

Task Results

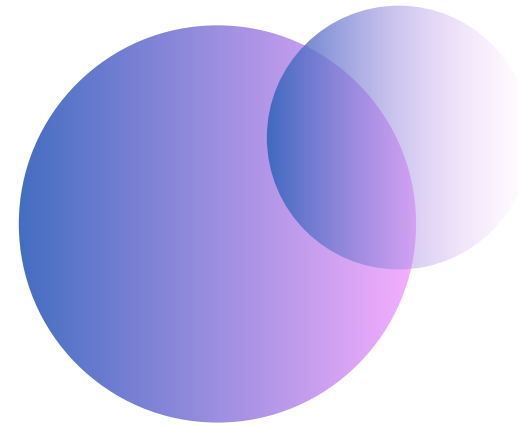
INDEX

1	Setup & Baseline
3	Turn Detector Plugin (Model-based)
5	SSML-Driven TTS Enhancements
7	Deliverables

2	VAD / Endpointing Parameter Tuning
4	Streaming STT & Partial Hypotheses
6	Latency & Quality Benchmarking

1

Setup & Baseline Background



Forked to <https://github.com/Ubaid-Ghante/voice-pipeline-agent-python>

Simple implementation using deepgram, cartesia and openai.
Created accounts in all the above services and got the API keys.

The example code give in Assignment seems for the older version.
Next slide contains a simple demo of the system working.



voice-pipeline-agent-python

EXPLORER

VOICE-PIPELINE-AGENT-PYTHON

task1

agent.py M .env U settings.py U logger_config.py U baseline.py U baseline.py task1 U X

task1 > baseline.py > entrypoint

6)

7 from livekit.agents import (

8 Agent,

9 AgentSession,

10 RoomInputOptions,

11 WorkerOptions,

12 cli,

13 JobContext,

14 AutoSubscribe

15)

16

17 from config.settings import settings

18 from config.logger_config import logger

19

20

21 async def entrypoint(ctx: JobContext):

22 await ctx.connect(auto_subscribe=AutoSubscribe.AUDIO_ONLY)

23

24 # Wait for the first participant to connect

25 participant = await ctx.wait_for_participant()

26 logger.info(f"Starting voice assistant for participant {participant.identity}")

27

PROBLEMS OUTPUT TERMINAL PORTS DEBUG CONSOLE

-US")

2025-07-30 12:56:56,991 - DEBUG livekit.agents - received user transcript {"user_transcript": "Alright. Perfect.", "language": "en-US"}

2025-07-30 12:57:07,298 - DEBUG livekit.agents - received user transcript {"user_transcript": "So how good is your speech to text?", "language": "en-US"}

2025-07-30 12:57:09,191 - DEBUG livekit.agents - received user transcript {"user_transcript": "It's", "language": "en-US"}

2025-07-30 12:57:21,385 - DEBUG livekit.agents - received user transcript {"user_transcript": "You", "language": "en-US"}

2025-07-30 12:57:32,001 - DEBUG livekit.agents - received user transcript {"user_transcript": "Okay. So this is a test without the bag.", "language": "en-US"}

2025-07-30 12:57:39,248 - INFO livekit.agents - shutting down worker {"id": "unregistered"}

2025-07-30 12:57:39,250 - DEBUG livekit.agents - shutting down job task {"reason": "", "user_initiated": false}

2025-07-30 12:57:39,250 - DEBUG livekit.agents - job exiting {"reason": "", "tid": 8753375, "job_id": "simulated-job-d7860149e69a"}

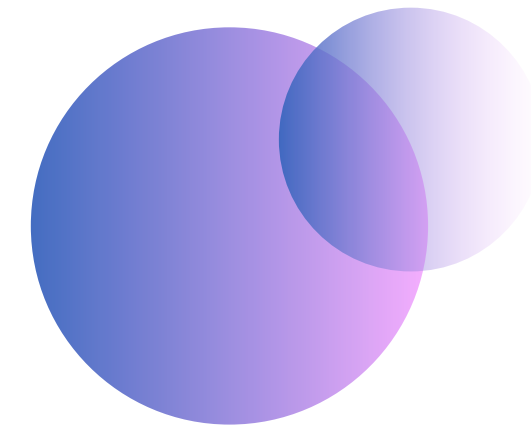
2025-07-30 12:57:39,254 - DEBUG livekit.agents - session closed {"reason": "job_shutdown", "error": null}

2025-07-30 12:57:39,254 - DEBUG livekit.agents - http_session(): closing the httpclient ctx

(voice-pipeline-agent-python) ubaidghante : 12:57 ~/Documents/Projects/VoicePipelineAgent/voice-pipeline-agent-python

Ln 31, Col 243 Spaces: 4 UTF-8 LF Python voice-pipeline-agent-python (3.12.6) Go Live

VAD / Endpointing Parameter Tuning



Min 0.2s and **Max** 3s -

Total latency :
2.0999720122199506 seconds

Min 2.0s and **Max** 6s -

Total latency :
3.34827495040372 seconds

Observation:

- Mainly added idle time if both min and max kept too high
- Starts interrupting if kept too low. High false triggers during natural pauses.

3

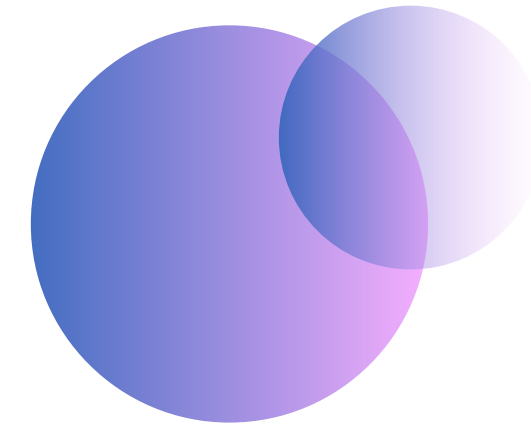
Turn Detector Plugin (Model-based)

Observation:

- Makes it way easier to talk without any false trigger.
- Allows **long pause** where it is needed.
- If eou probability threshold is not specified it takes from the pre configured hf "**livekit/turn-detector**" downloaded file.
- It waits for around 6 to 7 seconds after taking a pause on I think..

```
2025-07-30 16:29:18,355 - DEBUG livekit.agents - received user transcript {"user_transcript": "I think", "language": "en-US"}
2025-07-30 16:29:19,352 - DEBUG livekit.plugins.turn_detector - eou prediction {"eou_probability": 0.001277949777431786, "input": "me a sentence about artificial intelligence<|im_end|>\n<|im_start|>assistant\nhowdy artificial intelli", "duration": 0.094}
2025-07-30 16:29:26,680 - INFO voice-agent - TTS Metrics: type='tts_metrics' label='livekit.plugins.cartesia.tts.TTS' request_id=1753873166.6802402 ttfb=0.2818183330819011 duration=1.0776705420576036 audio_duration=2.4612916666666666 cancelled=False characters_count=31 streamed=True segment_id='584bfa35c9c0' speech_id='speech_7973bfdcad88'
```





Streaming STT & Partial Hypotheses

Observations

- Added stt_node overriding to log results
- Using preemptive generation reduces latency by a little.
- Experimented with basic examples. Major problem seems to be STT, WER is very high.
- LLM Hallucination is very less since most of them are generic test.

```
2025-07-30 17:55:41,612 - DEBUG livekit.agents - http_session(): creating a new httpclient ctx
2025-07-30 17:55:41,839 - DEBUG livekit.agents - using audio io: `ChatCLI` -> `AgentSession` -> `TranscriptSynchronizer` ->
2025-07-30 17:55:41,839 - DEBUG livekit.agents - using transcript io: `AgentSession` -> `TranscriptSynchronizer` -> `ChatCLI`
2025-07-30 17:55:48,533 - INFO voice-agent - #####-----]]

STT Text: Give me | Confidence: 0.9169922
2025-07-30 17:55:49,565 - INFO voice-agent - #####-----]]
-----]]

STT Text: Give me a | Confidence: 0.9248047
2025-07-30 17:55:50,567 - INFO voice-agent - -----]]
-----]]

STT Text: Give me a sentence | Confidence: 0.99902344
2025-07-30 17:55:51,542 - INFO voice-agent - #####-----]]
-----]]

STT Text: Give me a sentence on | Confidence: 0.9916992
2025-07-30 17:55:52,547 - DEBUG livekit.agents - received user transcript {"user_transcript": "Give me a sentence on", "lang
2025-07-30 17:55:52,588 - INFO voice-agent - ###-----]]
-----]]

STT Text: artificial income | Confidence: 0.7423096
2025-07-30 17:55:53,018 - DEBUG livekit.plugins.turn_detector - eou prediction {"eou_probability": 0.00021628230751957744, "

```

5

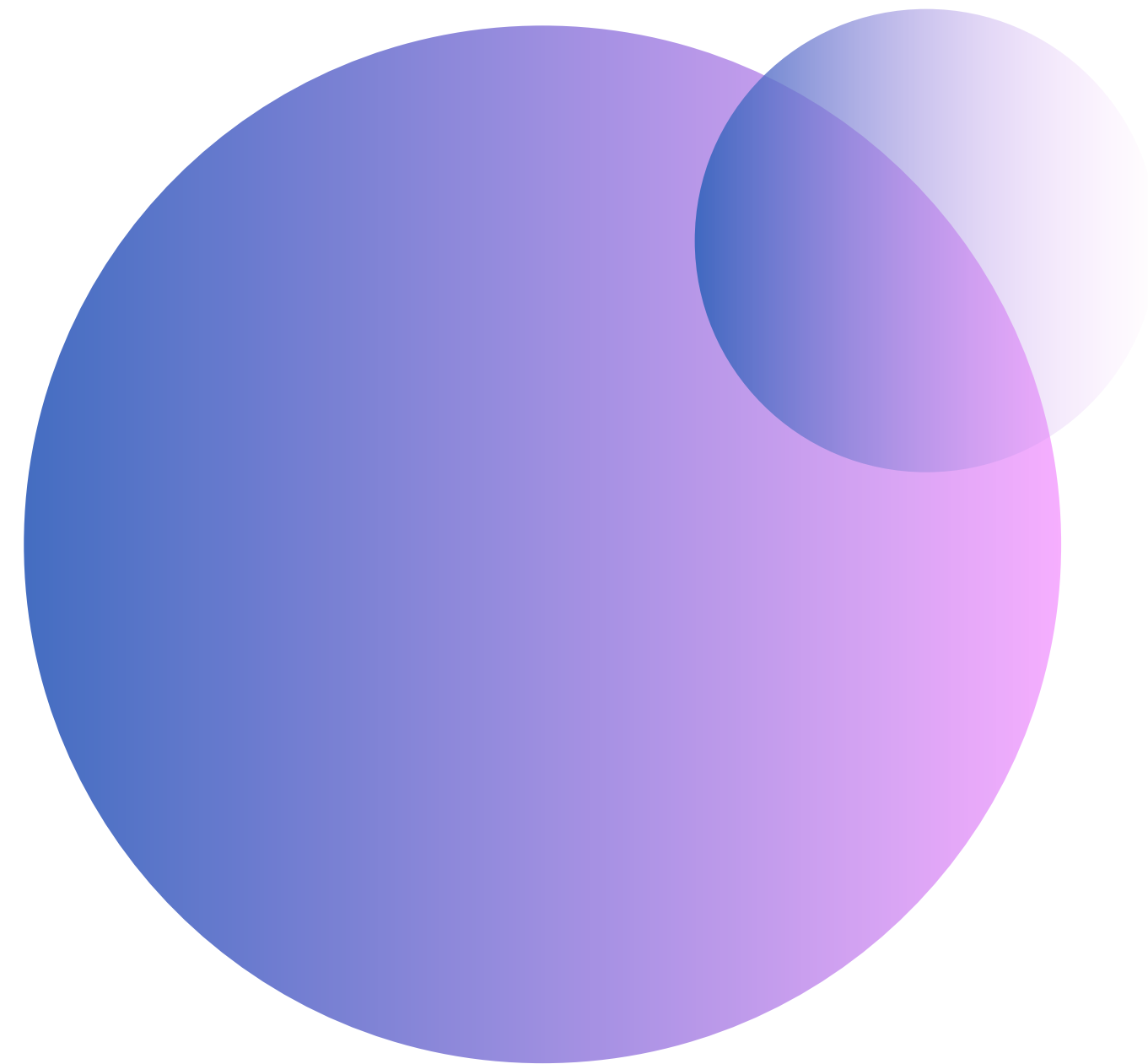
SSML-Driven TTS Enhancements

Observation:

- prosody works fine and allows pitch changes
- Spelling out words
- Having different pronunciations for words like nginx
- emphasis tag was not being responded too.

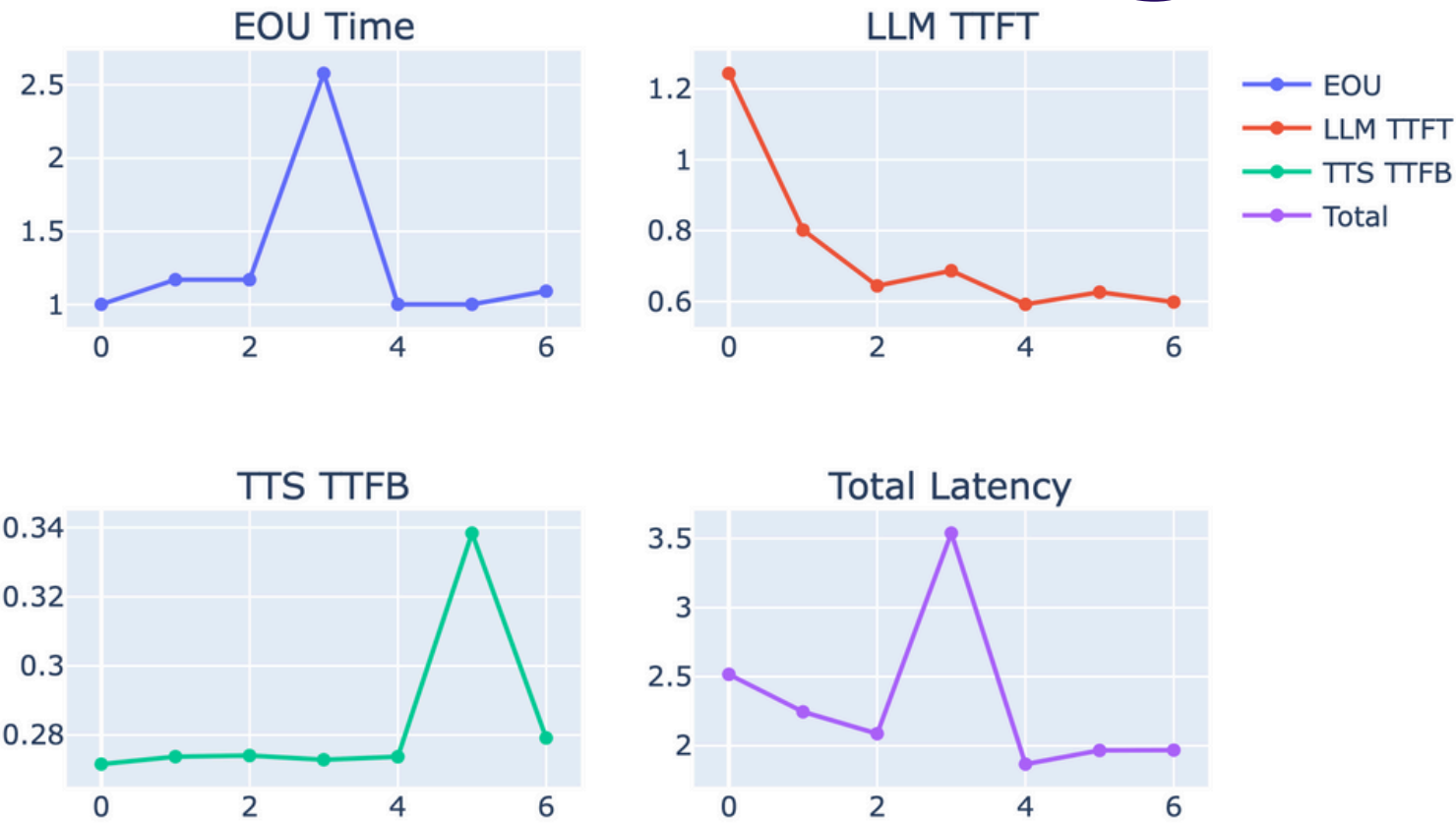
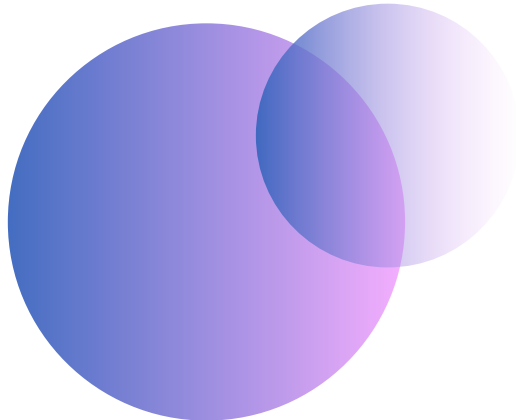
MOS Score - 3

- Rushes most part of the sentence.
- Punctuation pauses are not that great.

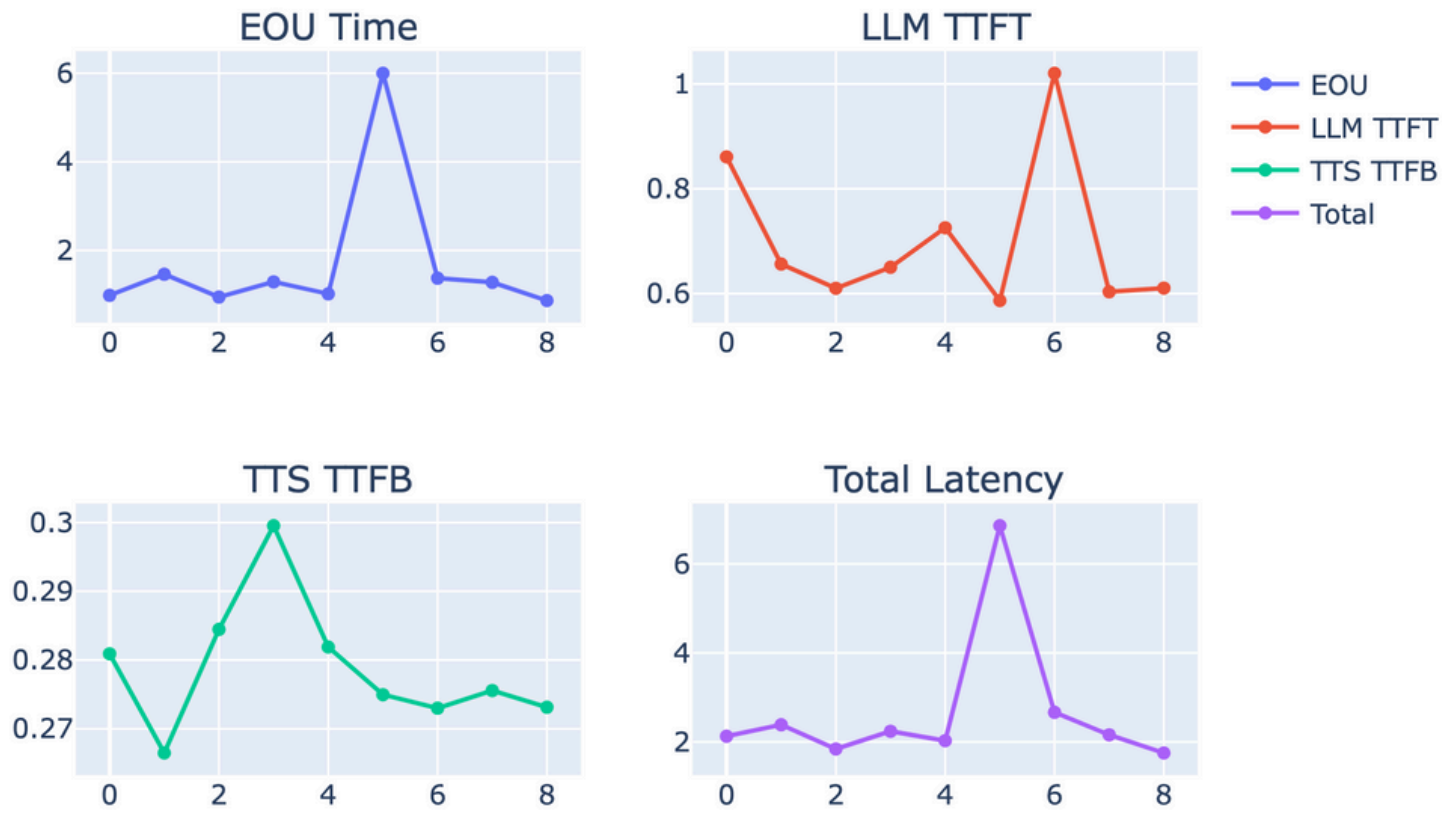


6

Latency & Quality Benchmarking



Baseline
MOS - 2

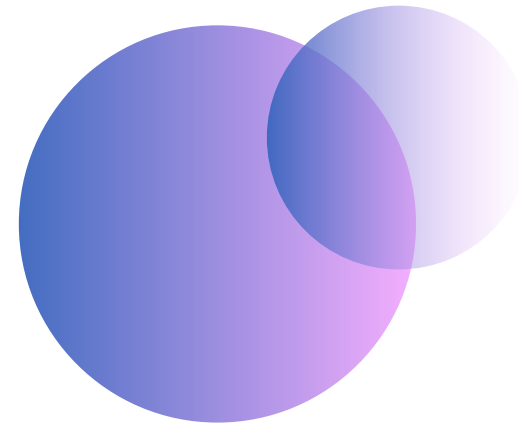


Advanced
MOS - 3



6

Latency & Quality Benchmarking



Observation:

- Baseline pipeline has fast responses from LLMs since the prompts are small and simple, but this reduces the accuracy of the response.
- Advance pipeline has a better prompt and does not generate stuff that cannot be spelt out, like having * for bold.
- TTS takes a little more time in Advance pipeline to do all the preprocessing compared to the baseline one.
- The overall latency is lower in Advance pipeline due to preemptive generation.
- Baseline model interrupts a lot while taking a pause.
- Since models are the same in both pipelines, the WER looks the same.



Code Repo

**Thank
You.**

Get In Touch With Us



+91 9284876115



ughante@gmail.com



LinkedIn