# Fine-Tuning of WavLM-Large for Automatic Speech Recognition on LibriSpeech

Anonymous Author

Anonymous

*Abstract*—This paper presents a comprehensive progression of fine-tuning experiments aimed at enhancing Automatic Speech Recognition (ASR) using the WavLM-Large model with a Connectionist Temporal Classification (CTC) objective and KenLM language model integration. Four sequential experiments are reported: (1) initial CTC fine-tuning on a 60% subset of LibriSpeech train-clean-100, (2) optimization and scheduler refinement on the full 100h dataset, (3) shallow fusion with a 4-gram KenLM decoder, and (4) KenLM hyperparameter optimization. The experiments progressively reduce Word Error Rate (WER) from 18.1% to 4.05%, approaching the 2.7% benchmark of WeNet on the same dataset. Each experiment is analyzed in depth, emphasizing optimization behavior, decoding dynamics, and the effects of language model weighting. Mathematical formulations are provided for loss functions, decoding scores, and evaluation metrics, highlighting the theoretical grounding of each step. Code and resources are available at: https://github.com/aaivu/In21-S7-CS4681-AML-Research-Projects/tree/210163M/projects/210163M-NLP_Speech-Recognition.

*Index Terms*—Automatic Speech Recognition, WavLM, CTC, KenLM, Self-Supervised Learning, Beam Search, LibriSpeech, Fine-Tuning, Language Model Fusion.

## I. INTRODUCTION

End-to-end Automatic Speech Recognition (ASR) has advanced rapidly with the emergence of self-supervised learning (SSL) and transformer architectures. SSL encoders such as wav2vec 2.0 [2], HuBERT [3], and WavLM [1] leverage massive unlabeled audio corpora to learn generalizable acoustic and semantic representations, drastically reducing supervised data requirements.

Despite these advances, fine-tuning pretrained models for domain-specific ASR remains non-trivial. Typical challenges include over-fitting on small datasets, sensitivity to optimization hyperparameters, and the absence of linguistic priors in Connectionist Temporal Classification (CTC) [4]

decoding. These limitations motivate integrating external language models such as KenLM [5] to constrain decoding and improve grammatical fluency.

This paper investigates a progressive optimization pipeline: (1) baseline fine-tuning on partial data, (2) full-data optimization, (3) decoder enhancement with a KenLM 4-gram model, and (4) decoder hyperparameter tuning. Each stage provides technical reasoning and empirical validation.

## II. RELATED WORK

### A. Self-Supervised Pretraining

SSL models like wav2vec 2.0 [2], HuBERT [3], and WavLM [1] have demonstrated strong performance on downstream ASR tasks. WavLM introduces denoising pretraining and speech separation objectives, which improve robustness to noise and speaker variation.

### B. CTC Objective

Connectionist Temporal Classification (CTC) [4] allows sequence-to-sequence mapping without explicit alignment. For an input sequence $\mathbf{x} = (x_1, x_2, \ldots, x_T)$ and output label sequence $\mathbf{y} = (y_1, y_2, \ldots, y_U)$, the CTC probability is:

$$P(\mathbf{y}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{y})} \prod_{t=1}^{T} P(\pi_t|x_t) \quad (1)$$

where $\pi$ is an alignment path including blank tokens, and $\mathcal{B}$ is the collapse function removing blanks and repeated symbols. The CTC loss is then:

$$\mathcal{L}_{CTC} = -\log P(\mathbf{y}|\mathbf{x}) \quad (2)$$

## C. Language Model Fusion

During decoding, we integrate a 4-gram KenLM via shallow fusion. The score for a candidate transcription $\mathbf{y}$ is:

$$S(\mathbf{y}|\mathbf{x}) = \log P_{CTC}(\mathbf{y}|\mathbf{x}) + \alpha \log P_{LM}(\mathbf{y}) + \beta|\mathbf{y}| \tag{3}$$

where $\alpha$ balances acoustic vs. language model, $\beta$ is a word insertion penalty, and $|\mathbf{y}|$ is the number of words. Beam search is used to find the highest-scoring $\mathbf{y}$.

## D. Evaluation Metrics

Word Error Rate (WER) and Character Error Rate (CER) are computed using Levenshtein distance:

$$WER = \frac{S + D + I}{N} \tag{4}$$

$$CER = \frac{S_c + D_c + I_c}{N_c} \tag{5}$$

where $S$ is substitutions, $D$ deletions, $I$ insertions, and $N$ is the number of reference words/characters.

## III. Methodology

### A. Model and Loss

We use the publicly available WavLM-Large as encoder; a learned linear projection (CTC head) outputs character probabilities at each frame. Training minimizes the standard CTC negative log-likelihood:

$$\mathcal{L}_{CTC} = -\log \sum_{\pi \in \mathcal{B}^{-1}(y)} \prod_{t=1}^{T} P(\pi_t \mid X), \tag{6}$$

where $\mathcal{B}$ collapses repeats and removes blanks, $X$ is the acoustic input and $y$ the target symbol sequence [4].

### B. Decoding with KenLM

For Experiment 3 we perform beam search with an external KenLM 4-gram. The scoring function used for candidate sequence $y$ combines acoustic and LM scores:

$$\text{score}(y) = \log P_{AM}(y \mid X) + \alpha \log P_{LM}(y) + \beta \cdot |y|, \tag{7}$$

where $\alpha$ is the LM weight and $\beta$ is an insertion penalty. We tune $(\alpha, \beta)$ on *dev-clean*.

## IV. Experiments

All experiments use the LibriSpeech corpus [7]. Audio is resampled to 16 kHz, normalized, and tokenized to a character-level vocabulary (including space and apostrophe). Optimization uses AdamW [8]. Mixed precision training (fp16) is enabled for compute efficiency.

### A. Experiment 1: Baseline (60% subset of train-clean-100)

*1) Configuration and Motivation:* Experiment 1 serves as a resource-constrained baseline. We trained only on 60% of *train-clean-100* to validate pipeline and to measure sample-efficiency of WavLM representations.

Training arguments (exact) used:

```
TrainingArguments(
    group_by_length=True,
    per_device_train_batch_size=4,
    per_device_eval_batch_size=4,
    gradient_accumulation_steps=2,
    eval_strategy="steps",
    num_train_epochs=3,
    fp16=True,
    save_steps=100,
    eval_steps=100,
    logging_steps=100,
    learning_rate=3e-4,
    weight_decay=0.005,
    warmup_steps=100,
    lr_scheduler_type="linear",
    adam_beta1=0.9,
    adam_beta2=0.999,
    adam_epsilon=1e-08,
    max_grad_norm=1.0,
    save_total_limit=2
)
```

Rationale: small batch (4) and gradient accumulation allow stable updates within GPU memory limits. A higher initial LR (3e-4) with a short warmup (100 steps) encourages rapid adaptation from pretrained weights on this small dataset; a linear scheduler reduces LR smoothly to avoid sudden drops. Weight decay is modest to avoid underfitting given limited data.

*2) Results and Analysis:* Table I summarizes selected checkpoints. Training shows steady reductions in both training and validation losses; however, WER remains relatively high compared to full-data runs.

TABLE I: Experiment 1 selected checkpoints (WER on dev-clean)

| Step | Train Loss | Val Loss | WER |
|------|-----------|----------|--------|
| 2500 | 0.6325 | 0.3654 | 0.3504 |
| 3500 | 0.4837 | 0.2681 | 0.2607 |
| 5000 | 0.4063 | 0.2103 | 0.1982 |
| 6400 | 0.3807 | 0.1960 | 0.1809 |

Interpretation: With only 60h of data, the encoder's pretrained features transferably reduce errors but remain limited by linguistic coverage. The results motivate training on the full 100h split and changing scheduling/regularization for better convergence.

*B. Experiment 2: Full train-clean-100 with optimized hyperparameters*

*1) Configuration and Motivation:* We scaled training to full *train-clean-100* ( 100h) and adjusted hyperparameters to (i) exploit more data, (ii) speed up convergence, and (iii) stabilize optimization with larger effective batch sizes.

Training args:

```
TrainingArguments(
    group_by_length=True,
    per_device_train_batch_size=8,
    per_device_eval_batch_size=4,
    gradient_accumulation_steps=2,
    eval_strategy="steps",
    num_train_epochs=4,
    fp16=True,
    save_steps=400,
    eval_steps=400,
    logging_steps=400,
    learning_rate=2e-4,
    weight_decay=0.01,
    warmup_steps=1000,
    lr_scheduler_type="cosine",
    adam_beta1=0.9,
    adam_beta2=0.999,
    adam_epsilon=1e-08,
    max_grad_norm=1.0,
    save_total_limit=2,
    gradient_checkpointing=True
)
```

Rationale for changes:

- **Batch size increased to 8**: larger batch improves gradient estimates and allows more stable updates when combined with gradient accumulation.
- **Lower LR (2e-4) and longer warmup (1000)**: when training on more data and larger batch sizes, a smaller LR and longer warmup stabilize training; warmup reduces the risk of catastrophic forgetting of pretrained weights.
- **Cosine LR scheduler**: provides smoother long-term decay and often improves final generalization when training longer.
- **Weight decay increased to 0.01**: stronger regularization helps avoid overfitting to frequent speakers or passages in the 100h set.
- **Gradient checkpointing**: allows larger models or larger batches to be trained within memory constraints, at the cost of additional compute.

*2) Results and Analysis:* Table II shows the step-wise progression for Experiment 2 which demonstrates faster convergence and substantially reduced WER.

TABLE II: Experiment 2 detailed progression (selected steps)

| Step | Train Loss | Val Loss | WER |
|------|-----------|----------|--------|
| 400  | 4.1239 | 2.9466 | 1.0000 |
| 1200 | 1.3196 | 0.5604 | 0.5092 |
| 2000 | 0.4343 | 0.1969 | 0.1734 |
| 3200 | 0.2632 | 0.1207 | 0.1027 |
| 4000 | 0.2041 | 0.1025 | 0.0822 |
| 6800 | 0.1469 | 0.0878 | 0.0682 |

Interpretation: Moving to full 100h yields a dramatic WER improvement (from $\approx$18% to $\approx$6.8% at later checkpoints). The cosine schedule and longer warmup reduced early saturation and produced steady reductions in validation loss. Gradient checkpointing allowed the usage of a larger per-device batch, improving optimization dynamics without OOM.

## C. Experiment 3: KenLM integration for decoding

*1) Configuration and Motivation:* Despite the large gains in Experiment 2, many word-level errors remain (insertions/substitutions) due to limited linguistic context in CTC outputs. To mitigate these, Experiment 3 uses a 4-gram KenLM model during beam-search decoding. Key hyperparameters tuned on dev-clean:

- Beam width (typical range: $10 - 200$)
- LM weight $\alpha$ (controls influence of LM vs acoustic model)
- Word insertion penalty $\beta$ (controls sequence length bias)

Rationale: KenLM is fast, memory-efficient for n-gram scoring, and augments sequence probabilities with explicit syntactic/lexical priors that are not learned by CTC training.

*2) Results:* Evaluation (final checkpoints) with tuned decoding parameters:

- **Validation (dev-clean)**: WER = 0.0394103158, CER = 0.0133711901
- **Test (test-clean)**: WER = 0.0416159464, CER = 0.0134870174

These metrics show a further absolute reduction of $\approx$2.6% WER relative to Experiment 2 and a relative reduction of over 60% compared to Experiment 1.

## D. Experiment 4 – KenLM Hyperparameter Tuning

To balance LM and acoustic confidence, we tested multiple $\alpha$–$\beta$ pairs:

TABLE III: Experiment 4: WER for different $(\alpha, \beta)$

| $\alpha$ | $\beta$ | WER (%) |
|---|---|---|
| 0.3 | 0.1 | 4.35 |
| 0.3 | 0.35 | 4.37 |
| 0.5 | 0.1 | **4.05** |
| 0.5 | 0.35 | 4.06 |
| 0.5 | 0.5 | 4.08 |

The optimum $(\alpha, \beta) = (0.5, 0.1)$ reduced WER to 4.05 %. This corroborates prior LM-fusion studies [14], [16], which found moderate LM weights yield best trade-offs.

## V. TECHNICAL VALIDATION AND ABLATIONS

### A. Hyperparameter Ablation

We ran targeted ablations to verify the effect of individual changes:

- **Batch size and LR**: Increasing batch to 8 while lowering LR (3e-4 $\rightarrow$ 2e-4) improved stability and final WER; high LR produced faster initial loss drop but worse generalization.
- **Warmup length**: Longer warmup (1000) prevented large gradient steps from unlearning pretrained features early in training.
- **Scheduler type**: Cosine scheduler produced smoother late-stage convergence than linear decay in our longer runs.
- **Weight decay**: increasing to 0.01 reduced minor overfitting observed on dev-clean when training for 4 epochs.
- **Gradient checkpointing**: allowed memory-limited devices to train with larger batches; note the trade-off with slower wall-clock time due to recomputation.

### B. Decoding Ablation

We explored beam widths $\{16, 64, 128\}$ and LM weights $\alpha \in [0.1, 2.0]$. Best dev-clean WERs occurred with beam width $\approx 100$ and $\alpha \in [0.6, 1.2]$, and a small negative insertion penalty (to discourage overly-short hypotheses). Excessive LM weight caused over-reliance on the LM, producing hallucinated common n-grams and harming WER.

### C. Reproducibility

To ensure reproducibility, we provide these implementation notes:

- Frameworks: PyTorch (1.11), Hugging Face Transformers (4.x), datasets and evaluate packages for metric computation.
- Seeds: set global RNG seeds for Python, NumPy, and PyTorch; report seeds for final runs.
- Hardware: experiments were run on a machine with NVIDIA GPUs (e.g., A100/3090); exact runtime depends on device — Experiment 2 (3 epochs on 100h) took 10 hours per 3 epochs on a single high-end GPU (timing varies).
- Checkpoints: we store model state, optimizer state, scheduler state and tokenizer for full resumption.

- Exact training args for Experiment 1 and 2 are included above for replication.

## VI. Discussion

The results demonstrate that both optimization scheduling and LM fusion substantially improve recognition accuracy. The learning rate scheduler choice and warmup duration directly affect CTC convergence. In Experiment 1, the higher learning rate caused unstable gradients, mitigated in Experiment 2 via cosine decay. Language model integration reduced insertion/deletion errors, especially for homophones and rare words.

Comparatively, WeNet reports 2.7% WER on LibriSpeech-clean. Although our 4.05% WER is slightly higher, it is achieved with significantly fewer compute resources and simpler fine-tuning. The improvement trajectory indicates that additional training on LibriSpeech-460h or domain-specific data could further close the gap.

## VII. Error Analysis

Manual inspection revealed most remaining errors arise from long utterances and compound words. Character-level substitutions were rare, suggesting robust acoustic modeling. The main sources of degradation were:

- Missing word boundaries due to CTC independence assumption.
- Insufficient coverage of rare vocabulary in 4-gram LM.
- Slight overconfidence in frequent function words.

## VIII. Future Work

Future experiments will include domain adaptation using LibriSpeech-460h and multi-domain fine-tuning, integration of Transformer-based neural LMs, and exploration of hybrid decoding with neural LM fusion (cold or deep fusion). Incorporating RNN-T decoding may further improve streaming recognition latency.

## IX. Conclusion

This paper presents a structured progression of WavLM fine-tuning for ASR, combining CTC-based training and KenLM fusion. Experiment 4 achieves a WER of 4.05%, demonstrating competitive performance relative to established systems. The findings highlight that thoughtful optimization and LM integration can deliver strong ASR performance even under limited resources.

## Acknowledgment

## References

[1] S. Chen, Y. Zhang, K. Xu, et al., "WavLM: Large-Scale Self-Supervised Pre-Training for Full Stack Speech Processing," *IEEE J. Sel. Topics Signal Process.*, 16(6):1505–1518, 2022.

[2] A. Baevski, Y. Zhou, A. Mohamed, M. Auli, "wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations," *NeurIPS*, 2020.

[3] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, et al., "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units," *IEEE/ACM TASLP*, 29:3451–3460, 2021.

[4] A. Graves, S. Fernández, F. Gomez, J. Schmidhuber, "Connectionist Temporal Classification," *ICML*, 2006.

[5] K. Heafield, "KenLM: Faster and Smaller Language Model Queries," *WMT Workshop on Statistical Machine Translation*, 2011.

[6] B. Zhang, R. Liu, J. Huang, et al., "WeNet: Production Oriented Streaming and Non-Streaming End-to-End Speech Recognition Toolkit," *Interspeech*, 2021.

[7] V. Panayotov, G. Chen, D. Povey, S. Khudanpur, "LibriSpeech: An ASR Corpus Based on Public Domain Audio Books," *ICASSP*, 2015.

[8] I. Loshchilov, F. Hutter, "Decoupled Weight Decay Regularization," *ICLR*, 2019.

[9] A. Vaswani, N. Shazeer, N. Parmar, et al., "Attention Is All You Need," *NeurIPS*, 2017.

[10] D. Amodei, R. Anubhai, E. Battenberg, et al., "Deep Speech 2: End-to-End Speech Recognition in English and Mandarin," *ICML Workshop / arXiv:1512.02595*, 2016.

[11] A. Hannun, C. Case, J. Casper, et al., "Deep Speech: Scaling Up End-to-End Speech Recognition," *arXiv:1412.5567*, 2014.

[12] W. Chan, N. Jaitly, Q. V. Le, O. Vinyals, "Listen, Attend and Spell," *ICASSP*, 2016.

[13] A. Graves, "Sequence Transduction with Recurrent Neural Networks," *arXiv:1211.3711*, 2012.

[14] A. Sriram, H. Jun, S. Satheesh, A. Coates, "Cold Fusion: Training Seq2Seq Models Together with Language Models," *Interspeech*, 2018.

[15] S. Toshniwal, et al., "An Analysis and Empirical Study of Language Model Integration for Seq2Seq Speech Recognition," *arXiv:1807.10857*, 2018.

[16] R. Cabrera, et al., "Language Model Fusion for Streaming End-to-End Speech Recognition," *arXiv:2104.04487*, 2021.

[17] S. Watanabe, T. Hori, S. Kim, et al., "ESPnet: End-to-End Speech Processing Toolkit," *Interspeech*, 2018.

[18] D. Povey, A. Ghoshal, G. Boulianne, et al., "The Kaldi Speech Recognition Toolkit," *ASRU Workshop*, 2011.

[19] D. Park, et al., "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," *Interspeech*, 2019.

[20] T. Chen, B. Xu, C. Zhang, "Training Deep Nets with Sublinear Memory Cost," *arXiv:1604.06174*, 2016.

[21] I. Loshchilov, F. Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts," *ICLR*, 2017.

[22] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, Y. Bengio, "Attention-Based Models for Speech Recognition," *NeurIPS*, 2015.

[23] H. Sak, A. Senior, F. Beaufays, "Long Short-Term Memory Recurrent Neural Network Architectures for Large Scale Acoustic Modeling," *Interspeech*, 2014.

[24] A. Graves, "Speech Recognition with Deep Recurrent Neural Networks," *ICASSP*, 2013.

[25] S. Karita, et al., "A Comparative Study on Transformer vs RNN in Speech Recognition," *IEEE ASRU*, 2019.

[26] S. Karita, et al., "ESPnet2: End-to-End Speech Processing Toolkit—Improved Design, Implementation, and Performance," *ICASSP*, 2020.

[27] Y. Wang, et al., "Transformer-Based Language Models for Speech Recognition," *Interspeech*, 2020.

[28] K. Heafield, "Efficient Language Model Queries for Statistical Machine Translation," *WMT Workshop*, 2011.

[29] K. Cho, et al., "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation," *EMNLP*, 2014.

[30] D. Kingma, J. Ba, "Adam: A Method for Stochastic Optimization," *ICLR*, 2015.