

# Lecture 2: Universal Tool Calling Demo

## AI Agent Development with Multi-Backend LLM Integration

(Jack) Yansong Li

October 7, 2025

## 1 Introduction and Course Overview

### 1.1 Learning Objectives

By the end of this lecture, students will be able to:

1. Understand the architecture of universal AI agents with automatic backend selection
2. Implement OpenAI-compatible tool calling mechanisms with proper error handling
3. Design and develop streaming architectures for real-time AI applications
4. Build cross-platform compatible AI systems with intelligent fallback strategies
5. Create extensible tool registries with dynamic registration capabilities
6. Apply modern software engineering patterns to AI agent development
7. Deploy and test AI agents across different hardware configurations

### 1.2 Project Overview

The Universal Tool Calling Demo is a sophisticated educational project that demonstrates advanced LLM tool calling functionality across multiple platforms. This project serves as a comprehensive example of modern AI agent development, showcasing:

- **Automatic Backend Selection:** Intelligent detection of system capabilities and optimal backend routing
- **OpenAI-Compatible Tool Calling:** Standard function calling format with extensible tool registry
- **Real-time Streaming:** Multi-phase streaming with visual feedback for different content types
- **Cross-Platform Compatibility:** Universal deployment across Windows, macOS, and Linux
- **Modern Development Practices:** Comprehensive testing, configuration management, and deployment strategies

### 1.3 Technology Stack

See Table 1.

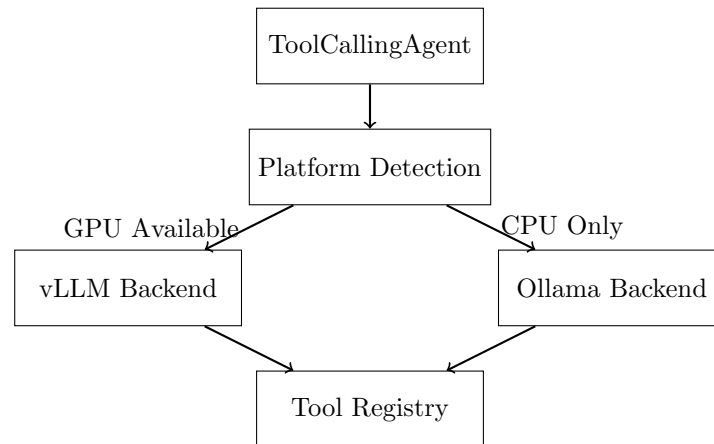
Component	Technology	Purpose
LLM Inference	vLLM ( $\iota=0.6.0$ )	High-performance GPU inference
Local Inference	Ollama ( $\iota=0.1.0$ )	CPU-based local inference
API Framework	FastAPI + Uvicorn	vLLM server management
Package Management	UV	Fast, reliable dependency management
Model	Qwen3-0.6B	Tool-optimized language model
Testing	pytest + coverage	Comprehensive test suite
Configuration	python-dotenv	Environment-based configuration

Table 1: Technology Stack Overview

## 2 Architecture and Design Patterns

### 2.1 High-Level Architecture

The system employs a **Universal Agent Pattern** with intelligent backend selection:



### 2.2 Core Design Patterns

#### 2.2.1 Strategy Pattern: Backend Selection

The system intelligently selects the optimal backend based on platform capabilities:

```

1 class ToolCallingAgent:
2     def _detect_best_backend(self) -> str:
3         """Intelligent backend selection based on system capabilities."""
4         system = platform.system()
5
6         # macOS - Always use Ollama for optimal compatibility
7         if system == "Darwin":
8             return "ollama"
9
10        # Windows/Linux - Check for CUDA availability
11        if torch.cuda.is_available():
12            gpu_memory = torch.cuda.get_device_properties(0).total_memory
13            if gpu_memory >= 4 * 1024**3: # 4GB minimum
14                return "vllm"
15
16        # Fallback to Ollama for CPU systems
17        return "ollama"

```

Listing 1: Backend Selection Strategy

### 2.2.2 Factory Pattern: Agent Creation

Automatic agent initialization based on detected capabilities:

```

1 def _initialize_backend(self):
2     """Factory method for backend initialization."""
3     if self.backend_type == "vllm":
4         self._init_vllm()
5     elif self.backend_type == "ollama":
6         self._init_ollama()
7     else:
8         raise ValueError(f"Unsupported backend: {self.backend_type}")

```

Listing 2: Factory Pattern Implementation

### 2.2.3 Observer Pattern: Streaming Architecture

Real-time streaming with multiple content types:

```

1 def chat_stream(self, message: str, use_tools: bool = True):
2     """Streaming chat with multiple chunk types."""
3     for chunk in self._process_stream(message):
4         chunk_type = chunk.get("type")
5         content = chunk.get("content", "")
6
7         if chunk_type == "thinking":
8             yield {"type": "thinking", "content": content} # Gray text
9         elif chunk_type == "tool_call":
10            yield {"type": "tool_call", "content": content} # Blue notification
11        elif chunk_type == "tool_result":
12            yield {"type": "tool_result", "content": content} # Green result
13        elif chunk_type == "content":
14            yield {"type": "content", "content": content} # Final response

```

Listing 3: Streaming Observer Pattern

### 2.2.4 Registry Pattern: Tool Management

Extensible tool system with dynamic registration:

```

1 class ToolRegistry:
2     def __init__(self):
3         self.tools = {}
4         self._register_built_in_tools()
5
6     def register_tool(self, name: str, function: callable,
7                      description: str, parameters: Dict):
8         """Register a new tool with OpenAI-compatible schema."""
9         self.tools[name] = {
10             "function": function,
11             "description": description,
12             "parameters": parameters
13         }
14
15     def execute_tool(self, name: str, arguments: Dict) -> str:
16         """Execute a registered tool with error handling."""
17         if name not in self.tools:
18             return json.dumps({"error": f"Tool {name} not found"})
19
20         try:
21             tool_func = self.tools[name]["function"]
22             result = tool_func(**arguments)
23             return json.dumps({"success": True, "result": result})
24         except Exception as e:
25             return json.dumps({
26                 "success": False,

```

```

27         "error": str(e),
28         "traceback": traceback.format_exc()
29     })

```

Listing 4: Tool Registry Pattern

## 3 Tool Calling Mechanisms

### 3.1 OpenAI-Compatible Function Calling

The system implements the standard OpenAI function calling format:

```

1  {
2      "tool_calls": [{
3          "id": "call_123",
4          "type": "function",
5          "function": {
6              "name": "get_weather",
7              "arguments": {
8                  "location": "Tokyo, Japan",
9                  "unit": "celsius"
10             }
11         }
12     }]
13 }

```

Listing 5: Function Calling Schema

### 3.2 ReAct Loop Implementation

The ReAct (Reasoning + Acting) pattern enables multi-step problem solving:

```

1  def _process_with_tools(self, message: str, max_iterations: int = 10):
2      """Multi-step reasoning with tool execution."""
3      iteration = 0
4
5      while iteration < max_iterations:
6          # Get model response
7          response = self._get_model_response(message)
8
9          # Parse tool calls from response
10         tool_calls = self._parse_tool_calls(response)
11
12         if not tool_calls:
13             # No tools needed, return final response
14             return response
15
16         # Execute tools and collect results
17         tool_results = []
18         for tool_call in tool_calls:
19             result = self.tool_registry.execute_tool(
20                 tool_call["name"],
21                 tool_call["arguments"]
22             )
23             tool_results.append({
24                 "tool_call_id": tool_call["id"],
25                 "result": result
26             })
27
28         # Add results to conversation history
29         self.conversation_history.extend(tool_results)
30
31         iteration += 1
32

```

```
33 return "Maximum iterations reached"
```

Listing 6: ReAct Loop Implementation

### 3.3 Built-in Tools Analysis

#### 3.3.1 Weather Tool

Real-time weather data retrieval with location parsing:

```
1 def get_current_temperature(location: str, unit: str = "celsius") -> str:
2     """Get current temperature for a location using Open-Meteo API."""
3     # Parse location to coordinates
4     lat, lon = parse_location(location)
5
6     # API call to Open-Meteo
7     url = f"https://api.open-meteo.com/v1/forecast"
8     params = {
9         "latitude": lat,
10        "longitude": lon,
11        "current_weather": True,
12        "temperature_unit": unit
13    }
14
15    response = requests.get(url, params=params)
16    data = response.json()
17
18    temperature = data["current_weather"]["temperature"]
19    return f"Current temperature in {location}: {temperature}{unit}"
```

Listing 7: Weather Tool Implementation

#### 3.3.2 Code Interpreter

Full Python environment with error handling:

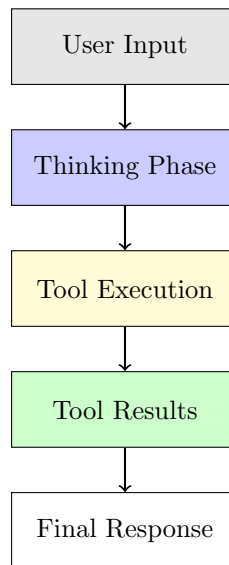
```
1 def code_interpreter(code: str) -> str:
2     """Execute Python code safely with output capture."""
3     try:
4         # Capture stdout and stderr
5         old_stdout = sys.stdout
6         old_stderr = sys.stderr
7         sys.stdout = captured_output = StringIO()
8         sys.stderr = captured_error = StringIO()
9
10        # Execute the code
11        exec(code, {"__builtins__": __builtins__})
12
13        # Restore streams
14        sys.stdout = old_stdout
15        sys.stderr = old_stderr
16
17        output = captured_output.getvalue()
18        error = captured_error.getvalue()
19
20        if error:
21            return f"Output:\\n{output}\\nError:\\n{error}"
22        return output if output else "Code executed successfully"
23
24    except Exception as e:
25        return json.dumps({
26            "success": False,
27            "error": str(e),
28            "traceback": traceback.format_exc()
29        })
```

Listing 8: Code Interpreter Tool

## 4 Streaming Architecture

### 4.1 Multi-Type Streaming System

The streaming architecture provides real-time feedback with different content types:



### 4.2 Streaming Implementation

```

1 CHUNK_TYPES = {
2     "thinking": "Internal reasoning process",
3     "tool_call": "Tool execution initiation",
4     "tool_result": "Tool execution completion",
5     "content": "Final response content",
6     "error": "Error notifications"
7 }
8
9 def chat_stream(self, message: str, use_tools: bool = True):
10     """Streaming chat with visual feedback."""
11     for chunk in self._process_stream(message):
12         chunk_type = chunk.get("type")
13         content = chunk.get("content", "")
14
15         if chunk_type == "thinking":
16             # Gray text for internal reasoning
17             print(f"\033[90m{content}\033[0m", end="", flush=True)
18             yield {"type": "thinking", "content": content}
19
20         elif chunk_type == "tool_call":
21             # Blue notification for tool execution
22             print(f"Tool: {content['name']}")
23             yield {"type": "tool_call", "content": content}
24
25         elif chunk_type == "tool_result":
26             # Green result display
27             print(f"Result: {content}")
28             yield {"type": "tool_result", "content": content}
29

```

```

30     elif chunk_type == "content":
31         # Standard text for final response
32         print(content, end="", flush=True)
33         yield {"type": "content", "content": content}

```

Listing 9: Streaming Chunk Processing

### 4.3 Performance Optimization

The streaming system includes several optimization strategies:

- **Chunk Batching:** Groups related content to reduce overhead
- **Buffer Management:** Intelligent buffering for smooth streaming
- **Backpressure Handling:** Prevents overwhelming the client with too many chunks
- **Timeout Protection:** Ensures streaming doesn't hang indefinitely

## 5 Cross-Platform Compatibility

### 5.1 Platform Detection Strategy

Intelligent platform detection with capability assessment:

```

1 def check_system_compatibility():
2     """Comprehensive platform and capability detection."""
3     system = platform.system()
4
5     # macOS Detection
6     if system == "Darwin":
7         # Check for Apple Silicon vs Intel
8         if platform.machine() == "arm64":
9             return "ollama", "Apple Silicon Mac (M1/M2)"
10        else:
11            return "ollama", "Intel Mac"
12
13    # Windows/Linux CUDA Detection
14    if torch.cuda.is_available():
15        gpu_name = torch.cuda.get_device_name(0)
16        gpu_memory = torch.cuda.get_device_properties(0).total_memory
17
18        if gpu_memory >= 4 * 1024**3: # 4GB minimum
19            return "vllm", f"NVIDIA GPU: {gpu_name} ({gpu_memory//1024**3}GB)"
20
21    # CPU-based fallback
22    cpu_count = multiprocessing.cpu_count()
23    memory = psutil.virtual_memory().total
24
25    return "ollama", f"CPU system ({cpu_count} cores, {memory//1024**3}GB RAM)"

```

Listing 10: Platform Detection Implementation

### 5.2 Backend Abstraction

Unified interface across different backends:

```

1 from abc import ABC, abstractmethod
2
3 class BaseAgent(ABC):
4     """Abstract base class for all backend implementations."""
5
6     @abstractmethod

```

```

7  def chat(self, message: str, use_tools: bool = True,
8      stream: bool = False) -> Union[str, Generator]:
9      """Process a chat message with optional tool usage."""
10     pass
11
12     @abstractmethod
13     def reset_conversation(self) -> None:
14         """Reset conversation history."""
15         pass
16
17     @abstractmethod
18     def get_system_info(self) -> Dict[str, Any]:
19         """Get system and model information."""
20         pass
21
22 class VLLMAgent(BaseAgent):
23     """vLLM-specific implementation with GPU acceleration."""
24
25     def __init__(self, model_name: str, host: str, port: int):
26         self.model_name = model_name
27         self.server = VLLMServer(host, port)
28         self.client = OpenAI(base_url=f"http://{host}:{port}/v1")
29
30     def chat(self, message: str, use_tools: bool = True,
31             stream: bool = False):
32         # vLLM-specific implementation
33         pass
34
35 class OllamaAgent(BaseAgent):
36     """Ollama implementation for local CPU inference."""
37
38     def __init__(self, model_name: str):
39         self.model_name = model_name
40         self.client = ollama.Client()
41
42     def chat(self, message: str, use_tools: bool = True,
43             stream: bool = False):
44         # Ollama-specific implementation
45         pass

```

Listing 11: Backend Abstraction Interface

## 5.3 Platform-Specific Setup

### 5.3.1 macOS Setup

```

1  # Install via Homebrew
2  brew install ollama
3
4  # Start Ollama service
5  ollama serve # Run in separate terminal
6
7  # Download default model
8  ollama pull qwen3:0.6b

```

Listing 12: macOS Installation

### 5.3.2 Windows Setup

```

1  # Download from official website
2  # https://ollama.com/download/windows
3
4  # Install and start service
5  # Ollama will run as system service

```



```

6
7 # Download model via command line
8 ollama pull qwen3:0.6b

```

Listing 13: Windows Installation

### 5.3.3 Linux Setup

```

1 # Install Ollama
2 curl -fsSL https://ollama.com/install.sh | sh
3
4 # Start service
5 sudo systemctl start ollama
6 sudo systemctl enable ollama
7
8 # Download model
9 ollama pull qwen3:0.6b
10
11 # For vLLM (requires CUDA)
12 # Install CUDA toolkit from NVIDIA
13 # vLLM will be automatically detected and used

```

Listing 14: Linux Installation

## 6 Testing and Development

### 6.1 Comprehensive Testing Framework

The project includes extensive testing with multiple categories:

```

1 tests/
2 |-- test_main.py           # Main agent functionality
3 |-- test_streaming.py      # Streaming implementation
4 |-- test_code_interpreter_full.py # Tool execution testing
5 '-- conftest.py           # Test fixtures and mocks

```

Listing 15: Test Structure

### 6.2 Unit Testing

```

1 def test_backend_detection():
2     """Test automatic backend selection logic."""
3     with patch('platform.system') as mock_system:
4         with patch('torch.cuda.is_available') as mock_cuda:
5             # Test macOS detection
6             mock_system.return_value = "Darwin"
7             backend, description = check_system_compatibility()
8             assert backend == "ollama"
9             assert "Mac" in description
10
11             # Test Windows with CUDA
12             mock_system.return_value = "Windows"
13             mock_cuda.return_value = True
14             backend, description = check_system_compatibility()
15             assert backend == "vllm"
16             assert "NVIDIA GPU" in description
17
18             # Test Linux without CUDA
19             mock_system.return_value = "Linux"
20             mock_cuda.return_value = False
21             backend, description = check_system_compatibility()

```

```
22 assert backend == "ollama"
```

Listing 16: Backend Detection Testing

### 6.3 Integration Testing

```
1 def test_tool_execution_workflow():
2     """Test complete tool execution workflow."""
3     agent = ToolCallingAgent(backend="ollama")
4
5     # Test weather tool execution
6     response = agent.chat("What's the weather in Tokyo?")
7
8     # Verify response contains weather information
9     assert "temperature" in response.lower() or "weather" in response.lower()
10    assert "Tokyo" in response
11
12    # Test code interpreter
13    response = agent.chat("Calculate 2 + 2")
14    assert "4" in response
```

Listing 17: Tool Execution Integration Test

### 6.4 Mock-Based Testing

```
1 def test_error_handling():
2     """Test error propagation and formatting."""
3     registry = ToolRegistry()
4
5     # Test tool not found
6     result = registry.execute_tool("nonexistent_tool", {})
7     result_dict = json.loads(result)
8     assert result_dict["success"] == False
9     assert "not found" in result_dict["error"]
10
11    # Test execution error
12    result = registry.execute_tool("code_interpreter",
13                                  {"code": "1/0"})
14    result_dict = json.loads(result)
15    assert result_dict["success"] == False
16    assert "ZeroDivisionError" in result_dict["error"]
17    assert "traceback" in result_dict
```

Listing 18: Mock Testing Strategy

### 6.5 Modern Development Workflow

```
1 # Install dependencies with UV
2 uv sync
3
4 # Run all tests
5 uv run pytest
6
7 # Run specific test file with verbose output
8 uv run pytest tests/test_main.py -v
9
10 # Run tests with coverage reporting
11 uv run pytest tests/ --cov=src/local_llm_serving
12
13 # Run single test by name pattern
14 uv run pytest -k "test_tool_execution"
15
```

```

16 # Format code with black
17 uv run black src/ tests/
18
19 # Sort imports with isort
20 uv run isort src/ tests/

```

Listing 19: Development Commands

## 7 Configuration Management

### 7.1 Environment-Based Configuration

```

1 # Load environment variables
2 from dotenv import load_dotenv
3 load_dotenv()
4
5 # Configuration with defaults
6 class Config:
7     MODEL_NAME = os.getenv("MODEL_NAME", "Qwen/Qwen3-0.6B")
8     VLLM_HOST = os.getenv("VLLM_HOST", "localhost")
9     VLLM_PORT = int(os.getenv("VLLM_PORT", "8000"))
10    LOG_LEVEL = os.getenv("LOG_LEVEL", "INFO")
11    MAX_ITERATIONS = int(os.getenv("MAX_ITERATIONS", "10"))
12    TOOL_TIMEOUT = int(os.getenv("TOOL_TIMEOUT", "30"))
13
14    # Backend-specific settings
15    OLLAMA_HOST = os.getenv("OLLAMA_HOST", "localhost")
16    OLLAMA_PORT = int(os.getenv("OLLAMA_PORT", "11434"))

```

Listing 20: Configuration Management

### 7.2 Logging Configuration

```

1 import logging
2 import sys
3
4 # Configure logging
5 logging.basicConfig(
6     level=getattr(logging, Config.LOG_LEVEL),
7     format='%(asctime)s - %(name)s - %(levelname)s - %(message)s',
8     handlers=[
9         logging.StreamHandler(sys.stdout),
10        logging.FileHandler('logs/agent.log')
11    ]
12 )
13
14 # Module-specific loggers
15 logger = logging.getLogger(__name__)

```

Listing 21: Logging Setup

## 8 Practical Exercises

### 8.1 Exercise 1: Custom Tool Development

**Objective:** Implement a custom tool that integrates with an external API.

**Requirements:**

- Create a tool that fetches current cryptocurrency prices

- Implement proper error handling and rate limiting
- Follow OpenAI function calling format
- Add comprehensive unit tests

**Solution Template:**

```

1 def get_crypto_price(symbol: str, currency: str = "USD") -> str:
2     """Get current cryptocurrency price."""
3     # Implementation here
4     pass
5
6 # Register the tool
7 registry.register_tool(
8     name="get_crypto_price",
9     function=get_crypto_price,
10    description="Get current price of a cryptocurrency",
11    parameters={
12        "type": "object",
13        "properties": {
14            "symbol": {"type": "string", "description": "Crypto symbol (BTC, ETH, etc.)"},
15            "currency": {"type": "string", "description": "Currency code (USD, EUR, etc.)"}
16        },
17        "required": ["symbol"]
18    }
19 )

```

Listing 22: Custom Tool Implementation

## 8.2 Exercise 2: Backend Implementation

**Objective:** Implement a new backend for a different LLM provider.

**Requirements:**

- Create a backend for Hugging Face Transformers
- Implement the BaseAgent interface
- Support both streaming and non-streaming modes
- Add tool calling capabilities
- Write integration tests

**Implementation Guidelines:**

```

1 class HuggingFaceAgent(BaseAgent):
2     """Hugging Face Transformers backend implementation."""
3
4     def __init__(self, model_name: str, device: str = "auto"):
5         self.model_name = model_name
6         self.device = self._get_device(device)
7         self.tokenizer = AutoTokenizer.from_pretrained(model_name)
8         self.model = AutoModelForCausalLM.from_pretrained(model_name)
9         self.model.to(self.device)
10
11     def chat(self, message: str, use_tools: bool = True,
12             stream: bool = False):
13         # Implementation here
14         pass

```

Listing 23: New Backend Template

### 8.3 Exercise 3: Cross-Platform Deployment

**Objective:** Deploy the application on different platforms with automated setup.

**Requirements:**

- Create Docker containers for different platforms
- Implement automated setup scripts
- Add health checks and monitoring
- Document deployment procedures
- Test on multiple cloud providers

### 8.4 Exercise 4: Performance Optimization

**Objective:** Optimize the streaming performance and reduce latency.

**Requirements:**

- Implement chunk batching for reduced overhead
- Add buffer management for smooth streaming
- Optimize tool execution parallelism
- Add caching for frequently used tools
- Measure and document performance improvements

## 9 Assessment Criteria

### 9.1 Learning Outcomes Assessment

Students will be evaluated based on their ability to:

**A. Architecture Understanding (25%)**

- Explain the Universal Agent Pattern and its benefits
- Describe the backend selection strategy and rationale
- Identify and explain the design patterns used

**B. Implementation Skills (35%)**

- Implement custom tools following best practices
- Create backend implementations with proper interfaces
- Handle errors gracefully with appropriate feedback

**C. Testing Proficiency (20%)**

- Write comprehensive unit tests with good coverage
- Implement integration tests for complex workflows
- Use mocking strategies effectively

**D. Cross-Platform Development (20%)**

- Deploy applications on different platforms
- Handle platform-specific requirements
- Document deployment procedures clearly

## 9.2 Practical Project Assessment

**Final Project:** Students must implement a complete AI agent application with:

- Custom domain-specific tools (minimum 3 tools)
- Multi-backend support (minimum 2 backends)
- Comprehensive test suite (minimum 80% coverage)
- Cross-platform deployment documentation
- Performance optimization and benchmarking
- User documentation and API reference

## 10 Conclusion

This lecture has covered the comprehensive architecture and implementation of the Universal Tool Calling Demo, demonstrating advanced concepts in:

- Modern AI agent development with multi-backend support
- OpenAI-compatible tool calling with extensible architecture
- Real-time streaming with visual feedback systems
- Cross-platform compatibility and intelligent fallback strategies
- Production-ready testing and deployment practices
- Software engineering patterns in AI application development

The project serves as an excellent foundation for understanding how to build sophisticated AI applications that work reliably across diverse computing environments while maintaining high performance and user experience standards.

### 10.1 Further Reading

1. OpenAI Function Calling Documentation: <https://platform.openai.com/docs/guides/function-calling>
2. ReAct Paper: "ReAct: Synergizing Reasoning and Acting in Language Models"
3. vLLM Documentation: <https://docs.vllm.ai/>
4. Ollama Documentation: <https://ollama.com/>
5. Python Packaging with pyproject.toml: <https://packaging.python.org/en/latest/guides/writing-pyproject-toml/>