

# Multi-Spec RAG System Setup Guide (Windows 11 Native)

## Prerequisites

### System Requirements

- Windows 11 (Windows 10 also works)
- 8GB+ RAM (16GB recommended)
- 10GB+ free disk space
- Internet connection for initial setup

### Enable Developer Features (Optional but Recommended)

1. Open Settings → Privacy & Security → For developers
2. Enable Developer Mode
3. Enable PowerShell execution policy:

```
powershell  
# Run PowerShell as Administrator  
Set-ExecutionPolicy -ExecutionPolicy RemoteSigned -Scope CurrentUser
```

## Step 1: Install Python

### Option A: Using Microsoft Store (Recommended)

1. Open Microsoft Store
2. Search for "Python 3.11" or "Python 3.12"
3. Install Python (by Python Software Foundation)

### Option B: Using Winget (Command Line)

```
powershell  
# Open PowerShell and run:  
winget install Python.Python.3.11
```

### Option C: Direct Download

1. Go to [python.org/downloads](https://python.org/downloads)
2. Download Python 3.11.x or 3.12.x
3. **Important:** Check "Add Python to PATH" during installation

## Verify Python Installation

```
powershell  
  
# Open Command Prompt or PowerShell  
python --version  
pip --version
```

---

## Step 2: Install Ollama (Local LLM Runtime)

### Download and Install Ollama

1. **Visit:** [ollama.ai/download](https://ollama.ai/download)
2. **Download:** Ollama for Windows
3. **Install:** Run the installer (default settings are fine)
4. **Verify:** Ollama should start automatically

### Or Use Winget

```
powershell  
  
winget install Ollama.Ollama
```

## Pull LLM Models

```
powershell  
  
# Open Command Prompt or PowerShell  
  
# For 8GB+ RAM systems  
ollama pull llama3.2:3b  
  
# For 16GB+ RAM systems (better performance)  
ollama pull llama3.1:8b  
  
# Lightweight for testing/slower systems  
ollama pull llama3.2:1b  
  
# Verify installation  
ollama list
```

---

## Step 3: Create Project Structure

### Create Project Directory

```
powershell
```

```
# Open Command Prompt or PowerShell
cd %USERPROFILE%
mkdir hardware-rag-system
cd hardware-rag-system
```

## Create Folder Structure

```
powershell
```

```
mkdir specs
mkdir src
mkdir data
mkdir logs
```

## Create Python Virtual Environment

```
powershell
```

```
# Create virtual environment
python -m venv rag-env
```

```
# Activate virtual environment
rag-env\Scripts\activate
```

```
# Upgrade pip
python -m pip install --upgrade pip
```

## Step 4: Install Python Dependencies

### Create requirements.txt

```
powershell
```

```
# Create requirements file
@"
# Core RAG dependencies
pymupdf==1.23.14
sentence-transformers==2.2.2
chromadb==0.4.18
ollama==0.1.7

# Data processing
numpy==1.24.3
pandas==2.0.3

# Utilities
python-dotenv==1.0.0
tqdm==4.66.1
colorama==0.4.6

# Optional: For advanced PDF processing
pdfplumber==0.9.0
pypdf==3.0.1
"@ | Out-File -FilePath requirements.txt -Encoding utf8
```

## Install Dependencies

```
powershell
```

```
# Make sure virtual environment is activated
rag-env\Scripts\activate

# Install all dependencies
pip install -r requirements.txt

# Install PyTorch (CPU version for Windows)
pip install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cpu
```

## Verify Installations

```
powershell
```

```
python -c "import sentence_transformers; print('✓ Sentence Transformers installed')"
python -c "import chromadb; print('✓ ChromaDB installed')"
python -c "import fitz; print('✓ PyMuPDF installed')"
python -c "import ollama; print('✓ Ollama client installed')"
```

## Step 5: Create Configuration Files

### Create config.py

```
powershell

@"
""""Configuration settings for the RAG system (Windows)"""
import os
from pathlib import Path

# Paths (Windows-compatible)
PROJECT_ROOT = Path(__file__).parent.parent
SPECS_DIR = PROJECT_ROOT / "specs"
DB_DIR = PROJECT_ROOT / "data" / "vector_db"
LOGS_DIR = PROJECT_ROOT / "logs"

# Embedding model settings
EMBEDDING_MODEL = "all-MiniLM-L6-v2" # Fast, good quality
# Alternative: "all-mpnet-base-v2" # Slower, better quality

# LLM settings
DEFAULT_LLM_MODEL = "llama3.2:3b"
# Alternative models: "llama3.2:1b", "llama3.1:8b"

# Chunking settings
MAX_CHUNK_SIZE = 1000
CHUNK_OVERLAP = 100

# Query settings
DEFAULT_RESULTS_PER_SPEC = 3
MAX_CONTEXT_LENGTH = 4000

# Create directories (Windows-safe)
DB_DIR.mkdir(parents=True, exist_ok=True)
LOGS_DIR.mkdir(parents=True, exist_ok=True)
"@ | Out-File -FilePath src\config.py -Encoding utf8
```

### Create spec download helper

```
powershell
```

```
@"
#!/usr/bin/env python3
"""

Script to help organize specification PDFs (Windows)
"""

import os
from pathlib import Path

def create_spec_structure():
    specs_dir = Path("specs")
    specs_dir.mkdir(exist_ok=True)

    # Create info file about where to get specs
    info_content = ""

    # Hardware Specification Sources

    ## SystemVerilog (IEEE 1800)
    - Source: IEEE Xplore Digital Library
    - Search: "IEEE 1800-2017" or "IEEE 1800-2023"
    - File: Save as 'ieee_1800_systemverilog.pdf'

    ## VHDL (IEEE 1076)
    - Source: IEEE Xplore Digital Library
    - Search: "IEEE 1076-2019"
    - File: Save as 'ieee_1076_vhdl.pdf'

    ## Ethernet (IEEE 802.3)
    - Source: IEEE Xplore Digital Library
    - Search: "IEEE 802.3-2022"
    - File: Save as 'ieee_802_3_ethernet.pdf'

    ## USB Specification
    - Source: USB Implementers Forum (usb.org)
    - Document: "Universal Serial Bus Specification Revision 3.2"
    - File: Save as 'usb_specification.pdf'

    ## PCIe Specification
    - Source: PCI-SIG (pcisig.com)
    - Document: "PCI Express Base Specification Revision 6.0"
    - File: Save as 'pcie_specification.pdf'
```

Note: Some specs require IEEE membership or purchase.  
For testing, you can use any technical PDF document.

m

```
with open(specs_dir / "README.md", "w", encoding='utf-8') as f:  
    f.write(info_content)  
  
print("Created specs directory structure")  
print("Please download your specification PDFs to the specs/ directory")  
print("See specs/README.md for download sources")  
  
if __name__ == "__main__":  
    create_spec_structure()  
"@ | Out-File -FilePath src\download_specs.py -Encoding utf8
```

## Step 6: Create the RAG System Files

### Create Windows-Optimized Multi-Spec RAG System

```
powershell
```

```
# This creates the main RAG system file
# Copy the multi_spec_rag.py code but with Windows path fixes
@"
#!/usr/bin/env python3
"""

Multi-Specification RAG System - Windows Optimized
"""

import fitz # PyMuPDF
import re
from typing import List, Dict, Tuple, Optional
import chromadb
from sentence_transformers import SentenceTransformer
from dataclasses import dataclass
import json
import os
from pathlib import Path
import ollama # For local LLM integration

@dataclass
class SpecChunk:
    text: str
    section: str
    title: str
    page: int
    chapter: str
    spec_type: str

class MultiSpecRAGSystem:
    def __init__(self, db_path: str = "./data/vector_db"):
        self.db_path = Path(db_path)
        self.db_path.mkdir(parents=True, exist_ok=True)

        print("Loading embedding model...")
        self.embedding_model = SentenceTransformer('all-MiniLM-L6-v2')

    # Initialize ChromaDB client
    self.chroma_client = chromadb.PersistentClient(path=str(self.db_path))

    # Define spec configurations
    self.spec_configs = {
        'systemverilog': {
            'collection_name': 'ieee_1800_systemverilog',
            'description': 'IEEE 1800 SystemVerilog Specification',
            'location': 'https://www.ieee.org/standards/mission-critical-computing/ieee-1800-system-verilog-specification.html'
        }
    }
```

```

'section_pattern': r'^(\d+\.\d+(?:\.\d+)?)\s+([^\n]+)',  

'chapter_pattern': r'^(\d+)\s+([A-Z][^\n]+)'  
},  

'vhdl': {  

    'collection_name': 'ieee_1076_vhdl',  

    'description': 'IEEE 1076 VHDL Specification',  

    'section_pattern': r'^(\d+\.\d+(?:\.\d+)?)\s+([^\n]+)',  

    'chapter_pattern': r'^(\d+)\s+([A-Z][^\n]+)'  
},  

'ethernet': {  

    'collection_name': 'ieee_802_3_ethernet',  

    'description': 'IEEE 802.3 Ethernet Specification',  

    'section_pattern': r'^(\d+\.\d+(?:\.\d+)?)\s+([^\n]+)',  

    'chapter_pattern': r'^(\d+)\s+([A-Z][^\n]+)'  
},  

'usb': {  

    'collection_name': 'usb_specification',  

    'description': 'USB Specification',  

    'section_pattern': r'^(\d+\.\d+(?:\.\d+)?)\s+([^\n]+)',  

    'chapter_pattern': r'^(\d+)\s+([A-Z][^\n]+)'  
},  

'pcie': {  

    'collection_name': 'pcie_specification',  

    'description': 'PCIe Specification',  

    'section_pattern': r'^(\d+\.\d+(?:\.\d+)?)\s+([^\n]+)',  

    'chapter_pattern': r'^(\d+)\s+([A-Z][^\n]+)'  
}  

}
}

```

```

# Initialize collections
self.collections = {}
for spec_type, config in self.spec_configs.items():
    self.collections[spec_type] = self.chroma_client.get_or_create_collection(
        name=config['collection_name'],
        metadata={"description": config['description'], "spec_type": spec_type})
)

```

```

def extract_text_with_structure(self, pdf_path: str, spec_type: str) -> List[SpecChunk]:  

    """Extract text while preserving document structure for any spec type"""
    doc = fitz.open(pdf_path)
    chunks = []
    current_chapter = ""
    config = self.spec_configs[spec_type]

    for page_num in range(min(len(doc), 100)): # Limit for testing
        page = doc[page_num]
        text = page.get_text()

```

```
# Look for chapter headers
chapter_match = re.search(config['chapter_pattern'], text, re.MULTILINE)
if chapter_match:
    current_chapter = f'{chapter_match.group(1)} {chapter_match.group(2)}'

# Look for section headers
sections = re.finditer(config['section_pattern'], text, re.MULTILINE)

prev_end = 0
prev_section_num = ""
prev_section_title = ""

for section_match in sections:
    section_num = section_match.group(1)
    section_title = section_match.group(2).strip()

    # Extract text from previous section
    if prev_end > 0:
        section_text = text[prev_end:section_match.start()].strip()
        if len(section_text) > 100:
            chunks.append(SpecChunk(
                text=section_text,
                section=prev_section_num,
                title=prev_section_title,
                page=page_num + 1,
                chapter=current_chapter,
                spec_type=spec_type
            ))
    prev_section_num = section_num
    prev_section_title = section_title
    prev_end = section_match.end()

# Handle remaining text on page
if prev_end > 0:
    remaining_text = text[prev_end:].strip()
    if len(remaining_text) > 100:
        chunks.append(SpecChunk(
            text=remaining_text,
            section=prev_section_num,
            title=prev_section_title,
            page=page_num + 1,
            chapter=current_chapter,
            spec_type=spec_type
        ))
```

```

doc.close()
return chunks

def intelligent_chunking(self, chunks: List[SpecChunk], max_chunk_size: int = 1000) -> List[SpecChunk]:
    """Break large chunks into smaller ones while preserving context"""
    refined_chunks = []

    for chunk in chunks:
        if len(chunk.text) <= max_chunk_size:
            refined_chunks.append(chunk)
            continue

        paragraphs = chunk.text.split('\n\n')
        current_text = ""

        for para in paragraphs:
            if len(current_text + para) > max_chunk_size and current_text:
                refined_chunks.append(SpecChunk(
                    text=current_text.strip(),
                    section=chunk.section,
                    title=chunk.title,
                    page=chunk.page,
                    chapter=chunk.chapter,
                    spec_type=chunk.spec_type
                ))
                # Add overlap
                current_text = current_text.split('\n')[-1] + '\n' + para
            else:
                current_text += '\n\n' + para if current_text else para

        if current_text.strip():
            refined_chunks.append(SpecChunk(
                text=current_text.strip(),
                section=chunk.section,
                title=chunk.title,
                page=chunk.page,
                chapter=chunk.chapter,
                spec_type=chunk.spec_type
            ))

    return refined_chunks

```

```

def load_spec(self, pdf_path: str, spec_type: str):
    """Load a single specification into the appropriate collection"""
    if spec_type not in self.spec_configs:
        raise ValueError(f"Unknown spec type: {spec_type}")

```

```

        raise ValueError(f"Unknown spec type: {spec_type!r}!")

print(f"Processing {spec_type.upper()} specification...")

# Extract and chunk
raw_chunks = self.extract_text_with_structure(pdf_path, spec_type)
refined_chunks = self.intelligent_chunking(raw_chunks)

# Prepare data
documents = [chunk.text for chunk in refined_chunks]
metadata = [
    {
        "section": chunk.section,
        "title": chunk.title,
        "page": chunk.page,
        "chapter": chunk.chapter,
        "spec_type": chunk.spec_type,
        "length": len(chunk.text)
    }
    for chunk in refined_chunks
]
ids = [f"{spec_type}_chunk_{i:05d}" for i in range(len(refined_chunks))]

# Generate embeddings
print(f"Generating embeddings for {len(documents)} chunks...")
embeddings = self.embedding_model.encode(documents, show_progress_bar=True)

# Store in appropriate collection
collection = self.collections[spec_type]
collection.add(
    documents=documents,
    metadata=metadata,
    embeddings=embeddings.tolist(),
    ids=ids
)

print(f"Successfully loaded {len(refined_chunks)} chunks for {spec_type}!")
return len(refined_chunks)

def query_single_spec(self, question: str, spec_type: str, n_results: int = 3) -> List[Dict]:
    """Query a specific specification"""
    if spec_type not in self.collections:
        raise ValueError(f"Spec type {spec_type} not loaded")

    collection = self.collections[spec_type]
    results = collection.query(
        query_texts=[question],

```

```

    n_results=n_results,
    include=['documents', 'metadatas', 'distances']
)

return self._format_results(results)

def smart_query(self, question: str) -> Dict[str, List[Dict]]:
    """Intelligently determine which specs to query based on question content"""
    question_lower = question.lower()

    # Keywords to identify relevant specs
    spec_keywords = {
        'systemverilog': ['systemverilog', 'sv', 'logic', 'module', 'interface', 'class', 'package'],
        'vhdl': ['vhdl', 'entity', 'architecture', 'signal', 'process', 'component'],
        'ethernet': ['ethernet', 'mac', 'phy', '802.3', 'csma', 'collision'],
        'usb': ['usb', 'endpoint', 'descriptor', 'enumeration', 'bulk', 'interrupt'],
        'pcie': ['pcie', 'pci express', 'tlp', 'transaction layer', 'link training']
    }

    # Determine relevant specs
    relevant_specs = []
    for spec_type, keywords in spec_keywords.items():
        if any(keyword in question_lower for keyword in keywords):
            relevant_specs.append(spec_type)

    # If no specific match, query all specs
    if not relevant_specs:
        return self.query_all_specs(question, n_results_per_spec=1)

    # Query only relevant specs
    results = {}
    for spec_type in relevant_specs:
        if spec_type in self.collections:
            results[spec_type] = self.query_single_spec(question, spec_type, n_results=3)

    return results

def query_all_specs(self, question: str, n_results_per_spec: int = 2) -> Dict[str, List[Dict]]:
    """Query all loaded specifications"""
    all_results = {}

    for spec_type in self.collections:
        try:
            results = self.query_single_spec(question, spec_type, n_results_per_spec)
            if results: # Only include if we got results
                all_results[spec_type] = results
        except Exception as e:
            print(f"Error querying {spec_type}: {e}")

```

```
except Exception as e:  
    print(f"Error querying {spec_type}: {e}")  
  
return all_results  
  
def _format_results(self, results) -> List[Dict]:  
    """Format ChromaDB results"""  
    formatted = []  
    for i in range(len(results['documents'][0])):  
        formatted.append({  
            'text': results['documents'][0][i],  
            'metadata': results['metadatas'][0][i],  
            'similarity_score': 1 - results['distances'][0][i]  
        })  
    return formatted
```

```
def generate_rag_response(self, question: str, model_name: str = "llama3.2:3b") -> str:  
    """Generate response using local LLM with RAG context"""  
    # Get relevant context  
    context_results = self.smart_query(question)
```

```
if not context_results:  
    return "I couldn't find relevant information in the loaded specifications."
```

```
# Build context string  
context_parts = []  
for spec_type, results in context_results.items():  
    context_parts.append(f"\n== {spec_type.upper()} SPECIFICATION ==")  
    for i, result in enumerate(results[:2], 1): # Limit to top 2 per spec  
        metadata = result['metadata']  
        context_parts.append(  
            f"Section {metadata['section']}: {metadata['title']}\n"  
            f"{result['text'][:800]}...\n")  
    )
```

```
context = "\n".join(context_parts)
```

```
# Create prompt  
prompt = f"""Based on the following hardware specification excerpts, please answer the question.
```

CONTEXT FROM SPECIFICATIONS:

{context}

QUESTION: {question}

Please provide a comprehensive answer based on the specification context above. If the context doesn't contain enough

```

try:
    # Use Ollama for local LLM inference
    response = ollama.chat(model=model_name, messages=[
        {'role': 'user', 'content': prompt}
    ])
    return response['message']['content']
except Exception as e:
    return f"Error generating response: {e}"

# Usage functions
def load_all_specs(rag_system: MultiSpecRAGSystem, specs_directory: str):
    """Load all specifications from a directory"""
    spec_files = {
        'systemverilog': 'ieee_1800_systemverilog.pdf',
        'vhdl': 'ieee_1076_vhdl.pdf',
        'ethernet': 'ieee_802_3_ethernet.pdf',
        'usb': 'usb_specification.pdf',
        'pcie': 'pcie_specification.pdf'
    }

    loaded_specs = []
    specs_path = Path(specs_directory)

    for spec_type, filename in spec_files.items():
        filepath = specs_path / filename
        if filepath.exists():
            try:
                chunks = rag_system.load_spec(str(filepath), spec_type)
                loaded_specs.append((spec_type, chunks))
                print(f"✓ Loaded {spec_type}: {chunks} chunks")
            except Exception as e:
                print(f"✗ Failed to load {spec_type}: {e}")
        else:
            print(f"⚠ File not found: {filepath}")

    return loaded_specs
"""

@ | Out-File -FilePath src\multi_spec_rag.py -Encoding utf8

```

## Create Main Application

```
powershell
```

```
@"
#!/usr/bin/env python3
"""

Enhanced Multi-Spec RAG System - Windows Version
"""

import sys
import os
from pathlib import Path

# Add src directory to path
sys.path.append(str(Path(__file__).parent))

from multi_spec_rag import MultiSpecRAGSystem, load_all_specs
from config import *
import logging

# Setup logging
logging.basicConfig(
    level=logging.INFO,
    format='%(asctime)s - %(levelname)s - %(message)s',
    handlers=[
        logging.FileHandler(LOGS_DIR / 'rag_system.log'),
        logging.StreamHandler()
    ]
)

def main():
    print("🚀 Hardware Specifications RAG System (Windows)")
    print("≡" * 60)

    # Initialize system
    try:
        rag_system = MultiSpecRAGSystem(str(DB_DIR))
        logging.info("RAG system initialized successfully")
    except Exception as e:
        logging.error(f"Failed to initialize RAG system: {e}")
        return

    # Check for spec files
    spec_files = list(SPECS_DIR.glob("*.pdf"))
    if not spec_files:
        print("⚠️ No PDF files found in {SPECS_DIR}")
        print("Please add your specification PDFs to the specs/ directory")
        print("Run: python src/download_specs.py for more info")
```

```
print(f"📁 Found {len(spec_files)} PDF files:")
for pdf_file in spec_files:
    print(f" - {pdf_file.name}")
```

```
# Load specifications
print("\n🔄 Loading specifications...")
loaded_specs = load_all_specs(rag_system, str(SPECS_DIR))
```

```
if not loaded_specs:
    print("❌ No specifications could be loaded")
    return
```

```
print(f"\n✅ Successfully loaded {len(loaded_specs)} specifications!")
```

```
# Interactive mode
interactive_mode(rag_system)
```

```
def interactive_mode(rag_system):
    """Interactive query interface"""
    print("\n" + "=" * 60)
    print("🎯 INTERACTIVE MODE - Ask questions about your specs!")
    print("=" * 60)
    print("Commands:")
    print(" help - Show help")
    print(" specs - List loaded specifications")
    print(" quit - Exit system")
    print(" Or just ask a question!")
```

```
while True:
```

```
    try:
```

```
        user_input = input("\n📝 Your question: ").strip()
```

```
        if not user_input:
```

```
            continue
```

```
        if user_input.lower() in ['quit', 'exit', 'q']:
```

```
            print("👋 Goodbye!")
```

```
            break
```

```
        elif user_input.lower() == 'help':
```

```
            show_help()
```

```
            continue
```

```
        elif user_input.lower() == 'specs':
```

```
            show_loaded_specs(rag_system)
```

continue

```
# Process question
print("\n🤔 Thinking...")
try:
    response = rag_system.generate_rag_response(user_input)
    print(f"\n🎯 Answer:\n{response}")
except Exception as e:
    print("❌ Error generating response: {e}")
    logging.error("Query error: {e}")

except KeyboardInterrupt:
    print("\n👋 Goodbye!")
    break
except Exception as e:
    print("❌ Unexpected error: {e}")
    logging.error("Unexpected error: {e}")
```

```
def show_help():
    """Show help information"""
    help_text = """
🔍 Question Examples:
"How do I declare a SystemVerilog interface?"
"What are USB endpoint types?"
"Explain Ethernet frame format"
"What is PCIe TLP?"
"VHDL process sensitivity list rules"
```

#### 💡 Tips:

- Be specific in your questions
- Mention the technology (SystemVerilog, USB, etc.) **for** better results
- Ask about syntax, concepts, or implementation details

.....

```
print(help_text)
```

```
def show_loaded_specs(rag_system):
    """Show information about loaded specifications"""
    print("\n🧾 Loaded Specifications:")
    for spec_type in rag_system.collections:
        collection = rag_system.collections[spec_type]
        count = collection.count()
        print(f" {spec_type.upper()}: {count} chunks")

if __name__ == "__main__":
    main()
@ | Out-File -FilePath src\rag_system.py -Encoding utf8
```

## Step 7: Create Convenience Scripts

### Create Startup Batch File

```
powershell

@"
@echo off
cd /d %~dp0
echo Starting Hardware RAG System...

rem Activate virtual environment
call rag-env\Scripts\activate.bat

rem Check if Ollama is running
tasklist /FI "IMAGENAME eq ollama.exe" 2>NUL | find /I /N "ollama.exe">NUL
if "%ERRORLEVEL%"=="1" (
    echo Starting Ollama...
    start "" ollama serve
    timeout /t 3 >nul
)

rem Start RAG system
python src\rag_system.py

pause
"@ | Out-File -FilePath start_rag.bat -Encoding ascii
```

### Create Test Script

```
powershell
```

```
@"
@echo off
cd /d %~dp0
call rag-env\Scripts\activate.bat

echo Testing system components...

echo.
echo 1. Testing Python packages...
python -c "import sentence_transformers; print('✓ Sentence Transformers')"
python -c "import chromadb; print('✓ ChromaDB')"
python -c "import fitz; print('✓ PyMuPDF')"
python -c "import ollama; print('✓ Ollama client')"

echo.
echo 2. Testing Ollama connection...
ollama list

echo.
echo 3. Testing spec directory...
python src\download_specs.py

echo.
echo System test complete!
pause
"@ | Out-File -FilePath test_system.bat -Encoding ascii
```

## Step 8: Test the Installation

### Run System Test

```
powershell
```

```
# Double-click test_system.bat or run:
.\test_system.bat
```

## Download Test PDFs

### 1. Add PDFs to specs folder:

- Place your specification PDFs in the `specs\` folder
- Name them according to the expected filenames
- For testing, any technical PDF will work

## First Run

```
powershell
```

```
# Double-click start_rag.bat or run:  
.\\start_rag.bat
```

---

## Step 9: Usage

### Starting the System

- **Option 1:** Double-click `start_rag.bat`
- **Option 2:** Open Command Prompt:

```
powershell
```

```
cd hardware-rag-system  
rag-env\\Scripts\\activate  
python src\\rag_system.py
```

### Example Questions

Once running, try asking:

- "How do I declare a SystemVerilog interface?"
- "What are USB descriptor types?"
- "Explain Ethernet collision detection"
- "What is PCIe link training?"

---

## Troubleshooting

### Common Windows Issues

#### 1. Python not found:

```
powershell
```

```
# Add Python to PATH manually  
setx PATH "%PATH%;C:\\Users\\%USERNAME%\\AppData\\Local\\Programs\\Python\\Python311"
```

#### 2. Ollama not starting:

```
powershell
```

```
# Check if Ollama service is running
```

```
tasklist | findstr ollama
```

```
# Restart Ollama
```

```
taskkill /F /IM ollama.exe
```

```
ollama serve
```

### 3. Virtual environment issues:

```
powershell
```

```
# Recreate virtual environment
```

```
rmdir /s rag-env
```

```
python -m venv rag-env
```

```
rag-env\Scripts\activate
```

```
pip install -r requirements.txt
```

### 4. ChromaDB permission errors:

```
powershell
```

```
# Run Command Prompt as Administrator
```

```
# Or move project to a folder with full permissions
```

### 5. PDF processing errors:

```
powershell
```

```
# Install Visual C++ Redistributable if needed
```

```
# Download from Microsoft website
```

## Performance Tips for Windows

### CPU Optimization

- **Use smaller models** on older systems: `ollama pull llama3.2:1b`
- **Close unnecessary programs** during initial spec loading
- **Use SSD storage** for better I/O performance

### Memory Management

- **Monitor Task Manager** during operation
- **Increase virtual memory** if needed (System → Advanced → Performance Settings)
- **Use 64-bit Python** (should be default on modern systems)

## Windows-Specific Optimizations

```
powershell

# Create performance config
@"
# Windows-optimized settings
EMBEDDING_BATCH_SIZE = 16 # Smaller batches for Windows
MAX_CHUNK_SIZE = 800      # Reduce memory usage
CACHE_DIR = "./cache"     # Local cache directory
"@ | Out-File -FilePath src\windows_config.py -Encoding utf8
```

## Step 10: Advanced Windows Features

### Windows Terminal Integration

1. **Install Windows Terminal** (from Microsoft Store)
2. **Create profile** for RAG system:

```
json

{
  "name": "RAG System",
  "commandline": "cmd.exe /k \"cd /d C:\\\\Users\\\\%USERNAME%\\\\hardware-rag-system && rag-env\\\\Scripts\\\\activate\"",
  "startingDirectory": "C:\\\\Users\\\\%USERNAME%\\\\hardware-rag-system",
  "icon": "🤖"
}
```

### Task Scheduler Automation

1. **Open Task Scheduler**
2. **Create Basic Task**:
  - Name: "Start Ollama Service"
  - Trigger: "At startup"
  - Action: "Start a program"
  - Program:
  - Arguments:

## Windows Defender Exclusions

```
powershell
```

```
# Run as Administrator to add exclusions for better performance
Add-MpPreference -ExclusionPath "C:\Users\$env:USERNAME\hardware-rag-system"
Add-MpPreference -ExclusionProcess "ollama.exe"
Add-MpPreference -ExclusionProcess "python.exe"
```

## Step 11: Desktop Integration

### Create Desktop Shortcuts

```
powershell
```

```
# Create desktop shortcut for RAG system
$WshShell = New-Object -comObject WScript.Shell
$Shortcut = $WshShell.CreateShortcut("$Home\Desktop\Hardware RAG System.lnk")
$Shortcut.TargetPath = "$Home\hardware-rag-system\start_rag.bat"
$Shortcut.WorkingDirectory = "$Home\hardware-rag-system"
$Shortcut.IconLocation = "shell32.dll,25"
$Shortcut.Description = "Hardware Specifications RAG System"
$Shortcut.Save()
```

```
Write-Host "Desktop shortcut created!"
```

### Windows Context Menu Integration

```
powershell
```

```
# Add "Query with RAG" to PDF context menu (optional)
$regPath = "HKEY_CLASSES_ROOT\pdf\shell\QueryRAG"
$regCommand = "$Home\hardware-rag-system\query_pdf.bat"

# This would require a separate query_pdf.bat script
```

## Step 12: Maintenance and Updates

### Create Update Script

```
powershell
```

```
@"
@echo off
echo Updating Hardware RAG System...
cd /d %~dp0
call rag-env\Scripts\activate.bat
```

```
echo Updating Python packages...
pip install --upgrade -r requirements.txt
```

```
echo Checking for Ollama updates...
ollama --version
```

```
echo Updating models...
ollama pull llama3:2:3b
```

```
echo Update complete!
pause
"@ | Out-File -FilePath update_system.bat -Encoding ascii
```

## Create Backup Script

```
powershell
```

```
@"
@echo off
set BACKUP_DIR=%USERPROFILE%\Documents\RAG_Backups
set TIMESTAMP=%date:~-4,4%%date:~-10,2%%date:~-7,2%
```

```
echo Creating backup of RAG system...
mkdir "%BACKUP_DIR%\%TIMESTAMP%" 2>nul
```

```
echo Backing up vector database...
xcopy "data\vector_db" "%BACKUP_DIR%\%TIMESTAMP%\vector_db" /E /I /Q
```

```
echo Backing up configuration...
copy "src\*.py" "%BACKUP_DIR%\%TIMESTAMP%\\" /Y
```

```
echo Backup complete: %BACKUP_DIR%\%TIMESTAMP%
pause
"@ | Out-File -FilePath backup_system.bat -Encoding ascii
```

## Step 13: GUI Integration (Optional)

# **Simple Tkinter GUI**

```
powershell
```

```
@"
#!/usr/bin/env python3
"""

Simple GUI for the RAG system (Windows)
"""

import tkinter as tk
from tkinter import scrolledtext, messagebox
import threading
import sys
from pathlib import Path

# Add src to path
sys.path.append(str(Path(__file__).parent / 'src'))

from multi_spec_rag import MultiSpecRAGSystem, load_all_specs
from config import *

class RAGGUI:

    def __init__(self):
        self.rag_system = None
        self.setup_gui()
        self.load_system()

    def setup_gui(self):
        self.root = tk.Tk()
        self.root.title("Hardware Specifications RAG System")
        self.root.geometry("800x600")

        # Question input
        tk.Label(self.root, text="Ask a question about hardware specifications:",
                font=("Arial", 12)).pack(pady=10)

        self.question_entry = tk.Entry(self.root, width=80, font=("Arial", 11))
        self.question_entry.pack(pady=5)
        self.question_entry.bind('<Return>', self.on_ask)

    # Ask button
    self.ask_button = tk.Button(self.root, text="Ask Question",
                               command=self.on_ask, font=("Arial", 11))
    self.ask_button.pack(pady=5)

    # Response area
    tk.Label(self.root, text="Response:", font=("Arial", 12)).pack(pady=(20,5))
```

```
self.response_text = scrolledtext.ScrolledText(  
    self.root, width=90, height=25, font=("Consolas", 10))  
self.response_text.pack(pady=5, padx=10, fill=tk.BOTH, expand=True)  
  
# Status bar  
self.status_label = tk.Label(self.root, text="Loading system...",  
    relief=tk.SUNKEN, anchor=tk.W)  
self.status_label.pack(side=tk.BOTTOM, fill=tk.X)  
  
def load_system(self):  
    """Load the RAG system in background"""  
def load_bg():  
    try:  
        self.rag_system = MultiSpecRAGSystem(str(DB_DIR))  
  
        # Check for specs  
        spec_files = list(SPECS_DIR.glob("*.pdf"))  
        if spec_files:  
            load_all_specs(self.rag_system, str(SPECS_DIR))  
            self.root.after(0, lambda: self.status_label.config(  
                text=f"Ready! Loaded {len(spec_files)} specifications"))  
        else:  
            self.root.after(0, lambda: self.status_label.config(  
                text="Ready! No specifications loaded - add PDFs to specs/ folder"))  
  
    except Exception as e:  
        self.root.after(0, lambda: self.status_label.config(  
            text=f"Error loading system: {e}"))  
  
    threading.Thread(target=load_bg, daemon=True).start()  
  
def on_ask(self, event=None):  
    """Handle question submission"""  
    question = self.question_entry.get().strip()  
    if not question:  
        return  
  
    if not self.rag_system:  
        messagebox.showerror("Error", "System not loaded yet!")  
        return  
  
    # Clear previous response  
    self.response_text.delete(1.0, tk.END)  
    self.response_text.insert(tk.END, "Thinking...\n")  
    self.ask_button.config(state=tk.DISABLED)  
    self.status_label.config(text="Processing question...")
```

```

def process_bg():
    try:
        response = self.rag_system.generate_rag_response(question)
        self.root.after(0, lambda: self.display_response(response))
    except Exception as e:
        self.root.after(0, lambda: self.display_response(f"Error: {e}"))

    threading.Thread(target=process_bg, daemon=True).start()

def display_response(self, response):
    """Display the response"""
    self.response_text.delete(1.0, tk.END)
    self.response_text.insert(tk.END, response)
    self.ask_button.config(state=tk.NORMAL)
    self.status_label.config(text="Ready")

def run(self):
    self.root.mainloop()

if __name__ == "__main__":
    app = RAGGUI()
    app.run()
"@ | Out-File -FilePath rag_gui.py -Encoding utf8

```

## Create GUI Launcher

```

powershell

@"
@echo off
cd /d %~dp0
call rag-env\Scripts\activate.bat

rem Install GUI dependencies if needed
pip install tkinter 2>nul

rem Start GUI
python rag_gui.py
"@ | Out-File -FilePath start_gui.bat -Encoding ascii

```

## Comparison: Windows vs WSL2

Feature	Native Windows	WSL2
<b>Setup Complexity</b>	★★★☆ (Simple)	★★★★☆ (Moderate)
<b>Performance</b>	★★★★★ (Good)	★★★★☆ (Slightly slower)
<b>NPU Access</b>	★★★☆☆ (Limited)	★★★★☆ (Better)
<b>File Integration</b>	★★★★★☆ (Excellent)	★★★★☆ (Bridge needed)
<b>GUI Support</b>	★★★★★☆ (Native)	★★★☆ (X11/WSLg)
<b>Maintenance</b>	★★★★☆☆ (Easy)	★★★★☆ (Moderate)

## Final Steps and Tips

### Quick Start Checklist

- Python 3.11+ installed
- Ollama installed and running
- Virtual environment created
- Dependencies installed
- PDF specifications in specs/ folder
- First run successful

### Daily Usage

1. **Double-click** `start_rag.bat`
2. **Wait** for "Ready!" message
3. **Ask questions** about your hardware specs
4. **Type** `quit` to exit

### Best Practices

- **Keep PDFs organized** in the specs/ folder
- **Use descriptive questions** for better results
- **Restart Ollama** if responses seem slow
- **Monitor RAM usage** during heavy queries
- **Backup your vector database** periodically

### Getting Help

- **Check logs** in `logs/rag_system.log`
- **Run test script** if issues arise
- **Restart system** for memory cleanup
- **Update regularly** using `update_system.bat`

# Success! 🎉

You now have a complete, Windows-native RAG system that can:

- Process multiple hardware specifications
- Provide instant, contextual answers
- Run entirely offline after setup
- Keep your proprietary specs private
- Scale to handle large document collections

The system will remember all loaded specifications between sessions, so subsequent startups are much faster. You can add new specs anytime by dropping PDFs into the specs/ folder and rerunning the system.

## Example usage:

- "*How do I implement a SystemVerilog interface for USB 3.0?*"
- "*What are the timing requirements for PCIe link training?*"
- "*Compare VHDL and SystemVerilog module instantiation syntax*"

Your personal hardware specification AI assistant is ready to use! 🚀