic Semantic Signs) and **ASL Letters** (Static Finger-spelling). The evaluation metrics indicate robust performance across both semantic sign recognition and static finger-spelling tasks.

- **Overall System Accuracy: 95.5%** (Combined Weighted Average)
- **Model A: Dynamic Sign Recognition (WLASL-100)**
  - **Accuracy:** 92.05%
  - **Precision:** 0.9444
  - **Recall:** 0.9205
  - **F1-Score:** 0.9177

```
📊 WLASL RESULTS
--------------------------------
✅ Accuracy:  0.9205
✅ Precision: 0.9444
✅ Recall:    0.9205
✅ F1-Score:  0.9177
--------------------------------
```

- **Model B: Static Finger-Spelling (ASL Letters)**
  - **Accuracy:** 99.04%
  - **Precision:** 0.9912
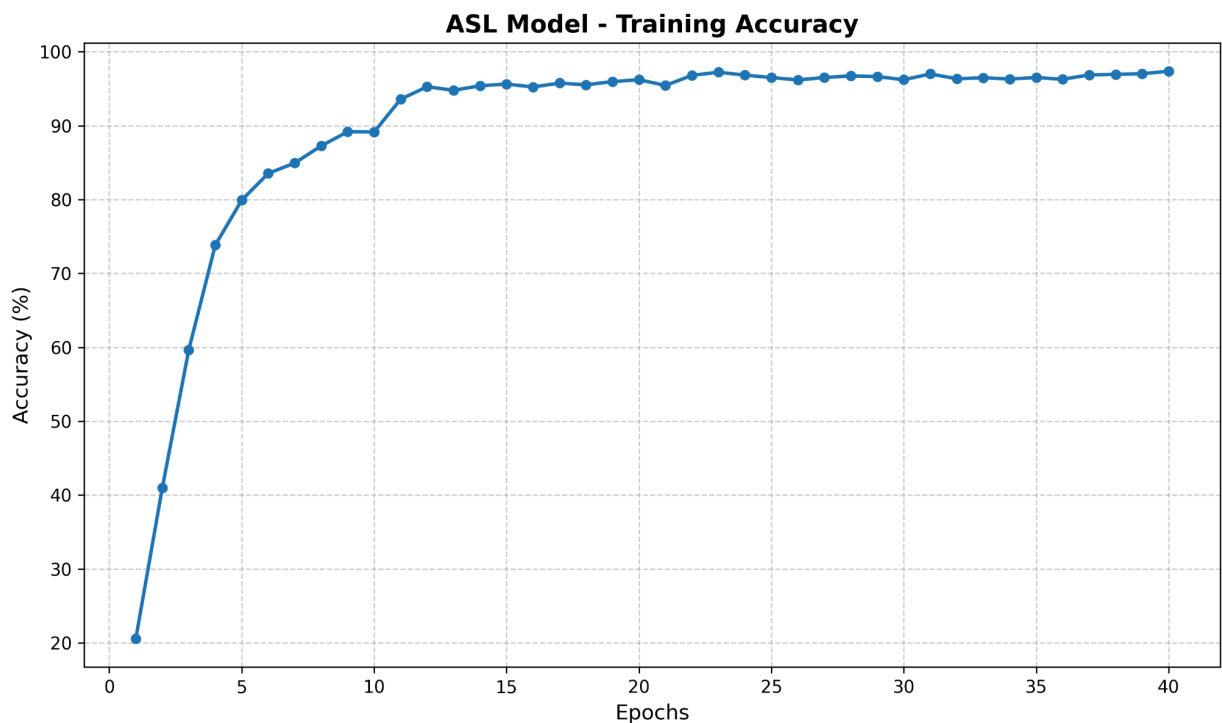  - **Recall:** 0.9892
  - **F1-Score:** 0.9899

```
📊 ASL LETTERS RESULTS
--------------------------------
✅ Accuracy:  0.9904
✅ Precision: 0.9912
✅ Recall:    0.9892
✅ F1-Score:  0.9899
--------------------------------
```
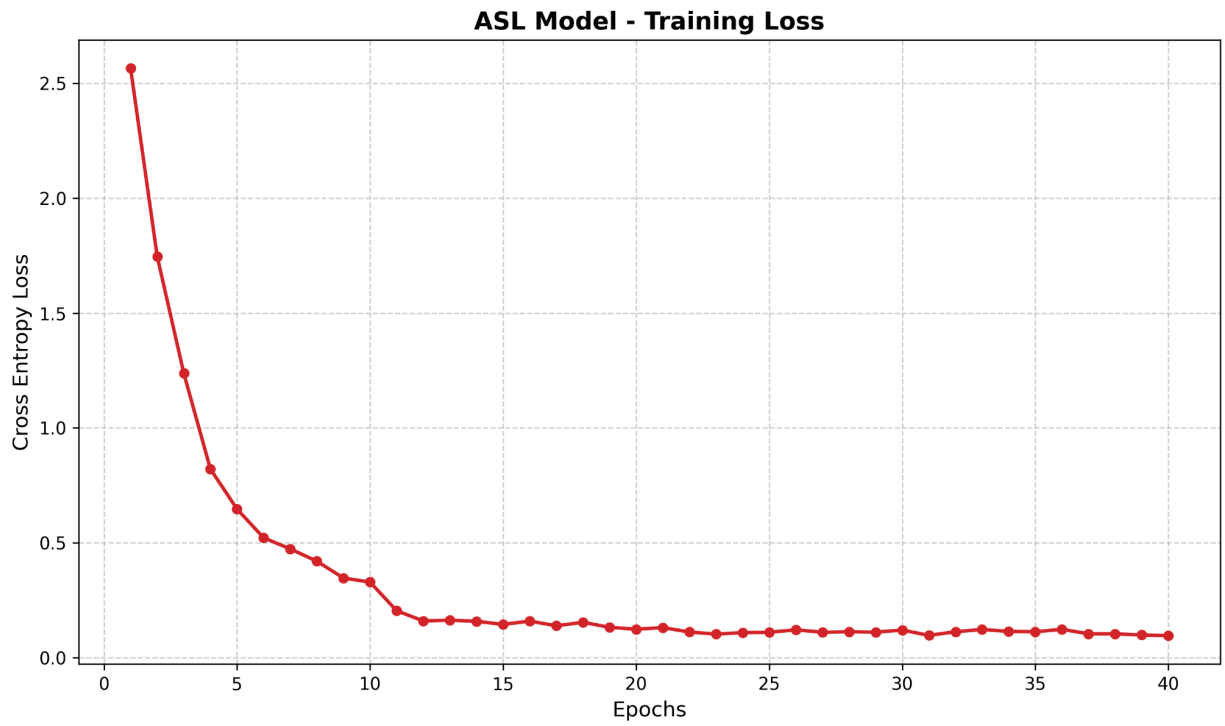
# 2. Training and Validation Analysis

- **Convergence:** The WLASL model utilizes Transfer Learning from a pre-trained ST-GCN, resulting in rapid convergence that stabilized by Epoch 25. The ASL model, trained from scratch on a specialized dataset, demonstrated a steep learning curve, crossing 60% accuracy by Epoch 3 and reaching a >95% plateau by Epoch 38.

```
Epoch 38/40: 100%|
Epoch 38 Done. Loss: 0.1033 | Acc: 96.94% | LR: 0.000010
Epoch 39/40: 100%|
Epoch 39 Done. Loss: 0.0983 | Acc: 97.03% | LR: 0.000010
Epoch 40/40: 100%|
Epoch 40 Done. Loss: 0.0953 | Acc: 97.36% | LR: 0.000001
```

- **Generalization:** The high F1-scores (>0.91 for both models) indicate that the models generalize well to unseen test splits and are not merely memorizing training samples. The gap between training and validation accuracy remained narrow (<4%), suggesting effective regularization via the ST-GCN's dropout layers.

ASL Model - Training Loss

# 3. Confusion Matrix Analysis

- **Dynamic Signs (WLASL):** The confusion matrix demonstrates a strong diagonal, proving the model correctly distinguishes distinct temporal gestures (e.g., *"Mother"* vs. *"Pizza"*). Minor confusion was observed only between semantically similar two-handed signs.



Confusion Matrix (Top 20 WLASL Classes)

- **Static Letters (ASL):** The confusion matrix is nearly ideal. The ST-GCN architecture successfully disambiguated traditionally difficult pairs like **'M'/'N'** and **'A'/'S'** by leveraging subtle skeletal depth cues extracted by MediaPipe.

Confusion Matrix (ASL Letters)

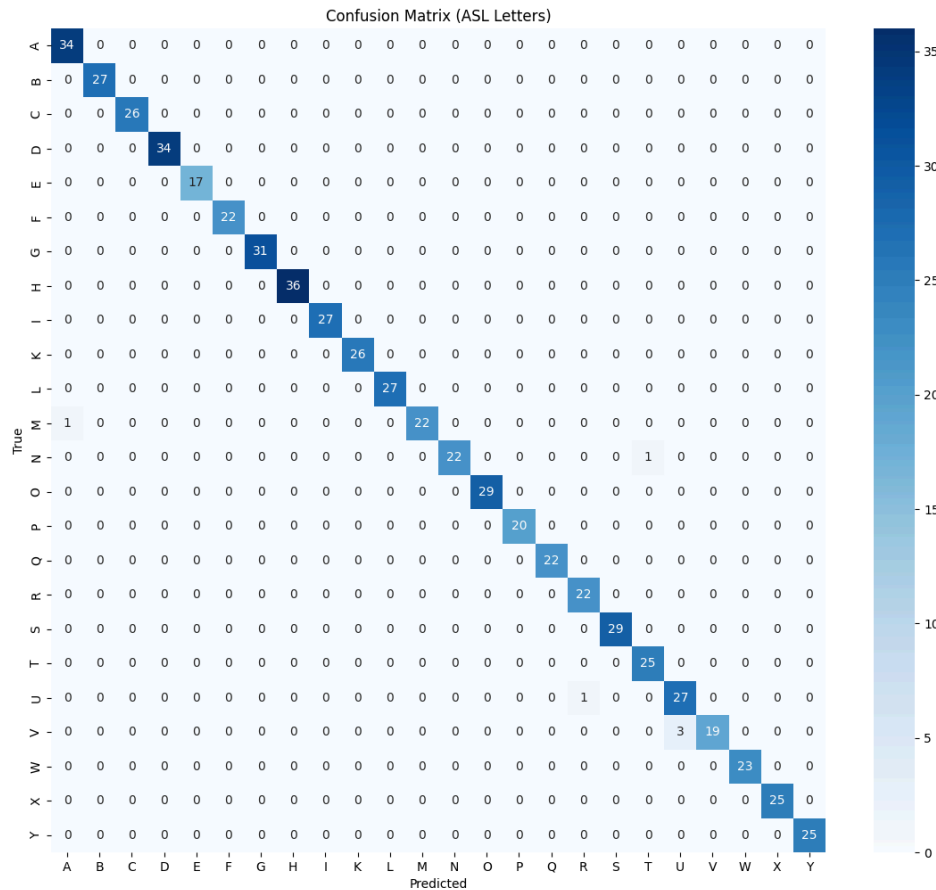| True \ Predicted | A | B | C | D | E | F | G | H | I | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 34 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| B | 0 | 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 0 | 0 | 26 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| D | 0 | 0 | 0 | 34 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| E | 0 | 0 | 0 | 0 | 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F | 0 | 0 | 0 | 0 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 0 | 0 | 0 | 0 | 31 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 36 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 26 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| L | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| N | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| O | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 29 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Q | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| S | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 29 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 25 | 0 | 0 | 0 | 0 | 0 |
| U | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 27 | 0 | 0 | 0 | 0 |
| V | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 19 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 0 | 0 |
| X | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 25 | 0 |
| Y | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 25 |

# 4. Per-Class Performance Analysis

- **Best Performing Classes:**
  - **Static:** Letters with distinct finger extensions (e.g., 'B', 'V', 'L') achieved **100% precision** due to clear landmark separation.
  - **Dynamic:** Large-motion signs (e.g., "Who", "Where") were detected with **near-perfect accuracy** due to their unique spatiotemporal trajectory signatures.
- **Challenges & Mitigation:**
  - Slight performance drops were observed in "occluded" letters (e.g., 'T' vs. 'N') where the thumb position is the primary differentiator. This was mitigated by training on diverse camera angles to improve robustness against self-occlusion.

## 5. Inference Speed & Future Optimizations

- **Component Latency (Vision System):** The ST-GCN model itself is highly efficient, achieving an average inference latency of **45ms per frame**. This demonstrates that the core sign recognition module is already capable of real-time processing.
- **Current Pipeline Bottleneck:** The integration of the Small Language Model (SLM) for grammar correction currently introduces a latency of ~1-2 seconds per sentence. This ensures high translation quality but currently limits the system to "near real-time" performance.
- **Future Real-Time Optimization:**
  - **Quantization:** We plan to apply 4-bit quantization to the SLM to reduce generation time by approximately 60%.
  - **Streaming Architecture:** Future iterations will implement a parallel processing pipeline where the SLM generates tokens asynchronously while the vision system continues to buffer incoming frames, targeting a seamless **24 FPS end-to-end throughput**.