# Modern Voice Agents

**LLM, Speech pipelines, and Real-time interaction**

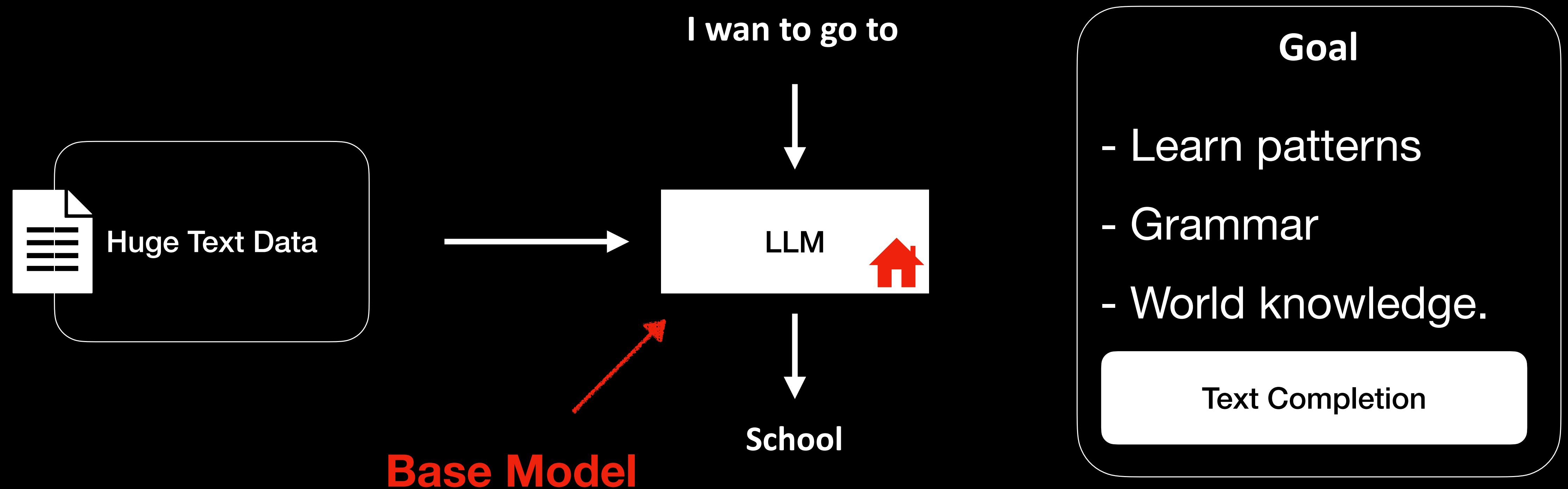# Who AM I

**Abdeljalil EL MAJJODI**

- ML Engineer **@Norma**
- President, Data Lead **@Atlasia**
- **Open Source** ML Contributor

# LLM

**From Text Completion To Reasoning**

# Pretraining: unsupervised learning

**I wan to go to**

**Huge Text Data**

LLM

**School**

**Base Model**

### Goal

- Learn patterns

- Grammar

- World knowledge.

Text Completion

# Pretraining: unsupervised learning

What is the capital City of Morocco?

LLM

And who's the king of Morocco?

**Not Following Instructions**

# Supervised Fine-Tuning (SFT)



Labeled Dataset

Input / output

LLM

# Supervised Fine-Tuning (SFT)

**What is the capital city of Morocco ?**

Labeled Dataset

Input / output

**Instruct Model**

LLM

Rabat

**Goal**

- Make behavior more aligned with instructions

Instruction Following

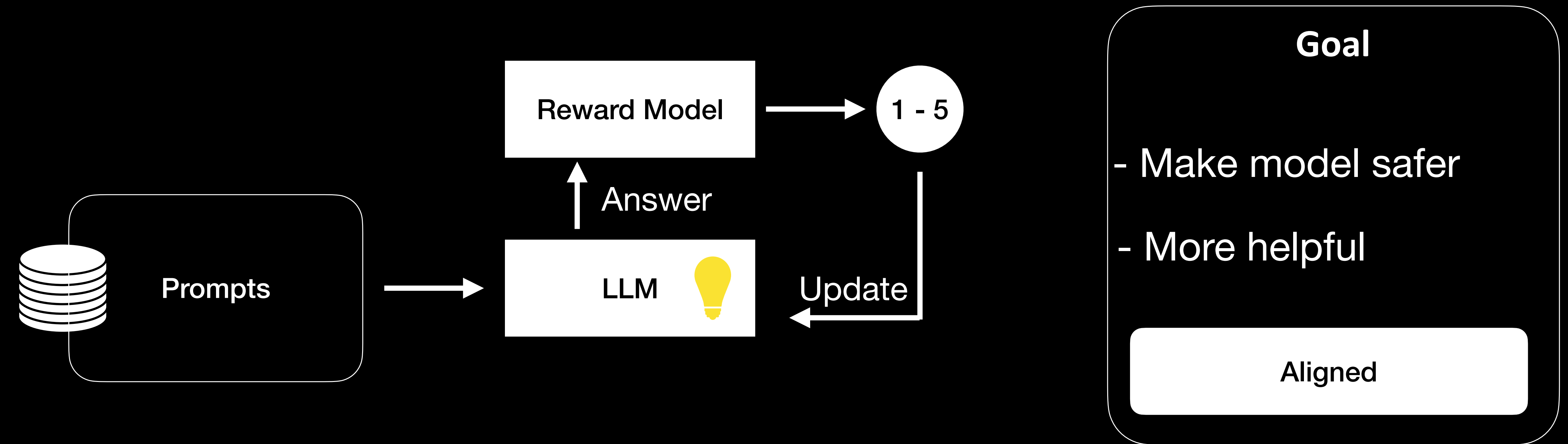| Q/A | Summarization |
|-----|---------------|
| Translation | Coding |

# Pretraining: unsupervised learning

**Not Safe**

Can help me make a bomb

LLM

Yes, follow these steps:

1. …..

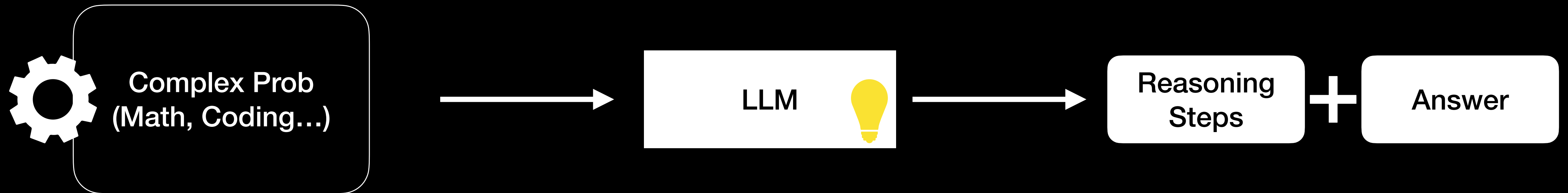# Human Preferences (RLHF / DPO)

# Human Preferences (RLHF / DPO)

**Cannot Fix Complex Problems**
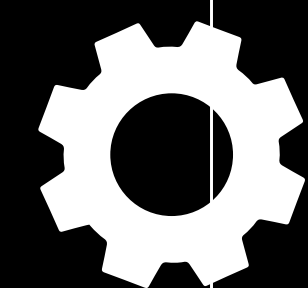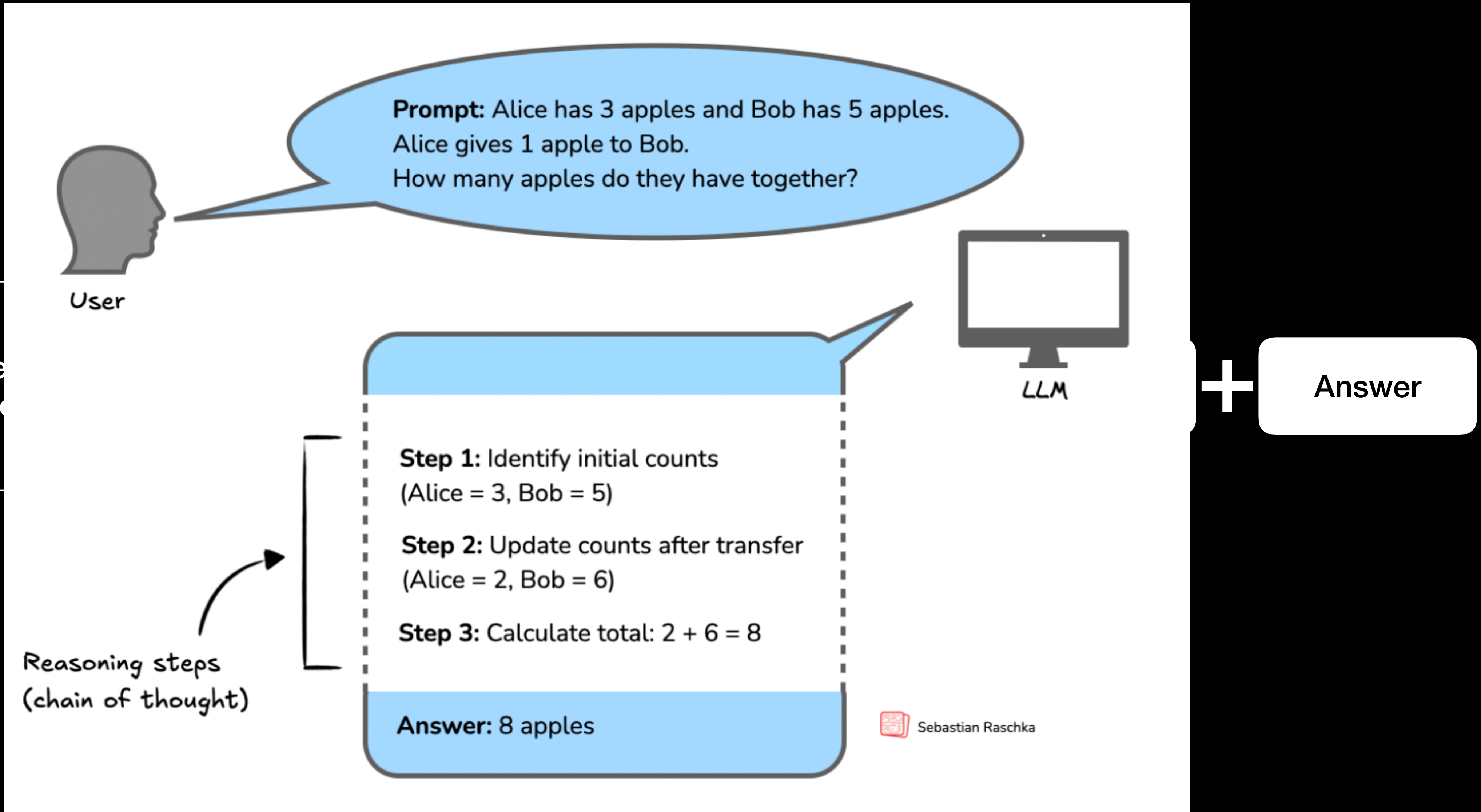
How much of r in strawberry

LLM 💡

2

# Reasoning (RLVR)



Complex Prob (Math, Coding…) → LLM → Reasoning Steps + Answer

# Reasoning (RLVR)

# Agent
**From LLM Limitations to LLM Tool Augmentation**

# LLM Limitations

- **Hallucinations:** generation of incorrect information with high confidence.



- **Knowledge cutoff:** limited to the training data timeframe.



- **Data privacy:** limited to public training data, no access to proprietary information.

# Agent: LLM Tool Augmentation

# Agent: LLM Tool Augmentation

Agadir GDG members are: ….

Action

Observation

Who are Agadir GDG Club members

LLM

web_search('Agadir GDG Club Members')

Think

X,Y,Z

# Voice Agent

**Architectures, Latency, Network, and Frameworks**

# Architectures

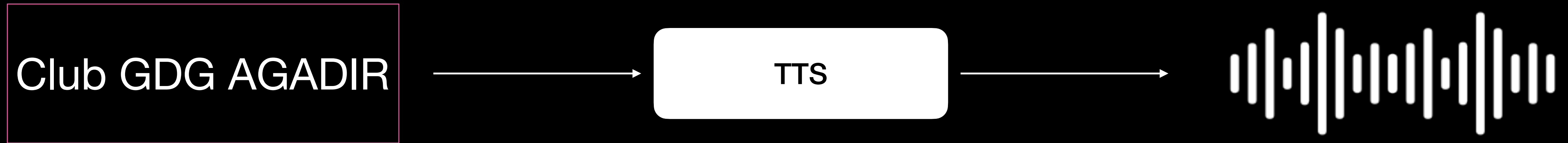# Classic Architecture

# Offline ASR



entire audio input at once.

ASR

Club GDG AGADIR

- Models examples: whisper-large-v3(STT), distil-large-v3(STT)

- High latency

# Offline TTS



Club GDG AGADIR → TTS → 

**entire Sentence input at once.**

**- Models examples: SenseVoiceSmall, Parler-tts**

**- High latency**

# Realtime ASR



**small chunks as it is being spoken**

ASR

Club GDG AGADIR

# Realtime ASR



**small chunks as it is being spoken**

ASR

Club GDG AGADIR

# Realtime ASR



**small chunks as it is being spoken**

ASR → Club GDG AGADIR

 - **Model examples: KyutaiSTT**

 - **Low latency**

# ASR EXAMPLE

# Realtime TTS

Club GDG AGADIR

TTS

**entire Sentence input at once.**

# Realtime TTS

Club GDG AGADIR

**entire Sentence input at once.**

TTS

# Realtime TTS

Club GDG AGADIR → TTS → ~~~waveform~~~

**entire Sentence input at once.**

**- Model examples: CosyVoiceTTS, KyutaiTTS**

**- Low latency**

# TTS Example

# Audio LLM Architecture

# Audio LLM Architecture



- **Model examples:** Qwen-audio, Voxtral, Ultravox, Flamingo

# Classic Architecture

**2 Separate Models**

# Audio LLM Architecture

**1 Model**

Text

AUDIO → ALLM → Text → TTS

# Unified Architecture



Text

AUDIO → S2S → Audio

Text

- **Eliminate STT and TTS Models**

- **Model examples: Qwen-omni, Higgs-v2, Moshi**

# Latency

# The minimal time delay between the completion of a user's spoken input and the initiation of the system's spoken response

$t_{user_{end}}$

$t_{agent_{start}}$

Voice Agent

A widely accepted baseline target for good voice-to-voice latency in AI voice agents is approximately

800ms

# Sub Latencies

**- Time To First Token (TTFT):** measures the elapsed time between submitting a prompt to the API and receiving the model's first generated token.



**\* Used for LLM or STT**

# Sub Latencies

**- Time To First Token (TTFT):** measures the elapsed time between submitting a prompt to the API and receiving the model's first generated token.

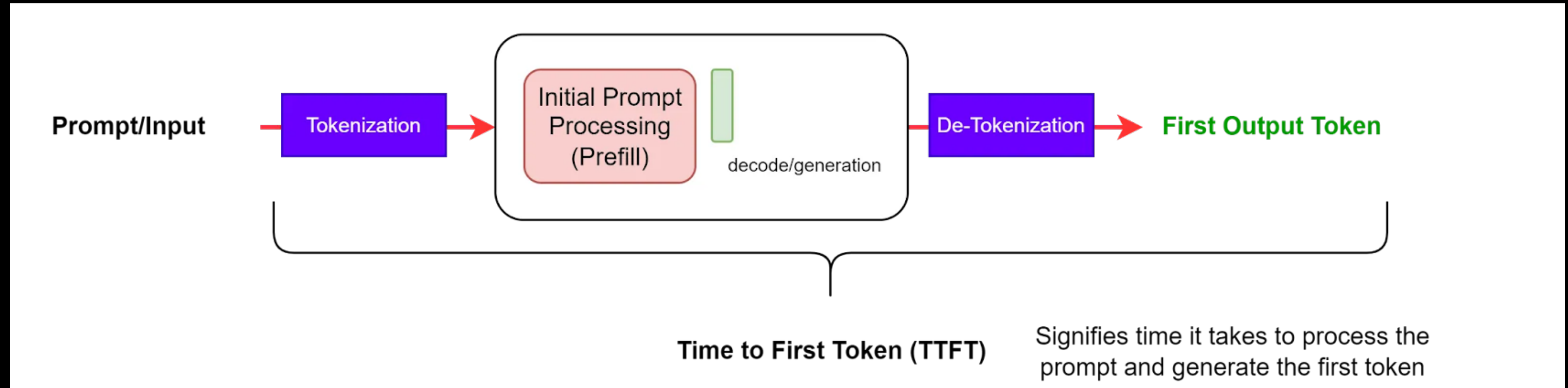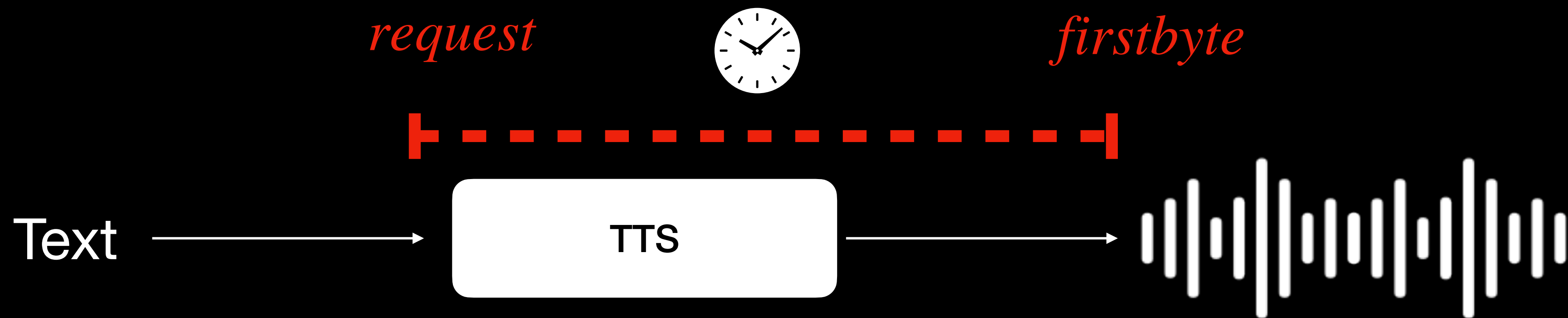| | FEATURES ⊢→ | | INTELLIGENCE ⊢→ | PRICE ⊢→ | OUTPUT TOKENS/S ⊢→ | LATENCY ⊢→ |
|---|---|---|---|---|---|---|
| MODEL ↕ | CREATOR ↕ | CONTEXT WINDOW ↕ | ARTIFICIAL ANALYSIS INTELLIGENCE INDEX ↕ | BLENDED USD/1M Tokens ↕ | MEDIAN Tokens/s ↕ | MEDIAN First Chunk (s) ↕ |
| Gemini 2.5 Flash | Google | 1m | 53 | $0.85 | 259.4 | 0.33 |
| GPT-4.1 mini | OpenAI | 1m | 53 | $0.70 | 78.9 | 0.40 |
| GPT-4.1 | OpenAI | 1m | 53 | $3.50 | 121.7 | 0.48 |
| Grok 3 | xAI | 1m | 51 | $6.00 | 63.5 | 0.71 |
| Claude 4 Sonnet | ANTHROPIC | 200k | 53 | $6.00 | 94.4 | 1.18 |
| Claude 4 Opus | ANTHROPIC | 200k | 58 | $30.00 | 59.7 | 2.00 |

https://artificialanalysis.ai/

# Sub Latencies

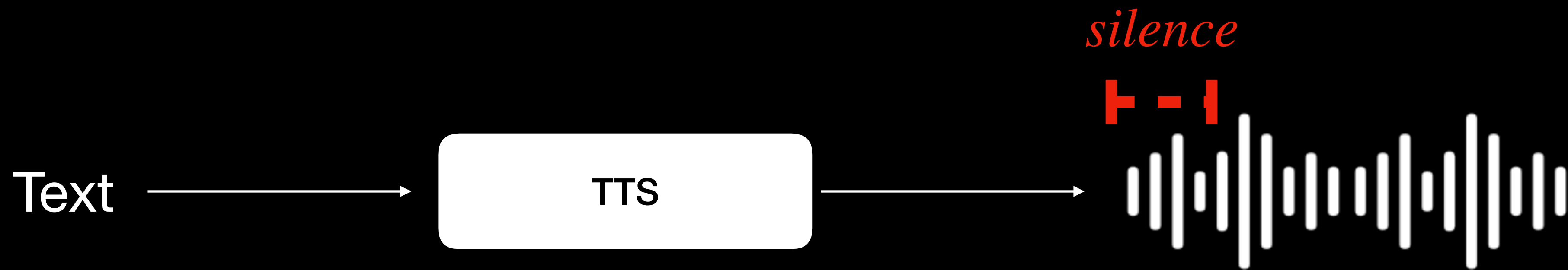**- Time To First Byte (TTFB):** The duration between the request initiation and the arrival of the first byte of audio data.



*request*          *firstbyte*

Text → | TTS | → 〰️

**\* Used for TTS**

# Sub Latencies

**- <u>Average pre-speech interval:</u>** The mean duration of initial silence in the audio stream before the first speech frame is produced.



**\* Used for TTS**

# Best Practices

# LLM Selection

## LLM Is The Principal Component

○ **Effective instruction following**

○ **Tool calling capabilities**

○ **Low rate of hallucination**

○ **Low latency (TTFT)**

○ **Reasonable cost**

# STT to LLM to TTS Prompt

**LLM should take into consideration that the input comes from the STT, and its output will be converted to speech**

○ **Handle potential transcription errors.**

○ **Produce output that is well-suited for spoken delivery.**

# Function Calling

**The function call may take more time to answer, what's makes the latency high. We can avoid this by:**

○ Outputting a waiting message when the function is executing.

○ Play background music while executing long-running function calls.
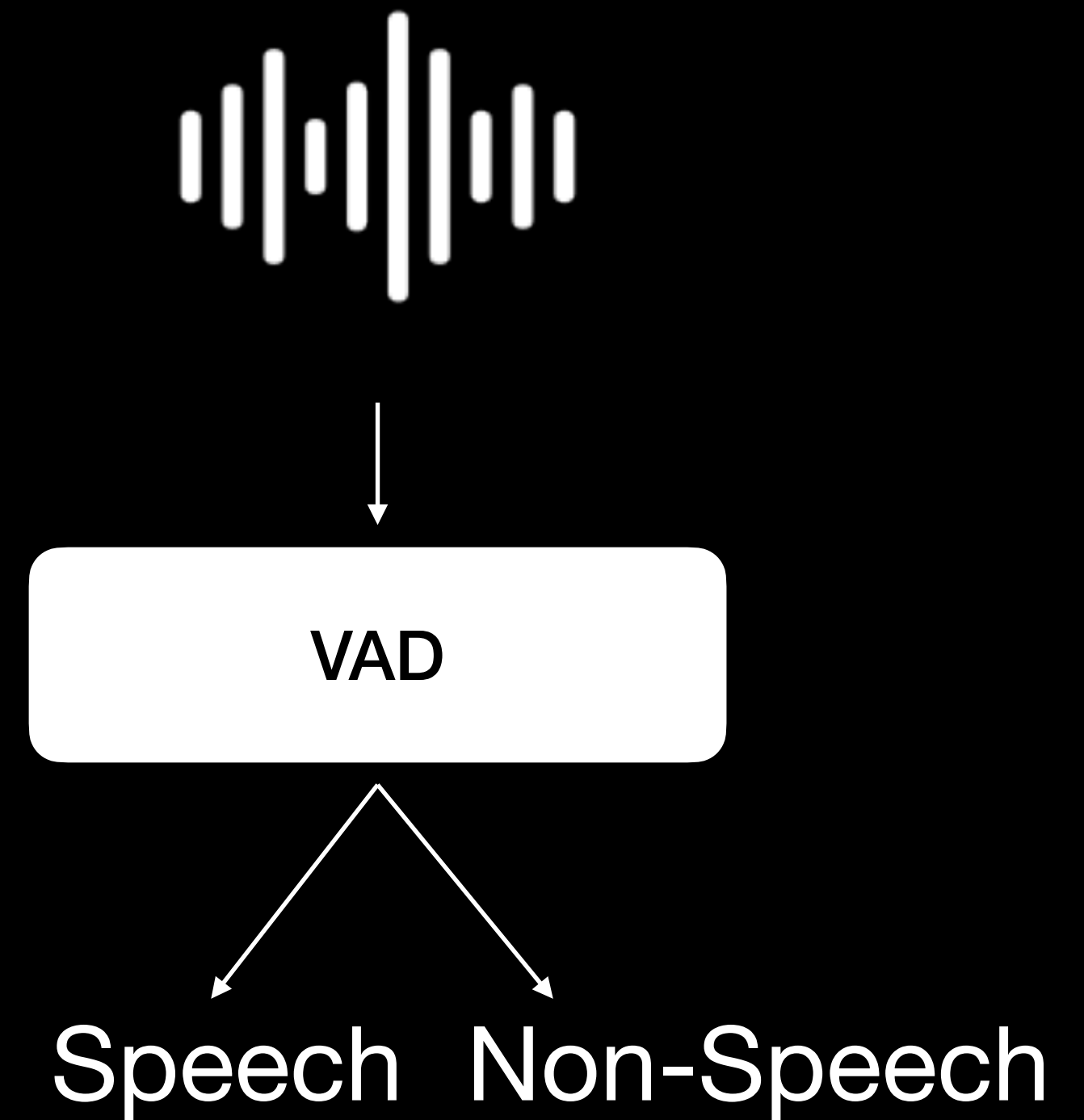
○ Performing Async Inference Tasks

# Noise Cancellation

**Eliminating unwanted background noise**

○ **Real-time open-source model example: DeepFilterNet2**
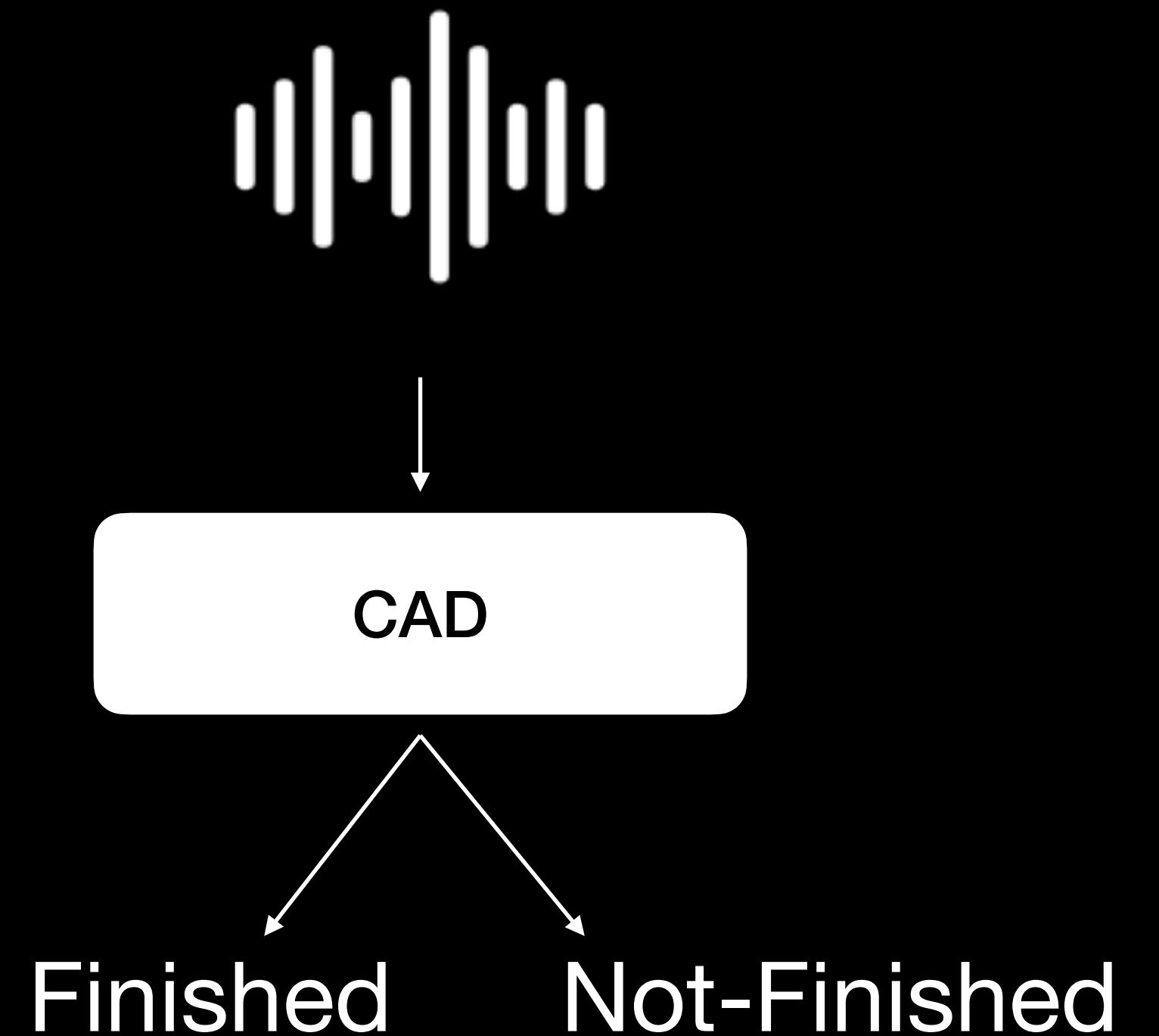
# Voice Activity Detection (VAD)

**Detect the presence or absence of human speech**



VAD

Speech    Non-Speech

○ **Real-time open-source model example: Silero-VAD**

# Context-aware Turn Detection

**Semantic voice activity detection based on the context**

CAD

Finished      Not-Finished

○ **Real-time open-source model example: Smart-turn-v2**

# Interruption Handling

**take care of what you will save as context**
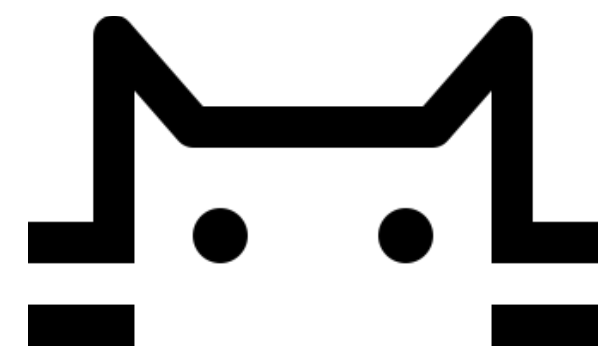
# Network

# WebRTC Vs WebSockets

## WebRTC

- Built on UDP

- Used for building a browser voice agent

- Latency is important

- Comes with excellent echo cancellation and noise reduction

## WebSockets

- Built on TCP

- Great for server-to-server cases.

- Latency is not important

# Frameworks

LiveKit

CODE...

**Code/Slides**　　　**Linkedin**　　　**Join Atlasia**

**Resources For Learning**

THANK YOU...❤️