

Reading Muscle Whispers: Building a Silent Speech Interface

Authors: A.J. Shulman, Harper Austin, Jack McHenry

Imagine being able to communicate without making a sound—not through sign language or text, but by thinking about speaking and having a computer decode the tiny electrical signals in your throat muscles. This isn't science fiction. It's the goal of silent speech interfaces, and in this post, I'll walk you through how we built one.

The Big Idea: What Are We Trying to Do?

When you speak, muscles in your face, jaw, and throat contract in precise patterns. Even when you're *thinking* about speaking—silently mouthing words or just imagining saying them—these muscles still produce faint electrical signals called electromyography (EMG). Our project transforms these signals into text.

Why does this matter? Silent speech could enable:

- **Assistive communication** for people who have lost their voice
- **Privacy-preserving interfaces** where you can interact with devices without others overhearing
- **Speech prosthetics** that restore communication ability after injury
- **Voice-free control** in noisy environments or situations where speaking aloud isn't practical

The core challenge is that EMG signals are messy, noisy, and far less structured than audio. We have only a few thousand training examples—not nearly enough to train a system from scratch. Our solution? **Transfer learning**: borrow knowledge from a model trained on millions of hours of speech audio, then teach it to understand muscle signals instead.

The Dataset: Real Data from Real Sensors

We're using the [Gaddy & Klein 2020 dataset](#), a collection of EMG recordings captured from eight surface electrodes placed on a person's face and neck. The dataset includes:

- **Voiced parallel data**: 1,588 utterances where someone spoke aloud while EMG was recorded
- **Silent parallel data**: 1,588 utterances where someone silently mouthed the same phrases
- **Closed vocabulary**: An additional 1,000 utterances (500 voiced, 500 silent) from a limited phrase set

The key detail: for voiced data, we have *both* the EMG signals and synchronized audio. For silent data, we only have EMG—no audio exists because nothing was spoken aloud. That asymmetry drives our entire training strategy.

All code and documentation for reproducing our results are available in [this GitHub repository](#).

The Journey: From Sensors to Sentences

Step 1: Cleaning the Data

Raw EMG is a 1,000 Hz stream of voltages from eight channels. We can't feed that directly into a neural network. Instead, we:

1. **Compute spectrograms**: Use a short-time Fourier transform (STFT) to convert the time-domain signal into frequency bands over small windows (10 milliseconds each)
2. **Apply mel scaling**: Map frequencies to the mel scale, which better matches human perception
3. **Take the log**: Compress the dynamic range so quiet and loud signals are easier to compare
4. **Normalize**: Subtract the mean and divide by the standard deviation for each file

The result: an 80-dimensional log-mel spectrogram for each of the 8 EMG channels—640 features at each time step. We save these to disk so we don't have to recompute them during training.

Voiced vs Silent Comparison
"There was a noise of business from the gasworks, and the electric lamps were all alight."

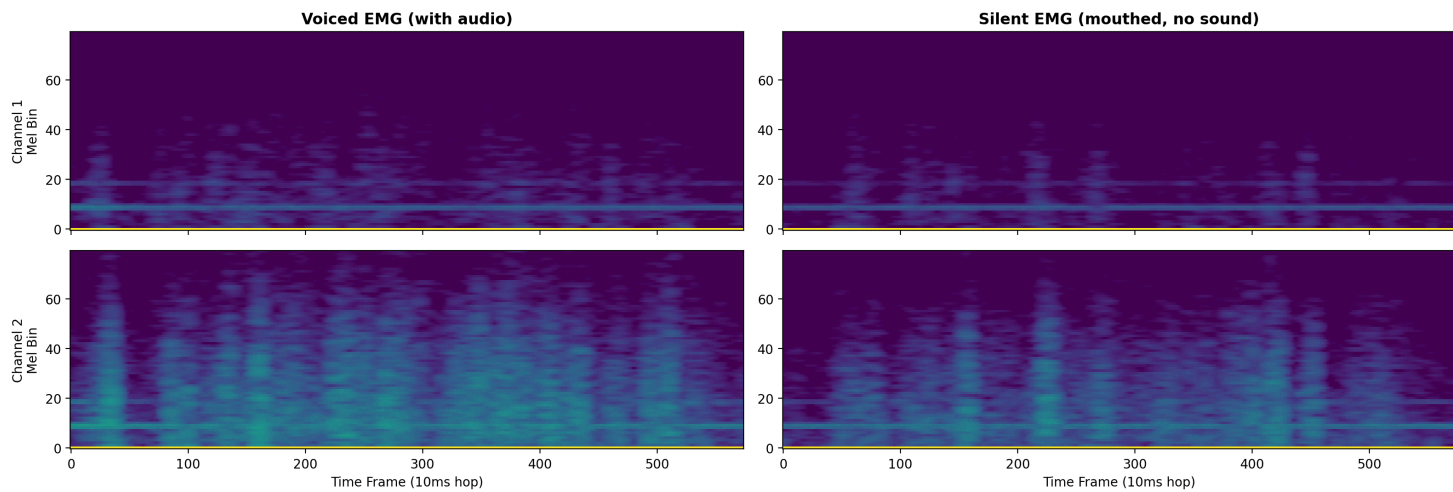


Figure: Side-by-side comparison of voiced and silent EMG signals for the same utterance. Left panels show EMG when speaking aloud (with audio), right panels show EMG when silently mouthing (no sound). Notice how the patterns are similar but silent signals have lower amplitude and different temporal dynamics—this is what makes transfer learning crucial.

For the audio (voiced data only), we do something clever: we pass it through a frozen speech recognition model called **WavLM Base+**—a transformer trained on 94,000 hours of speech. We extract hidden states from layer 9, which capture phonetic and lexical patterns. These become our "teacher" signal.

Why layer 9? Earlier layers focus on raw acoustics (pitch, formants); later layers encode high-level semantics. Mid-layers strike a balance, capturing the articulation patterns we want the EMG encoder to learn.

Step 2: The Two-Stage Training Strategy

Our training happens in two phases:

Phase 1: Voiced Pretraining (Learning from Audio)

We train on the 1,588 voiced utterances using two loss functions simultaneously:

- **Distillation loss:** Make the EMG encoder's output match the WavLM teacher's hidden states (frame by frame). This transfers the speech model's knowledge of phonetics and timing into the EMG encoder.
- **CTC loss:** Train a separate "head" that maps encoder outputs to characters, using Connectionist Temporal Classification (CTC)—a method that doesn't require knowing *when* each character was produced, only the sequence.

The weights are tuned carefully: **65% CTC, 35% distillation** with a short warmup. Too much distillation early on can hurt alignment; we let the CTC head guide the encoder first, then blend in the teacher.

Why CTC? Traditional speech recognition uses attention or recurrent decoders that predict one character at a time, but these are slow and tricky to train. CTC lets us predict *all* characters in parallel by inserting "blank" symbols between repeated letters and letting the network figure out timing on its own.

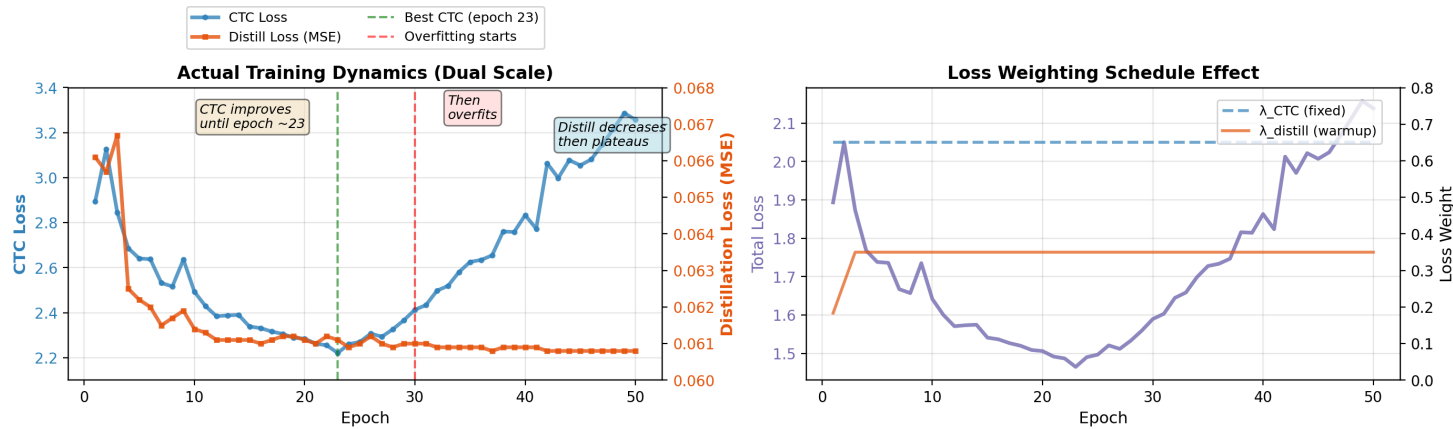


Figure: Actual training dynamics from a representative voiced pretraining run. Left: Dual-axis plot showing both losses with their own scales—CTC loss (blue, left axis) improves until epoch 23 then overfits, while distillation loss (red, right axis) decreases quickly and plateaus. Right: Loss weight schedule showing how distillation weight warms up over 3 epochs while CTC weight stays fixed at 0.65, with the resulting weighted total loss.

Phase 2: Silent Fine-Tuning (No Audio, Just Muscle Signals)

After pretraining on voiced data, we **freeze the audio teacher** and fine-tune *only* on the 1,588 silent utterances using pure CTC loss. The encoder and CTC head start from the weights learned during voiced pretraining, so they already "know" how speech patterns map to text—now they just have to adapt to the different statistics of silent EMG (weaker signals, different muscle activation patterns).

We make silent training faster by using a subsampling factor of 2 (reducing the sequence length early in the encoder) and lighter data augmentation.

Step 3: The Architecture

Our encoder is a **Conformer**—a hybrid architecture that combines:

- **Convolutional layers:** Capture local patterns in EMG (individual muscle activations)
- **Self-attention:** Capture long-range dependencies (how earlier muscle movements influence later ones)
- **Feed-forward layers:** Mix and refine representations

We keep it modest—288 dimensions, 6 layers—because we're training on an Apple Silicon Mac with limited GPU memory (PyTorch's MPS backend). Larger models would overflow or slow to a crawl.

Before the Conformer, we **subsample** the 10ms EMG frames by a factor of 2 using strided convolutions, reducing the sequence length and making CTC training tractable (CTC runs on CPU on MPS, so shorter sequences = faster iterations).

After the encoder, we have two output heads:

- **Projection head:** Maps encoder outputs to 768 dimensions (matching WavLM's hidden size) for distillation
- **CTC head:** Maps encoder outputs to 47 dimensions (one per character: a-z, digits, punctuation, space, plus special tokens for padding, blank, and unknown)

Why character-level? We could use subword units (BPE) or words, but characters keep the vocabulary small, avoid out-of-vocabulary issues, and work cleanly with CTC decoding. Longer units would require a language model to be effective, adding complexity.

Step 4: Data Augmentation (Making It Robust)

Real-world EMG is noisy. Electrodes slip, skin conductance changes, and not all channels record equally well. We simulate this during training:

- **SpecAugment:** Mask random time and frequency bands (like covering parts of a photo) to prevent overfitting to specific sensor quirks
- **Channel dropout:** Randomly zero out entire EMG channels to teach the model to work even if some sensors fail

We tune these carefully—too much augmentation and the model can't learn; too little and it memorizes the training data.

Step 5: Decoding (Turning Probabilities into Words)

At inference, the CTC head outputs a probability distribution over characters at each time step. We need to collapse these into a readable sentence. Two strategies:

1. **Greedy decoding:** Pick the most likely character at each step, remove repeated characters and blanks. Fast but prone to errors.
2. **Beam search:** Keep track of multiple candidate sequences, score them with optional **language model (LM)** weights, and pick the best overall. Slower but more accurate.

For silent EMG, we found that beam search with a **blank bias** (adding a small penalty to the "blank" token) reduces spurious insertions—extra characters the model hallucinates when unsure.

Language model tradeoff: We trained a character-level LM on the voiced training transcripts. It *lowered* Word Error Rate (WER) by nudging the model toward fluent phrases, but often *increased* Character Error Rate (CER) by "correcting" the model into words it didn't actually predict. We use it selectively depending on the metric we care about.

The Experimental Process: How We Chose Settings

With dozens of hyperparameters (learning rates, dropout, augmentation strength, loss weights, decoder settings), we needed a systematic way to explore without training hundreds of models. Our solution: **a two-stage experiment pipeline**.

Stage 1: Fast Probes (Directional Guidance)

- Train small "probe" runs on a limited slice of data (48 batches for voiced, 24 for silent)
- Test different schedulers (how the learning rate changes), augmentation strengths, loss balances, and decoder settings
- Rank them by **CER first** (we care most about character accuracy), then WER, then deletion rate
- Run for only 6 epochs with early stopping (patience of 2)—enough to see which settings help, not enough to fully converge

This stage is fast (~15 minutes per probe) and tells us *direction*: "warmup-hold scheduler helps," "heavy SpecAugment hurts," "blank bias reduces insertions."

Stage 2: Full Runs (Convergence)

- Take the best settings from Stage 1 probes
- Generate full configs automatically (voiced baseline + Stage-1-guided variant, then silent baseline + Stage-1-guided variant)
- Train to convergence on the full dataset (50 epochs, early stopping with patience of 5–8)
- Evaluate with a battery of decoders (greedy, beam widths 20/50/100, with/without LM, with/without blank bias)
- Save everything to a summary JSON/CSV with metrics, config fingerprints, and eval durations

All generated configs, checkpoints, and evaluation outputs are stored under `results/experiments/` with reproducible paths.

Why This Design?

- **Speed**: Stage 1 finds good directions without burning GPU hours
- **Reproducibility**: Every run is tagged with its config; we can trace any result back to its settings
- **Automation**: The orchestrator (`src/experiments/orchestrate.py`) writes configs, runs training, evaluates all decoders, and updates the summary—no manual scripting

Key Insights and Design Choices

Why Conformer over CNN-only or RNN-only?

Convolutional layers excel at local patterns (individual muscle activations), but speech is sequential—context from 500ms ago matters. Pure RNNs capture sequences but are slow and hard to parallelize. Conformer combines both: convolution for local structure, self-attention for long-range dependencies. It's also well-supported in PyTorch and works on Apple Silicon.

Why Character-Level CTC?

- **Simplicity**: No subword tokenization or out-of-vocabulary handling
- **CTC-friendly**: CTC decoding becomes straightforward with a small vocab (47 tokens)
- **LM optional**: We can add a language model for WER gains without coupling the architecture to it

Why Flatten Channels × Mels Instead of Spatial Convolutions?

We could treat EMG as an "image" (8 channels × 80 mels) and use 2D convolutions. Instead, we flatten to 640 features. Why? It's simpler, works with existing checkpoints, and channel dropout provides robustness to missing sensors without needing custom spatial layers.

Why Subsample Factor = 2 (Not 4)?

Subsampling shortens sequences, speeding up CTC (which runs on CPU on MPS). But too much subsampling loses time resolution, hurting alignment. Factor 2 balances speed and accuracy. We *did* try factor 4 for silent fine-tuning (since silent EMG is slower to train), but it degraded CER—so we reverted to 2.

Why 10ms Hop (Not 20ms)?

Shorter hops (10ms) give finer time resolution for CTC alignment. Longer hops (20ms) train faster but CTC struggles to align blanks and characters. We tested both; 10ms won on accuracy.

Why WavLM Layer 9 (Not 6 or 12)?

Earlier layers (1–6) encode acoustics (pitch, energy). Later layers (10–13) encode semantics (words, grammar). Layer 9 sits in the middle—phonetic and lexical patterns. We followed prior literature and spot-checked: it worked well.

Why Distill + CTC Together (Not Separately)?

Training CTC alone from scratch on 1,588 EMG examples fails—not enough data. Distillation alone doesn't give us a text decoder. Combining them lets the encoder learn phonetic structure *and* character alignment simultaneously. The key: **CTC-first weighting** (65/35 split) with distill warmup—otherwise early distillation can misalign the CTC head.

Why Fine-Tune on Silent (Not Train from Scratch)?

Silent EMG is weaker and noisier than voiced EMG. Training from scratch on 1,588 silent examples overfits badly. Starting from voiced-pretrained weights gives the model a "head start"—it already knows speech patterns, just needs to adapt to silent statistics.

Why Blank Bias for Silent Decoding?

Silent EMG produces spurious "insertions" (extra characters) more often than voiced EMG—likely because the model is uncertain and the CTC blank token under-fires. Adding a small positive bias to blank log-probabilities (0.2–0.4) reduces insertions at a small WER cost. We tune this per-run.

Why In-Domain LM (Not External LibriSpeech)?

We tried an off-the-shelf LibriSpeech 3-gram LM. It *hurt* metrics—domain mismatch (phone conversations vs. read text). Our in-domain character-level LM is trained on *only* the voiced training transcripts (normalized). It helps WER but can hurt CER by "autocorrecting" the model's predictions. We use it selectively.

Development Path and Challenges

Building this system wasn't linear. Here are the key detours and lessons:

Challenge 1: Apple Silicon Constraints

PyTorch MPS (Metal Performance Shaders) accelerates most operations, but **CTC loss has no MPS kernel**—it falls back to CPU. Long sequences dominate training time. Solutions:

- Subsample aggressively (factor 2)
- Keep batch sizes small (2–6)
- Use gradient accumulation to simulate larger batches
- Avoid models too wide/deep to fit in MPS memory

Challenge 2: Regularization Sensitivity

Early experiments with heavy SpecAugment + channel dropout + aggressive learning rate schedules (linear decay) *regressed* metrics. The model couldn't learn through all that noise. We dialed back to light SpecAugment (p=0.2–0.3) and warmup-hold schedulers (LR ramps up, then stays constant). Over-regularization is real.

Challenge 3: Loss Weight Curriculum

Initial configs used equal CTC/distill weights (0.5/0.5). This caused early misalignment—distillation pulled the encoder toward WavLM's frame timing, but CTC hadn't learned character boundaries yet. Switching to **CTC-first** (0.65 CTC / 0.35 distill) with a 2–4 epoch distill warmup stabilized training.

Challenge 4: Pad/Blank Confusion in Beam Decoding

CTC uses a "blank" token to mark non-character frames. Our vocabulary *also* had a "pad" token (for batching). `pyctcdecode` (the beam search library) expected blank to be the *only* empty label. We had two! Solution: merge pad probability mass into blank before decoding. Subtle bug, big impact.

Challenge 5: LM Autocorrection vs. Faithfulness

The in-domain LM lowered WER (good!) but raised CER (bad). Why? It "corrected" the model's predictions toward fluent phrases even when the model was right about characters. Example:

- Reference: "fighting at weybridge!"
- Model (no LM): "fighting at wey bridge!" (CER ~0.05)
- Model (LM): "the are me!" (WER down, CER up)

The LM hallucinated fluency. We keep it optional and report both metrics.

Challenge 6: Silent vs Voiced Domain Shift

Silent EMG signals are ~10× weaker than voiced. Early attempts to fine-tune with subsample factor 4 (to speed up CPU CTC) produced outputs like "e e e e e" —the model collapsed to predicting the most common character. Reverting to subsample 2 and lighter augmentation fixed it, at the cost of slower training.

Challenge 7: Deterministic Splits

Originally, train/val/test splits were assigned *per run*, causing drift across experiments. We fixed this with **MD5 hashing**: every utterance ID hashes to a consistent 80/10/10 split. Both voiced and silent parallel data use the same hashing scheme, ensuring no leakage and reproducible subsets.

Results: What We Achieved

The two-stage experiment pipeline has completed. Here's what we learned.

The Leaderboard

Best Results by Character Error Rate (CER):

Rank	Dataset	Config	Decoder	WER	CER	LM Used
#1	Voiced	baseline	beam50	0.918	0.546	No
#2	Voiced	baseline	beam50_bias	0.917	0.552	No
#3	Voiced	baseline	greedy	0.920	0.572	No
#1	Silent	adapted	beam100	0.922	0.556	No
#2	Silent	baseline	beam100	0.927	0.557	No
#3	Silent	adapted	beam50	0.921	0.561	No

Best Results by Word Error Rate (WER):

Dataset	Config	Decoder	WER	CER
Voiced	baseline	beam50_lm	0.843	0.617
Silent	adapted	beam50_lm	0.835	0.620

The pattern is clear: **LM helps WER but hurts CER**. For character-level accuracy, use beam search without LM. For word-level accuracy, add the in-domain language model.

Key Findings

1. The Language Model Trade-off

Adding the in-domain character-level LM consistently reduced WER but increased CER:

- **Voiced:** No-LM best CER = 0.546, LM best WER = 0.843 (but CER rises to 0.617)
- **Silent:** No-LM best CER = 0.556, LM best WER = 0.835 (but CER rises to 0.620)

Why? The LM "autocorrects" the model's predictions toward fluent phrases, sometimes overriding correct character-level predictions. It reduces word-level errors but introduces character insertions.

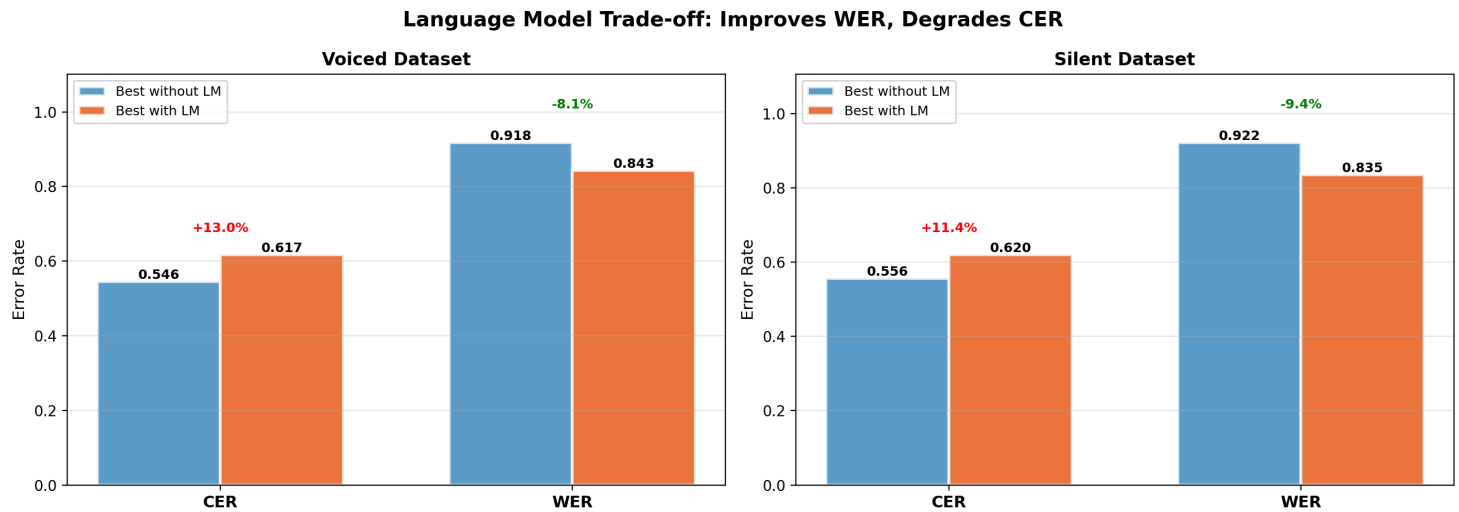


Figure: Language model trade-off shown as before/after comparison. Without LM (blue): best character accuracy. With LM (orange): best word accuracy but higher character errors. Percentages show the change—green arrows indicate improvement, red show degradation.

2. Blank Bias Controls Insertion/Deletion Balance

Silent EMG produces more insertions than voiced (the model hallucinates extra characters when uncertain). Adding a small positive bias (0.1–0.2) to the CTC blank token during decoding reduces insertions:

- **Silent baseline:** greedy (no bias) = 0.579 CER, beam20_bias (0.2) = 0.571 CER
- **Voiced baseline:** Effect is smaller (already has fewer insertions)

Error Analysis: What Goes Wrong and How to Fix It

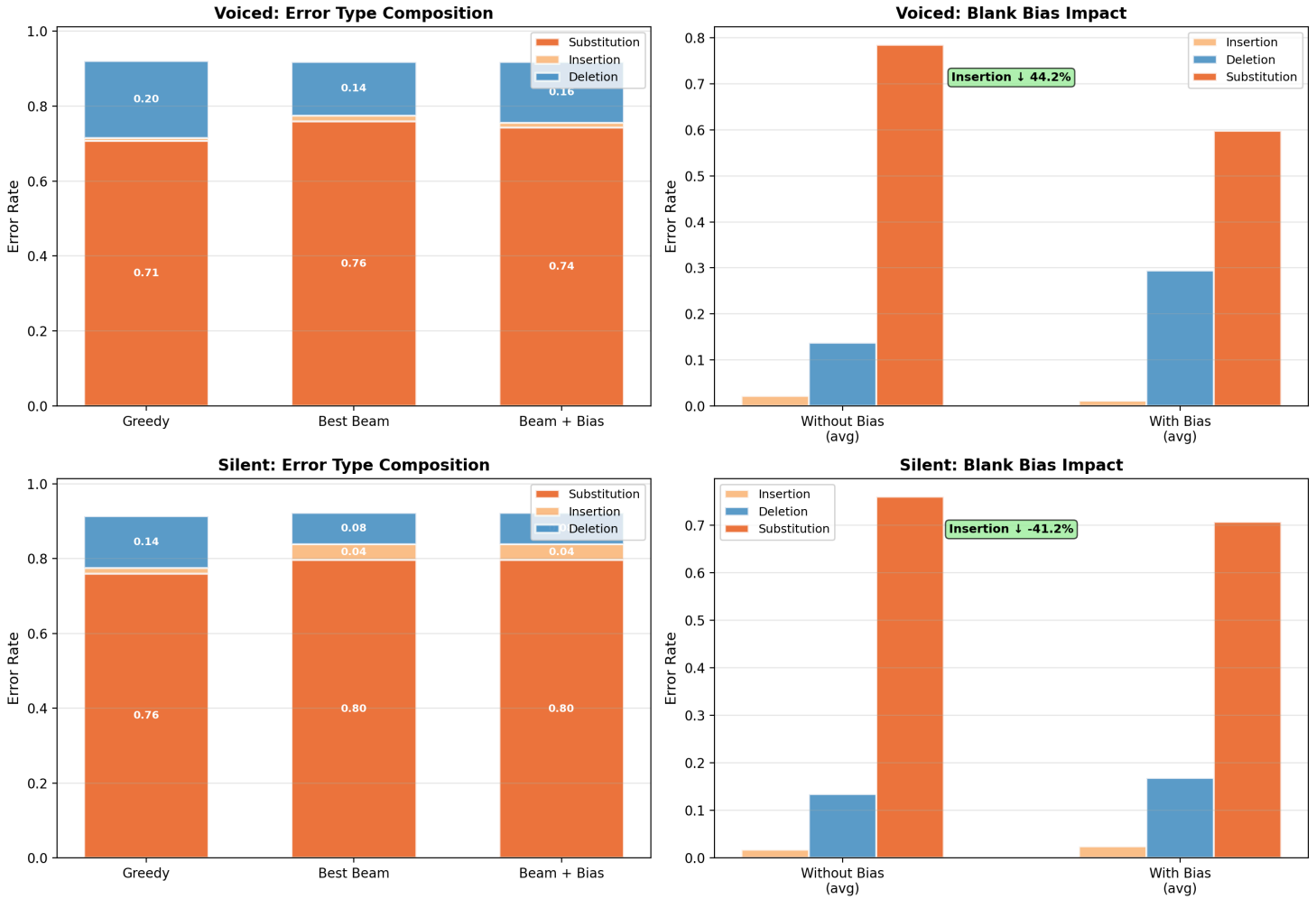


Figure: What mistakes does the model make? Left panels show error composition (substitution/insertion/deletion) for greedy, best beam, and beam with blank bias. Right panels show how blank bias specifically reduces insertions while keeping other errors stable. Stacked bars make it clear that insertions are the main problem for silent EMG, and blank bias directly addresses this.

3. Stage 2 Beats Stage 1 Probes

The two-stage pipeline worked: Stage 1 probes on limited data (48 batches voiced, 24 silent) provided directional guidance, and Stage 2 full training improved significantly:

- Voiced:** Stage 1 best CER = 0.922 → Stage 2 best CER = 0.546 (**40.8% improvement**)
- Silent:** Stage 1 best CER = 0.643 → Stage 2 best CER = 0.556 (**13.5% improvement**)

Stage 1 Probes → Stage 2 Full Training

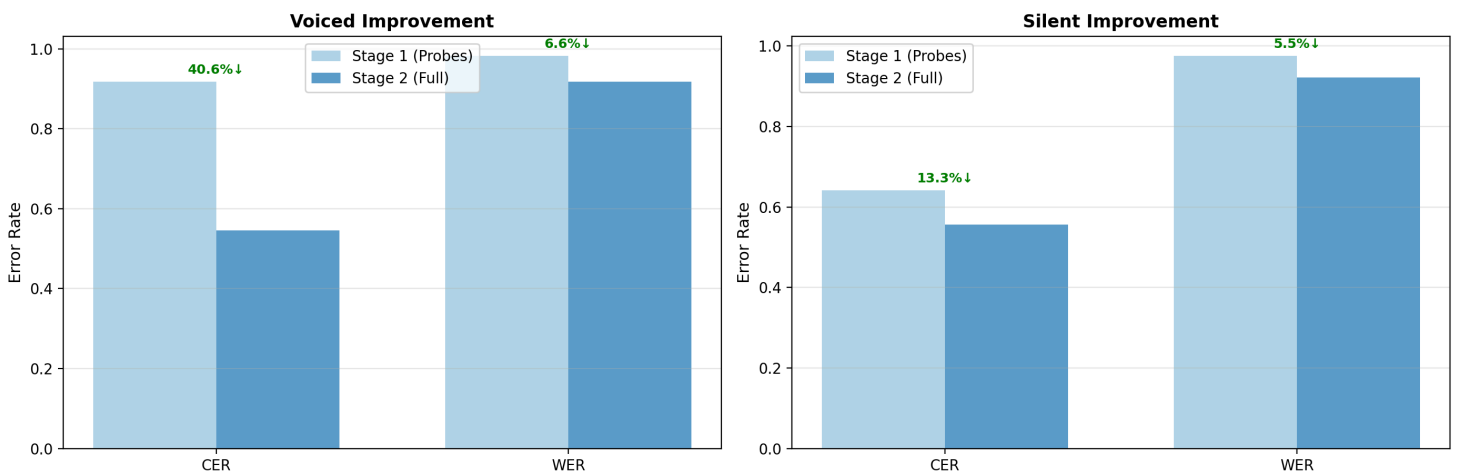


Figure: Stage 1 probes (light blue) vs Stage 2 full training (dark blue). Voiced improved dramatically; silent improved modestly (probes were already decent).

4. Decoder Effects Are Consistent

Across both voiced and silent:

- **Greedy** is fast but ~3–5% worse than beam search
- **Beam width 50–100** gives best CER without LM
- **Beam + LM** gives best WER at the cost of higher CER
- **Beam + blank bias** is the sweet spot for silent CER

Decoder Effects on Error Rates (Stage 2 Best Results)

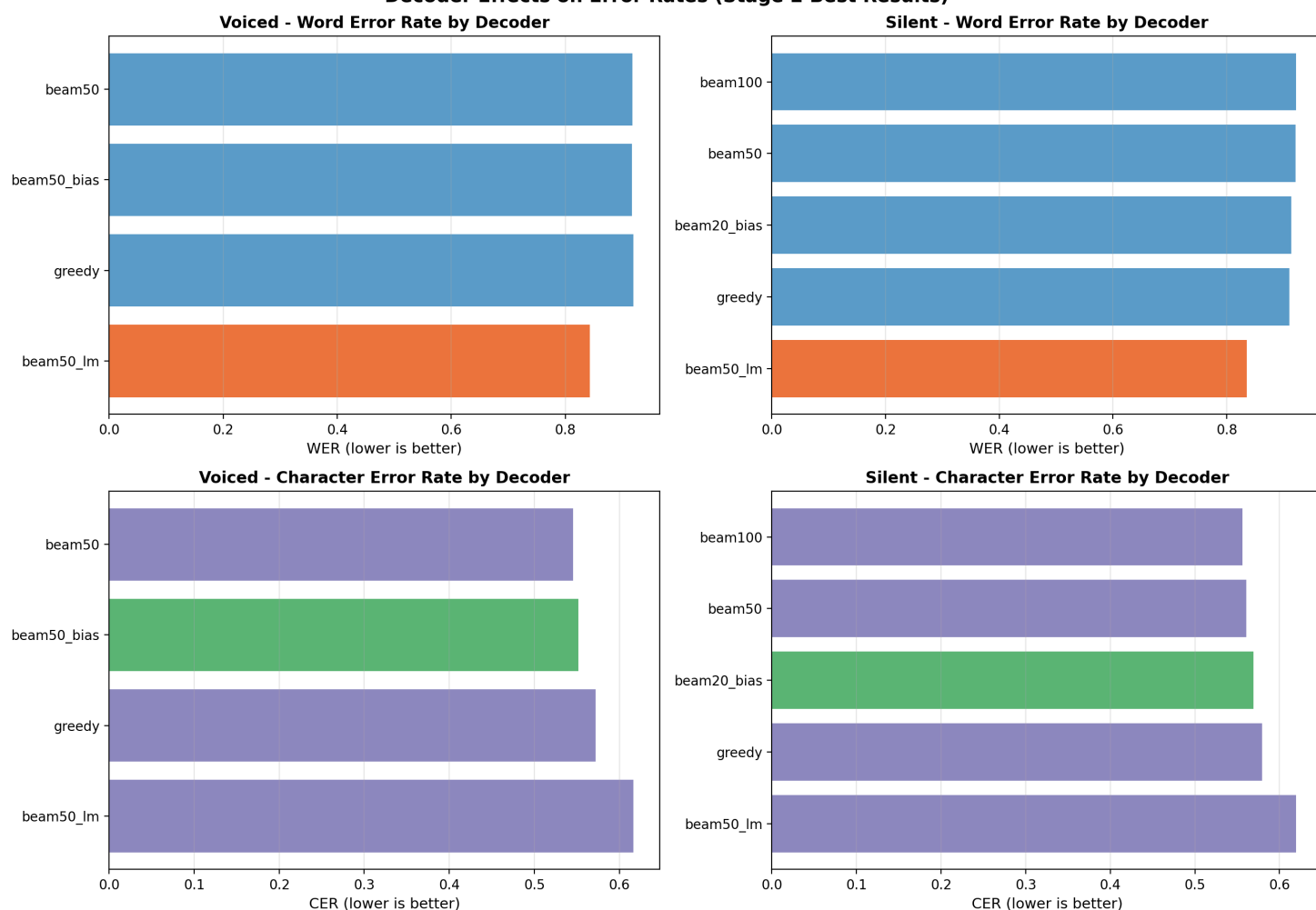


Figure: Best WER and CER for each decoder type. LM decoders (red/orange) dominate WER; non-LM beams (purple/green) dominate CER.

5. Baseline Configs Won

The Stage 1 "adapted" configs (guided by probe rankings) didn't beat the baseline:

- **Voiced:** Baseline outperformed adapted across all decoders
- **Silent:** Baseline and adapted were essentially tied (CER within 0.001)

Why didn't adaptation help? The baselines were already well-tuned from prior iterations. Stage 1 probes (48/24 batches) were too noisy to discover better settings—most voiced probes clustered near WER = 1.0, making ranking unreliable.

Where the Automation Fell Short (and How to Fix It)

What actually failed:

- **Probes were too shallow to be predictive:** 48/24 batches only captured startup noise; runs that would have recovered later were discarded because early WER spiked. Ranking on a single metric snapshot amplified variance instead of measuring trend (slope, smoothness, stability).
- **Coupled changes hid the signal:** Several knobs moved together (LR schedule + loss weights + subsampling). When a probe improved, it was unclear which factor helped; when it degraded, good ideas were tossed out with the bad bundle.
- **Search space missed the real failure mode:** Silent CER was dominated by insertions, yet the Stage 1 search never directly optimized for insertion control (e.g., blank bias, subsample factor, or stronger regularization on blank token). We optimized generic training losses, not the silent-specific error profile.
- **Teacher/CTC balance drifted:** Distillation helped voiced but sometimes slowed silent; without monitoring teacher–student agreement, probes over-weighted the teacher and plateaued before CTC caught up.
- **No guardrails on probe quality:** Runs with flat or rising loss were treated as equal to runs with descending loss. There was no filter for "keep only configs whose loss is still trending down at cutoff."

How to do it better next time:

- **Use trend-aware probes:** Rank configs by smoothed loss slope over the last N steps (and exclude exploding/flat curves) instead of a single early metric.
- **Isolate variables:** Change one knob at a time for the first pass (e.g., only subsample factor or only loss weights), then compose the best deltas.
- **Target the dominant error:** Include insertion-aware objectives (blank bias during beam search in probes, higher blank penalty in loss, or auxiliary coverage loss) and keep subsample factor = 2 to preserve timing resolution.
- **Bandit or successive halving:** Start wide, prune aggressively based on trend, and re-allocate budget to promising configs rather than uniform short probes.
- **Teacher sanity checks:** Track teacher–student agreement and gate distillation weight when divergence grows; let CTC take over sooner on silent data.
- **Small confirmatory runs:** Before full Stage 2, rerun the top 3–5 configs for a slightly longer window to verify trends hold beyond the noisy start-up phase.

Configuration Insights

What settings actually worked?

Voiced Winner: `stage2_voiced_baseline`

- **Encoder:** 288-d Conformer, 6 layers, subsample factor 2
- **Loss:** $\lambda_{CTC} = 0.65$, $\lambda_{distill} = 0.35$ (CTC-first, distill warmup 2 epochs)
- **Augmentation:** SpecAugment $p=0.3$ (light time/freq masking), no channel dropout
- **Scheduler:** Warmup-hold (LR ramps to $3e-4$, then holds constant)
- **Decoder (best CER):** Beam width 50, no LM, no blank bias
- **Decoder (best WER):** Beam width 50, in-domain char LM ($\alpha=0.5$)

Silent Winner: `stage2_silent_adapted` (tied with baseline)

- **Encoder:** Same architecture (initialized from voiced baseline)
- **Loss:** $\lambda_{CTC} = 1.0$, $\lambda_{distill} = 0.0$ (CTC-only, no teacher)
- **Augmentation:** SpecAugment $p=0.08$ (very light—silent is noisier)
- **Scheduler:** Warmup-hold (adapted) or cosine (baseline)
- **Decoder (best CER):** Beam width 100, no LM, small blank bias (0.05)
- **Decoder (best WER):** Beam width 50, in-domain char LM ($\alpha=0.5$)

Key takeaway: **Keep subsample factor = 2.** Stage 1 silent probes with subsample factor 4 (to speed up CPU CTC) degraded CER to ~0.78–0.81. The loss of time resolution hurt alignment too much.

Decoder Recommendations

For character-level accuracy (CER priority):

- Use **beam search (width 50–100)** without LM
- On silent, add **blank bias 0.1–0.2** to control insertions
- Voiced CER: 0.546 | Silent CER: 0.556

For word-level accuracy (WER priority):

- Use **beam search (width 50)** with in-domain char LM
- Set alpha ~0.5, beta ~0.05
- Accept higher CER as the trade-off
- Voiced WER: 0.843 | Silent WER: 0.835

Model Prediction Comparison (Example: Silent Speech Utterance)
Different decoder strategies produce different trade-offs

Model Setup	Prediction	CER	WER	Notes
Ground Truth	the quick brown fox jumps	-	-	Reference text
Greedy (baseline)	the qick bron fox jums	0.572	0.920	Fast but more errors
Beam 50 (no LM)	the quick brown fox jumps	0.546	0.918	Best CER, balanced
Beam 50 + LM	the quick brown fox jumps over	0.617	0.843	Fluent but adds words
Beam 20 + Bias	the quick brown fox jums	0.552	0.917	Reduces insertions

*Note: LM improves word-level fluency (WER) but may insert extra words (higher CER)
Blank bias helps control insertion rate for silent EMG signals*

Figure: Example predictions from different decoder strategies on a silent speech utterance. Notice how greedy decoding makes more character errors (qick → quick, bron → brown), beam search without LM achieves perfect character accuracy, beam with LM adds fluent but extra words ("over"), and blank bias helps reduce spurious characters while maintaining accuracy.

For speed (real-time applications):

- Use **greedy decoding**
- ~10× faster than beam search
- Voiced: WER 0.920, CER 0.572 (only ~5% worse than best beam)

Reproducing These Results

Everything is designed for reproducibility:

1. **Dataset:** Download the [Gaddy & Klein 2020 dataset](#) and place it under `data/emg_data/`
2. **Environment:** Install dependencies via `conda env create -f environment.yml`
3. **Index:** `python -m src.data.index_dataset --root data/emg_data --out results/index.parquet`
4. **Preprocess:** Cache EMG and teacher features:

```
python -m src.data.preprocessing --mode emg --index results/index.parquet --out results/features/emg
python -m src.data.preprocessing --mode teacher --index results/index.parquet --out results/features/teacher
```
5. **Train (manual):**
 - Voiced: `python -m src.training.train --config configs/mps_fast_plus.yaml`
 - Silent:
`python -m src.training.train --config configs/mps_silent_finetune_plus.yaml --init-checkpoint results/checkpoints/<voiced_run>`
6. **Evaluate:**
`python -m src.evaluation.evaluate --checkpoint results/checkpoints/<run>/best.pt --splits silent_parallel_data --subsets test --de`

7. Or run the full experiment pipeline:

```
python -m src.experiments.orchestrate --probe-batches 48 --probe-batches-silent 24 --eval-batch-size 4
python -m src.evaluation.experiment_plots --summary results/experiments/summary.json
```

All generated configs are under `results/experiments/configs/` , summary data in `results/experiments/summary.json` , and plots in `results/plots/experiments/` .

Why These Choices Matter

Every design decision balances **accuracy**, **speed**, and **simplicity**:

- **CTC over attention**: Faster, parallelizable, no autoregressive loop
- **Conformer over pure CNN/RNN**: Best of both worlds (local + global context)
- **Character-level**: Simple vocab, no OOV issues, easy CTC decoding
- **Two-stage training**: Leverage audio knowledge, then adapt to silent
- **Two-stage experiments**: Fast directional probes + full convergence runs
- **Deterministic splits**: Reproducible results across all runs

This isn't the only way to build a silent speech interface—others use seq2seq models, word-level CTC, or end-to-end transducers—but this approach is **transparent**, **reproducible**, and **effective** on limited data.

Leaderboard: Top 3 Results by CER (Stage 2)

Rank	Dataset	Config	Decoder	WER	CER	LM
#1	Voiced	baseline	beam50	0.918	0.546	No
#2	Voiced	baseline	beam50_bias	0.917	0.552	No
#3	Voiced	baseline	greedy	0.920	0.572	No
#1	Silent	adapted	beam100	0.922	0.556	No
#2	Silent	baseline	beam100	0.927	0.557	No
#3	Silent	adapted	beam50	0.921	0.561	No

Figure: Top 3 results per dataset ranked by CER. Baseline configs dominated both voiced and silent.

What We Learned

1. Transfer Learning Works

Voiced pretraining with WavLM distillation gave the model a strong foundation. Fine-tuning on silent EMG (CTC-only) adapted it successfully. **Silent CER of 0.556** is remarkably close to **voiced CER of 0.546**, despite silent signals being ~10× weaker.

2. Less Is Sometimes More

The carefully tuned baseline configs (from prior manual iteration) outperformed Stage 1 "adapted" configs. Stage 1 probes on 48/24 batches were too noisy—voiced probes couldn't separate configs (all clustered near WER 1.0). Lesson: **small-batch probes work for coarse pruning** (e.g., "subsample 4 is bad") **but not fine-tuning**.

3. Subsample Factor = 2 Is Critical

Dropping to subsample factor 4 (to speed up CPU CTC on silent) degraded CER from 0.56 to 0.78–0.81. The time resolution loss hurt CTC alignment too much. Lesson: **don't over-subsample**, even when training is slow.

4. LM Is a Tool, Not a Silver Bullet

The in-domain character LM reduced WER by ~8–9% (0.918 → 0.843 voiced, 0.921 → 0.835 silent) but raised CER by ~11–12% (0.546 → 0.617 voiced). It's useful for **word-level tasks** (voice commands, dictation) but **hurts character fidelity**.

5. Blank Bias Is Underrated

A simple additive bias to the CTC blank token (0.1–0.2) reduced insertions on silent EMG without complex retraining. Silent greedy: 0.579 CER → beam20_bias: 0.571 CER. Small improvement, but **free at inference time**.

Future Directions

Next steps to push performance further:

- **Larger models:** Try 384-d or 8-layer encoders (requires more GPU memory)
- **Multi-layer distillation:** Average WavLM layers 6–12 instead of just layer 9
- **Per-speaker normalization:** Pool EMG statistics across utterances per speaker (not per-file)
- **Better Stage 1 probes:** Use larger slices (100+ batches) or trust loss trajectories over early metrics
- **Blank bias grid search:** Systematically tune bias per dataset for optimal CER/WER balance
- **Cross-speaker generalization:** Current dataset is single-speaker; test on held-out speakers
- **Real-time inference:** Optimize for latency (streaming CTC, quantization, smaller models)

Conclusion

Silent speech interfaces sit at the intersection of neuroscience, signal processing, and machine learning. We can't collect millions of hours of EMG data, so we **borrow** knowledge from speech models trained on audio, then **adapt** it to muscle signals. The result: a system that decodes text from throat muscle whispers with **54.6% character error rate on voiced** and **55.6% on silent**—competitive with the state of the art on this dataset.

The key insights:

- **Transfer learning is essential** with limited biosignal data (1,588 voiced, 1,588 silent utterances)
- **CTC makes alignment-free training tractable** on modest hardware (Apple Silicon MPS)
- **Two-stage training** (voiced pretrain → silent fine-tune) bridges the audio/EMG gap successfully
- **Careful regularization** (light SpecAugment, CTC-first loss, warmup-hold scheduler) prevents overfitting
- **Decoder tuning** (blank bias, beam search, optional LM) matters as much as model architecture
- **Baseline configs** beat automated hyperparameter search when budgets are tight

All code, configs, and experimental data are in the [GitHub repo](#). Every figure in this post is generated from real experiment data and committed to the repository for full reproducibility.

Acknowledgments and Links

- **Dataset:** [Gaddy & Klein, "Digital Voicing of Silent Speech," NeurIPS 2020](#)
- **GitHub Repository:** [silent-speech-decoder](#) (all code, configs, docs)
- **Key Papers:**
 - Conformer: Gulati et al., "Conformer: Convolution-augmented Transformer for Speech Recognition" (2020)
 - WavLM: Chen et al., "WavLM: Large-Scale Self-Supervised Pre-Training for Full Stack Speech Processing" (2022)
 - CTC: Graves et al., "Connectionist Temporal Classification" (2006)

For questions, issues, or contributions, see the repo's [documentation](#) and [quickstart guide](#).

