# Comprehensive Guide: Agentic AI, Web Sockets, Scrapy, and Apache NiFi

### **Table of Contents**

- 1. Agentic AI
- 2. Web Sockets
- 3. Scrapy
- 4. Apache NiFi
- 5. Key Comparisons
- 6. Practice Questions

## **Agentic AI**

### **Definition and Core Concepts**

Agentic AI refers to artificial intelligence systems designed to act autonomously on behalf of users to achieve specified goals. Unlike traditional AI systems that respond to direct commands, agentic AI can:

- Operate with varying degrees of autonomy
- Take initiative to complete tasks
- Make decisions based on context and prior knowledge
- Plan and execute multi-step processes
- Learn from outcomes and adapt strategies

### **Key Components of Agentic AI**

#### 1. Goal Setting and Understanding

- $\,^{\circ}$  Ability to interpret user intentions
- Converting vague instructions into concrete objectives
- Prioritizing competing goals

#### 2. Planning and Reasoning

- Breaking down complex tasks into manageable steps
- Anticipating potential obstacles
- Developing contingency plans
- Causal reasoning about actions and consequences

#### 3. Tool Usage

- API integration
- Accessing external databases
- Using web search
- Operating system/application control
- Calling specialized models for specific subtasks

#### 4. Memory and Context Management

Short-term working memory

- Long-term knowledge storage
- Context preservation across multiple steps
- Episodic memory of previous interactions

#### 5. Execution and Monitoring

- Taking actions based on plans
- Progress tracking
- Error detection
- Course correction

#### 6. Reflection and Learning

- Self-assessment of actions
- Improvement through experience
- Documentation of processes

### Architectures for Agentic AI

#### 1. LLM-Based Agents

- ReAct: Reasoning and Acting
- LangChain/AutoGPT frameworks
- Chain-of-thought prompting

### 2. Cognitive Architectures

- o ACT-R
- ∘ SOAR
- o GOMS

#### 3. Multi-Agent Systems

- Specialist agents with different roles
- Communication protocols between agents
- Consensus mechanisms

## **Ethical and Safety Considerations**

### 1. Alignment

- Ensuring agent goals match human intentions
- Preventing goal distortion or misinterpretation

#### 2. Oversiaht

- Human-in-the-loop monitoring
- Tiered autonomy levels
- Interruptibility

#### 3. Responsibility and Accountability

- Clear attribution of decision-making
- Explainability of agent actions
- Liability frameworks

#### **Current Limitations**

- 1. Hallucination and Reasoning Gaps
- 2. Context Window Constraints
- 3. Tool Integration Challenges
- 4. Long-Term Planning Deficiencies
- 5. Robustness Against Adversarial Inputs

### **Applications**

- 1. **Personal Assistants** Email management, calendar scheduling
- 2. Research Assistants Literature review, data analysis
- 3. **Software Development** Code generation, debugging, optimization
- 4. **Business Process Automation** Customer service, inventory management
- 5. Educational Tutoring Personalized learning paths, feedback

### **Web Sockets**

### **Definition and Core Concepts**

WebSocket is a communication protocol that provides full-duplex communication channels over a single TCP connection. Unlike HTTP, WebSockets:

- Maintain persistent connections
- Allow real-time bidirectional data transfer
- Have minimal overhead after initial handshake
- Support both text and binary data formats

#### WebSocket Protocol

#### 1. Handshake Process

- Client initiates with HTTP upgrade request
- Server responds with 101 Switching Protocols
- Connection transitions to WebSocket protocol

Client request example:

GET /chat HTTP/1.1

Host: server.example.com

Upgrade: websocket
Connection: Upgrade

Sec-WebSocket-Key: dGhlIHNhbXBsZSBub25jZQ==

Sec-WebSocket-Version: 13

Server response example:

HTTP/1.1 101 Switching Protocols

Upgrade: websocket
Connection: Upgrade

Sec-WebSocket-Accept: s3pPLMBiTxaQ9kYGzzhZRbK+x0o=

#### 1. Frame Structure

- FIN bit (indicates final fragment)
- RSV bits (protocol extensions)
- Opcode (text, binary, close, ping, pong)
- Mask bit
- Payload length

- Masking key (if mask bit set)
- Payload data

#### 2. Connection Management

- Ping/pong frames for keeping connections alive
- Close frames with status codes
- Reconnection strategies

### **Implementation in Different Languages**

1. JavaScript (Browser)

```
const socket = new WebSocket('ws://example.com/socket');
  socket.onopen = function(event) {
    console.log('Connection established');
    socket.send('Hello Server!');
  };
  socket.onmessage = function(event) {
    console.log('Message from server:', event.data);
  };
  socket.onclose = function(event) {
    console.log('Connection closed, code:', event.code);
  };
  socket.onerror = function(error) {
    console.error('WebSocket error:', error);
  };
2. Python (Server with websockets library)
  import asyncio
  import websockets
  async def echo(websocket):
      async for message in websocket:
          await websocket.send(f"Echo: {message}")
  async def main():
      async with websockets.serve(echo, "localhost", 8765):
          await asyncio.Future() # Run forever
  asyncio.run(main())
```

3. Java (Server with javax.websocket)

```
@ServerEndpoint("/echo")
public class EchoServer {
    @0n0pen
    public void onOpen(Session session) {
        System.out.println("Connected: " + session.getId());
    }
    @OnMessage
    public String onMessage(String message, Session session)
        return "Echo: " + message;
    }
    @OnClose
    public void onClose(Session session) {
        System.out.println("Disconnected: " +
        session.getId());
    }
}
```

### WebSocket vs. HTTP

Feature	WebSocket	НТТР
Connection	Persistent	Request/ Response (stateless)
Overhead	Low after handshake	Headers with each request
Communication	Bidirectional	Client- initiated only
Real-time	Yes	No (requires polling)
Use Cases	Live updates, chat, gaming	Document transfer, API calls

### **WebSocket Security Considerations**

#### 1. Origin Verification

- Server should validate Origin header
- Prevents cross-site WebSocket hijacking

#### 2. Authentication and Authorization

- Token-based authentication
- Session management

#### 3. Data Validation

- Input sanitization
- Frame size limits

#### 4. Transport Security

- Using WSS (WebSocket Secure) over TLS/SSL
- Certificate validation

#### Common Use Cases

#### 1. Real-time Applications

- Chat applications
- Collaborative editing
- Live dashboards

#### 2. Gaming

- Multiplayer online games
- Real-time game state synchronization

#### 3. Financial Applications

- Live stock tickers
- Trading platforms

#### 4. **IoT**

- Device monitoring
- Remote control applications

### **Popular WebSocket Libraries and Frameworks**

#### 1. JavaScript

- Socket.IO
- ∘ ws (Node.js)
- SockIS

#### 2. Python

- websockets
- Flask-SocketIO
- Django Channels

#### 3. **Java**

- JSR 356 (Java API for WebSockets)
- Spring WebSockets
- Jetty WebSocket

## **Scrapy**

### **Definition and Core Concepts**

Scrapy is a high-level Python web crawling and scraping framework designed for extracting structured data from websites. Key aspects include:

- Asynchronous networking for high-performance crawling
- Built-in support for selecting and extracting data
- Pipeline architecture for data processing
- Extensible design with middleware components
- Robust request/response handling

### **Architecture and Components**

#### 1. Engine

- Coordinates all components
- Controls data flow between components
- Handles exceptions

#### 2. Scheduler

- Receives requests from engine
- Queues requests for future crawling
- Implements priority and duplicate filtering

#### 3. Downloader

- Fetches web pages
- Returns responses to the engine
- Handles HTTP protocol details

#### 4. Spiders

- User-defined classes that:
  - Define start URLs
  - Parse responses
  - Extract data
  - Follow links

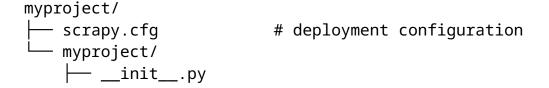
#### 5. Item Pipeline

- Processes scraped items
- Validates extracted data
- Removes duplicates
- Stores data in databases or files

#### 6. Middleware

- Downloader Middleware: Processes requests before downloading and responses after
- Spider Middleware: Processes spider input (responses) and output (items and requests)

### **Project Structure**



```
    items.py  # project items definition
    middlewares.py  # project middlewares
    ipipelines.py  # project pipelines
    ipipelines.py  # project settings
    ipipelines  # project settings
    ipipelines  # directory with spiders
    ipipelines  # project settings
    ipipelines  # project settings
    ipipelines  # project settings
    ipipelines  # project pipelines
    ipipelines  # project middlewares
    ipipelines  # project pipelines
    ipipelines  # project pipeli
```

### Creating a Spider

```
import scrapy
class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
    1
    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text': quote.css('span.text::text').get(),
                'author': quote.css('small.author::text').get(),
                'tags': quote.css('div.tags
        a.tag::text').getall(),
            }
        next_page = response.css('li.next a::attr(href)').get()
        if next_page is not None:
            yield response.follow(next_page, self.parse)
```

### **Data Extraction Techniques**

#### 1. CSS Selectors

```
# Select all <h1> elements
response.css('h1')

# Extract text from <h1>
response.css('h1::text').get()

# Extract attribute value
response.css('a::attr(href)').get()
```

```
# Select with class
    response.css('div.quote')
  2. XPath Selectors
    # Select all <h1> elements
    response.xpath('//h1')
    # Extract text from <h1>
    response.xpath('//h1/text()').get()
    # Extract attribute value
    response.xpath('//a/@href').get()
    # Select with class
    response.xpath('//div[@class="quote"]')
  3. Item Loaders
    from scrapy.loader import ItemLoader
    from myproject.items import ProductItem
    def parse(self, response):
        loader = ItemLoader(item=ProductItem(),
             response=response)
        loader.add_css('name', 'h1.product-name::text')
        loader.add_xpath('price', '//span[@class="price"]/
            text()')
        loader.add_value('url', response.url)
        yield loader.load item()
Item Pipeline
class PricePipeline:
    def process_item(self, item, spider):
        if 'price' in item:
            # Convert price string to float
            value = item['price'].replace('$', '')
            item['price'] = float(value)
        return item
class DatabasePipeline:
    def open_spider(self, spider):
        self.conn = database_connection()
    def process_item(self, item, spider):
```

```
self.conn.save(item)
return item

def close_spider(self, spider):
    self.conn.close()
```

### **Settings Configuration**

```
# scrapy settings file
# Configure maximum concurrent requests
CONCURRENT_REQUESTS = 16
# Configure delay between requests
DOWNLOAD_DELAY = 1
# Enable or disable user agents
USER_AGENT = 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
        AppleWebKit/537.36'
# Configure item pipelines
ITEM_PIPELINES = {
   'myproject.pipelines.PricePipeline': 300,
   'myproject.pipelines.DatabasePipeline': 800,
}
# Enable cache
HTTPCACHE ENABLED = True
HTTPCACHE_EXPIRATION_SECS = 86400
```

### **Handling JavaScript and Dynamic Content**

#### 1. Splash Integration

```
# Parse rendered HTML
pass
```

#### 2. Selenium Integration

```
from scrapy import Spider
from selenium import webdriver
class SeleniumSpider(Spider):
    name = "selenium_example"
    def __init__(self):
        self.driver = webdriver.Chrome()
    def start_requests(self):
        self.driver.get('https://example.com')
        # Wait for JavaScript to execute
        yield scrapy.Request(
            self.driver.current url,
            self.parse,
            dont_filter=True
        )
    def parse(self, response):
        # Parse content
        pass
    def closed(self, reason):
        self.driver.quit()
```

### **Running Scrapy**

#### 1. Command Line

```
# Create new project
scrapy startproject myproject

# Create spider
scrapy genspider example example.com

# Run spider
scrapy crawl quotes

# Save output to file
scrapy crawl quotes -o quotes.json
```

```
# Run shell for testing selectors
scrapy shell "http://quotes.toscrape.com"
```

#### 2. From Python Script

```
from scrapy.crawler import CrawlerProcess
from myproject.spiders.quotes import QuotesSpider

process = CrawlerProcess({
    'USER_AGENT': 'Mozilla/5.0',
    'FEED_FORMAT': 'json',
    'FEED_URI': 'output.json'
})

process.crawl(QuotesSpider)
process.start()
```

### **Best Practices and Ethical Scraping**

#### 1. Respect robots.txt

```
# In settings.py
ROBOTSTXT_OBEY = True
```

### 2. Rate Limiting

```
# In settings.py
DOWNLOAD_DELAY = 2 # 2 seconds between requests
CONCURRENT_REQUESTS_PER_DOMAIN = 8
```

#### 3. User Agents

```
# In settings.py
USER_AGENT = 'MyBot/1.0 (contact@example.com)'
```

### 4. Proxy Rotation

```
# In settings.py
DOWNLOADER_MIDDLEWARES = {
    'scrapy_proxies.RandomProxy': 100,
    'scrapy.downloadermiddlewares.httpproxy.HttpProxyMiddleware': 110,
}
PROXY_LIST = '/path/to/proxy/list.txt'
```

#### 5. Handling Authentication

## **Apache NiFi**

### **Definition and Core Concepts**

Apache NiFi is an open-source data integration and dataflow automation tool designed to automate the movement of data between disparate systems. Key aspects:

- Visual command center for data routing
- Data provenance tracking
- Built for high throughput
- Web-based user interface
- Highly configurable
- Data transformation capabilities

#### Core Architecture

#### 1. Flow-Based Programming Model

- Visual dataflow representation
- Data-driven processing
- Component-based development

#### 2. Key Components

- FlowFile: Core data unit with content and attributes
- **Processor**: Components that perform actions on FlowFiles
- **Connection**: Links between processors with queues
- **Process Group**: Container for organizing processors
- Controller Service: Shared services for processors
- **Reporting Task**: Components for monitoring the system

#### 3. NiFi Cluster Architecture

- **Zero-Master Clustering**: All nodes perform same tasks
- **ZooKeeper**: Coordinates cluster state
- **Primary Node**: Handles cluster coordination

### **User Interface and Key Features**

- 1. Canvas: Visual workspace for building dataflows
- 2. Components Toolbar: Available processors and components
- 3. **Status Bar**: System status information
- 4. Global Menu: Configuration and administrative options
- 5. **Search**: Finding components in complex flows
- 6. Parameter Context: Managing flow parameters
- 7. **Templates**: Reusable flow designs

#### **Processors and Data Flow**

#### 1. Common Processor Types

- Data Ingestion: GetFile, GetHTTP, ListenTCP, ConsumeKafka
- Transformation: SplitText, JoltTransformJSON, ExecuteScript
- Routing: RouteOnAttribute, RouteOnContent, ControlRate
- **Database**: ExecuteSQL, PutDatabaseRecord, QueryDatabaseTable
- o Distribution: PutFile, PutS3Object, PublishKafka
- **Monitoring**: MonitorActivity, NotifyEmail, PostSlack

#### 2. FlowFile Lifecycle

- Creation from source data
- Transformation through processors
- Routing based on attributes or content
- Delivery to destination
- Provenance events recorded throughout

#### 3. Data Provenance

- Tracks complete lineage of all data
- Records every action performed on data
- Enables auditing and troubleshooting
- Replay capability for failed flows
- Configurable retention policy

### **Example NiFi Data Flow**

#### 1. File Ingestion to Database

- GetFile: Monitor directory for new files
- SplitCSV: Split file into individual records
- ConvertRecord: Convert CSV to ISON
- ValidateRecord: Ensure data meets schema
- PutDatabaseRecord: Insert into database
- PutFile: Archive processed files

#### 2. API Integration Flow

- InvokeHTTP: Call external API
- EvaluateIsonPath: Extract fields
- JoltTransformJSON: Reshape data
- AttributesToJSON: Convert attributes to JSON
- PublishKafka: Send to messaging system

### **Configuration and Management**

#### 1. System Properties

- Web UI settings
- Security configurations
- Thread pool sizes
- FlowFile repository settings

#### 2. Controller Services

- Database connection pools
- SSL context services
- Schema registries
- Distributed cache services

#### 3. Parameter Contexts

- Named sets of parameters
- Reusable across process groups
- Environment-specific values
- Secure sensitive parameters

#### 4. State Management

- Local vs. cluster state
- Persistent component state
- Stateful processors

#### 5. Backpressure and Flow Control

- Connection queue thresholds
- Pressure-sensitive prioritization
- Upstream notification

### **Security Features**

#### 1. Authentication

- LDAP/Active Directory integration
- Kerberos
- OpenID Connect
- Client certificates

#### 2. Authorization

- Fine-grained access control
- Policy-based permissions
- User/group management
- Component-level restrictions

#### 3. Data Protection

- Encryption at rest
- TLS for transit
- Sensitive property encryption
- Secure processor configuration

#### 4. Auditing

- User action logging
- Data access tracking
- Configuration changes

### **Performance Optimization**

#### 1. Threading Model

- Event-driven architecture
- Concurrent task scheduling
- Processor-specific concurrency

#### 2. Sizing and Scaling

- JVM memory allocation
- Disk space planning
- Network capacity
- Cluster sizing

#### 3. Tuning Techniques

- Batch processing configuration
- Connection queue tuning
- Back pressure threshold adjustment
- Processor scheduling strategies

#### 4. Performance Testing

- Load simulation
- Resource utilization monitoring
- Bottleneck identification

### **Integration with Other Technologies**

#### 1. Apache Ecosystem

- Hadoop: HDFS integration
- Kafka: Producing and consuming messages
- Spark: Data processing integration
- HBase: NoSOL database interaction

#### 2. Cloud Platforms

- AWS services (S3, DynamoDB, etc.)
- Azure services (Blob Storage, Event Hub)
- Google Cloud Platform services

#### 3. Databases

- Relational databases (MySQL, PostgreSQL, etc.)
- NoSQL databases (MongoDB, Cassandra)
- Data warehouses (Snowflake, Redshift)

#### 4. Messaging Systems

- Apache Kafka
- · RabbitMQ
- JMS implementations

### NiFi Extensions and Registry

#### 1. Custom Processors

- Java-based processor development
- $\,{}^{\circ}$  Maven archetypes for processor generation
- Component lifecycle management

#### 2. NiFi Registry

- Version control for flows
- Flow deployment
- Environment management

Change tracking

#### 3. NiFi Toolkit

- CLI operations
- Flow analysis
- Automated deployment
- Configuration scripts

### **Monitoring and Operations**

### 1. Built-in Monitoring

- Bulletins
- Status history
- System diagnostics
- Statistics

#### 2. External Monitoring

- JMX metrics
- Prometheus integration
- Reporting tasks
- Log analysis

### 3. Operational Considerations

- Backup and recovery
- Disaster planning
- Upgrade strategies
- Capacity planning

## **Key Comparisons**

### **Agentic AI vs. Traditional Automation Systems**

Feature	Agentic AI	Traditional Automation
Adaptability	High - can handle novel situations	Limited - follows predefined rules
Decision Making	Autonomous based on context	Programmed logic only
Planning	Can create multi-step plans dynamically	Follows fixed workflows
Learning		Requires manual updates

Feature	Agentic AI	Traditional Automation
	Can improve from experience	
Complexity Handling	Handles ambiguity well	Requires explicit instructions

## WebSockets vs. REST API

Feature	WebSockets	REST API
Connection Type	Persistent, bidirectional	Request-response, stateless
Use Case	Real-time updates, chat	CRUD operations, resource access
Protocol	ws:// or wss://	http:// or https://
Data Transfer	Push and pull	Pull only
Server Load	Higher connection maintenance	Higher request processing
Complexity	More complex to implement	Simpler, standardized

# Scrapy vs. Beautiful Soup

Feature	Scrapy	Beautiful Soup
Type	Full-featured framework	Library
Performance	Fast (asynchronous)	Slower (synchronous)

Feature	Scrapy	Beautiful Soup
Complexity	Higher learning curve	Simple to use
Features	Request handling, pipelines, scheduling	Parsing only
Concurrency	Built-in	Manual implementation
Use Case	Large-scale production scraping	Simple scripts, one-off tasks

# Apache NiFi vs. Apache Airflow

Feature	Apache NiFi	Apache Airflow
Paradigm	Data flow	Task scheduling
UI	Rich, real-time flow visualization	DAG representation
Focus	Data movement and transformation	Workflow orchestration
State	Stateful data tracking	Task state tracking
Scaling	Horizontal scaling for throughput	Worker scaling for tasks
Use Case	Continuous data processing	Batch processing, ETL

## **Practice Questions**

### **Agentic AI**

- 1. What are the key components that distinguish agentic AI from traditional AI systems?
- 2. Explain the importance of tool usage in agentic AI systems.
- 3. What are the main challenges in implementing reliable planning in agentic AI?
- 4. How does the ReAct framework improve agentic AI reasoning?
- 5. What ethical considerations are most critical when deploying agentic AI systems?

#### Web Sockets

- 1. Describe the WebSocket handshake process and its relationship to HTTP.
- 2. What are the main advantages of WebSockets over HTTP polling for real-time applications?
- 3. How do WebSocket frames handle binary vs. text data?
- 4. Explain how WebSocket connections are maintained and why ping/pong frames are important.
- 5. What security considerations should be addressed when implementing WebSockets?

### **Scrapy**

- 1. Describe the main components of Scrapy's architecture and how they interact.
- 2. Compare and contrast CSS selectors and XPath selectors in Scrapy.
- 3. How does Scrapy handle JavaScript-rendered content?
- 4. Explain the purpose of middleware in Scrapy and provide an example.
- 5. What ethical considerations should be taken into account when web scraping?

### **Apache NiFi**

- 1. What is a FlowFile in Apache NiFi and how does it relate to data provenance?
- 2. Explain the difference between a Processor and a Controller Service in NiFi
- 3. How does NiFi's clustering architecture ensure high availability?
- 4. Describe the purpose of Parameter Contexts and how they improve flow management.
- 5. What strategies can be used to optimize NiFi performance for high-throughput data flows?