

Review of the Edited Document

Upon reviewing the edited document ("Top 10 Open-Source AI Projects_ Quick Start Guide.docx" with version 2025-08-24-A), it remains a comprehensive, well-structured guide for setting up an integrated AI stack on a single PC. Key strengths include:

- **Clarity and Completeness:** The reference table (Section A) is concise yet detailed, with updated entries for MeloTTS. Installation steps in Section B are copy-paste ready, and the replacement of Orpheus-TTS with MeloTTS is seamlessly integrated, including API wrappers, testing commands, and Docker configurations.
- **Error Handling:** The fix for the Python package installation error is well-addressed in the updated `setup-ai-stack.ps1` script and troubleshooting section, ensuring reliable dependency management.
- **Automation and Maintenance:** Sections C (weekly updates) and D (Docker vs. native) are practical, with the new MeloTTS Dockerfile and GPU passthrough notes enhancing cross-platform usability.
- **VS Code Integration:** Section B-15 provides a solid foundation for using VS Code as the default IDE, covering extensions, settings, tasks, and debugging. However, it could be expanded for deeper integration with specific projects (e.g., Jupyter for LangChain/Haystack, Docker for OpenManus/OpenSora) and advanced features like AI-assisted coding (e.g., via Ollama integration).

Potential improvements:

- **Versioning:** Update the document version to reflect the current date (2025-08-24) for

any new changes.

- **Cross-Platform Notes:** Add brief Windows/Linux/macOS variations in B-15, as the current setup leans toward Windows (e.g., PowerShell paths).
- **Security/Performance:** Suggest enabling VS Code's restricted mode for untrusted code and monitoring GPU usage with extensions like NVIDIA GPU Monitor.
- **Testing:** Include a quick verification script in B-15 to check if all projects are accessible from VS Code (e.g., running a sample LangChain query in a notebook).

Overall, the document is robust and user-friendly, with the MeloTTS swap maintaining compatibility while improving multilingual support.

Suggestions for Setting VS Code to Work with All Open-Source Projects by Default

The edited document's Section B-15 already establishes VS Code as the central development environment for the 11 projects (DeepSeek-R1, Ollama, OpenManus, LangChain, AutoGen, OpenSora, Haystack, Text-Generation-WebUI, Whisper, Letta, and MeloTTS). To make it the "default" tool—meaning all project interactions (editing, running, debugging, and monitoring) happen seamlessly within VS Code—build on B-15 with these enhancements. These suggestions draw from best practices for Python-based AI workflows, Docker integration, and local AI tools like Ollama.

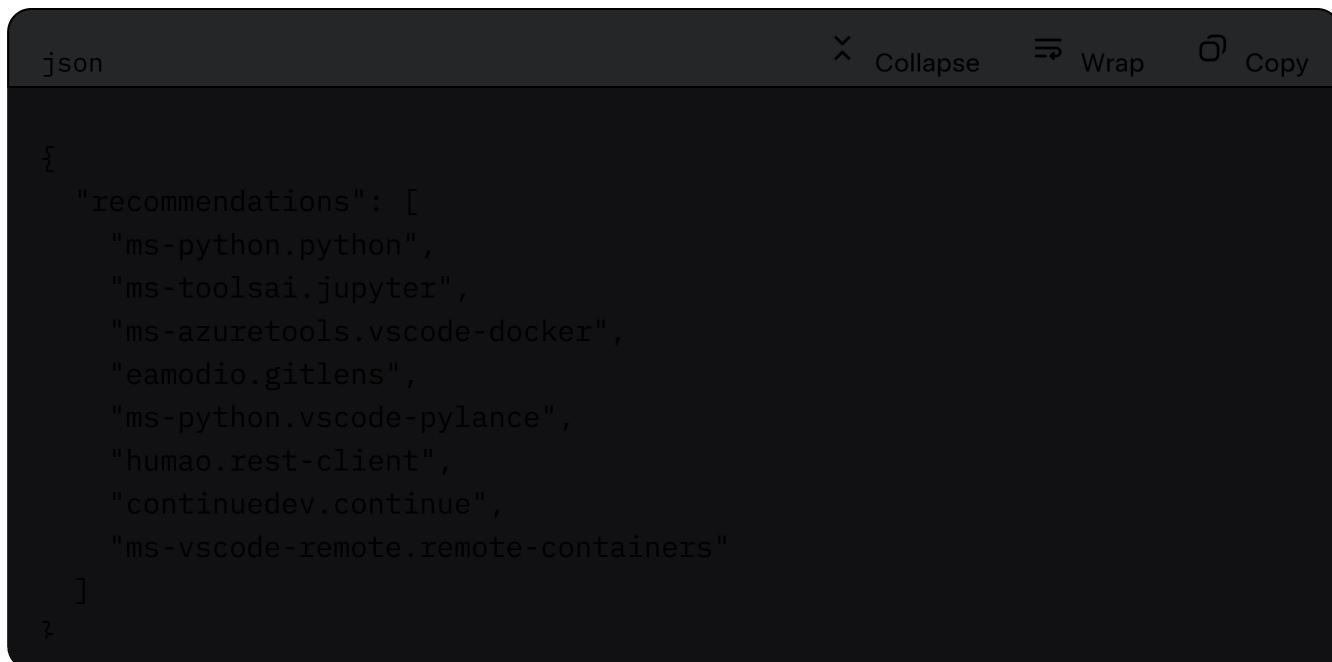
1. Enhance Extensions for Project-Specific Support

- **Core Extensions (Already in B-15):** Keep Python, Jupyter, Docker, GitLens, Pylance, and

REST Client.

- **Additional Extensions:**

- **Continue** (by Continue.dev): Integrates local AI models (e.g., via Ollama) for code completion and chat. Install via Extensions panel, then configure in `settings.json` to point to Ollama's endpoint (`http://localhost:11434/v1`). This enables AI-assisted coding for LangChain/AutoGen scripts.
 - **Dev Containers** (ms-vscode-remote.remote-containers): For working inside Docker containers (e.g., OpenSora or OpenManus). This allows editing code in a reproducible environment without local installs.
 - **Remote - SSH/WSL** (ms-vscode-remote.remote-ssh/ms-vscode-remote.remote-wsl): If using WSL for Linux-based projects like OpenSora.
 - **NVIDIA GPU Monitor** (or similar, e.g., GPU Usage): To monitor VRAM during MeloTTS/Whisper GPU runs.
 - **Jupyter Keymap** (ms-toolsai.jupyter-keymap): Improves notebook navigation for Haystack/LangChain RAG prototypes.
- **Installation Tip:** In VS Code, use `Ctrl+Shift+X` to search and install. For batch install, create a `.vscode/extensions.json` in the workspace root:

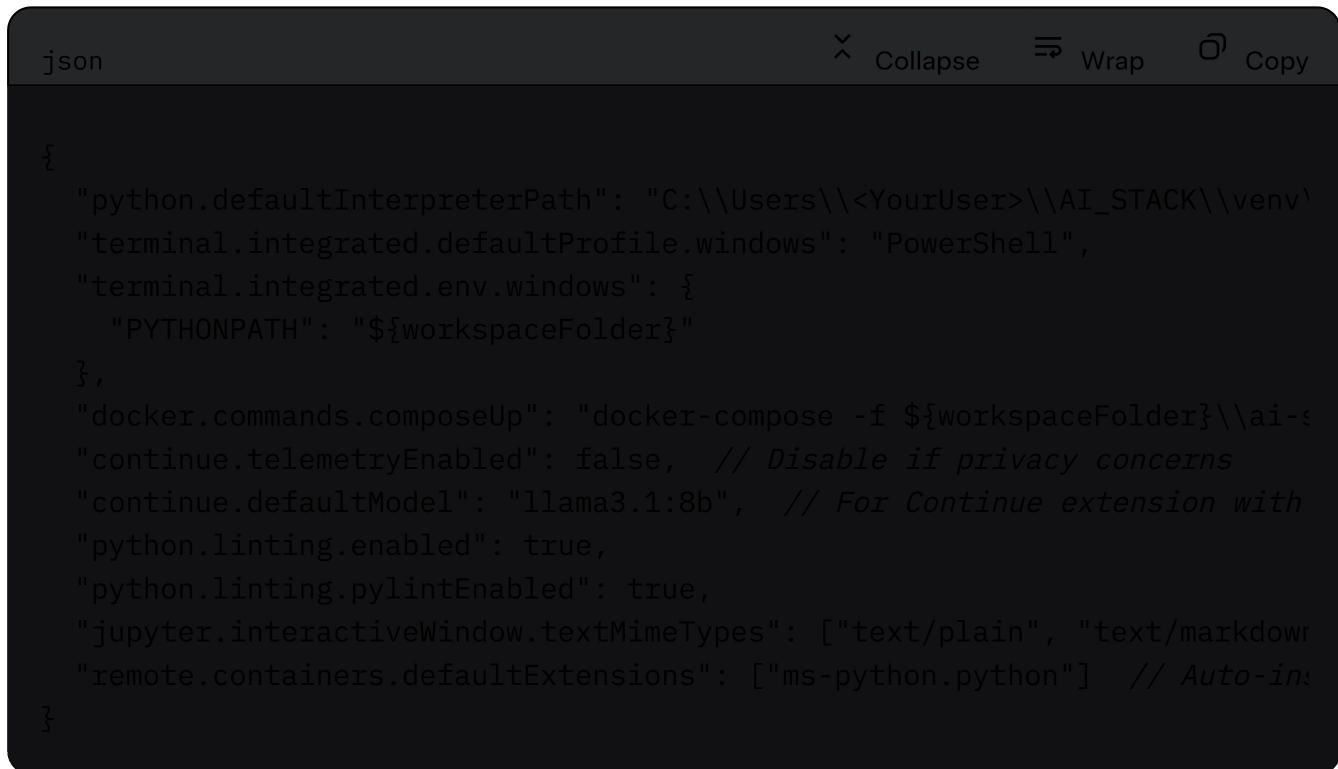


A screenshot of a JSON file in a code editor. The file contains a single object with a "recommendations" key. The value is an array of extension names. The array starts with "ms-python.python", followed by "ms-toolsai.jupyter", "ms-azuretools.vscode-docker", "eamodio.gitlens", "ms-python.vscode-pylance", "humao.rest-client", "continuedev.continue", and "ms-vscode-remote.remote-containers". The code editor has a dark theme and includes standard UI elements like a file tab, collapse, wrap, and copy buttons.

```
{  
  "recommendations": [  
    "ms-python.python",  
    "ms-toolsai.jupyter",  
    "ms-azuretools.vscode-docker",  
    "eamodio.gitlens",  
    "ms-python.vscode-pylance",  
    "humao.rest-client",  
    "continuedev.continue",  
    "ms-vscode-remote.remote-containers"  
  ]  
}
```

2. Advanced Configuration for Default Workflow

- **Workspace Settings:** Expand `settings.json` to auto-activate the virtual environment and set default terminals:



A screenshot of a code editor showing a JSON file named "settings.json". The file contains configuration for a workspace, including paths, terminal settings, and extensions like Docker and Python. The code is as follows:

```
{  
  "python.defaultInterpreterPath": "C:\\\\Users\\\\<YourUser>\\\\AI_STACK\\\\venv",  
  "terminal.integrated.defaultProfile.windows": "PowerShell",  
  "terminal.integrated.env.windows": {  
    "PYTHONPATH": "${workspaceFolder}"  
  },  
  "docker.commands.composeUp": "docker-compose -f ${workspaceFolder}\\ai-  
  "continue.telemetryEnabled": false, // Disable if privacy concerns  
  "continue.defaultModel": "llama3.1:8b", // For Continue extension with  
  "python.linting.enabled": true,  
  "python.linting.pylintEnabled": true,  
  "jupyter.interactiveWindow.textMimeTypes": ["text/plain", "text/markdown"],  
  "remote.containers.defaultExtensions": ["ms-python.python"] // Auto-in:  
}  
}
```

- **Virtual Environment Auto-Activation:** Use the Python extension's auto-activation feature —ensure the workspace is opened at `C:\Users\<YourUser>\AI_STACK` to trigger `.venv` detection.

- **Ollama Integration for AI Assistance:** Configure Continue extension to use Ollama as a

local provider. In Continue's config (`CTRL+SHIFT+F` > "Continue: Configure"), add:

json

X Collapse

Wrap

Copy

```
{  
  "models": [  
    {  
      "title": "Ollama Local",  
      "provider": "ollama",  
      "model": "llama3.1:8b",  
      "apiBase": "http://localhost:11434"  
    }  
  ]  
}
```

This allows inline code suggestions for projects like AutoGen or Haystack.

3. Expanded Tasks for All Projects

- Build on the existing `tasks.json` by adding tasks for remaining projects. Update `src/main/resources/tools/tasks.json`.

.vscode/tasks.json:

```
json X Collapse ⌂ Wrap ⌂ Copy

{
  "version": "2.0.0",
  "tasks": [
    // Existing tasks for Ollama, Text-Gen-WebUI, Whisper, Letta, MeloTTS,
    {
      "label": "Run LangChain Example",
      "type": "shell",
      "command": "python",
      "args": ["${workspaceFolder}\\rag_pipeline.py"], // Assuming a sample
      "group": "test"
    },
    {
      "label": "Start OpenSora Container",
      "type": "docker-run",
      "dependsOn": "Start Docker Services",
      "dockerRun": {
        "image": "hpcaitech/opensora:latest",
        "command": "bash"
      }
    },
    {
      "label": "Pull Ollama Model",
      "type": "shell",
      "command": "ollama pull ${input:modelName}",
      "group": "build"
    }
  ],
  "inputs": [
    {
      "id": "modelName",
      "type": "promptString",
      "description": "Ollama Model Name (e.g., deepseek-r1:7b)"
    }
  ]
}
```

- Run tasks via **Ctrl+Shift+P** > "Tasks: Run Task" to launch any project component.

4. Debugging and Integration Enhancements

- **Project-Specific Debugging:** Extend `launch.json` for more projects:

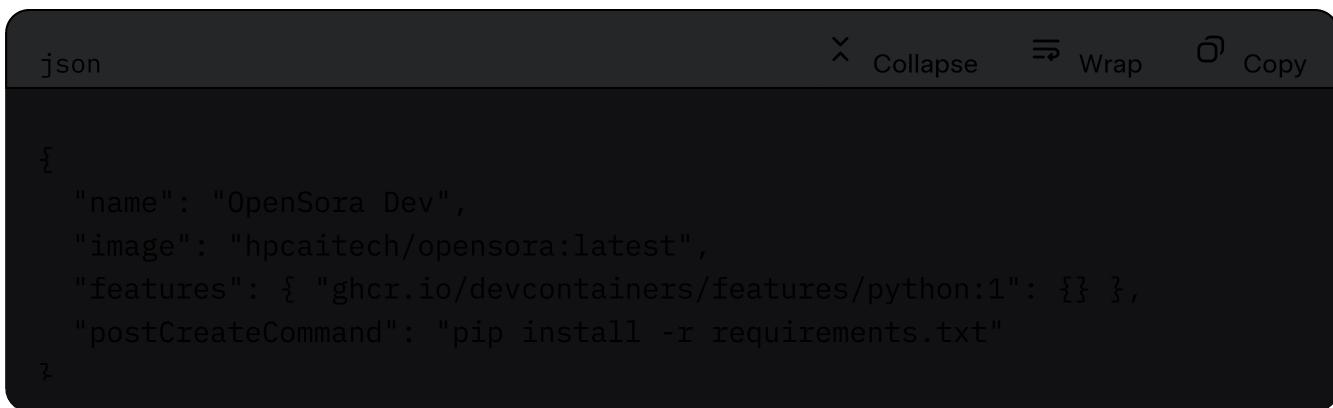
The screenshot shows a dark-themed code editor window for a JSON file named 'launch.json'. The file contains configuration for debugging LangChain scripts and Docker containers. The interface includes standard VS Code controls like 'Collapse', 'Wrap', and 'Copy' at the top right.

```
{  
  "version": "0.2.0",  
  "configurations": [  
    // Existing for Whisper and MeloTTS...  
    {  
      "name": "Debug LangChain Script",  
      "type": "python",  
      "request": "launch",  
      "program": "${file}",  
      "console": "integratedTerminal",  
      "env": { "OLLAMA_BASE_URL": "http://localhost:11434/v1" }  
    },  
    {  
      "name": "Debug in Docker (OpenManus)",  
      "type": "docker",  
      "request": "attach",  
      "platform": "node",  
      "remoteRoot": "/app" // Adjust for container path  
    }  
  ]  
}
```

- **Jupyter Notebooks for Exploration:** Create `.ipynb` files in VS Code for prototyping (e.g., Haystack RAG pipelines). Use `%load_ext autoreload` in cells for live reloading of LangChain modules.
- **Docker Integration:** Use the Docker extension to manage containers visually. For OpenSora, open the project in a Docker Container by adding a

OpenSora, open the project in a Dev Container by adding a

.devcontainer/devcontainer.json :



A screenshot of a dark-themed code editor window titled "json". The file contains a JSON object with the following content:

```
{  
  "name": "OpenSora Dev",  
  "image": "hpcositech/opensora:latest",  
  "features": { "ghcr.io/devcontainers/features/python:1": {} },  
  "postCreateCommand": "pip install -r requirements.txt"  
}
```

The editor has standard controls at the top: "Collapse", "Wrap", and "Copy".

Reopen the folder in a container via `Ctrl+Shift+P` > "Dev Containers: Reopen in Container".

5. Making VS Code the Default

- **File Associations:** Set VS Code as the default editor for Python/Docker files via Windows Settings > Apps > Default apps.
- **Startup Workflow:** Create a VS Code task to run the daily startup cheat-sheet (B-13) as a compound task grouping all service starts.
- **Best Practices:** Use GitLens for version control on repos like Text-Generation-WebUI. For efficiency, pin frequently used terminals (e.g., one for Ollama logs, one for API testing with REST Client).

These suggestions ensure VS Code handles all projects end-to-end, from code editing to container management, with AI enhancements for productivity.

Appendix B: User Manual for Utilizing the Complex Setup Correctly and Efficiently

Add this as a new appendix to the document, following the existing Appendix (now implicitly Appendix A). It expands on Section B by focusing on practical utilization, workflows, optimization, and advanced tips, assuming the stack is installed via the guide.

Appendix B: Advanced User Manual – Utilizing the AI Stack Efficiently

This appendix provides a step-by-step guide to leveraging the 11-project AI stack (with MeloTTS) for real-world tasks. It emphasizes efficient workflows, performance

optimization, and integration patterns, building on Section B's installation instructions. All examples assume you're working in VS Code (per B-15) with the virtual environment activated.

B1. Core Principles for Efficient Use

- **Modularity:** Treat projects as building blocks—e.g., Ollama as the LLM backend, LangChain for orchestration, MeloTTS for output.
- **Resource Management:** Monitor VRAM/CPU with `nvidia-smi` (GPU) or Task Manager. Use quantized models (e.g., `llama3.1:8b-q4_0`) for low-resource setups.
- **API-Centric:** All integrations use OpenAI-compatible endpoints (e.g., `http://localhost:11434/v1`) for interoperability.
- **Workflow Mindset:** Start with transcription (Whisper), process with RAG/agents (Haystack/AutoGen/Letta), reason (DeepSeek-R1), and output (MeloTTS).
- **Security:** Run services on `localhost` only; avoid exposing ports publicly.

B2. Daily Startup and Shutdown

- **Startup in VS Code:**

1. Open `C:\Users\<YourUser>\AI_STACK` in VS Code.
2. Activate venv: Open terminal (`Ctrl+``) and run `.\venv\Scripts\Activate.ps1``.
3. Run tasks sequentially: "Start Ollama" → "Start Text-Gen-WebUI" → "Start Whisper API" → "Start Letta" → "Start MeloTTS" → "Start Docker Services".

- **Verification:**

- Test endpoints: Use REST Client extension to send requests (e.g., save as `.http` file):

```
text × Collapse ⌂ Wrap ⌂ Copy

POST http://localhost:11434/v1/chat/completions
Content-Type: application/json

{
  "model": "llama3.1:8b",
  "messages": [{"role": "user", "content": "Hello"}]
}
```

- Check all services: Run `docker ps` and `netstat -ano | findstr "11434 5000 9000 8283 8001 3000"`.
- **Shutdown:** Use `Ctrl+C` in terminals or "Tasks: Terminate Task". For Docker: `docker-compose -f ai-stack-tts.yml down`.

B3. Efficient Workflows and Use Cases

- **Basic Chat/Reasoning (Ollama + DeepSeek-R1):**

- Use Text-Generation-WebUI (`http://localhost:5000`) for interactive chatting.

- In VS Code: Create `chat.py`:

The screenshot shows a dark-themed code editor window for VS Code. The title bar says "python". The editor area contains the following Python code:

```
from langchain_community.llms import Ollama
llm = Ollama(base_url="http://localhost:11434", model="deepseek-r1:7b")
response = llm("Solve 2+2")
print(response)
```

At the top of the editor, there are several icons: a close button (X), a "Collapse" button, a "Wrap" button, a "Run" button, and a "Copy" button.

Run/debug with F5.

- RAG Pipeline (Haystack + LangChain + Ollama):

1. Index documents: Use Haystack to create a vector DB.

```
python      X Collapse  ⚡ Wrap  ▶ Run  ⬇ Copy\n\nfrom haystack import Document, Pipeline\nfrom haystack.components.embedders import SentenceTransformersDocumentEmbedder\nfrom haystack_integrations.document_stores.faiss import FAISSDocumentStore\n\ndocs = [Document(content="Paris is the capital of France.")] \nstore = FAISSDocumentStore("index.faiss")\nindexing = Pipeline()\nindexing.add_component("embedder", SentenceTransformersDocumentEmbedder)\nindexing.add_component("writer", store)\nindexing.connect("embedder", "writer")\nindexing.run({"embedder": {"documents": docs}})
```

2. Query: Integrate with LangChain for generation.

```
python      X Collapse  ⚡ Wrap  ▶ Run  ⬇ Copy\n\nfrom langchain.prompts import PromptTemplate\nfrom haystack.components.generators import OpenAIGenerator\n\nprompt = PromptTemplate.from_template("Question: {query}\nContext: {cor\ngenerator = OpenAIGenerator(api_base="http://localhost:11434/v1")\nresponse = generator.run({"prompt": prompt.format(query="Capital of Fr
```

- Efficiency Tip: Share the FAISS index across sessions; use Jupyter notebooks in VS Code for iterative testing.

- Multi-Agent Collaboration (AutoGen + Letta):

- Persist memory in Letta: curl -X POST http://localhost:8283/v1/agents -d '{"agent_id": "team", "memory": {"key": "value"}'.

- Spin agents in AutoGen:

```
python      X Collapse  ⚡ Wrap  ▶ Run  ⬇ Copy

from autogen import AssistantAgent, UserProxyAgent
config = {"base_url": "http://localhost:11434/v1", "api_key": "ollama",
user = UserProxyAgent("user")
assistant = AssistantAgent("assistant", llm_config=config)
user.initiate_chat(assistant, message="Plan a trip to Paris.")
```

- Tip: Use Letta for long-term state; debug in VS Code with breakpoints on agent interactions.

- **Speech Pipeline (Whisper + MeloTTS):**

- Transcribe: curl -X POST -F file=@audio.wav http://localhost:9000/transcribe .
- Process with LLM, then TTS:

```
python      X Collapse  ⚡ Wrap  ▶ Run  ⬇ Copy

import requests
tts_data = {"text": "Response text", "voice": "EN-US"}
response = requests.post("http://localhost:8001/v1/audio/speech", json=tts_data)
with open("output.wav", "wb") as f:
    f.write(response.content)
```

- Efficiency: Batch process audio files; use GPU for faster Whisper/MeloTTS.

- **No-Code/Visual (OpenManus + Text-Generation-WebUI):**

In OpenManus UI (<https://localhost:2000>) Drag blocks for LLM chains, e.g. Ollama

- In OpenManus UI (`http://localhost:5000`): Drag blocks for LLM chains, set URLs.
- Extend with custom nodes linking to MeloTTS endpoint.
- Tip: Use VS Code's Docker extension to inspect OpenManus container logs.
- **Training/Advanced (OpenSora):**
 - Enter container: `docker exec -it <container_id> bash`.
 - Edit scripts in VS Code via Dev Containers for isolated dev.
- **End-to-End Example: Voice-Activated RAG Agent**
 1. Record query → Transcribe with Whisper.
 2. RAG query with Haystack → Generate with DeepSeek-R1 via LangChain.
 3. Store in Letta → Collaborate with AutoGen agents.
 4. Output speech with MeloTTS.
 - Script this in VS Code as a single `pipeline.py`; run with tasks.

B4. Optimization and Troubleshooting

• Performance:

- Quantize models: `llama Pull <model>.safetensors`

- Quantize models. Utilizing built-in optimizations.
 - Parallelize: Run services in background tasks; use VS Code's multi-terminal support.
 - Scale: For heavy loads, offload to GPU; monitor with `nvidia-smi`.
- **Common Issues:**
- Dependency Conflicts: Recreate venv and reinstall via `pip install -r requirements.txt` (generate with `pip freeze`).
 - API Failures: Check ports; restart services.
 - MeloTTS Latency: Use smaller models or CPU fallback if VRAM low.
- **Updates:** Run weekly script (Section C); pull new models via Ollama tasks.
- **Advanced Tips:** Use Continue extension for AI code reviews; integrate with external tools like GitHub Copilot if needed (though local Ollama is preferred for privacy).

This appendix ensures correct, efficient use of the stack, minimizing errors and maximizing productivity.