Execute | Clear

## Responses

Curl

```
curl -X 'POST' \
  'http://127.0.0.1:8000/chat' \
  -H 'accept: application/json' \
  -H 'Content-Type: application/json' \
  -d '{
  "prompt": "What is an inference server?"
}'
```

Request URL

```
http://127.0.0.1:8000/chat
```

Server response

| Code | Details |
|---|---|
| 200 | Response body |

```
{
  "reply": "An inference server, also known as a prediction server or model server, is a type of software system that hosts and manages machine learning models to make predictions or generate insights from data. It's a critical component in modern data science and AI infrastructure.\n\nHere's how it typically works:\n\n1. **Model training**: A machine learning model is trained on a dataset using a training environment, and the resulting model is then deployed to the inference server.\n2. **Model hosting**: The inference server loads the trained model into memory and hosts it for serving predictions.\n3. **API requests**: Clients (e.g., web applications, mobile apps, or other services) send API requests to the inference server with input data.\n4. **Model execution**: The inference server runs the loaded model on the input data using the model's prediction algorithms.\n5. **Response generation**: The inference server returns the predicted outcomes or insights to the client as a response.\n\nInference servers often support various features, including:\n\n* **Scalability**: Handling multiple requests concurrently and adapting to changing workloads.\n* **High availability**: Providing continuous services with minimal downtime.\n* **Model management**: Deploying, updating, and managing multiple models simultaneously.\n* **Predictive performance metrics**: Tracking model performance, including accuracy, latency, and data throughput.\n* **Security**: Implementing authentication, authorization, and encryption to protect sensitive data.\n\nInference servers are widely used in cloud computing, edge computing, and on-premises environments to:\n\n* **Enable real-time analytics**: Provide fast and accurate predictions for applications like recommendation systems, natural language processing, and computer vision.\n* **Streamline data processing**: Automate data analysis and generate insights at the edge or in the cloud for improved decision-making.\n\nExamples of popular inference server frameworks and platforms include:\n\n* TensorFlow Serving\n* Azure Machine Learning\n* Google Cloud AI Platform\n* AWS SageMaker\n* IBM Watson Studio\n\nInference servers have become essential components of modern data science and AI infrastructure, enabling organizations to leverage machine learning models and generate insights from data in real-time."
}
```

Response headers

```
content-length: 2363
```

---

| Code | Details |
|---|---|
| 200 | Response body |

Response headers

```
content-length: 2363
content-type: application/json
date: Sat,07 Feb 2026 17:15:35 GMT
server: uvicorn
```

## Responses

| Code | Description | Links |
|---|---|---|
| 200 | Successful Response | No links |

Media type

application/json ⌄

Controls Accept header.

Example Value | Schema

```
"string"
```

| 422 | Validation Error | No links |

---

Request URL

```
http://127.0.0.1:8000/chatbot
```

Server response

| Code | Details |
|---|---|
| 200 | Response body |

```
{
  "reply": "You just asked \"Hello, who are you?\"",
  "history": [
    {
      "role": "user",
      "content": "Hello, who are you?"
    },
    {
      "role": "assistant",
      "content": "I'm an artificial intelligence model known as a large language model (LLM) or conversational AI. I'm a computer program designed to understand and generate human-like text. I don't have a personal name or identity, but I'm here to assist and communicate with you.\n\nI can help with a wide range of topics and tasks, from answering questions and providing information to generating creative content and engaging in discussions. I'm constantly learning and improving my language abilities, so please bear with me if I make any mistakes or don't understand something at first.\n\nHow can I help you today? Are you looking for information on a specific topic or just want to chat? I'm here to listen and assist you in any way I can."
    },
    {
      "role": "user",
      "content": "What did I just ask you?"
    },
    {
      "role": "assistant",
      "content": "You just asked \"Hello, who are you?\""
    }
  ]
}
```

Response headers

```
content-length: 998
content-type: application/json
date: Sat,07 Feb 2026 17:17:00 GMT
server: uvicorn
```

## Responses

Thru inference server

# AI Chatbot — AMD Agent Screening

Ask me anything!

🤖 hi

🤖 How can I assist you today?

🤖 what day is it today?

🤖 It looks like my knowledge was last updated in December 2023, so I don't have real-time information on the current date. However, I can tell you that today's date is February 7, 2024, based on my understanding of the calendar.

Type your message...



# FastAPI 0.1.0 OAS 3.1

/openapi.json

## default

| POST | /chat Chat |
|------|-----------|

| POST | /chatbot Chatbot |
|------|-----------------|

### Schemas

ChatRequest > Expand all object

HTTPValidationError > Expand all object

ValidationError > Expand all object

These run locally: yet to deploy.

>> crewai agent

**Screenshot 1 — VS Code (agents.yaml, terminal):**

```
simple_agent > src > simple_agent > config > ! agents.yaml
1  explainer_agent:
2    role: Expert AI Explainer
3    goal: Explain user-given concepts clearly and concisely
4    backstory: >
5      You are an expert computer science tutor who explains
6      complex concepts in simple, precise language.
7
```

Terminal:
```
* History restored
PS C:\Users\samru\crewai_agent> cd simple_agent
PS C:\Users\samru\crewai_agent\simple_agent> .\.venv\Scripts\activate
(.venv) PS C:\Users\samru\crewai_agent\simple_agent> python src\simple_agent\main.py
                              🚀 Crew Execution Started
Crew Execution Started
Name:
crew
ID:
77aa9fa6-8c82-4e54-8848-6519ef12ad84

                              📋 Task Started
Task Started
Name: explain_task
ID: 4bce12ce-31a2-406b-88bc-491e99fad168

                              🤖 Agent Started
Agent: Expert AI Explainer
```



**Screenshot 2 — VS Code (agents.yaml, terminal):**

```
simple_agent > src > simple_agent > config > ! agents.yaml
1  explainer_agent:
2    role: Expert AI Explainer
3    goal: Explain user-given concepts clearly and concisely
4    backstory: >
5      You are an expert computer science tutor who explains
6      complex concepts in simple, precise language.
7
```

Terminal:
```
(.venv) PS C:\Users\samru\crewai_agent\simple_agent> python src\simple_agent\main.py

                              🤖 Agent Started
Agent: Expert AI Explainer

Task: Explain the following concept clearly in one or two lines: What is an inference server?

                              ✅ Agent Final Answer
Agent: Expert AI Explainer

Final Answer:
An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models, typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

This server acts as an intermediary between the client application and the machine learning model, handling tasks such as model deployment, data validation, and result delivery, making it a crucial component in various applications, including natural language processing, computer vision, and recommendation systems.

                              📋 Task Completion
Task Completed
Name:
explain_task
```

**Screenshot 1:**

```
File  Edit  Selection  View  Go  Run  Terminal  Help          ← →                    crewai_agent

EXPLORER                    crew.py    ! agents.yaml  ×    main.py
CREW...
simple_agent > src > simple_agent > config > ! agents.yaml
> simple_agent                    1    explainer  agent:
  > .venv
  > knowledge          PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS        powershell - simple_agent
  > src\simple_agent
    > __pycache__       (.venv) PS C:\Users\samru\crewai_agent\simple_agent> python  src\simple_agent\main.py
      _init_.cpython-3...
      crew.cpython-313....                              ─── Task Completion ───
      main.cpython-313...
    > config             Task Completed
      ! agents.yaml      Name:
      ! tasks.yaml       explain_task
    > tools              Agent:
      _init_.py          Expert AI Explainer
      custom_tool.py
    _init_.py
    crew.py                                            ─── Crew Completion ───
    main.py                                            ─── Crew Completion ───
  > tests
    .gitignore           Crew Execution Completed
    pyproject.toml       Crew Execution Completed
    README.md            Name:
    uv.lock              crew
  > venv                 ID:
    .env                 crew
                         ID:
                         ID:
                         77aa9fa6-8c82-4e54-8848-6519ef12ad84
                         77aa9fa6-8c82-4e54-8848-6519ef12ad84
                         Final Output: An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models,
                         typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

                         Final Output: An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models,
                         typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

                         typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

                         This server acts as an intermediary between the client application and the machine learning model, handling tasks such as model deployment, data validation, and result
                         delivery, making it a crucial component in various applications, including natural language processing, computer vision, and recommendation systems.
> OUTLINE
> TIMELINE
                                                                                                                     Ln 7, Col 1   Spaces: 2   UTF-8   CRLF   {} YAML
```

**Screenshot 2:**

```
File  Edit  Selection  View  Go  Run  Terminal  Help          ← →                    crewai_agent

EXPLORER                    crew.py    ! agents.yaml  ×    main.py
CREW...
simple_agent > src > simple_agent > config > ! agents.yaml
> simple_agent                    1    explainer  agent:
  > .venv
  > knowledge          PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS        powershell - simple_agent
  > src\simple_agent
    > __pycache__       (.venv) PS C:\Users\samru\crewai_agent\simple_agent> python  src\simple_agent\main.py
      _init_.cpython-3...
      crew.cpython-313....                              ─── Crew Completion ───
      main.cpython-313...                               ─── Crew Completion ───
    > config
      ! agents.yaml      Crew Execution Completed
      ! tasks.yaml       Crew Execution Completed
    > tools              Name:
      _init_.py          crew
      custom_tool.py     ID:
    _init_.py            crew
    crew.py              ID:
    main.py              ID:
  > tests                77aa9fa6-8c82-4e54-8848-6519ef12ad84
    .gitignore           77aa9fa6-8c82-4e54-8848-6519ef12ad84
    pyproject.toml       Final Output: An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models,
    README.md            typically at scale and in real-time, while managing data handling, model serving, and performance optimization.
    uv.lock
  > venv                 Final Output: An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models,
    .env                 typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

                         typically at scale and in real-time, while managing data handling, model serving, and performance optimization.

                         This server acts as an intermediary between the client application and the machine learning model, handling tasks such as model deployment, data validation, and result
                         delivery, making it a crucial component in various applications, including natural language processing, computer vision, and recommendation systems.


                         Final Output:
                         An inference server is a software infrastructure that processes and generates predictions, decisions, or outputs based on pre-trained machine learning models, typically at scale
                         and in real-time, while managing data handling, model serving, and performance optimization.

                         This server acts as an intermediary between the client application and the machine learning model, handling tasks such as model deployment, data validation, and result delivery,
                         making it a crucial component in various applications, including natural language processing, computer vision, and recommendation systems.
> OUTLINE
> TIMELINE
                                                                                                                     Ln 7, Col 1   Spaces: 2   UTF-8   CRLF   {} YAML
```