# Trax v2 Research and Architecture Analysis: A Deep Dive into Performance, Speaker Diarization, and Advanced Features

## The Next Frontier of Accuracy: Multi-Pass Processing and Domain-Specific Enhancement

The pursuit of transcription accuracy beyond the 95% baseline achieved by Whisper distil-large-v3 is a primary driver for Trax v2. The research indicates that this can be accomplished through two distinct yet complementary strategies: multi-pass processing to refine transcriptions iteratively, and domain-specific enhancement to tailor the model's understanding to specialized content. These approaches move beyond a single-pass inference model, embracing a more sophisticated pipeline architecture where outputs from one stage serve as inputs or context for subsequent stages, leading to significant quality improvements.

Multi-pass processing represents a paradigm shift in ASR systems, designed to bridge the gap between real-time responsiveness and offline-quality accuracy. This strategy involves chaining multiple models together, often with different strengths, to progressively improve the final transcript. One of the most compelling examples is the two-pass end-to-end speech recognition system developed by Google researchers [8]. This architecture pairs a fast, streaming Recognizer Neural Transducer (RNN-T) model with a slower, non-streaming Listen, Attend and Spell (LAS) model that shares an encoder network [4] [8]. The RNN-T acts as a first pass, providing a preliminary hypothesis quickly, while the LAS model performs a deeper rescoring of the top-K hypotheses from the first pass, leveraging its attention-based mechanism to capture longer-range dependencies [8]. This approach has been shown to achieve a 17%-22% relative Word Error Rate (WER) reduction compared to the RNN-T alone, effectively closing the quality gap with traditional non-streaming models, all while keeping the latency increase under 200ms—a trade-off that is highly favorable for many applications [2] [4] [8]. Further innovation comes from cascaded systems that optimize for efficiency; one study demonstrated reducing the frame rate of the second pass by half resulted in a 20% reduction in Real-Time Factor (RTF) and 13% power savings without impacting final accuracy [22].

Another powerful technique within the multi-pass framework is iterative refinement, which leverages the output of a transcription model to directly improve its own performance. Research shows that self-supervised speech models like HuBERT become better at representing linguistic features with each training iteration, improving their correlation with canonical phoneme and word identities while de-correlating from speaker identity [23]. This suggests that using a model's own pseudo-labels to create new training data for subsequent iterations enhances its core capabilities. An even more advanced concept is mutual enhancement, where ASR and Voice Conversion (VC) models are trained in a loop, with the ASR model generating text to train the VC model, and the VC model

generating synthesized audio to augment the ASR training set [1]. While complex, this approach demonstrates how models can learn from each other without requiring massive annotated datasets, pointing towards a future where Trax could continuously improve its own transcription engine.

Domain-specific enhancement addresses the challenge of maintaining high accuracy across diverse content types, such as technical lectures, medical consultations, or noisy conference calls. The industry standard for this is full-scale fine-tuning on domain-specific data, but this is computationally expensive and risks catastrophic forgetting of general knowledge [32]. Fortunately, the field of Parameter-Efficient Fine-Tuning (PEFT) offers elegant solutions. Low-Rank Adaptation (LoRA) is a standout method that freezes the vast majority of the pre-trained model's weights and injects small, trainable "rank decomposition" matrices into the transformer layers [31] [33]. This drastically reduces memory requirements and training time while achieving near-full fine-tuning performance. For instance, LoRA was used to adapt Whisper for domain adaptation with less than 0.1% of the total model parameters, resulting in WER reductions of over 3 points [28]. Other PEFT methods include prompt tuning, which learns a small, trainable "prompt" vector to steer the model's behavior, and speech prefix tuning (SPT), which appends fixed-length vectors to input features and has been shown to outperform LoRA on certain tasks [34] [57]. A particularly innovative approach involves text-only fine-tuning, where only the Language Model (LLM) component of a Speech LLM is adapted using unpaired target-domain text, preserving the original speech encoder's integrity and avoiding performance degradation on general domains [29]. This allows Trax to build a highly accurate, specialized model for a niche domain like financial reporting or legal proceedings without bloating the overall system.

| Feature | Multi-Pass Processing | Domain-Specific Enhancement (PEFT) |
|---|---|---|
| Core Principle | Chaining multiple models (e.g., fast-first, slow-second) to refine results iteratively [8]. | Modifying a small fraction of a pre-trained model's weights to adapt it to a specific domain [31]. |
| Primary Benefit | Achieves near-offline accuracy with low latency [4]. Enables high accuracy in specialized contexts [51]. | |
| Key Techniques | Shared-encoder architectures [8], N-best rescoring [8], adaptive beam search [4], cascaded systems [22]. | Low-Rank Adaptation (LoRA) [36], Prompt Tuning [41], Text-Only Fine-Tuning [29], Adapter Tuning [30]. |
| Performance Impact | 17-22% relative WER reduction vs. single-pass models [2]. 3-11 absolute point WER reduction [28] [29]. | |
| Implementation Cost | Increased complexity in pipeline architecture. Requires managing multiple models. | Low computational cost for adaptation. Minimal storage overhead for adapter modules. |

| Feature | Multi-Pass Processing | Domain-Specific Enhancement (PEFT) |
|---|---|---|
| Use Case for Trax | Ideal for creating a "quality" processing path that users can select for critical transcripts. Enables creation of lightweight, specialized "modules" for different verticals. | |

In summary, the path to 99.5%+ accuracy for Trax v2 is not through a single architectural leap but through a layered, intelligent approach. By implementing a multi-pass processing pipeline, Trax can offer superior accuracy as a selectable feature. By integrating PEFT techniques like LoRA, Trax can provide deep domain specialization without sacrificing its core generality or performance, positioning itself as a versatile and powerful tool for a wide range of professional use cases.

## Mastering Conversational Audio: State-of-the-Art Speaker Diarization and Voice Profiling

Speaker diarization—the process of identifying who spoke when in a conversation—is a critical feature for any modern transcription platform, transforming a monolithic transcript into an actionable document. For Trax v2, achieving robust and accurate speaker diarization is paramount, especially given the user's interest in handling conversations. The current state of the art offers a spectrum of solutions, from established open-source frameworks to cutting-edge deep learning models, each with distinct trade-offs in accuracy, latency, and resource consumption.

The most effective speaker diarization systems today are typically modular, combining several components: a Voice Activity Detector (VAD) to segment the audio into speech turns, an embedding extractor to generate a compact speaker representation ("d-vector" or "x-vector") for each turn, and a clustering algorithm to group these embeddings by speaker [9]. Frameworks like Pyannote.audio have become a de facto standard in this space, offering a well-engineered implementation of this pipeline [10] [12]. However, recent advancements in end-to-end (E2E) neural speaker diarization promise to simplify this process. Models like EEND (End-to-End Neural Speaker Diarization) replace the separate clustering step with a single neural network that predicts the number of speakers and assigns a label to each frame of the input [7] [13]. While promising, E2E models often face challenges with latency and require fixed speaker limits, making them less suitable for real-time applications or scenarios with an unknown number of participants [13].

Comparative studies provide crucial insights into the performance of these systems. In a direct comparison on a hobbyist-grade project, Taishin Maeda evaluated Pyannote.audio against NVIDIA NeMo on two different audio files [12]. On a 5-minute, two-speaker file, NeMo achieved a lower Diarization Error Rate (DER) of 16.1% compared to Pyannote's 25.2%, albeit at the cost of double the execution time (63.9s vs 31.3s). On a more challenging 9-minute, nine-speaker file, Pyannote performed slightly better with a DER of 8.3% versus NeMo's 9.7% (with pre-identified speakers) [12]. Another study benchmarked various systems on the Voxconverse dataset and found that DIART, based on `pyannote/segmentation` and `pyannote/embedding`, had the lowest latency at

just 0.057 seconds per chunk on a CPU, whereas another E2E model, UIS-RNN-SML, became impractically slow on long recordings [10]. These findings suggest that for a project like Trax, which values both accuracy and performance, a hybrid approach might be optimal: using a highly efficient, lightweight system like DIART or a custom-built module for initial, real-time processing, and reserving heavier, more accurate models like Pyannote for post-processing or user-selected "enhanced" analysis modes.

Latency is a critical factor, especially for real-time applications. Most modern systems operate with some degree of look-ahead, analyzing a few hundred milliseconds of future audio to make more confident decisions about speaker changes and endpointing. Bilal Rahou et al. proposed a causal segmentation model that uses a multi-latency look-ahead during training, allowing it to dynamically adjust its latency to balance performance with speed, nearly matching offline model accuracy with just 500ms of look-ahead [16]. In contrast, AssemblyAI's diarization model is currently limited to asynchronous transcription, not real-time streams, highlighting the technical hurdles involved [18]. For Trax v2, a configurable latency setting would be a powerful feature, allowing users to trade a slight delay for significantly improved accuracy in detecting overlapping speech and speaker turns.

Beyond simple diarization, voice profiling and privacy-preserving methods represent the next frontier. Speaker identification, the task of labeling who a speaker is, can be enhanced by adapting models like ECAPA-TDNN with speaker embeddings, which has been shown to improve DER on children's speech data [19]. For privacy-sensitive applications, techniques like zero-party authentication, where no actual voice samples are stored, become essential. Furthermore, a novel and highly effective technique involves using a Large Language Model (LLM) as a post-processing step to correct diarization errors. Researchers fine-tuned a Mistral 7b model on the Fisher corpus to analyze transcripts from various ASR systems and correct speaker labels, demonstrating an ASR-agnostic correction capability that significantly improved accuracy [20,50]. This opens up a fascinating possibility for Trax v2: after producing a raw transcript with speaker labels, it could run the text through a specialized LLM-based diarization-corrector to produce a polished, expertly labeled version. This approach decouples the core transcription task from the complex, context-dependent task of speaker attribution, potentially leading to higher overall accuracy and greater flexibility.

| System / Technique | Key Strengths | Key Weaknesses | Latency Profile | Source(s) |
|---|---|---|---|---|
| Pyannote.audio | High accuracy, mature ecosystem, extensive documentation. | Slower than some alternatives, especially on long audio. | ~31s for 5min audio on RTX 3090. | [12] |
| NVIDIA NeMo | Lower DER than Pyannote on short, clean audio. | Slower execution time, requires more GPU memory. | ~64s for 5min audio on RTX 3090. | [12] |
| DIART (**pyannote**) | Extremely low latency (~57ms/chunk on CPU), scalable. | May be less accurate than fully optimized models on very challenging data. | Very low latency (0.057s per chunk). | [10] |

| System / Technique | Key Strengths | Key Weaknesses | Latency Profile | Source(s) |
|---|---|---|---|---|
| UIS-RNN-SML | High accuracy, but latency increases dramatically with audio length. | Becomes impractical for long recordings (>9s for 9min audio). | High latency that scales with audio duration. | [10] |
| LLM Post-Processing | ASR-agnostic, can fix contextual errors, improves accuracy. | Adds computational overhead, requires a fine-tuned LLM. | Not specified, but adds to overall processing time. | [20] [50] |
| SelfVC Framework | No explicit speaker labels needed, works on unlabeled data. | Focuses on voice conversion, not diarization. | Not applicable. | [24] [26] |

Ultimately, the choice of speaker diarization technology for Trax v2 depends on the desired user experience. A pure hobby project might prioritize a quick and easy integration of a library like Pyannote.audio. A more ambitious v2 could implement a dual-pathway architecture: a fast, low-latency pathway for real-time transcription and a slower, more accurate pathway for post-processing that employs advanced techniques like E2E models or LLM-based correction. This would give users control over the trade-off between immediacy and precision, delivering a truly state-of-the-art conversational transcription experience.

# Architectural Evolution: From Iterative Refinement to Scalable Cloud-Native Systems

To realize the ambitions of Trax v2—achieving 99.5%+ accuracy, supporting thousands of concurrent users, and enabling advanced features like multi-pass processing and domain-specific models—the underlying architecture must evolve significantly from the current production-ready, protocol-based design. The existing architecture, centered on a batch processor with parallel workers, excels at deterministic, sequential tasks. However, the new requirements demand a more dynamic, distributed, and service-oriented structure. This evolution will involve decomposing the monolith into microservices, adopting a message-driven communication pattern, and embracing cloud-native principles for scalability and resilience.

The most fundamental architectural change required is the transition from a synchronous, blocking batch processor to an asynchronous, event-driven workflow. The current system processes files in a tight loop, which is simple but inefficient for the complex pipelines envisioned for v2. An event-driven architecture (EDA) is far better suited. In this model, the process begins when a user uploads a file. The system creates a `TranscriptionJob` event containing metadata (file ID, source language, requested enhancements) and publishes it to a message broker like Apache Kafka or RabbitMQ. This immediately returns control to the user, fulfilling the low-latency requirement for

initiating a job. Multiple, independent worker services then subscribe to this topic. One worker might handle the initial audio preprocessing, another the primary transcription, a third the speaker diarization, and so on. Each service performs its task and, upon completion, publishes a new event (e.g., `PrimaryTranscriptionComplete`) with its output, triggering the next service in the chain. This decouples the processing stages, allowing them to scale independently and fail without bringing down the entire system. It also naturally enables the multi-pass and multi-model processing flows discussed previously, where the output of one model becomes the input for another.

This EDA forms the basis for a microservice architecture. Instead of a single, large application, Trax v2 would consist of a collection of small, focused services: 1. API Gateway: The single entry point for all client requests. It authenticates users, routes requests to the appropriate backend service, and aggregates responses. 2. Transcription Service: Manages the lifecycle of transcription jobs, interacting with the message broker to trigger and coordinate workflows. 3. Worker Services: Specialized services for different processing tasks (e.g., `WhisperWorker`, `DeepSeekEnhancer`, `DiarizationWorker`). These can be scaled independently based on their computational intensity. 4. Model Management Service: Handles the loading, caching, and versioning of machine learning models. This is crucial for efficiently swapping in different PEFT-adapted models for various domains. 5. Storage Service: Manages access to the PostgreSQL database and object storage for audio files and processed transcripts. 6. Metrics & Logging Service: Collects telemetry data to monitor system health, performance, and error rates.

This modular design offers immense benefits. It allows teams to develop and deploy services independently, facilitates experimentation with new models, and improves maintainability. For example, if a new, more accurate diarization model is released, only the Diarization Worker needs to be updated and redeployed, without touching the rest of the system.

Scalability is a key success metric, targeting 1000+ concurrent transcriptions. This is best achieved through containerization and orchestration. Docker should be used to package each service, ensuring consistency across development and production environments. Kubernetes would then serve as the orchestrator, managing the deployment, scaling, and operation of these containers. Kubernetes' Horizontal Pod Autoscaler (HPA) can automatically increase the number of replicas for a service like the `WhisperWorker` when CPU utilization or the length of the message queue exceeds a threshold, and decrease them when load is low. This ensures resources are used efficiently. To meet the <1GB memory per worker target, careful selection of container base images and optimization of the Python environment (e.g., using `uv` as planned) is critical. Additionally, using smaller, more efficient models, such as quantized versions of Whisper, can further reduce memory footprint [17].

Finally, the architecture must incorporate mechanisms for cost optimization and reliability. Caching is a powerful tool. Transcripts of frequently used content or common phrases can be cached in a system like Redis to avoid redundant processing and reduce API costs. Intelligent caching of intermediate results from multi-pass processing can also yield significant performance gains. For reliability, the system must be designed for failure. This includes implementing idempotency keys for API requests to prevent duplicate processing, using dead-letter queues in the message broker to handle failed messages for later inspection, and ensuring all services are stateless so they can be restarted or replaced without losing data. With a cloud-native architecture built on these principles, Trax v2 can confidently scale to meet demanding workloads while remaining performant, cost-effective, and resilient.

# Optimizing for Scale and Speed: Strategies for Concurrent Transcription and Resource Efficiency

Achieving the ambitious targets of 1000+ concurrent transcriptions and <$0.005 per transcript requires a multi-faceted approach to optimization, focusing on workload distribution, resource management, and computational efficiency. The foundation of this effort lies in moving beyond the current single-machine, multi-worker setup to a distributed, cloud-native architecture capable of horizontal scaling. This involves leveraging containerization, message queues, and efficient model deployment strategies to maximize throughput and minimize operational costs.

The first step toward high concurrency is to eliminate bottlenecks in the processing pipeline. As previously discussed, transitioning to an event-driven architecture with a message broker is central. This decouples the frontend from the backend processing, allowing the system to accept thousands of new transcription jobs instantly without being blocked by the processing capacity. The message broker acts as a buffer, smoothing out spikes in demand. The workers that consume from this queue can then be deployed as a scalable Kubernetes deployment. When the volume of jobs increases, Kubernetes can automatically spin up more worker pods to consume messages from the queue in parallel, distributing the load across multiple machines or cores. This horizontal scaling is the most direct way to handle thousands of concurrent users.

Efficient resource management within each worker is equally critical. The goal is to keep the memory usage below 1GB per worker. This can be achieved through several techniques. First, selecting leaner base images for Docker containers (e.g., `python:slim` instead of a full OS image) and carefully managing dependencies is important. Second, and most critically, is the use of Post-Training Quantization (PTQ). PTQ is a technique that converts the floating-point weights of a trained model into lower-bitwidth integers (e.g., 8-bit or 4-bit) without retraining, significantly reducing the model's memory footprint and accelerating computation. Research has shown that w8-a8 quantization (8-bit weights, 8-bit activations) generally preserves accuracy, while w4-a8 can cause significant degradation in smaller models but is surprisingly robust in larger ones like Whisper Small [17]. Methods like GPTQ and SpQR have demonstrated strong robustness across configurations [17]. By applying PTQ, Trax can deploy multiple instances of the Whisper model on a single GPU, drastically increasing batch processing capacity and reducing the overall hardware cost per transcription.

Further computational efficiency can be gained by optimizing the processing logic itself. For multi-pass systems, as explored in the accuracy section, one study successfully reduced the frame rate of the second pass by 50% without affecting final accuracy, leading to a 20% reduction in Real-Time Factor (RTF) and 13% power savings [22]. This principle can be applied to Trax's v2 processing pipeline, where the initial, faster pass can be executed with a higher frame rate, and a more computationally intensive second pass can run at a lower frame rate on a subset of the audio. This targeted optimization focuses compute resources where they are most needed, improving overall efficiency.

Cost optimization is intrinsically linked to resource efficiency. The target of <$0.005 per transcript is aggressive and will only be met by minimizing every component of the cost equation: compute, storage, and data transfer. Beyond model quantization and efficient pipelines, strategic use of spot instances or preemptible VMs in the cloud can dramatically reduce compute costs, provided the

system is designed to gracefully handle interruptions. Storage costs can be managed by using tiered storage, where raw audio files are stored in cheaper archival storage and only moved to high-performance storage when actively being processed. Data transfer costs, particularly if using a cloud provider, can be minimized by running the API gateway, message broker, and worker services within the same availability zone or region. Finally, intelligent caching, as mentioned earlier, can reduce the need for reprocessing identical content, saving both time and compute cycles.

The following table summarizes key optimization strategies and their potential impact:

| Optimization Strategy | Description | Potential Impact | Source(s) |
|---|---|---|---|
| Event-Driven Architecture | Use a message broker (e.g., Kafka) to decouple job submission from processing. | Enables thousands of concurrent job submissions and independent scaling of worker services. | Analytical Reasoning |
| Horizontal Scaling | Deploy worker services as scalable Kubernetes deployments. | Directly supports handling thousands of concurrent transcription tasks. | [22] |
| Post-Training Quantization (PTQ) | Reduce model size by converting weights to lower-bitwidth integers (e.g., w4-a8). | Reduces memory usage, accelerates inference, and increases batch size, lowering cost per transcript. | [17] |
| Efficient Multi-Pass Pipelines | Reduce computational load in later passes (e.g., by lowering frame rate). | Decreases overall latency and computational cost without sacrificing accuracy. | [22] |
| Cloud-Native Cost Management | Utilize spot/preemptible instances, regional deployments, and tiered storage. | Drastically reduces compute and data transfer costs, meeting aggressive pricing targets. | [5] |
| Parameter-Efficient Fine-Tuning (PEFT) | Use LoRA or similar methods for domain adaptation. | Avoids deploying large, full-scale fine-tuned models, saving memory and storage. | [28] [31] |

By systematically applying these strategies, Trax v2 can build a highly performant, scalable, and cost-effective platform. The architectural shift to a distributed, event-driven system provides the necessary foundation for concurrency. Within that system, optimizations in model quantization, pipeline design, and cloud resource management will ensure that the performance and cost targets are not just met, but exceeded.

# The User Experience Imperative: Designing a Modern Interface and Workflow

While backend performance and feature set are crucial, the ultimate success of Trax v2 hinges on a seamless and intuitive user experience. For a hobby project aimed at becoming a serious tool, the user interface must evolve from a functional but dated Command Line Interface (CLI) to a modern, web-based platform that simplifies complex workflows and provides powerful, accessible tools for reviewing and editing transcripts. The focus should be on streamlining the path from audio upload to final, usable text, catering to the needs of researchers, journalists, and other professionals who rely on accurate transcription.

The immediate priority is the development of a comprehensive web interface. This interface should be built using a modern front-end framework like React or Vue.js, ensuring a responsive and mobile-friendly design [25]. The core workflow should be clear and logical. Upon visiting the site, a user should see a prominent "Upload Audio" button. After uploading a file, the system should present a clean dashboard displaying the status of the transcription job. Once the job is complete, the interface should display the transcript in a readable format. Crucially, this is not just a static display. The transcript should be interactive, allowing users to click on any word to hear the corresponding audio snippet, a feature that greatly aids verification and correction.

One of the most valuable additions to enhance the user experience is real-time collaboration. While the user indicated no critical integrations, the ability for multiple users to review, edit, and comment on a single transcript simultaneously is a powerful productivity tool. This feature, however, presents a significant engineering challenge, especially regarding performance. To support this, Trax v2 must be architected with real-time capabilities from the ground up. This likely involves using WebSockets or a similar persistent connection technology to facilitate low-latency updates. The system must be designed to handle concurrent edits efficiently, perhaps using Operational Transformation (OT) or Conflict-Free Replicated Data Types (CRDTs) to merge changes from multiple users without conflict. The target of <500ms latency for updates is achievable but will require a highly optimized backend and a well-designed front-end architecture [3].

The interface should also provide advanced export options and a flexible workflow. Users should be able to download transcripts in a variety of formats, including SRT for subtitles, DOCX for editable documents, and JSON for programmatic access. Integration with popular note-taking platforms like Obsidian or Notion, though not a "critical" partner, would be a significant value-add and could be implemented via a browser extension or bookmarklet that allows users to send highlighted text directly to their preferred tool. The workflow should be streamlined, minimizing clicks and cognitive load. For example, after correcting an error, the user should be able to immediately request a new translation of a specific sentence or paragraph without having to re-run the entire transcription job.

Finally, the design must be tailored to different user types. A researcher may need to tag specific sections of the transcript with metadata, while a journalist may prioritize quick searching and quoting. The interface should be adaptable, perhaps through user profiles or customizable dashboards, to accommodate these varying needs. The goal is to create an interface that feels both powerful and intuitive, empowering users to leverage the advanced capabilities of the Trax engine without being overwhelmed by its complexity. By investing in a modern, collaborative, and user-

centric web interface, Trax v2 can transform from a powerful engine into an indispensable tool for anyone working with spoken language.

# Synthesizing the Future: A Roadmap for Trax v2 Implementation and Success

The journey from Trax v1.0.0 to a next-generation transcription platform is an exciting opportunity to build a system that is not only faster and more accurate but also architecturally robust and user-centric. The research clearly indicates that the path forward involves a deliberate and phased implementation, starting with foundational architectural upgrades before layering on advanced features. This synthesis outlines a practical roadmap to guide the development of Trax v2, ensuring that the project remains focused, manageable, and aligned with the user's goals of performance and ambition.

Phase 1: Foundation and Core Pipeline (4-6 Weeks)

The initial phase must establish the groundwork for all future enhancements. The primary objective is to overhaul the current architecture. 1. Architectural Decomposition: Begin by breaking down the monolithic CLI and batch processor into a suite of microservices. Develop the core services: an API Gateway, a Transcription Service, and a set of generic Worker Services. Integrate a message broker (e.g., Kafka) to enable the event-driven workflow. 2. Implement Multi-Pass Engine: Build the core engine for multi-pass processing. This involves designing a workflow definition language or configuration system that allows for the chaining of different models (e.g., a fast Whisper variant followed by a DeepSeek enhancement pass). This phase should focus on the orchestration logic rather than building entirely new models. 3. Establish Baseline Accuracy: Implement the current best-in-class single-pass accuracy pipeline using the existing Whisper and DeepSeek models. This serves as a stable baseline against which the improvements from Phase 2 can be measured. Document performance metrics (e.g., 95%+ WER on test sets).

Phase 2: Advanced Features and Domain Adaptation (6-8 Weeks)

With the new architecture in place, this phase focuses on adding the advanced features that differentiate Trax v2. 1. Integrate Speaker Diarization: Implement a robust speaker diarization system. A pragmatic approach would be to integrate a lightweight, efficient library like DIART for real-time processing and pair it with a more accurate model like Pyannote.audio for post-processing. Add user-facing controls to toggle diarization on/off and select the level of detail. 2. Deploy Parameter-Efficient Fine-Tuning (PEFT): Implement a system for applying PEFT methods like LoRA. This involves creating a Model Management Service that can handle the storage and loading of adapter modules. Develop a user interface or API endpoint for selecting a domain-specific model for a particular job. 3. Develop Confidence Scoring: Implement a confidence scoring mechanism. Given the mixed results in the literature, a practical approach would be to extract available scores from the underlying models and provide them as a supplementary tool, clearly documenting their limitations. Do not treat them as a reliable error detection mechanism.

Phase 3: Scalability and Optimization (4-6 Weeks)

This phase is dedicated to ensuring the platform can handle high loads and remain cost-effective. 1. Containerize the Application: Package all microservices into Docker containers. This ensures portability and lays the groundwork for deployment. 2. Orchestrate with Kubernetes: Deploy the application on a Kubernetes cluster. Configure Horizontal Pod Autoscalers (HPA) to automatically scale the worker services based on message queue depth or CPU load. 3. Optimize for Performance and Cost: Apply Post-Training Quantization (PTQ) to the core Whisper model to reduce its memory footprint and accelerate inference. Benchmark the system to verify that the <1GB memory per worker and <$0.005 per transcript targets are met. Optimize the multi-pass pipeline for computational efficiency.

Phase 4: User Interface and Polishing (2-4 Weeks)

The final phase brings the platform to life for the user. 1. Build the Web Interface: Develop a modern, responsive web interface using a framework like React or Vue.js. This interface should manage the entire workflow: uploading files, monitoring job status, viewing and editing transcripts, and downloading results. 2. Implement Real-Time Collaboration: Develop the back-end and front-end logic for real-time collaboration. Start with basic functionality (e.g., shared cursor, simultaneous highlighting) and iterate based on usability testing. 3. Final Testing and Documentation: Conduct comprehensive testing, including performance benchmarks against the v1.0.0 baseline. Generate detailed documentation for developers and end-users.

In conclusion, the development of Trax v2 is a feasible and highly rewarding endeavor. By following this structured roadmap, the project can systematically address the key challenges of architecture, accuracy, scalability, and user experience. The most critical decision is the architectural pivot to a distributed, event-driven system. This choice will unlock the ability to implement multi-pass processing, PEFT, and high concurrency, transforming Trax from a competent tool into a powerful and scalable platform for the future of speech recognition.

## Reference

1. Iteratively Improving Speech Recognition and Voice Conversion https://arxiv.org/abs/2305.15055

2. Two-Pass End-to-End Speech Recognition - Google Research https://research.google/pubs/two-pass-end-to-end-speech-recognition/

3. Two-pass endpoint detection for speech recognition - arXiv https://arxiv.org/html/2401.08916v1

4. [1908.10992] Two-Pass End-to-End Speech Recognition - ar5iv - arXiv https://ar5iv.labs.arxiv.org/html/1908.10992

5. Align-Refine: Non-autoregressive speech recognition via iterative ... https://www.amazon.science/publications/align-refine-non-autoregressive-speech-recognition-via-iterative-realignment

6. Two-Pass Endpoint Detection for Speech Recognition - IEEE Xplore https://ieeexplore.ieee.org/document/10389743/

7. [PDF] End-to-End Neural Speaker Diarization with an Iterative Refinement ... https://www.isca-archive.org/interspeech_2022/rybicka22_interspeech.pdf

8. Two-Pass End-to-End Speech Recognition - ResearchGate https://www.researchgate.net/publication/335830044_Two-Pass_End-to-End_Speech_Recognition

9. An enhanced deep learning approach for speaker diarization using ... https://www.nature.com/articles/s41598-025-09385-1

10. Systematic Evaluation of Online Speaker Diarization Systems ... - arXiv https://arxiv.org/html/2407.04293v1

11. [Literature Review] Two-pass Endpoint Detection for Speech ... https://www.themoonlight.io/en/review/two-pass-endpoint-detection-for-speech-recognition

12. Pyannote.audio vs Nvidia Nemo, and Post-Processing Approach ... https://docs.voice-ping.com/voiceping-corporation-company-profile/apr-2024-speaker-diarization-performance-evaluation-pyannoteaudio-vs-nvidia-nemo-and-post-processing-approach-using-openais-gpt-4-turbo-1

13. A Review of Common Online Speaker Diarization Methods - arXiv https://arxiv.org/html/2406.14464v1

14. Exploring the trade-off between speed and accuracy in real-time ... https://blog.speechmatics.com/latency_accuracy

15. [PDF] Latency and Quality Trade-offs for Simultaneous Speech-to-Speech ... https://www.isca-archive.org/interspeech_2023/dugan23_interspeech.pdf

16. [PDF] Multi-latency look-ahead for streaming speaker segmentation https://www.isca-archive.org/interspeech_2024/rahou24_interspeech.pdf

17. Edge-ASR: Towards Low-Bit Quantization of Automatic Speech ... https://arxiv.org/html/2507.07877v2

18. What is speaker diarization and how does it work? (Complete 2025 ... https://assemblyai.com/blog/what-is-speaker-diarization-and-how-does-it-work

19. Optimizing Speaker Diarization for the Classroom https://jedm.educationaldatamining.org/index.php/JEDM/article/download/841/240

20. LLM-based speaker diarization correction: A generalizable approach https://www.sciencedirect.com/science/article/abs/pii/S0167639325000391

21. A review of the best ASR engines and the models powering them in ... https://www.gladia.io/blog/a-review-of-the-best-asr-engines-and-the-models-powering-them-in-2024

22. Efficient Cascaded Streaming ASR System Via Frame Rate Reduction https://ieeexplore.ieee.org/document/10389645/

23. Iterative refinement, not training objective, makes HuBERT behave ... https://arxiv.org/html/2508.08110v1

24. SelfVC: Voice Conversion With Iterative Refinement using Self ... https://research.nvidia.com/labs/conv-ai/publications/2024/2024-selfvc/

25. [PDF] Comparative Analysis of Personalized Voice Activity Detection ... https://www.isca-archive.org/interspeech_2024/buddi24_interspeech.pdf

26. (PDF) SelfVC: Voice Conversion With Iterative Refinement using ... https://www.researchgate.net/publication/381121265_SelfVC_Voice_Conversion_With_Iterative_Refinement_using_Self_Transformations

27. Thinking aloud.. LoRA & Prompt Tuning - DeepLearning.AI https://community.deeplearning.ai/t/thinking-aloud-lora-prompt-tuning/465150

28. A Domain Adaptation Framework for Speech Recognition Systems ... https://arxiv.org/html/2501.12501v1

29. Low-Resource Domain Adaptation for Speech LLMs via Text-Only ... https://arxiv.org/html/2506.05671v1

30. A Comparison of Parameter-Efficient ASR Domain Adaptation ... https://ieeexplore.ieee.org/document/10445894/

31. [PDF] Smarter Fine-Tuning: How LoRA Enhances Large Language Models https://hal.science/hal-04983079/document

32. Fine-tuning ASR Models: Boosting Accuracy and Adaptability https://lamarr-institute.org/blog/fine-tuning-asr-models/

33. Fine-Tuning Transformers Efficiently: A Survey on LoRA and Its Impact https://www.preprints.org/manuscript/202502.1637/v1

34. Parameter-efficient adaptation with multi-channel adversarial ... https://asmp-eurasipjournals.springeropen.com/articles/10.1186/s13636-025-00406-5

35. [PDF] Low Rank Adaptation for Multilingual Speech Emotion Recognition https://www.isca-archive.org/interspeech_2024/goncalves24_interspeech.pdf

36. [PDF] The Role of LoRA in Parameter-Efficient Adaptation | TechRxiv https://www.techrxiv.org/users/887510/articles/1269329/master/file/data/Revolutionizing_Large_Model_Fine_Tuning__The_Role_of_LoRA_in_Parameter_Efficient_Adaptation/Revolutionizing_Large_Model_Fine_Tuning__The_Role_of_LoRA_in_Parameter_Efficient_Adaptation.pdf

37. Machine Learning Confidence Scores — All You Need to Know as a ... https://medium.com/voice-tech-global/machine-learning-confidence-scores-all-you-need-to-know-as-a-conversation-designer-8babd39caae7

38. How to Use Confidence Scores in Machine Learning Models - Mindee https://www.mindee.com/blog/how-use-confidence-scores-ml-models

39. Evaluating ASR Confidence Scores for Automated Error Detection in ... https://arxiv.org/html/2503.15124v1

40. Using transcription confidence scores to improve slot filling in ... - AWS https://aws.amazon.com/blogs/machine-learning/using-transcription-confidence-scores-to-improve-slot-filling-in-amazon-lex/

41. [PDF] Prompt-tuning in ASR systems for efficient domain-adaptation https://assets.amazon.science/cf/6f/65b75c8544fabc2e2adab334140c/prompt-tuning-in-asr-systems-for-efficient-domain-adaptation.pdf

42. Modular Domain Adaptation for Conformer-Based Streaming ASR https://www.researchgate.net/publication/373248113_Modular_Domain_Adaptation_for_Conformer-Based_Streaming_ASR

43. [PDF] Improving Speech Recognition with Prompt-based Contextualized ... https://www.isca-archive.org/interspeech_2024/manhtienanh24_interspeech.pdf

44. Prompting Large Language Models for Zero-Shot Domain ... https://www.researchgate.net/publication/377538976_Prompting_Large_Language_Models_for_Zero-Shot_Domain_Adaptation_in_Speech_Recognition

45. What is the significance of confidence scores in speech recognition? https://zilliz.com/ai-faq/what-is-the-significance-of-confidence-scores-in-speech-recognition

46. What is the significance of confidence scores in speech recognition? https://milvus.io/ai-quick-reference/what-is-the-significance-of-confidence-scores-in-speech-recognition

47. What do confidence scores mean in speech recognition? https://stackoverflow.com/questions/61331681/what-do-confidence-scores-mean-in-speech-recognition

48. [PDF] Using Automatically Created Confidence Measures - LSEG https://www.lseg.com/content/dam/data-analytics/en_us/documents/white-papers/lseg-itg-automatic-transcript-research-paper.pdf

49. Ensuring Transcription Accuracy: Techniques and Best Practices https://waywithwords.net/resource/transcription-accuracy-best-practices/

50. LLM-based speaker diarization correction: A generalizable approach https://arxiv.org/html/2406.04927v3

51. How accurate is speech-to-text in 2025? - AssemblyAI https://www.assemblyai.com/blog/how-accurate-speech-to-text

52. Survey of End-to-End Multi-Speaker Automatic Speech Recognition ... https://arxiv.org/html/2505.10975v1

53. Speech-to-Text APIs: Key Players and Innovations in 2024 - Krisp https://krisp.ai/blog/speech-to-text-apis-key-players-and-innovations-in-2024/

54. Moving beyond word error rate to evaluate automatic speech ... https://www.sciencedirect.com/science/article/pii/S0165178125003385

55. Top Real-Time Speech-to-Text Tools in 2024 - Galileo AI https://galileo.ai/blog/best-real-time-speech-to-text-tools

56. Method and system for correcting speech-to-text auto-transcription ... https://patents.google.com/patent/US20200160866A1/en

57. Prompt-tuning in ASR systems for efficient domain-adaptation - arXiv https://arxiv.org/abs/2110.06502