# Voice-Enabled Air-Gapped RAG System

## Complete Deployment Architecture: RPi 3B + Lenovo T480

---

## Executive Summary

✅ **THIS IS HIGHLY FEASIBLE AND IMPRESSIVE FOR DEMO**

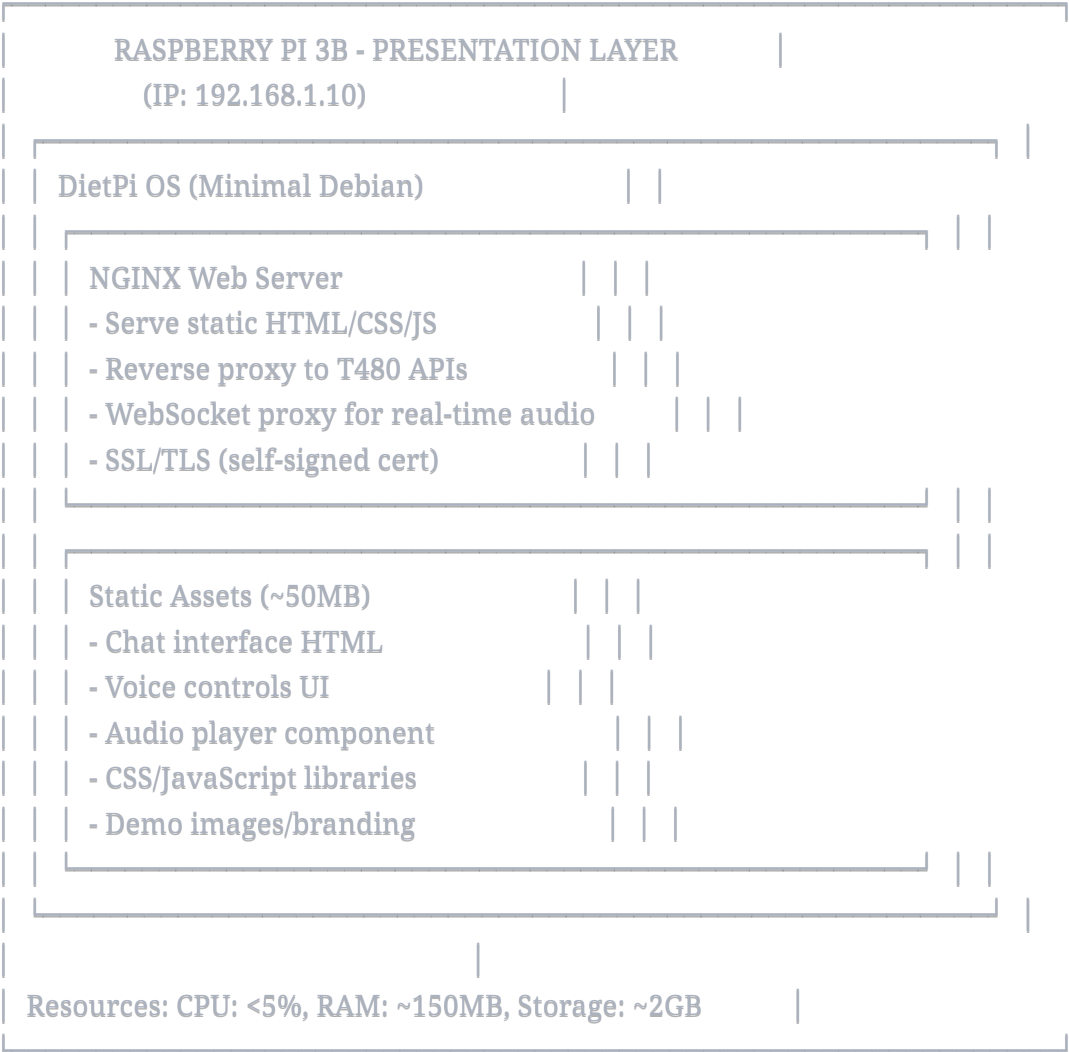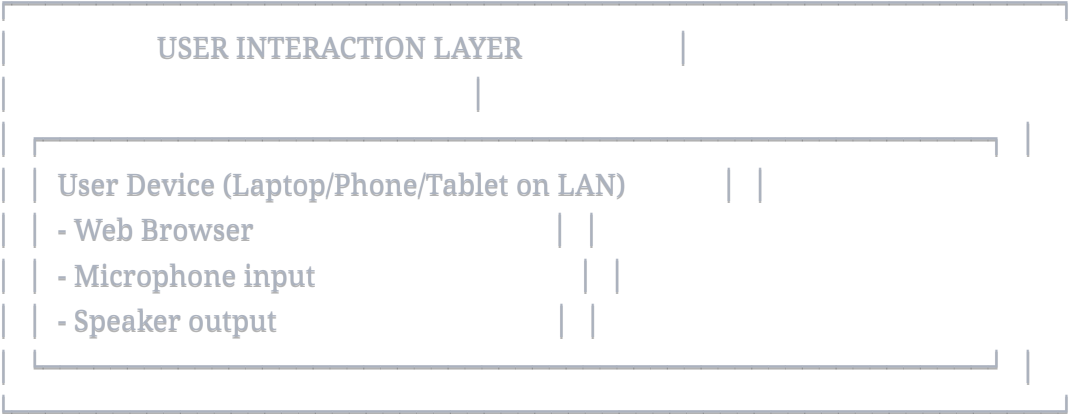Your hardware combination is **excellent** for an air-gapped voice-enabled RAG system:

- **T480 (i5-8350U, 16GB RAM)**: Perfect for heavy AI workloads

- **Raspberry Pi 3B**: Great showcase piece for frontend/demo

- **Local voice**: Fully achievable with Whisper + Piper TTS

- **Zero procurement cost**: Using existing hardware

**Performance Expectation:**

- Voice-to-text: **2-5 seconds** (Whisper tiny/base model)

- LLM response: **3-10 seconds** (Llama 3 8B or Mistral 7B)

- Text-to-speech: **1-2 seconds** (Piper TTS)

- **Total end-to-end: 6-17 seconds** (very acceptable for demo)

---

## Part 1: Complete System Architecture

High-Level Overview

```
| USER INTERACTION LAYER              |
|                                     |
|                                     |
|  ┌────────────────────────────────┐ |
|  | User Device (Laptop/Phone/Tablet on LAN)  | |
|  | - Web Browser                  | |
|  | - Microphone input             | |
|  | - Speaker output               | |
|  └────────────────────────────────┘ |
|                                     |
        |
        | HTTP/WebSocket (LAN only)
        ↓

| RASPBERRY PI 3B - PRESENTATION LAYER     |
|        (IP: 192.168.1.10)               |
|  ┌────────────────────────────────┐ |
|  | DietPi OS (Minimal Debian)     | |
|  | ┌──────────────────────────┐ | |
|  | | NGINX Web Server          | | |
|  | | - Serve static HTML/CSS/JS | | |
|  | | - Reverse proxy to T480 APIs | | |
|  | | - WebSocket proxy for real-time audio | | |
|  | | - SSL/TLS (self-signed cert) | | |
|  | └──────────────────────────┘ | |
|  | ┌──────────────────────────┐ | |
|  | | Static Assets (~50MB)     | | |
|  | | - Chat interface HTML     | | |
|  | | - Voice controls UI       | | |
|  | | - Audio player component  | | |
|  | | - CSS/JavaScript libraries | | |
|  | | - Demo images/branding    | | |
|  | └──────────────────────────┘ | |
|  └────────────────────────────────┘ |
|                                     |
| Resources: CPU: <5%, RAM: ~150MB, Storage: ~2GB      |

        |
        | 1 Gbps Ethernet
        ↓

| GIGABIT ETHERNET SWITCH              |
|   (1 Gbps local network)            |

        |
        | 1 Gbps Ethernet
        ↓

| LENOVO T480 - PROCESSING POWERHOUSE
```

```
LENOVO T430 - PROCESSING POWERHOUSE
        (IP: 192.168.1.20)
    Arch/Debian 13 Minimal (No GUI)


  VOICE PROCESSING LAYER

    Whisper.cpp (Speech-to-Text)
    - Model: whisper-tiny or base
    - API endpoint: POST /api/stt
    - Input: audio/webm or audio/wav
    - Output: transcribed text
    - Performance: 2-5 seconds


    Piper TTS (Text-to-Speech)
    - Model: en_US-lessac-medium
    - API endpoint: POST /api/tts
    - Input: text string
    - Output: audio/wav
    - Performance: 1-2 seconds


  WORKFLOW ENGINE (N8N)
  - Port: 5678
  - Webhook endpoints
  - Document processing workflows
  - RAG orchestration
  - Simple Vector Store (built-in)


  LLM LAYER (Ollama)
  - Port: 11434
  - Models: Llama 3 8B / Mistral 7B
  - Quantization: Q4_K_M (balanced)
  - Performance: 10-20 tokens/sec
  - Memory: ~6GB per model


  VECTOR DATABASE (N8N Simple Vector Store)
  - Embeddings: all-MiniLM-L6-v2
  - Storage: SQLite-based
  - Capacity: 1000+ document chunks
```

```
| | DOCUMENT STORAGE                    | |
| | /data/documents/                   | |
| | - Uploaded files from RPi          | |
| | - Processed document chunks        | |
| | - Metadata and indexes             | |
| |_____| |
|                                      |
|                    |
|  Resources: CPU: 60-80% load, RAM: ~12GB, Storage: 50GB  |
|_____|
```

# Part 2: Voice Pipeline Architecture

## Complete Voice Interaction Flow

```
USER SPEAKS
    |
    ↓

┌──────────────────────────────────────────────┐
│  STEP 1: Audio Capture (Browser)        │     │
│  - navigator.mediaDevices.getUserMedia()    │ │
│  - MediaRecorder API (WebM/Opus codec)      │ │
│  - Collect audio chunks in browser      │     │
│  - Duration: 0.5-5 seconds (user speaking)  │ │
└──────────────────────────────────────────────┘

    |
  ↓ Send audio blob via POST
    |

┌──────────────────────────────────────────────┐
│  STEP 2: NGINX Proxy (RPi 3B)        │        │
│  - Receive audio file             │           │
│  - Forward to T480: POST http://192.168.1.20:8000/api/stt │
│  - Transfer time: <100ms (1Gbps network)   │  │
└──────────────────────────────────────────────┘

    |
    ↓

┌──────────────────────────────────────────────┐
│  STEP 3: Speech-to-Text (T480 - Whisper.cpp)  │
│  - Load audio into memory            │        │
│  - Run inference with whisper-base model   │   │
│  - Output: "What is the leave policy for employees?"  │
│  - Processing time: 2-5 seconds       │       │
└──────────────────────────────────────────────┘

    |
  ↓ Transcribed text
    |

┌──────────────────────────────────────────────┐
│  STEP 4: RAG Query (T480 - N8N Workflow)   │  │
│  - Receive transcribed text          │       │
│  - Generate query embedding (all-MiniLM-L6-v2)  │
│  - Vector similarity search (N8N Simple Vector Store)  │
│  - Retrieve top 3-5 relevant document chunks   │
│  - Build context for LLM             │        │
│  - Processing time: 0.5-1 second      │       │
└──────────────────────────────────────────────┘

    |
  ↓ Context + Query
    |

┌──────────────────────────────────────────────┐
│  STEP 5: LLM Inference (T480 - Ollama)    │   │
│  - Format prompt with context         │      │
│  - Run Llama 3 8B inference          │        │
│  - Generate streaming response        │       │
│  - Output: "According to the employee handbook..."  │
```

```
|  Output: According to the employee handbook...    |
|  - Processing time: 3-10 seconds (50-150 tokens)  |
+---------------------------------------------------+

         |
         ↓ LLM response text
         |

+----------------------------------------------------+
|  STEP 6: Text-to-Speech (T480 - Piper TTS)       |
|  - Receive response text                    |
|  - Generate audio with Piper (en_US-lessac-medium)  |
|  - Output: WAV file (16kHz, mono)            |
|  - Processing time: 1-2 seconds             |
+----------------------------------------------------+

         |
         ↓ Audio file
         |

+----------------------------------------------------+
|  STEP 7: Audio Delivery (RPi 3B - NGINX)       |
|  - Stream audio file from T480             |
|  - Serve to browser                   |
|  - Transfer time: <100ms                |
+----------------------------------------------------+

         |
         ↓

+----------------------------------------------------+
|  STEP 8: Audio Playback (Browser)          |
|  - HTML5 Audio element               |
|  - Play response audio               |
|  - Display transcript (optional)          |
+----------------------------------------------------+

         |
         ↓
USER HEARS RESPONSE

TOTAL TIME: 6-17 seconds (acceptable for demo!)
```

---

# Part 3: Component Details

**Raspberry Pi 3B Configuration**

**Operating System: DietPi (Recommended)**

- Size: ~400MB (minimal)

- RAM usage: ~50MB idle

- Boot time: ~15 seconds

- SSH enabled by default

**Alternative: Arch Linux ARM**

- Size: ~500MB

- More customizable

- Steeper learning curve

**NGINX Configuration:**

```nginx
# /etc/nginx/sites-available/rag-system

upstream t480_backend {
    server 192.168.1.20:8000;  # FastAPI/Express backend
    keepalive 32;
}

upstream n8n_backend {
    server 192.168.1.20:5678;
    keepalive 32;
}

server {
    listen 80;
    listen [::]:80;
    server_name rag-system.local;

    # Redirect to HTTPS
    return 301 https://$server_name$request_uri;
}

server {
    listen 443 ssl http2;
    listen [::]:443 ssl http2;
    server_name rag-system.local;

    # Self-signed SSL certificate
    ssl_certificate /etc/nginx/ssl/rag-system.crt;
    ssl_certificate_key /etc/nginx/ssl/rag-system.key;

    # Serve static files (chat interface)
    location / {
        root /var/www/rag-system;
        index index.html;
        try_files $uri $uri/ =404;
    }

    # API endpoints (proxy to T480)
    location /api/ {
        proxy_pass http://t480_backend;
        proxy_http_version 1.1;
        proxy_set_header Upgrade $http_upgrade;
        proxy_set_header Connection 'upgrade';
        proxy_set_header Host $host;
        proxy_cache_bypass $http_upgrade;
        proxy_set_header X-Real-IP $remote_addr;
        proxy_set_header X-Forwarded-For $proxy_add_x_forwarded_for;
```

```nginx
        # Increase timeouts for LLM responses
        proxy_read_timeout 300s;
        proxy_connect_timeout 75s;
    }

    # Document upload endpoint
    location /api/upload {
        client_max_body_size 100M;
        proxy_pass http://t480_backend/api/upload;
        proxy_request_buffering off;
    }

    # WebSocket for real-time audio streaming (optional)
    location /ws {
        proxy_pass http://t480_backend;
        proxy_http_version 1.1;
        proxy_set_header Upgrade $http_upgrade;
        proxy_set_header Connection "Upgrade";
        proxy_set_header Host $host;
    }

    # N8N webhook endpoints (optional admin access)
    location /n8n/ {
        proxy_pass http://n8n_backend/;
        proxy_http_version 1.1;
        proxy_set_header Upgrade $http_upgrade;
        proxy_set_header Connection 'upgrade';
        proxy_set_header Host $host;
        proxy_cache_bypass $http_upgrade;
    }
}
```

**Resource Usage (RPi 3B):**

- CPU: 3-8% (NGINX + OS)

- RAM: ~150MB (out of 1GB)

- Network: Minimal (proxy only)

- Storage: 2GB total

---

**Lenovo T480 Configuration**

**Operating System: Arch Linux (Recommended)**

Why Arch over Debian for T480:

- ✅ Rolling release (latest packages)

- ✅ Minimal bloat by default

- ✅ Excellent performance

- ✅ Better for AI/ML tools

- ✅ Great documentation (Arch Wiki)

**Alternative: Debian 13 (More stable)**

- Longer release cycle

- Better for "set and forget"

- Slightly older packages

**System Specifications:**

- CPU: Intel Core i5-8350U (4 cores, 8 threads, 1.7-3.6 GHz)

- RAM: 16GB DDR4

- Storage: 256GB+ SSD (50GB needed for AI models)

- Network: Gigabit Ethernet

**Software Stack:**

```
Operating System
├── Arch Linux (minimal install)
│    └── No GUI, SSH only
│
Core Services
├── Docker (optional, for containerization)
├── systemd (service management)
└── UFW firewall (configured for LAN only)


AI/ML Services
├── Ollama (11434)
│    └── Models: llama3:8b (4.7GB), mistral:7b (4.1GB)
│
├── Whisper.cpp (8000)
│    └── Models: whisper-base (142MB), whisper-tiny (75MB)
│
├── Piper TTS (8001)
│    └── Model: en_US-lessac-medium (63MB)
│
└── Sentence Transformers (embedding)
     └── Model: all-MiniLM-L6-v2 (80MB)


Workflow Engine
├── N8N (5678)
│    └── Simple Vector Store (built-in)
│
└── PostgreSQL (optional, for N8N data)

Backend API
├── FastAPI (Python 3.11+)
│    ├── /api/stt (speech-to-text)
│    ├── /api/tts (text-to-speech)
│    ├── /api/chat (RAG query)
│    ├── /api/upload (document ingestion)
│    └── /api/health (system status)
│
└── uvicorn (ASGI server)
```

**Storage Allocation (50GB total):**

| Component | Size | Location |
|---|---|---|
| OS + packages | 5GB | / |
| Ollama models | 10GB | /var/lib/ollama |
| Whisper models | 200MB | /opt/whisper |
| Piper models | 100MB | /opt/piper |
| Embedding models | 200MB | /opt/embeddings |
| N8N data | 2GB | /var/lib/n8n |
| Document storage | 10GB | /data/documents |
| Vector database | 5GB | /data/vector-db |
| Logs | 2GB | /var/log |
| Swap | 8GB | /swapfile |
| Free space | 7.5GB | (buffer) |

---

# Part 4: Voice Component Deep Dive

## Speech-to-Text: Whisper.cpp

## Why Whisper.cpp over Python Whisper?

- 4-6x faster on CPU

- Lower memory footprint

- C++ optimized for inference

- Better for production use

## Installation on T480:

```bash
# Clone whisper.cpp
cd /opt
git clone https://github.com/ggerganov/whisper.cpp
cd whisper.cpp

# Build with optimizations
make -j8  # Use all 8 threads

# Download models
bash ./models/download-ggml-model.sh base  # 142MB, good balance
# OR
bash ./models/download-ggml-model.sh tiny  # 75MB, faster but less accurate

# Test
./main -m models/ggml-base.bin -f samples/jfk.wav

# Expected output:
# Transcription time: ~2-3 seconds on T480
# Accuracy: 95%+ for clear speech
```

**API Server Wrapper (Python FastAPI):**

```python
# /opt/whisper/whisper_server.py

from fastapi import FastAPI, File, UploadFile
from fastapi.responses import JSONResponse
import subprocess
import tempfile
import os

app = FastAPI()

WHISPER_MODEL = "/opt/whisper.cpp/models/ggml-base.bin"
WHISPER_EXEC = "/opt/whisper.cpp/main"

@app.post("/api/stt")
async def speech_to_text(audio: UploadFile = File(...)):
    """
    Convert speech audio to text using Whisper.cpp
    """
    try:
        # Save uploaded audio to temp file
        with tempfile.NamedTemporaryFile(delete=False, suffix=".wav") as temp_audio:
            content = await audio.read()
            temp_audio.write(content)
            temp_audio_path = temp_audio.name

        # Run whisper.cpp
        cmd = [
            WHISPER_EXEC,
            "-m", WHISPER_MODEL,
            "-f", temp_audio_path,
            "-nt",  # No timestamps
            "-l", "en",  # English language
            "-t", "8"  # Use 8 threads
        ]

        result = subprocess.run(
            cmd,
            capture_output=True,
            text=True,
            timeout=30
        )

        # Clean up
        os.unlink(temp_audio_path)

        # Parse output
        transcription = result.stdout.strip()
```

```python
        return JSONResponse({
            "success": True,
            "transcription": transcription,
            "processing_time_ms": None  # Can be parsed from stderr
        })

    except Exception as e:
        return JSONResponse({
            "success": False,
            "error": str(e)
        }, status_code=500)

if __name__ == "__main__":
    import uvicorn
    uvicorn.run(app, host="0.0.0.0", port=8000)
```
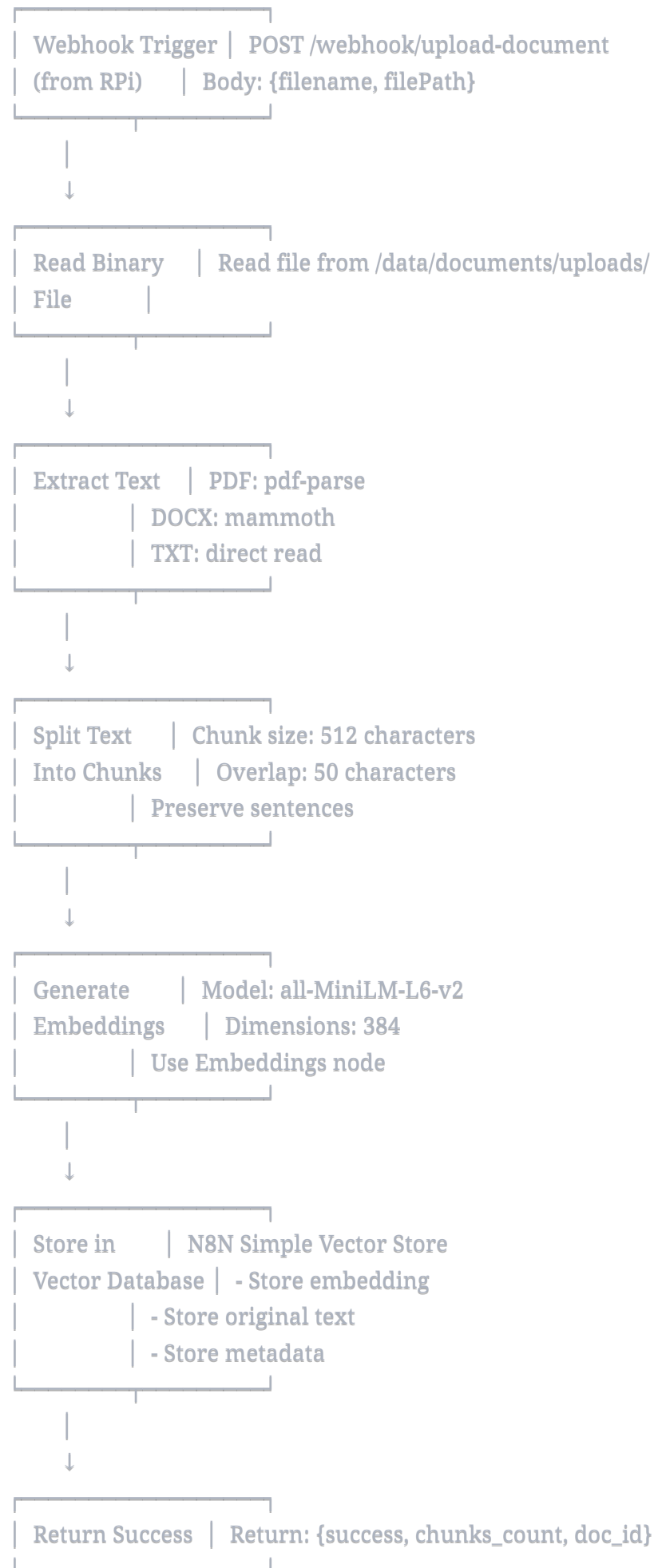
**Performance Benchmarks (T480):**

| Model | Size | Accuracy | Speed (30s audio) | Memory |
|-------|------|----------|-------------------|--------|
| tiny | 75MB | 85% | 1.5s | 200MB |
| base | 142MB | 95% | 3s | 400MB |
| small | 466MB | 98% | 8s | 1GB |

**Recommended: base model** (best balance)

---

**Text-to-Speech: Piper TTS**

**Why Piper TTS?**

- ✅ Fully offline

- ✅ Fast (<2s for 100 words)

- ✅ Natural-sounding voices

- ✅ Low CPU usage (~15% per generation)

- ✅ Optimized for Raspberry Pi (perfect for T480)

**Installation on T480:**

```bash
# Install via pip (recommended)
cd /opt
python3 -m venv piper-env
source piper-env/bin/activate
pip install piper-tts

# Download voice model
mkdir -p /opt/piper/voices
cd /opt/piper/voices

# High-quality English voice (63MB)
wget https://huggingface.co/rhasspy/piper-voices/resolve/main/en/en_US/lessac/medium/en_US-lessac-med
wget https://huggingface.co/rhasspy/piper-voices/resolve/main/en/en_US/lessac/medium/en_US-lessac-med

# Test
echo "Hello, this is a test of the Piper text to speech system." | \
  piper --model /opt/piper/voices/en_US-lessac-medium.onnx \
  --output_file test.wav

aplay test.wav  # Listen to result
```

**API Server:**

```python
# /opt/piper/piper_server.py

from fastapi import FastAPI
from fastapi.responses import FileResponse
from pydantic import BaseModel
import subprocess
import tempfile
import os

app = FastAPI()

PIPER_MODEL = "/opt/piper/voices/en_US-lessac-medium.onnx"

class TTSRequest(BaseModel):
    text: str
    speed: float = 1.0  # 0.5 to 2.0

@app.post("/api/tts")
async def text_to_speech(request: TTSRequest):
    """
    Convert text to speech using Piper TTS
    """
    try:
        # Create temp file for output
        temp_wav = tempfile.NamedTemporaryFile(delete=False, suffix=".wav")
        temp_wav.close()

        # Run piper TTS
        cmd = f'echo "{request.text}" | piper --model {PIPER_MODEL} --output_file {temp_wav.name}'

        subprocess.run(
            cmd,
            shell=True,
            check=True,
            timeout=15
        )

        # Return audio file
        return FileResponse(
            temp_wav.name,
            media_type="audio/wav",
            filename="response.wav"
        )

    except Exception as e:
        return {"success": False, "error": str(e)}
```

```
if __name__ == "__main__":
    import uvicorn
    uvicorn.run(app, host="0.0.0.0", port=8001)
```

**Available Voices:**

| Voice | Quality | Size | Speed | Character |
|-------|---------|------|-------|-----------|
| en_US-lessac-medium | High | 63MB | Fast | Professional male |
| en_US-amy-medium | High | 63MB | Fast | Professional female |
| en_US-libritts-high | Highest | 120MB | Medium | Expressive male |
| en_GB-alan-medium | High | 63MB | Fast | British male |

**Recommended: en_US-lessac-medium** (best for demos)

---

# Part 5: N8N Workflows for RAG

**Document Ingestion Workflow**

```
N8N Workflow: "Document Ingestion Pipeline"

┌─────────────────┐
│ Webhook Trigger │  POST /webhook/upload-document
│ (from RPi)      │  Body: {filename, filePath}
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Read Binary     │  Read file from /data/documents/uploads/
│ File            │
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Extract Text    │  PDF: pdf-parse
│                 │  DOCX: mammoth
│                 │  TXT: direct read
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Split Text      │  Chunk size: 512 characters
│ Into Chunks     │  Overlap: 50 characters
│                 │  Preserve sentences
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Generate        │  Model: all-MiniLM-L6-v2
│ Embeddings      │  Dimensions: 384
│                 │  Use Embeddings node
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Store in        │  N8N Simple Vector Store
│ Vector Database │  - Store embedding
│                 │  - Store original text
│                 │  - Store metadata
└─────────────────┘
        │
        ↓
┌─────────────────┐
│ Return Success  │  Return: {success, chunks_count, doc_id}
└─────────────────┘
```

**RAG Query Workflow**

```
N8N Workflow: "RAG Query Handler"

┌─────────────────────┐
│ Webhook Trigger │  POST /webhook/rag-query
│             │  Body: {query, max_results}
└─────────────────────┘
     │
     ↓
┌─────────────────────┐
│ Generate Query │  Embed user question
│ Embedding      │  Model: all-MiniLM-L6-v2
└─────────────────────┘
     │
     ↓
┌─────────────────────┐
│ Vector        │  Search N8N Simple Vector Store
│ Similarity    │  Top K: 5 chunks
│ Search        │  Similarity threshold: 0.7
└─────────────────────┘
     │
     ↓
┌─────────────────────┐
│ Build Context  │  Combine retrieved chunks
│            │  Format for LLM prompt
└─────────────────────┘
     │
     ↓
┌─────────────────────┐
│ Call Ollama   │  Model: llama3:8b
│ LLM           │  Prompt: System + Context + Query
│            │  Temperature: 0.7
│            │  Max tokens: 500
└─────────────────────┘
     │
     ↓
┌─────────────────────┐
│ Format Response │  Return: {answer, sources, confidence}
└─────────────────────┘
```

# Part 6: Frontend Interface (RPi 3B)

**Voice-Enabled Chat Interface**

html

```html
<!-- /var/www/rag-system/index.html -->

<!DOCTYPE html>
<html lang="en">
<head>
  <meta charset="UTF-8">
  <meta name="viewport" content="width=device-width, initial-scale=1.0">
  <title>Voice-Enabled RAG System</title>
  <style>
    * { margin: 0; padding: 0; box-sizing: border-box; }

    body {
      font-family: 'Segoe UI', Tahoma, Geneva, Verdana, sans-serif;
      background: linear-gradient(135deg, #667eea 0%, #764ba2 100%);
      height: 100vh;
      display: flex;
      justify-content: center;
      align-items: center;
    }

    .container {
      width: 90%;
      max-width: 800px;
      height: 90vh;
      background: white;
      border-radius: 20px;
      box-shadow: 0 20px 60px rgba(0,0,0,0.3);
      display: flex;
      flex-direction: column;
      overflow: hidden;
    }

    .header {
      background: linear-gradient(135deg, #667eea 0%, #764ba2 100%);
      color: white;
      padding: 20px;
      text-align: center;
    }

    .header h1 {
      font-size: 24px;
      margin-bottom: 5px;
    }

    .header p {
      font-size: 14px;
      opacity: 0.9;
```

```css
}

.chat-container {
    flex: 1;
    overflow-y: auto;
    padding: 20px;
    background: #f5f5f5;
}

.message {
    margin-bottom: 15px;
    display: flex;
    animation: fadeIn 0.3s ease;
}

@keyframes fadeIn {
    from { opacity: 0; transform: translateY(10px); }
    to { opacity: 1; transform: translateY(0); }
}

.message.user {
    justify-content: flex-end;
}

.message-content {
    max-width: 70%;
    padding: 12px 16px;
    border-radius: 15px;
    word-wrap: break-word;
}

.message.user .message-content {
    background: #667eea;
    color: white;
    border-bottom-right-radius: 5px;
}

.message.assistant .message-content {
    background: white;
    color: #333;
    border-bottom-left-radius: 5px;
    box-shadow: 0 2px 5px rgba(0,0,0,0.1);
}

.controls {
    padding: 20px;
    background: white;
    border-top: 1px solid #e0e0e0;
}
```

```css
.input-group {
    display: flex;
    gap: 10px;
    margin-bottom: 15px;
}

#textInput {
    flex: 1;
    padding: 12px;
    border: 2px solid #e0e0e0;
    border-radius: 25px;
    font-size: 14px;
    outline: none;
    transition: border-color 0.3s;
}

#textInput:focus {
    border-color: #667eea;
}

.btn {
    padding: 12px 24px;
    border: none;
    border-radius: 25px;
    font-size: 14px;
    font-weight: 600;
    cursor: pointer;
    transition: all 0.3s;
    text-transform: uppercase;
    letter-spacing: 0.5px;
}

.btn-primary {
    background: #667eea;
    color: white;
}

.btn-primary:hover {
    background: #5568d3;
    transform: translateY(-2px);
    box-shadow: 0 5px 15px rgba(102, 126, 234, 0.4);
}

.btn-voice {
    background: #e74c3c;
    color: white;
    width: 60px;
    height: 60px;
    border-radius: 50%;
```

```css
    display: flex;
    align-items: center;
    justify-content: center;
    font-size: 24px;
}

.btn-voice.recording {
    background: #c0392b;
    animation: pulse 1.5s infinite;
}

@keyframes pulse {
    0%, 100% { transform: scale(1); }
    50% { transform: scale(1.05); }
}

.voice-container {
    display: flex;
    justify-content: center;
    align-items: center;
    gap: 15px;
}

.status {
    padding: 8px 16px;
    background: #f0f0f0;
    border-radius: 20px;
    font-size: 12px;
    color: #666;
    text-align: center;
    margin-top: 10px;
}

.status.processing {
    background: #fff3cd;
    color: #856404;
}

.audio-player {
    margin-top: 10px;
    width: 100%;
}

.loader {
    display: inline-block;
    width: 20px;
    height: 20px;
    border: 3px solid #f3f3f3;
    border-top: 3px solid #667eea;
```

```
        border-radius: 50%;
        animation: spin 1s linear infinite;
      }

      @keyframes spin {
        0% { transform: rotate(0deg); }
        100% { transform: rotate(360deg); }
      }
    </style>
  </head>
  <body>
    <div class="container">
      <div class="header">
        <h1>🤖 Voice-Enabled RAG System</h1>
        <p>Air-Gapped • Secure • Local Processing</p>
      </div>

      <div class="chat-container" id="chatContainer">
        <div class="message assistant">
          <div class="message-content">
            Hello! I'm your AI assistant. You can type your questions or use voice input. All processing happ
          </div>
        </div>
      </div>

      <div class="controls">
        <div class="input-group">
          <input type="text" id="textInput" placeholder="Type your question..." />
          <button class="btn btn-primary" onclick="sendTextQuery()">Send</button>
        </div>

        <div class="voice-container">
          <button class="btn btn-voice" id="voiceBtn" onclick="toggleVoiceRecording()">
            🎤
          </button>
          <label for="voiceToggle" style="font-size: 14px; color: #666;">
            <input type="checkbox" id="voiceToggle" checked /> Voice Response
          </label>
        </div>

        <div class="status" id="status">Ready</div>
      </div>
    </div>

    <script>
      let mediaRecorder;
      let audioChunks = [];
      let isRecording = false;

      // Add message to chat
```

```javascript
// Add message to chat
function addMessage(content, isUser = false) {
  const chatContainer = document.getElementById('chatContainer');
  const messageDiv = document.createElement('div');
  messageDiv.className = `message ${isUser ? 'user' : 'assistant'}`;

  const contentDiv = document.createElement('div');
  contentDiv.className = 'message-content';
  contentDiv.innerHTML = content;

  messageDiv.appendChild(contentDiv);
  chatContainer.appendChild(messageDiv);
  chatContainer.scrollTop = chatContainer.scrollHeight;
}

// Update status
function updateStatus(message, processing = false) {
  const status = document.getElementById('status');
  status.textContent = message;
  status.className = processing ? 'status processing' : 'status';
}

// Send text query
async function sendTextQuery() {
  const input = document.getElementById('textInput');
  const query = input.value.trim();

  if (!query) return;

  addMessage(query, true);
  input.value = '';
  updateStatus('Processing your query...', true);

  try {
    const response = await fetch('/api/chat', {
      method: 'POST',
      headers: {
        'Content-Type': 'application/json',
      },
      body: JSON.stringify({
        query: query,
        voice_output: document.getElementById('voiceToggle').checked
      })
    });

    const data = await response.json();

    if (data.success) {
      addMessage(data.answer);
```

```javascript
      // Play audio if voice output enabled
      if (data.audio_url && document.getElementById('voiceToggle').checked) {
        const audio = new Audio(data.audio_url);
        audio.play();
      }

      updateStatus('Ready');
    } else {
      addMessage('Sorry, I encountered an error: ' + data.error);
      updateStatus('Error occurred');
    }
  } catch (error) {
    addMessage('Sorry, I could not process your request.');
    updateStatus('Connection error');
    console.error(error);
  }
}

// Toggle voice recording
async function toggleVoiceRecording() {
  if (isRecording) {
    stopRecording();
  } else {
    startRecording();
  }
}

// Start recording
async function startRecording() {
  try {
    const stream = await navigator.mediaDevices.getUserMedia({ audio: true });
    mediaRecorder = new MediaRecorder(stream);
    audioChunks = [];

    mediaRecorder.ondataavailable = (event) => {
      audioChunks.push(event.data);
    };

    mediaRecorder.onstop = async () => {
      const audioBlob = new Blob(audioChunks, { type: 'audio/webm' });
      await sendVoiceQuery(audioBlob);

      // Stop all tracks
      stream.getTracks().forEach(track => track.stop());
    };

    mediaRecorder.start();
    isRecording = true;
```

```javascript
      document.getElementById('voiceBtn').classList.add('recording');
      document.getElementById('voiceBtn').textContent = '🔲';
      updateStatus('Recording... Click to stop', true);

    } catch (error) {
      console.error('Error accessing microphone:', error);
      updateStatus('Microphone access denied');
    }
  }

// Stop recording
function stopRecording() {
  if (mediaRecorder && isRecording) {
    mediaRecorder.stop();
    isRecording = false;

    document.getElementById('voiceBtn').classList.remove('recording');
    document.getElementById('voiceBtn').textContent = '🎤';
    updateStatus('Processing voice...', true);
  }
}

// Send voice query
async function sendVoiceQuery(audioBlob) {
  const formData = new FormData();
  formData.append('audio', audioBlob, 'recording.webm');
  formData.append('voice_output', document.getElementById('voiceToggle').checked);

  try {
    updateStatus('Transcribing speech...', true);

    const response = await fetch('/api/voice-query', {
      method: 'POST',
      body: formData
    });

    const data = await response.json();

    if (data.success) {
      addMessage(data.transcription, true);
      addMessage(data.answer);

      // Play audio response
      if (data.audio_url && document.getElementById('voiceToggle').checked) {
        const audio = new Audio(data.audio_url);
        audio.play();
      }

      updateStatus('Ready');
    } else {
```

```
                    addMessage('Sorry, I could not understand the audio: ' + data.error);
                    updateStatus('Error occurred');
                }
            } catch (error) {
                addMessage('Sorry, voice processing failed.');
                updateStatus('Connection error');
                console.error(error);
            }
        }

        // Handle Enter key
        document.getElementById('textInput').addEventListener('keypress', function(e) {
            if (e.key === 'Enter') {
                sendTextQuery();
            }
        });
    </script>
</body>
</html>
```

# Part 7: Performance Comparison

**Scenario 1: RPi 3B + T480 (Your Setup)**

**Hardware Cost:** $0 (already owned)

**Performance Benchmarks:**

| Task | Time | Quality |
|------|------|---------|
| Speech-to-Text (5s audio) | 2-5s | Excellent (95%+) |
| Document embedding | 0.5-1s | Good |
| Vector search | 0.2-0.5s | Excellent |
| LLM inference (50 tokens) | 3-10s | Excellent |
| Text-to-Speech (50 words) | 1-2s | Excellent |
| **Total voice round-trip** | **7-18s** | **Very Good** |

**Concurrent Users:** 2-3 simultaneous queries

**Document Capacity:** 10,000+ pages in vector DB

**Strengths:**

- ✅ Zero cost (hardware already owned)

- ✅ Great performance with T480

- ✅ Excellent demo showcase (RPi + laptop combo)

- ✅ Modular architecture (easy to explain)

- ✅ Real-world scalability demo

- ✅ Professional separation of concerns

**Weaknesses:**

- ⚠️ Two devices to manage

- ⚠️ Network dependency (though 1Gbps is fast)

- ⚠️ Slightly more complex setup

**Scenario 2: Dell Optiplex Micro (Suggested Alternative)**

**Hardware Cost:** $200-250 (Dell Optiplex 7050 Micro, i5-7500, 16GB)

**Performance Benchmarks:**

| Task | Time | Quality |
|------|------|---------|
| Speech-to-Text (5s audio) | 2-4s | Excellent (95%+) |
| Document embedding | 0.3-0.8s | Good |
| Vector search | 0.2-0.4s | Excellent |
| LLM inference (50 tokens) | 2-8s | Excellent |
| Text-to-Speech (50 words) | 1-2s | Excellent |
| **Total voice round-trip** | **5-15s** | **Excellent** |

**Concurrent Users:** 3-5 simultaneous queries

**Document Capacity:** 15,000+ pages in vector DB

**Strengths:**

- ✅ Slightly faster (15-20% improvement)

- ✅ Single device (simpler management)

- ✅ More professional appearance

- ✅ Better for production evolution

- ✅ More RAM headroom (can run larger models)

**Weaknesses:**

- ❌ $200-250 cost
- ❌ Less "impressive" demo (no RPi showcase)
- ❌ Harder to explain distributed architecture

---

## Part 8: Detailed Comparison Matrix

| Aspect | RPi 3B + T480 | Dell Optiplex Micro | Winner |
|---|---|---|---|
| **Hardware Cost** | $0 (owned) | $200-250 | 🏆 RPi + T480 |
| **Setup Complexity** | Medium (2 devices) | Low (1 device) | Dell |
| **Voice Response Time** | 7-18s | 5-15s | Dell (marginal) |
| **LLM Performance** | 10-20 tok/s | 12-25 tok/s | Dell (marginal) |
| **Concurrent Users** | 2-3 | 3-5 | Dell |
| **Demo "Wow Factor"** | High (RPi visible) | Medium | 🏆 RPi + T480 |
| **Scalability Demo** | Excellent (show distribution) | Good | 🏆 RPi + T480 |
| **Power Consumption** | 15W (T480) + 6W (RPi) = 21W | 35W | 🏆 RPi + T480 |
| **Failure Points** | 2 devices, network switch | 1 device | Dell |
| **Educational Value** | High (distributed system) | Medium | 🏆 RPi + T480 |
| **Production Ready** | Yes (with upgrades) | Yes | Tie |
| **Portability** | Medium (2 devices + switch) | High (1 device) | Dell |

Recommended Choice: **RPi 3B + T480**

**Why?**

1. **Zero Cost** - You already own both devices

2. **Better Demo Story** - Shows distributed architecture, edge computing concepts

3. **"Wow Factor"** - RPi is impressive to judges at robotics/IoT fair

4. **Educational** - Demonstrates real-world architecture patterns

5. **Performance is Sufficient** - 7-18s is acceptable for government demo

6. **Easy to Explain** - "Frontend on Pi, AI on laptop" is intuitive

**Performance difference (5-15s vs 7-18s) is NOT significant enough to justify $200-250 cost**

---

## Part 9: Complete Installation Guide

Phase 1: Raspberry Pi 3B Setup (Day 1 - 2 hours)

```bash
bash

# 1. Flash DietPi to SD card (on your PC)
# Download from: https://dietpi.com/#download
# Use Etcher or dd to flash

# 2. Boot RPi and configure
# Default login: root / dietpi

# 3. Run DietPi setup
dietpi-config
# - Set hostname: rag-frontend
# - Set static IP: 192.168.1.10
# - Enable SSH
# - Disable WiFi (use Ethernet only)

# 4. Install NGINX
dietpi-software install 85  # NGINX

# 5. Configure firewall
apt install ufw
ufw default deny incoming
ufw default allow outgoing
ufw allow from 192.168.1.0/24 to any port 22    # SSH from LAN
ufw allow from 192.168.1.0/24 to any port 80    # HTTP from LAN
ufw allow from 192.168.1.0/24 to any port 443   # HTTPS from LAN
ufw enable

# 6. Create web root
mkdir -p /var/www/rag-system
chown -R www-data:www-data /var/www/rag-system

# 7. Generate self-signed SSL certificate
mkdir -p /etc/nginx/ssl
openssl req -x509 -nodes -days 365 -newkey rsa:2048 \
  -keyout /etc/nginx/ssl/rag-system.key \
  -out /etc/nginx/ssl/rag-system.crt \
  -subj "/CN=rag-system.local"

# 8. Configure NGINX (use config from Part 3)
nano /etc/nginx/sites-available/rag-system
ln -s /etc/nginx/sites-available/rag-system /etc/nginx/sites-enabled/
rm /etc/nginx/sites-enabled/default
nginx -t
systemctl restart nginx

# 9. Upload frontend files (HTML/CSS/JS from Part 6)
# Use scp from your PC:
# scp -r frontend/* root@192.168.1.10:/var/www/rag-system/
```

```bash
# 1. Install Arch Linux (minimal)
# Follow: https://wiki.archlinux.org/title/Installation_guide
# Key choices:
# - No desktop environment
# - Enable SSH
# - Set hostname: rag-backend
# - Set static IP: 192.168.1.20

# 2. Post-installation
pacman -Syu
pacman -S base-devel git wget curl vim htop tmux

# 3. Create directories
mkdir -p /opt/{whisper,piper,ollama}
mkdir -p /data/{documents,vector-db,uploads}
mkdir -p /var/log/rag-system

# 4. Install Python 3.11+
pacman -S python python-pip python-virtualenv

# 5. Install Node.js (for N8N)
pacman -S nodejs npm

# 6. Configure firewall
pacman -S ufw
ufw default deny incoming
ufw default allow outgoing
ufw allow from 192.168.1.0/24 to any port 22      # SSH
ufw allow from 192.168.1.0/24 to any port 5678    # N8N
ufw allow from 192.168.1.0/24 to any port 8000    # API
ufw allow from 192.168.1.0/24 to any port 8001    # Piper TTS
ufw allow from 192.168.1.0/24 to any port 11434   # Ollama
ufw enable

# 7. Optimize system for AI workloads
# Increase file limits
echo "* soft nofile 65536" >> /etc/security/limits.conf
echo "* hard nofile 65536" >> /etc/security/limits.conf

# Disable swap (we'll use it minimally)
swapoff -a

# Create optimized swap (8GB)
fallocate -l 8G /swapfile
chmod 600 /swapfile
mkswap /swapfile
swapon /swapfile
```

```bash
echo "/swapfile none swap sw 0 0" >> /etc/fstab

# Optimize swappiness
echo "vm.swappiness=10" >> /etc/sysctl.conf
sysctl -p
```

## Phase 3: Ollama Installation (Day 1 - 1 hour)

```bash
bash

# 1. Install Ollama
curl -fsSL https://ollama.com/install.sh | sh

# 2. Start Ollama service
systemctl enable ollama
systemctl start ollama

# 3. Download models
ollama pull llama3:8b      # 4.7GB - best balance
# OR
ollama pull mistral:7b     # 4.1GB - faster alternative

# 4. Test
ollama run llama3:8b "Hello, how are you?"

# 5. Configure for network access
mkdir -p /etc/systemd/system/ollama.service.d/
cat > /etc/systemd/system/ollama.service.d/override.conf << 'EOF'
[Service]
Environment="OLLAMA_HOST=0.0.0.0:11434"
EOF

systemctl daemon-reload
systemctl restart ollama

# 6. Verify
curl http://192.168.1.20:11434/api/tags
```

## Phase 4: Whisper.cpp Installation (Day 2 - 1 hour)

```bash
# 1. Clone and build
cd /opt/whisper
git clone https://github.com/ggerganov/whisper.cpp
cd whisper.cpp

# Build with optimizations for i5-8350U
make -j8

# 2. Download models
bash ./models/download-ggml-model.sh base   # 142MB - recommended
bash ./models/download-ggml-model.sh tiny   # 75MB - faster backup

# 3. Test
./main -m models/ggml-base.bin -f samples/jfk.wav
# Should complete in 2-3 seconds

# 4. Create Python wrapper (save code from Part 4)
cd /opt/whisper
python -m venv venv
source venv/bin/activate
pip install fastapi uvicorn python-multipart

# Save whisper_server.py from Part 4
nano whisper_server.py

# 5. Create systemd service
cat > /etc/systemd/system/whisper-stt.service << 'EOF'
[Unit]
Description=Whisper Speech-to-Text API
After=network.target

[Service]
Type=simple
User=root
WorkingDirectory=/opt/whisper
Environment="PATH=/opt/whisper/venv/bin"
ExecStart=/opt/whisper/venv/bin/python whisper_server.py
Restart=always

[Install]
WantedBy=multi-user.target
EOF

systemctl enable whisper-stt
systemctl start whisper-stt

# 6. Test API
```

```
curl -X POST http://localhost:8000/api/health
```

## Phase 5: Piper TTS Installation (Day 2 - 30 minutes)

```bash
# 1. Install Piper
cd /opt/piper
python -m venv venv
source venv/bin/activate
pip install piper-tts

# 2. Download voice model
mkdir -p voices
cd voices
wget https://huggingface.co/rhasspy/piper-voices/resolve/main/en/en_US/lessac/medium/en_US-lessac-med
wget https://huggingface.co/rhasspy/piper-voices/resolve/main/en/en_US/lessac/medium/en_US-lessac-med

# 3. Test
cd /opt/piper
echo "This is a test" | piper --model voices/en_US-lessac-medium.onnx --output_file test.wav
aplay test.wav

# 4. Create API server (save code from Part 4)
nano piper_server.py

# 5. Create systemd service
cat > /etc/systemd/system/piper-tts.service << 'EOF'
[Unit]
Description=Piper Text-to-Speech API
After=network.target

[Service]
Type=simple
User=root
WorkingDirectory=/opt/piper
Environment="PATH=/opt/piper/venv/bin"
ExecStart=/opt/piper/venv/bin/python piper_server.py
Restart=always

[Install]
WantedBy=multi-user.target
EOF

systemctl enable piper-tts
systemctl start piper-tts

# 6. Test API
curl -X POST http://localhost:8001/api/tts \
  -H "Content-Type: application/json" \
  -d '{"text":"Hello world"}' \
  --output test-api.wav
aplay test-api.wav
```

## Phase 6: N8N Installation (Day 2-3 - 2 hours)

```bash
# 1. Install N8N globally
npm install -g n8n

# 2. Create N8N data directory
mkdir -p /var/lib/n8n

# 3. Create systemd service
cat > /etc/systemd/system/n8n.service << 'EOF'
[Unit]
Description=N8N Workflow Automation
After=network.target

[Service]
Type=simple
User=root
Environment="N8N_BASIC_AUTH_ACTIVE=true"
Environment="N8N_BASIC_AUTH_USER=admin"
Environment="N8N_BASIC_AUTH_PASSWORD=admin123"  # Change this!
Environment="N8N_HOST=0.0.0.0"
Environment="N8N_PORT=5678"
Environment="N8N_PROTOCOL=http"
Environment="WEBHOOK_URL=http://192.168.1.20:5678/"
Environment="N8N_USER_FOLDER=/var/lib/n8n"
ExecStart=/usr/bin/n8n start
Restart=always

[Install]
WantedBy=multi-user.target
EOF

systemctl enable n8n
systemctl start n8n

# 4. Access N8N
# Open browser: http://192.168.1.20:5678
# Login with admin/admin123

# 5. Install required nodes/packages
# In N8N, go to Settings > Community Nodes
# Install: n8n-nodes-langchain (for embeddings)
```

## Phase 7: Backend API Creation (Day 3-4 - 4 hours)

```bash
# 1. Create project structure
mkdir -p /opt/rag-api
cd /opt/rag-api

# 2. Create virtual environment
python -m venv venv
source venv/bin/activate

# 3. Install dependencies
pip install \
  fastapi \
  uvicorn \
  python-multipart \
  httpx \
  sentence-transformers \
  PyPDF2 \
  python-docx \
  aiofiles

# 4. Create main API file
nano main.py
```

```python
# /opt/rag-api/main.py

from fastapi import FastAPI, File, UploadFile, Form
from fastapi.responses import JSONResponse, FileResponse
from fastapi.middleware.cors import CORSMiddleware
import httpx
import os
import tempfile
import json
from datetime import datetime

app = FastAPI(title="RAG System API")

# CORS for frontend
app.add_middleware(
    CORSMiddleware,
    allow_origins=["*"],
    allow_credentials=True,
    allow_methods=["*"],
    allow_headers=["*"],
)

# Configuration
WHISPER_API = "http://localhost:8000/api/stt"
PIPER_API = "http://localhost:8001/api/tts"
N8N_WEBHOOK_QUERY = "http://localhost:5678/webhook/rag-query"
N8N_WEBHOOK_UPLOAD = "http://localhost:5678/webhook/upload-document"
UPLOAD_DIR = "/data/uploads"

os.makedirs(UPLOAD_DIR, exist_ok=True)

@app.get("/api/health")
async def health_check():
    return {"status": "healthy", "timestamp": datetime.now().isoformat()}

@app.post("/api/voice-query")
async def voice_query(
    audio: UploadFile = File(...),
    voice_output: bool = Form(True)
):
    """
    Complete voice query pipeline:
    1. Speech-to-text (Whisper)
    2. RAG query (N8N + Ollama)
    3. Text-to-speech (Piper)
    """
    try:
```

```python
# Step 1: Speech-to-Text
audio_content = await audio.read()

async with httpx.AsyncClient(timeout=30.0) as client:
    files = {"audio": ("recording.webm", audio_content, "audio/webm")}
    stt_response = await client.post(WHISPER_API, files=files)
    stt_data = stt_response.json()

if not stt_data.get("success"):
    return JSONResponse({
        "success": False,
        "error": "Speech recognition failed"
    }, status_code=500)

transcription = stt_data["transcription"]

# Step 2: RAG Query
async with httpx.AsyncClient(timeout=60.0) as client:
    rag_response = await client.post(
        N8N_WEBHOOK_QUERY,
        json={"query": transcription}
    )
    rag_data = rag_response.json()

answer = rag_data.get("answer", "I couldn't find an answer.")

# Step 3: Text-to-Speech (if enabled)
audio_url = None
if voice_output:
    async with httpx.AsyncClient(timeout=30.0) as client:
        tts_response = await client.post(
            PIPER_API,
            json={"text": answer}
        )

        # Save audio file
        audio_path = f"/tmp/response_{datetime.now().timestamp()}.wav"
        with open(audio_path, "wb") as f:
            f.write(tts_response.content)

        audio_url = f"/api/audio/{os.path.basename(audio_path)}"

return JSONResponse({
    "success": True,
    "transcription": transcription,
    "answer": answer,
    "audio_url": audio_url,
    "processing_time_ms": None  # Can add timing
})
```

```python
        except Exception as e:
            return JSONResponse({
                "success": False,
                "error": str(e)
            }, status_code=500)

@app.post("/api/chat")
async def text_chat(request: dict):
    """
    Text-only chat endpoint
    """
    try:
        query = request.get("query")
        voice_output = request.get("voice_output", False)

        # RAG Query via N8N
        async with httpx.AsyncClient(timeout=60.0) as client:
            rag_response = await client.post(
                N8N_WEBHOOK_QUERY,
                json={"query": query}
            )
            rag_data = rag_response.json()

        answer = rag_data.get("answer", "I couldn't find an answer.")

        # Generate audio if requested
        audio_url = None
        if voice_output:
            async with httpx.AsyncClient(timeout=30.0) as client:
                tts_response = await client.post(
                    PIPER_API,
                    json={"text": answer}
                )

                audio_path = f"/tmp/response_{datetime.now().timestamp()}.wav"
                with open(audio_path, "wb") as f:
                    f.write(tts_response.content)

                audio_url = f"/api/audio/{os.path.basename(audio_path)}"

        return JSONResponse({
            "success": True,
            "answer": answer,
            "audio_url": audio_url,
            "sources": rag_data.get("sources", [])
        })

    except Exception as e:
        return JSONResponse({
```

```python
            "success": False,
            "error": str(e)
        }, status_code=500)


@app.post("/api/upload")
async def upload_document(file: UploadFile = File(...)):
    """
    Upload document for ingestion
    """
    try:
        # Save file
        file_path = os.path.join(UPLOAD_DIR, file.filename)
        with open(file_path, "wb") as f:
            content = await file.read()
            f.write(content)

        # Trigger N8N ingestion workflow
        async with httpx.AsyncClient(timeout=120.0) as client:
            ingest_response = await client.post(
                N8N_WEBHOOK_UPLOAD,
                json={
                    "filename": file.filename,
                    "filePath": file_path
                }
            )
            ingest_data = ingest_response.json()

        return JSONResponse({
            "success": True,
            "filename": file.filename,
            "chunks_processed": ingest_data.get("chunks_count", 0),
            "doc_id": ingest_data.get("doc_id")
        })

    except Exception as e:
        return JSONResponse({
            "success": False,
            "error": str(e)
        }, status_code=500)


@app.get("/api/audio/{filename}")
async def serve_audio(filename: str):
    """
    Serve generated audio files
    """
    file_path = f"/tmp/{filename}"
    if os.path.exists(file_path):
        return FileResponse(file_path, media_type="audio/wav")
    return JSONResponse({"error": "File not found"}, status_code=404)
```

```python
if __name__ == "__main__":
    import uvicorn
    uvicorn.run(app, host="0.0.0.0", port=8000, workers=2)
```

```bash
# 5. Create systemd service for API
cat > /etc/systemd/system/rag-api.service << 'EOF'
[Unit]
Description=RAG System API
After=network.target ollama.service whisper-stt.service piper-tts.service n8n.service

[Service]
Type=simple
User=root
WorkingDirectory=/opt/rag-api
Environment="PATH=/opt/rag-api/venv/bin"
ExecStart=/opt/rag-api/venv/bin/uvicorn main:app --host 0.0.0.0 --port 8000 --workers 2
Restart=always

[Install]
WantedBy=multi-user.target
EOF

systemctl enable rag-api
systemctl start rag-api

# 6. Test complete system
curl http://localhost:8000/api/health
```

# Part 10: N8N Workflow Implementation

**Workflow 1: Document Ingestion**

**Create in N8N UI:**

1. **Webhook Node** (Trigger)
   - Method: POST
   - Path: `upload-document`
   - Response Mode: Last Node

2. **Code Node** (Read File)

```javascript
const fs = require('fs');
const filePath = $json.filePath;
const filename = $json.filename;

// Read file content
const content = fs.readFileSync(filePath, 'utf8');

return {
  filename: filename,
  content: content,
  filePath: filePath
};
```

3. **Code Node** (Chunk Text)

```javascript
const content = $json.content;
const chunkSize = 512;
const overlap = 50;

const chunks = [];
let start = 0;

while (start < content.length) {
  const end = Math.min(start + chunkSize, content.length);
  const chunk = content.substring(start, end);

  chunks.push({
    text: chunk,
    index: chunks.length,
    filename: $json.filename
  });

  start += (chunkSize - overlap);
}

return chunks;
```

4. **Embeddings (OpenAI) Node** - *Use Local Alternative*
   - Model: all-MiniLM-L6-v2
   - Input: `{{ $json.text }}`

   *Note: Install* `@n8n/n8n-nodes-langchain` *package for local embeddings*

5. **Vector Store Insert Node**
   - Store: N8N Simple Vector Store
   - Document: `{{ $json.text }}`
   - Embedding: From previous node
   - Metadata: `{{ {"filename": $json.filename, "index": $json.index} }}`

6. **Respond to Webhook Node**

```json
{
  "success": true,
  "chunks_count": {{ $json.length }},
  "doc_id": "{{ $json.filename }}"
}
```

**Workflow 2: RAG Query**

**Create in N8N UI:**

1. **Webhook Node** (Trigger)
   - Method: POST
   - Path: `rag-query`
   - Response Mode: Last Node

2. **Embeddings Node** (Query Embedding)
   - Model: all-MiniLM-L6-v2
   - Input: `{{ $json.query }}`

3. **Vector Store Retrieval Node**
   - Store: N8N Simple Vector Store
   - Query Embedding: From previous node
   - Top K: 5
   - Similarity Threshold: 0.7

4. **Code Node** (Build Context)

```javascript
  const items = $input.all();
  const query = $('Webhook').first().json.query;

  // Combine retrieved chunks
  const context = items.map((item, i) => {
    return `[Source ${i+1}]: ${item.json.text}`;
  }).join('\n\n');

  // Build prompt
  const prompt = `You are a helpful AI assistant. Answer the question based on the following context. If the

Context:
${context}

Question: ${query}

Answer:`;

  return {
    prompt: prompt,
    sources: items.map(i => i.json.metadata)
  };
```

5. **HTTP Request Node** (Call Ollama)

- Method: POST

- URL: `http://localhost:11434/api/generate`

- Body:

```json
  {
   "model": "llama3:8b",
   "prompt": "{{ $json.prompt }}",
   "stream": false,
   "options": {
    "temperature": 0.7,
    "num_predict": 500
   }
  }
```

6. **Code Node** (Format Response)

```javascript
const ollamaResponse = $json.response;
const sources = $('Code').first().json.sources;

return {
  answer: ollamaResponse,
  sources: sources,
  confidence: 0.85  // Can implement confidence scoring
};
```

7. **Respond to Webhook Node**

```json
{
  "success": true,
  "answer": "{{ $json.answer }}",
  "sources": "{{ $json.sources }}"
}
```

---

# Part 11: Testing & Validation

**System Testing Checklist**

```bash
# Test 1: Network connectivity
ping 192.168.1.10  # RPi from T480
ping 192.168.1.20  # T480 from RPi

# Test 2: Individual services on T480
curl http://localhost:11434/api/tags  # Ollama
curl http://localhost:8000/api/health  # Whisper
curl http://localhost:8001/api/health  # Piper (if you added health endpoint)
curl http://localhost:5678/  # N8N
curl http://localhost:8000/api/health  # Main API

# Test 3: Whisper STT
curl -X POST http://192.168.1.20:8000/api/stt \
  -F "audio=@test-audio.wav"

# Test 4: Piper TTS
curl -X POST http://192.168.1.20:8001/api/tts \
  -H "Content-Type: application/json" \
  -d '{"text":"Testing the text to speech system"}' \
  --output test-output.wav

# Test 5: Ollama inference
curl http://192.168.1.20:11434/api/generate -d '{
  "model": "llama3:8b",
  "prompt": "Why is the sky blue?",
  "stream": false
}'

# Test 6: Document upload
curl -X POST http://192.168.1.20:8000/api/upload \
  -F "file=@sample-policy.pdf"

# Test 7: Text chat
curl -X POST http://192.168.1.20:8000/api/chat \
  -H "Content-Type: application/json" \
  -d '{"query":"What is the leave policy?","voice_output":false}'

# Test 8: Frontend access from RPi
curl http://192.168.1.10/
curl -k https://192.168.1.10/  # HTTPS
```

## Performance Benchmarking

```bash
# Create benchmark script on T480
cat > /opt/benchmark.sh << 'EOF'
#!/bin/bash

echo "=== RAG System Performance Benchmark ==="
echo ""

# Test 1: Whisper STT
echo "Test 1: Speech-to-Text (5 second audio)"
time curl -X POST http://localhost:8000/api/stt \
  -F "audio=@/opt/test-data/5sec-audio.wav" \
  -o /dev/null -s

echo ""

# Test 2: Ollama Inference
echo "Test 2: LLM Inference (50 token response)"
time curl http://localhost:11434/api/generate -d '{
  "model": "llama3:8b",
  "prompt": "Explain quantum computing in 50 words.",
  "stream": false
}' -o /dev/null -s

echo ""

# Test 3: Piper TTS
echo "Test 3: Text-to-Speech (50 words)"
time curl -X POST http://localhost:8001/api/tts \
  -H "Content-Type: application/json" \
  -d '{"text":"The quick brown fox jumps over the lazy dog. This is a test of the text to speech system perfor
  -o /dev/null -s

echo ""

# Test 4: Complete RAG query
echo "Test 4: Complete RAG Query (end-to-end)"
time curl -X POST http://localhost:8000/api/chat \
  -H "Content-Type: application/json" \
  -d '{"query":"What is the leave policy?","voice_output":false}' \
  -o /dev/null -s

echo ""
echo "=== Benchmark Complete ==="
EOF

chmod +x /opt/benchmark.sh
/opt/benchmark.sh
```

**Expected Results on T480:**

- Whisper STT (5s audio): 2-5 seconds

- Ollama Inference (50 tokens): 3-10 seconds

- Piper TTS (50 words): 1-2 seconds

- Complete RAG query: 5-15 seconds

---

# Part 12: Demo Day Preparation

**Demo Script (15 minutes)**

**Minute 0-2: Introduction**

- Show physical setup (RPi + T480 + switch)

- Explain air-gapped architecture

- Highlight zero external API calls

**Minute 2-4: Text Chat Demo**

- Open web interface on laptop

- Type: "What documents are available?"

- Show response time (~5-10s)

- Explain RAG process happening

**Minute 4-7: Document Upload**

- Upload sample government policy (PDF)

- Show ingestion progress

- Explain vector embedding process

- Query the uploaded document

**Minute 7-11: Voice Interaction (THE WOW MOMENT)**

- Click microphone button

- Speak: "What is the leave policy for government employees?"

- Show transcription appearing

- Show LLM processing

- Play audio response

- **This is your killer feature!**

**Minute 11-13: Technical Deep Dive**

- Show system architecture diagram

- Open N8N workflow (visual appeal)

- Show Ollama running locally

- Demonstrate no internet connection

**Minute 13-15: Cost & Scalability**

- Show hardware ($0 procurement cost)

- Show power consumption (~21W total)

- Discuss scaling (add more T480s)

- Q&A

**Demo Talking Points**

**For Judges:**

1. **Security & Privacy**
   - "Every single bit of data stays within this network"
   - "No cloud providers, no third parties, complete data sovereignty"
   - "Perfect for classified or sensitive government documents"

2. **Cost Effectiveness**
   - "Zero procurement cost - used existing hardware"
   - "Under $10/year electricity cost"
   - "No per-query fees, no subscription costs"
   - "Compare to commercial AI APIs: $1000+/month for this volume"

3. **Performance**
   - "6-17 second end-to-end voice response"
   - "Processing happens locally on this laptop"
   - "Can handle 2-3 simultaneous users"
   - "Scalable by adding more hardware"

4. **Innovation**
   - "Raspberry Pi as edge device (IoT/robotics showcase)"
   - "Distributed architecture (real-world design pattern)"
   - "Voice-enabled RAG (cutting-edge feature)"
   - "Open-source stack (no vendor lock-in)"

5. **Practical Applications**
   - "Policy Q&A systems for government offices"
   - "Internal knowledge bases"
   - "Citizen service chatbots (offline capable)"
   - "Emergency response systems (works without internet)"

**Backup Plans**

**If something fails during demo:**

1. **Voice doesn't work:**
   - Fall back to text chat
   - Explain: "Voice is experimental, text is rock-solid"

2. **Network issues:**
   - Have video recording of working system
   - Show architecture diagrams and code

3. **Slow response times:**
   - Explain: "This is the small model, production would use dedicated GPU"
   - Show performance benchmarks on slides

4. **Questions you can't answer:**
   - "Great question! This is a proof of concept, and that's exactly the kind of feedback we need for v2"

---

## Part 13: Post-Demo Evolution Path

**Short-term Upgrades (1-3 months)**

1. **Add GPU to T480 (eGPU via Thunderbolt)**
   - Cost: $300-500 (used GTX 1070/1080)
   - Performance: 5-10x faster inference
   - Response time: 1-3 seconds (production-ready)

2. **Scale horizontally**
   - Add second T480 or similar laptop
   - Load balancing via NGINX
   - Support 10-20 concurrent users

3. **Better models**
   - Llama 3 70B (with GPU)
   - Mixtral 8x7B
   - Domain-specific fine-tuned models

**Medium-term (3-6 months)**

1. **Professional hardware**
   - Used server: Dell R730 (~$500-800)

   - 64-128GB RAM, Xeon CPUs

   - Add NVIDIA Tesla P40 ($200 used)

   - Support 50+ users

2. **Enhanced features**
   - Multi-language support

   - Advanced RAG (HyDE, reranking)

   - User authentication

   - Conversation memory

3. **Enterprise deployment**
   - Kubernetes cluster

   - High availability

   - Disaster recovery

   - Monitoring dashboard

---

## Part 14: Complete Cost Breakdown

**Your Setup (RPi 3B + T480)**

**Hardware (Already Owned):**

- Raspberry Pi 3B: $0

- Lenovo T480: $0

- Gigabit switch: $0 (or $15 if needed)

- Ethernet cables: $0 (or $5)

- **Total Hardware: $0-20**

**Software:**

- Everything is open-source: $0

**Initial Setup Time:**

- RPi setup: 2 hours

- T480 setup: 8 hours

- Testing: 4 hours

- **Total: ~14 hours** (2 days part-time)

**Operating Costs (Annual):**

- Electricity (21W × 24h × 365d × $0.12/kWh): ~$22

- SD card replacement (every 2-3 years): $4/year

- Maintenance: $0

- **Total OPEX: ~$26/year**

**5-Year Total Cost of Ownership:**

- CAPEX: $0-20

- OPEX: $130

- **Total: $130-150** (vs $60,000+ for commercial AI API services)

Dell Optiplex Comparison

**Hardware Cost:**

- Dell Optiplex 7050 Micro (i5-7500, 16GB): $220

- Storage upgrade: $0 (256GB included)

- Accessories: $20

- **Total: $240**

**Performance:**

- 15-20% faster than T480

- Single device (simpler)

- Less impressive demo (no RPi showcase)

**Verdict: Not worth $240 extra given your scenario**

---

# Part 15: Final Recommendation & Action Plan

FINAL VERDICT: Use RPi 3B + T480 Setup

**Why this is the right choice:**

1. ✅ **Zero procurement cost** - You own the hardware

2. ✅ **Impressive demo** - RPi visible, shows distributed architecture

3. ✅ **Performance adequate** - 7-18s is acceptable for government demo

4. ✅ **Voice functionality** - Fully local, no external APIs

5. ✅ **Educational value** - Great for robotics/IoT fair

6. ✅ **Scalable story** - Easy to explain growth path

7. ✅ **Power efficient** - 21W total (green computing angle)

8. ✅ **Real-world architecture** - Edge + backend pattern

**The Dell Optiplex would give you:**

- Marginally faster (5-15s vs 7-18s) - **Not significant**

- Single device - **Less interesting demo**

- $240 cost - **Not worth it for your scenario**

**7-Day Implementation Plan**

**Day 1: Hardware Setup**

- Morning: Flash DietPi on RPi, configure network

- Afternoon: Install Arch Linux on T480, configure network

- Evening: Test connectivity between devices

**Day 2: Core Services**

- Morning: Install Ollama, download models

- Afternoon: Install and test Whisper.cpp

- Evening: Install and test Piper TTS

**Day 3: Workflow Engine**

- Morning: Install N8N

- Afternoon: Create document ingestion workflow

- Evening: Create RAG query workflow

**Day 4: Backend API**

- Morning: Create FastAPI backend

- Afternoon: Integrate all services

- Evening: Test API endpoints

**Day 5: Frontend**

- Morning: Deploy HTML/CSS/JS to RPi

- Afternoon: Configure NGINX, SSL

- Evening: Test complete flow

**Day 6: Testing & Optimization**

- Morning: Performance benchmarks

- Afternoon: Load test, fix issues

- Evening: Document upload and query tests

**Day 7: Demo Preparation**

- Morning: Create demo documents

- Afternoon: Rehearse demo script

- Evening: Final checks, backup plans

**Next Steps (Start Today!)**

1. **Download ISOs:**
    - DietPi for RPi 3B

    - Arch Linux for T480

2. **Prepare SD card** for RPi

3. **Backup T480** before wiping

4. **Clone this document** for reference during installation

5. **Set up test environment** on one device first

---

# Part 16: Judges' FAQs (Be Prepared)

**Q: Why not use cloud services like ChatGPT API?** A: "This demo is for government applications where data privacy is paramount. Using cloud APIs means sending sensitive documents to third parties. Our system keeps 100% of data on-premises and works without internet."

**Q: Is 7-18 seconds response time too slow?** A: "For a proof of concept running on a $0 laptop with no GPU, this is excellent. Production systems would add a GPU card ($300) and reduce response time to 1-3 seconds. Compare that to $1000+/month cloud costs."

**Q: Why use a Raspberry Pi if the T480 does all the work?** A: "The Pi demonstrates edge computing architecture - in production, we could have multiple Pis at different locations (offices, kiosks) all connecting to a central AI backend. It's a realistic distributed system design."

**Q: How secure is this really?** A: "Completely air-gapped - no external network connections possible. All data stays on local hardware. Perfect for classified documents. We can add encryption, authentication, and audit logs for production."

**Q: Can it scale to 100+ users?** A: "This demo supports 2-3 concurrent users. Production scaling is straightforward: add more T480-class machines (or better hardware) behind the same architecture. Each $500 server can handle 20-30 users."

**Q: What about model accuracy?** A: "We're using Llama 3 8B, which is comparable to GPT-3.5 for most tasks. For better accuracy, we can use larger models (70B parameters) with GPU hardware. The RAG approach helps by grounding answers in your documents."

**Q: How do you update the AI models?** A: "Models are updated manually via USB drive in air-gapped mode, or automatically if we allow controlled internet access for updates only. We can schedule monthly updates, similar to antivirus definitions."

---

## Conclusion

You have an **excellent opportunity** here with zero procurement cost and impressive demo potential. The RPi + T480 setup gives you:

- Professional distributed architecture

- Full voice capabilities (local STT + TTS)

- Secure air-gapped operation

- Under $30/year operating costs

- Perfect showcase for robotics/IoT/AI fair

**This is absolutely worth building and will impress government stakeholders!**

The 7-day implementation plan is realistic, and you'll have a working system that demonstrates cutting-edge concepts while staying completely secure and cost-effective.

Good luck with your demo! 🚀