# The Problem

- Major challenge for restaurant industry

   **High Labor Costs**

- Solution

   **Let AI solve it**

# The Solution

- An order taking bot with speech capabilities
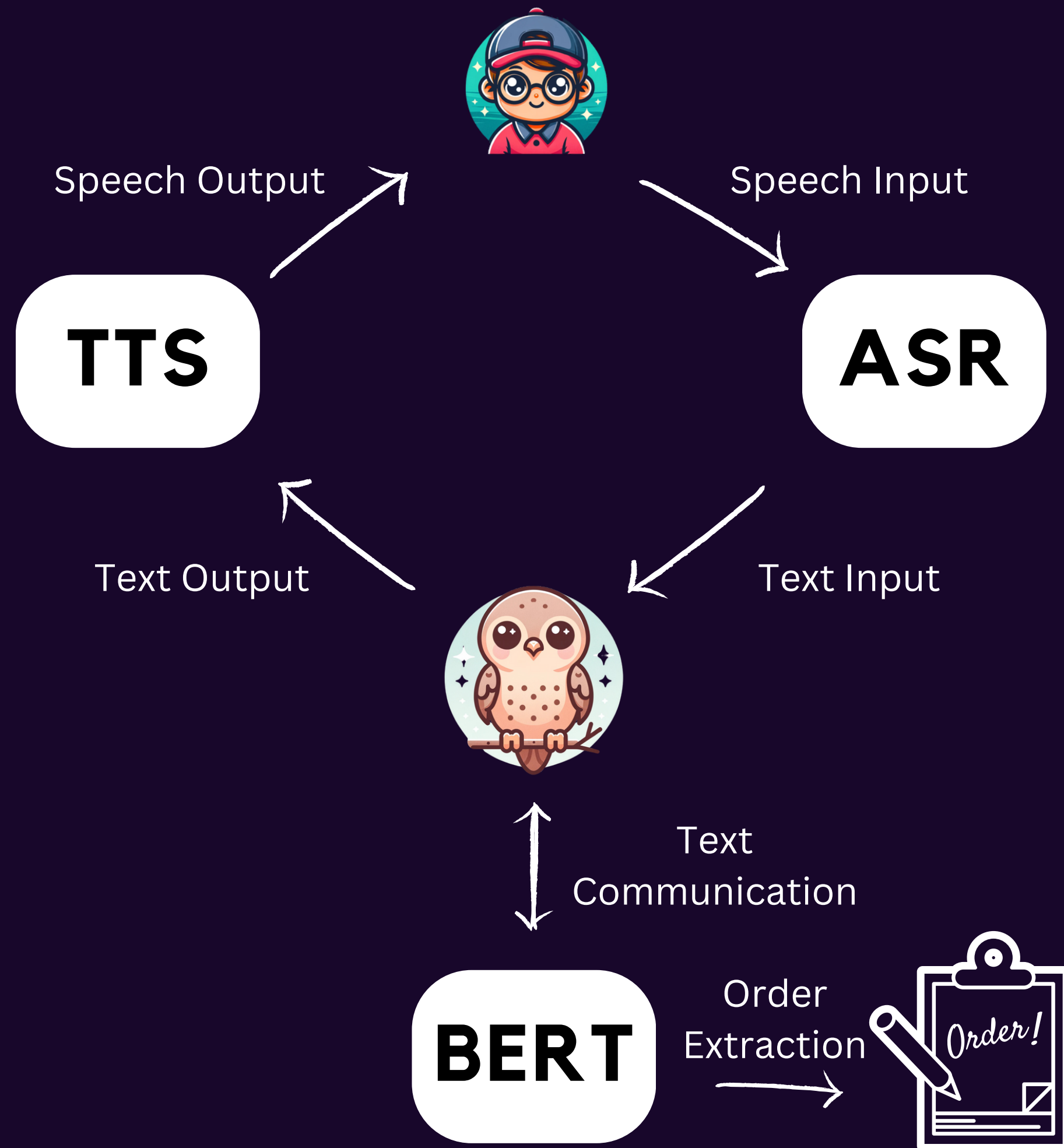
### FALCON BARISTA

- No manpower required for order taking
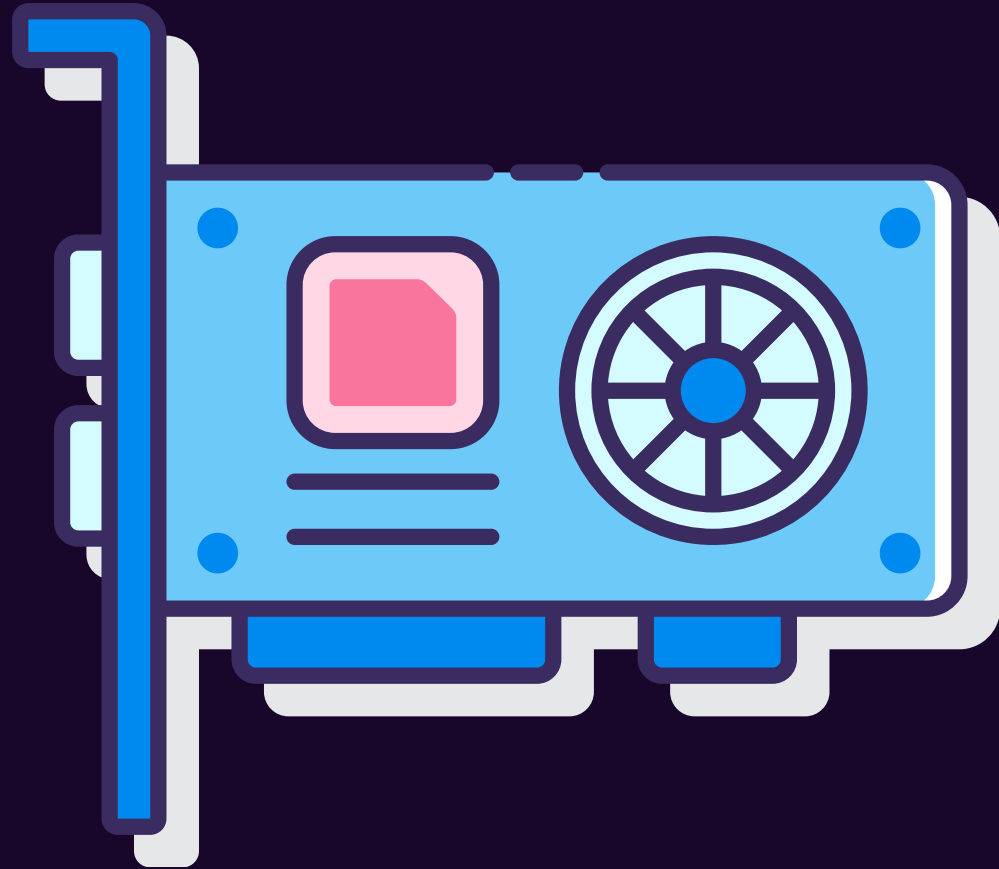- For order taking on counters, drive-throughs, and over the phone

**DEEP DREAM**

# Salient Features

- Has ability to hear and speak

- Falcon-7B LLM for conversation

- Fine-tuned BERT for information extraction

Speech Output

Speech Input

**TTS**

**ASR**

Text Output

Text Input

Text Communication

**BERT**

Order Extraction

Order!

# Our Innovation

Works on single
16 GB RAM GPU

- **Falcon Barista has minimal compute needs.**

- **Most chatbots use LLMs with over 100 billion parameters.**

- **Falcon Barista uses the fine-tuned Falcon-7B LLM with just 7 billion parameters along with smaller fine-tuned BERT models to give best performance / compute ratio**

# DEMO VIDEO

## Falcon Barista - Proof of Concept (POC)

Note: This is just a POC and has several issues

1. High Latency

2. Models are still in fine tuning phase and get confused easily
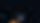
Look how cute is Falcon Barista

Chat with Falcon Barista

Audio

Record from microphone

Restart Chat

End Chat and Confirm Order

Use via API · Built with Gradio

[Video link](#)

# THANK YOU