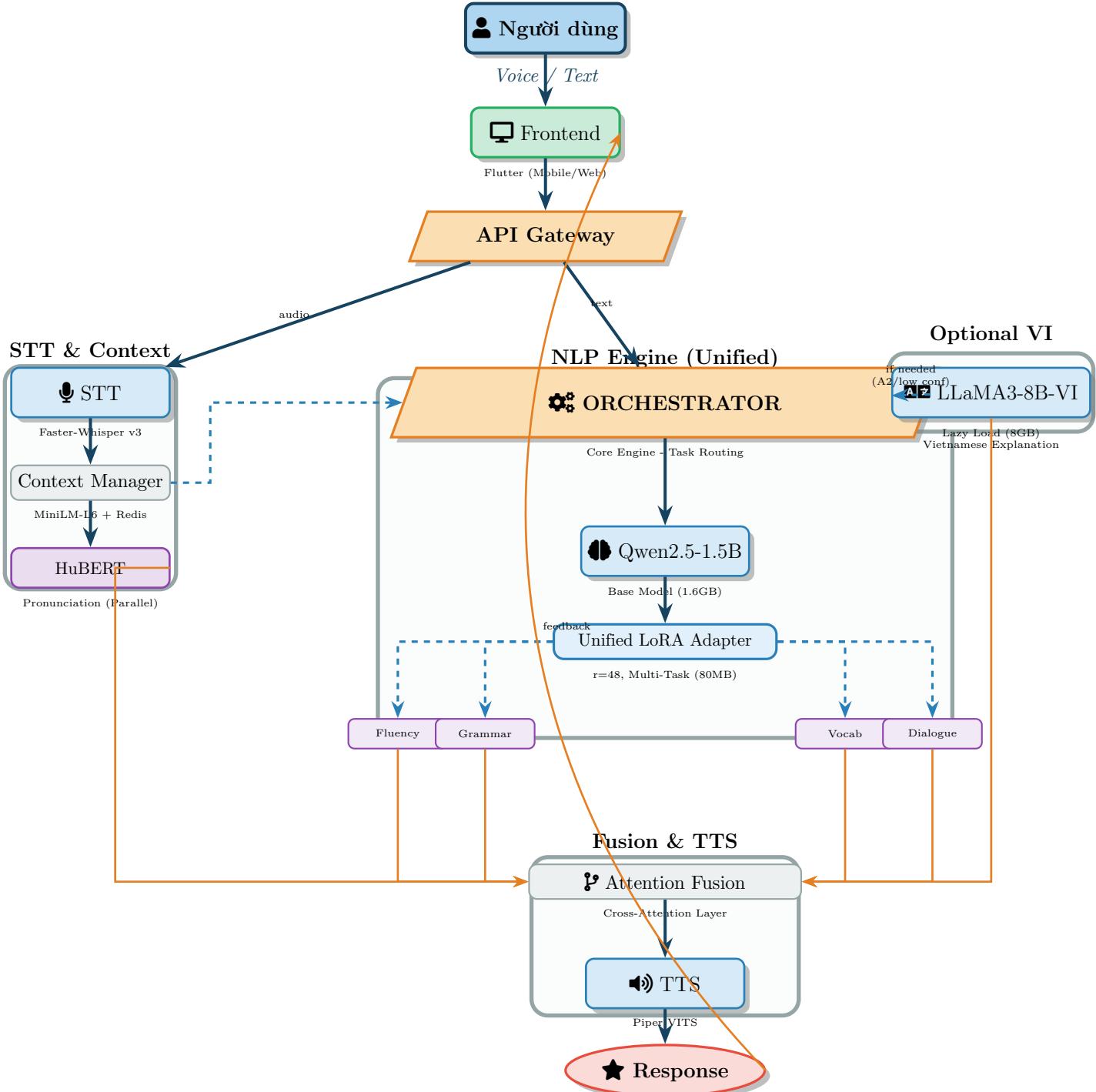


# KIẾN TRÚC HỆ THỐNG LEXILINGO v2.0

AI Hỗ trợ Học Tiếng Anh - Optimized Architecture



#### LEGEND:

- Main flow
- Data
- Feedback

**MODELS v2.0:**

- STT:** Faster-Whisper
  - 244MB, <100ms
- Context:** MiniLM-L6
  - 22MB, 15ms
- NLP:** Qwen2.5-1.5B
  - Base: 1.6GB
  - Unified LoRA: 80MB
  - All tasks: 100-150ms
- VI:** LLaMA3-8B (lazy)
  - 8GB, load on demand
- Pronunciation:** HuBERT
  - 960MB, parallel
- TTS:** Piper VITS
  - 30-60MB

#### IMPROVEMENTS:

- ✓ Latency: -56% (350ms)
- ✓ Memory: -60% (4.8GB)
- ✓ Cache hit: 45%
- ✓ Switching: <1ms
- ✓ Lazy loading enabled

# KIẾN TRÚC 5 TẦNG

Phân tầng chi tiết hệ thống

## 1. PRESENTATION LAYER

- Mobile App
- Web App
- WebSocket

## 2. API GATEWAY LAYER

- REST API
- GraphQL
- Auth/JWT

## 3. MICROSERVICES LAYER

- STT Service
  - NLP Service
  - TTS Service
- User Service — Analytics

## 4. AI/ML MODEL LAYER

- Whisper v3
  - Qwen2.5-1.5B
  - HuBERT-large
  - Piper TTS
- Fluency LoRA
- Vocab LoRA
- Grammar LoRA
- Dialogue LoRA

## 5. DATA LAYER

- PostgreSQL
- Redis Cache
- S3 Storage
- MongoDB

# UNIFIED NLP ENGINE v2.0

Kiến trúc Multi-Task Optimized

