# Voice to Text Chat with Lampi Fine-Tuned Llama 2 7B and Whisper Models

GCP T2A VMs, powered by Ampere® Altra CPU and high performance Ampere® AI inference engine, deliver best-in-class GPU-Free AI inference performance on standard AI frameworks, including PyTorch, TensorFlow, and ONNX-RT.

## Ampere Altra Powered Generative AI Inference

Ampere® **Cloud Native Processors** satisfy the performance requirements of widely used generative AI models such as Whisper, originally developed by Open AI and used especially for ASR (automatic speech recognition), and Meta's Llama 2, especially for the smaller versions of the Llama model. Ampere CPUs provide the **best price-performance and optimized power draw** for all kinds of ML deployments.
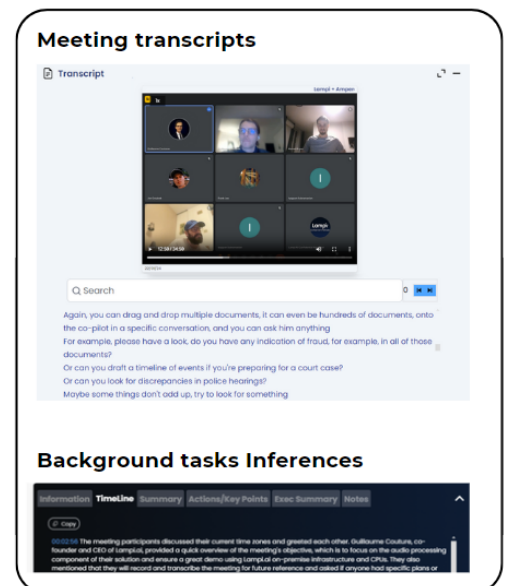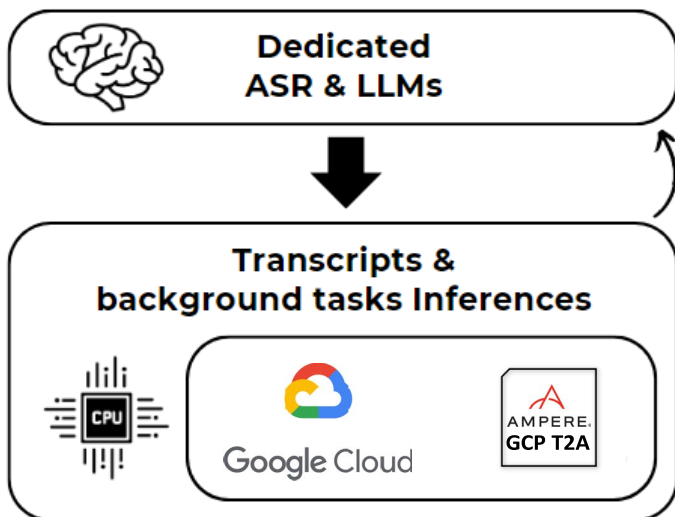
## About Lampi

Lampi provides AI workflow copilot oriented to business use cases - agnostic and fully modular adapted on the fly to client's business needs. Lampi deploys multimodal AI on the fly allowing for a consistency in results to reach the desired quality and compliance levels.

## Setup

Deployment of the dynamic Lampi fine-tuned versions of the open-source generative AI **Llama 2 7B** and **Whisper** models running on Ampere-based **GCP T2A** cloud instance with **Ampere® Optimized AI Frameworks**. The demo showcases real-time chatbot interaction with the training set of Ampere AI and energy efficiency content chosen to inform the audience visiting the Ampere booth at Google Cloud Next.

## Key Benefits Demonstrated

- Meets or exceeds the necessary **low latency** requirements for real-time generative AI chatbot interaction. (voice-to-text, voice-to-voice, text-to-voice, and text-to-text)

- Ampere delivers best throughput for Whisper model inference outcompeting some of the most commonly used GPUs.

- Delivers the best **price-performance** in GPU-Free AI inference in cloud deployment scenarios.

- Ampere Optimized AI Frameworks work right out-of-the-box and can be downloaded free of charge directly from Google Cloud Marketplace or, as docker images, from the Ampere AI solutions page.

- Ampere Cloud Native Processors enable **easy scaling** and can be **dynamically provisioned** based on the performance requirements of the user's applications.
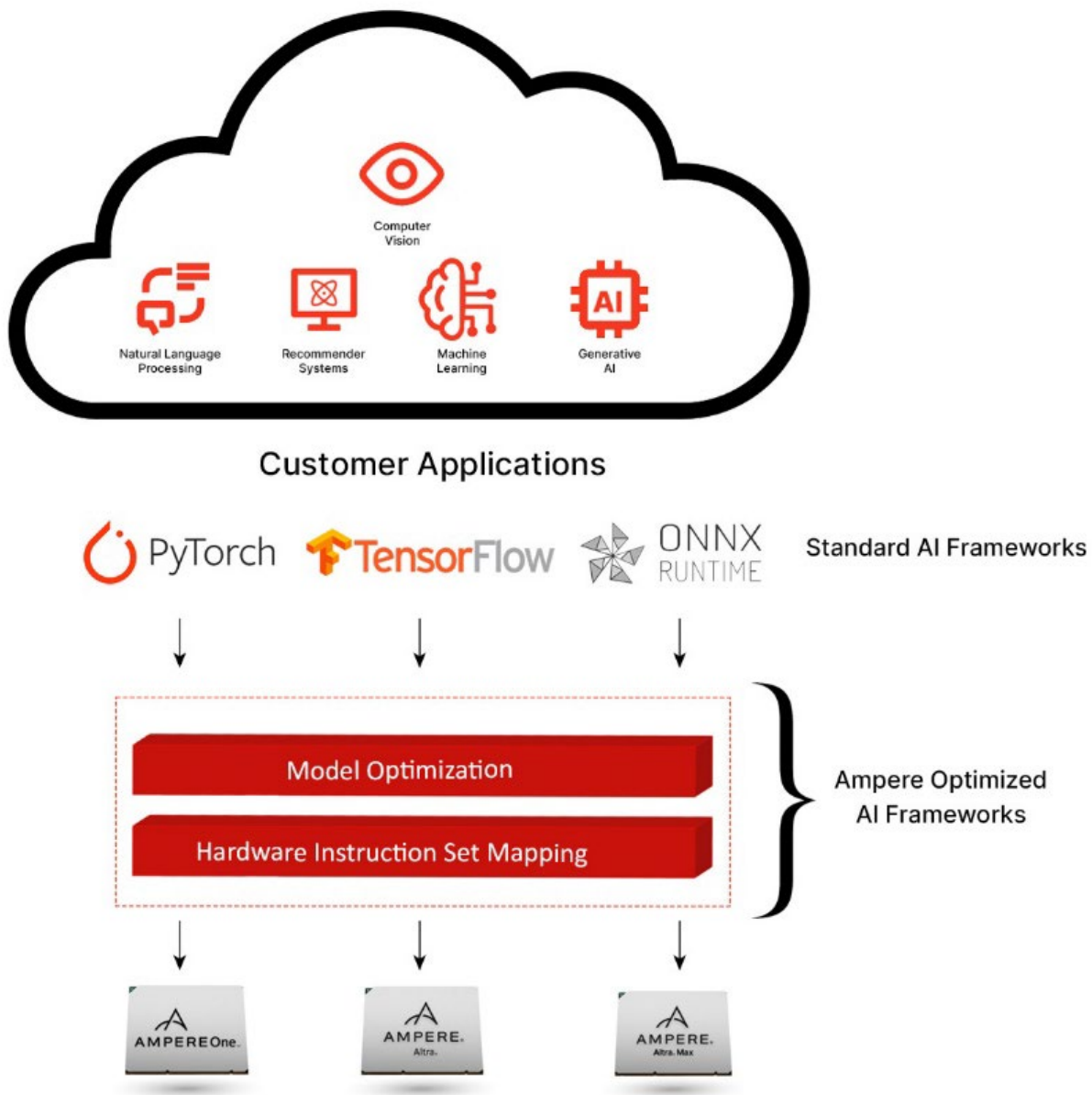
## Real-time Voice-to-Text Chat

This demo shows a partner solution of a **generative AI chatbot deployment**. It processes incoming real-time input from a dedicated external microphone available to the visitors of the Ampere booth at Google Cloud Next. The demo runs at a **real-time performance level** with sufficient **low latency** and per second token generation rate to satisfy the user needs.

## Resources

Ampere Optimized Pytorch, and other Ampere Optimized AI Frameworks, can be accessed directly from Google Cloud Marketplace.

The docker images are also available in the downloads section of Ampere AI Solutions web page. All software is available free of charge and runs straight out-of-the-box with no additional coding required.

**Figure 2:** The integration of Ampere Optimized AI Frameworks with Ampere Altra Cloud Native Processors