



xAI Guides Compilation

This document combines key guidance from multiple xAI documentation pages: **Reasoning**, **Chat**, **The Hitchhiker's Guide to Grok (Tutorial)** and **Tools Overview**. Each section summarizes the original content, preserves headings and code examples, and provides citations to the relevant documentation lines or screenshots.

Reasoning

The **Reasoning** guide explains how xAI's reasoning models, such as `grok-4`, think through problems and expose their internal thought processes. Variants labelled `grok-4-non-reasoning` are based on `grok-4` but disable reasoning. Parameters like `presencePenalty`, `frequencyPenalty` and `stop` are not supported by reasoning models ¹.

Key Features

- **Think Before Responding:** The model reasons step-by-step before delivering an answer ².
- **Math & Quantitative Strength:** Strong at numerical challenges and logic puzzles ³.
- **Reasoning Trace:** The model's thoughts are available via the `reasoning_content` or `encrypted_content` fields in the response object ⁴. You can access this trace through the `message.reasoning_content` attribute of the chat completion response ⁵. For `grok-4` the trace may be encrypted if `use_encrypted_content=true` ⁶.

Control How Hard the Model Thinks

The `reasoning_effort` parameter (not supported by `grok-4`) controls how much time the model spends thinking. It accepts two values ⁷:

- `low` – minimal thinking time for quick responses.
- `high` – maximum thinking time for complex problems.

Usage Example

Below is a complete example showing how to use `grok-4` to multiply `101` by `3` and report both the completion and reasoning token counts. Note that reasoning models require a longer timeout to allow for their internal reasoning:

```
import os

from xai_sdk import Client
from xai_sdk.chat import system, user
```

```

client = Client(
    api_key=os.getenv("XAI_API_KEY"),
    timeout=3600, # Override default timeout with longer timeout for reasoning
    models
)

chat = client.chat.create(
    model="grok-4",
    messages=[system("You are a highly intelligent AI assistant.")],
)
chat.append(user("What is 101*3?"))

response = chat.sample()

print("Final Response:")
print(response.content)
print("\nNumber of completion tokens:")
print(response.usage.completion_tokens)
print("\nNumber of reasoning tokens:")
print(response.usage.reasoning_tokens)

```

Sample output:

```

Final Response:
The result of 101 multiplied by 3 is 303.

Number of completion tokens:
14

Number of reasoning tokens:
310

```

(Code and output from the Reasoning page ⁸.)

Notes on Consumption

Reasoning tokens are counted in your final consumption, so using a higher `reasoning_effort` will increase token usage ⁹.

Chat

The **Chat** guide introduces the xAI chat completions API. Chat accepts “text in, text out” and is the most popular feature of xAI API ¹⁰. It supports summarizing articles, creative writing, Q&A, customer support and even coding tasks ¹⁰.

Prerequisites

- **xAI Account:** You need an xAI account to access the API ¹¹.
- **API Key:** Ensure your API key has access to the chat endpoint and the chat model ¹².

If you don't know how to create these, follow the Hitchhiker's Guide to Grok (see tutorial section) ¹³. You can create an API key via the xAI console ¹⁴.

Basic Chat Completions Example

Chat requests send a list of messages to the API; the model processes them and returns a response. You can also stream the response, covered in the Streaming Response guide ¹⁵.

Here is a simple chat example using the xAI Python SDK ¹⁶:

```
import os

from xai_sdk import Client
from xai_sdk.chat import user, system

client = Client(
    api_key=os.getenv("XAI_API_KEY"),
    timeout=3600,
)

chat = client.chat.create(model="grok-4")
chat.append(system("You are a PhD-level mathematician."))
chat.append(user("What is 2 + 2?"))

response = chat.sample()
print(response.content)
```

The response would be:

```
'2 + 2 equals 4.'
```

Conversations

xAI's API is stateless; each chat request is processed independently without automatically using previous history. To preserve context, you must include prior messages in the `messages` list ¹⁷. Messages have roles:

- `system` – defines how the model should respond (instructions) ¹⁸.
- `user` – represents user requests or data ¹⁸.
- `assistant` – used for the model's responses or to inject previous answers ¹⁹.

For example, a conversation might look like the following (JSON format) ²⁰ :

```
[
  {
    "role": "system",
    "content": [{ "type": "text", "text": "You are a helpful and funny
assistant." }]
  },
  {
    "role": "user",
    "content": [{ "type": "text", "text": "Why don't eggs tell jokes?" }]
  },
  {
    "role": "assistant",
    "content": [{ "type": "text", "text": "They'd crack up!" }]
  },
  {
    "role": "user",
    "content": [{ "type": "text", "text": "Can you explain the joke?" }]
  }
]
```

Message Role Order Flexibility

Unlike some providers, xAI allows mixing `system`, `user` and `assistant` messages in any order ²¹ . You can have multiple system messages at different points in the conversation or start with a user message first. Example structures include multiple system messages or user-first messages ²² .

The Hitchhiker's Guide to Grok (Tutorial)

This tutorial walks you through getting started with the xAI API ²³ .

Step 1: Create an xAI Account

You need to create an xAI account to access the API ²⁴ . After creating an account, add credits before using the API ²⁵ .

Step 2: Generate an API Key

Create an API key from the **API Keys Page** in the xAI console ²⁶ . Store it securely, ideally as an environment variable or in a `.env` file ²⁷ .

Step 3: Make Your First Request

With your API key available, you can make your first API request. The tutorial uses `curl` to send a request to the `/chat/completions` endpoint ²⁸:

```
curl https://api.x.ai/v1/chat/completions \
-H "Content-Type: application/json" \
-H "Authorization: Bearer $XAI_API_KEY" \
-m 3600 \
-d '{
  "messages": [
    {
      "role": "system",
      "content": "You are Grok, a highly intelligent, helpful AI assistant."
    },
    {
      "role": "user",
      "content": "What is the meaning of life, the universe, and everything?"
    }
  ],
  "model": "grok-4",
  "stream": false
}'
```

Step 4: Make a Request from Python or Javascript

You can also perform the same request via the xAI Python SDK (or compatible OpenAI/Anthropic SDKs) ²⁹. For Python ³⁰:

```
import os

from xai_sdk import Client
from xai_sdk.chat import user, system

client = Client(
    api_key=os.getenv("XAI_API_KEY"),
    timeout=3600,
)

chat = client.chat.create(
    model="grok-4",
    messages=[system("You are Grok, a highly intelligent, helpful AI
assistant.")],
)
chat.append(user("What is the meaning of life, the universe, and everything?"))
```

```
response = chat.sample()
print(response.content)
```

Step 5: Use Grok to Analyze Images

Certain Grok models accept both text and images. For example ³¹ :

```
import os

from xai_sdk import Client
from xai_sdk.chat import user, image

client = Client(
    api_key=os.getenv("XAI_API_KEY"),
    timeout=3600,
)

chat = client.chat.create(model="grok-4")
chat.append(
    user(
        "What's in this image?",
        image("https://science.nasa.gov/wp-content/uploads/2023/09/web-first-
images-release.png")
    )
)

response = chat.sample()
print(response.content)
```

The agent will respond with a description such as:

"This image is a photograph of a region in space, specifically a part of the Carina Nebula, captured by the James Webb Space Telescope..." ³² .

Monitoring Usage

As you use your API key, you are charged based on the number of tokens used. You can monitor usage on the xAI Console Usage Page ³³ . For per-request tracking, each API response includes a `usage` object detailing prompt (input) and completion (output) token usage ³⁴ . An example usage object ³⁵ :

```
{
  "usage": {
    "prompt_tokens": 37,
    "completion_tokens": 530,
    "total_tokens": 800,
```

```

    "prompt_tokens_details": {
      "text_tokens": 37,
      "audio_tokens": 0,
      "image_tokens": 0,
      "cached_tokens": 8
    },
    "completion_tokens_details": {
      "reasoning_tokens": 233,
      "audio_tokens": 0,
      "accepted_prediction_tokens": 0,
      "rejected_prediction_tokens": 0
    },
    "num_sources_used": 0
  }
}

```

If you send requests too frequently or with long prompts, you might hit rate limits; see the [Consumption and Rate Limits guide](#) for details ³⁶.

Next Steps

After learning the basics, explore the **Models** page to start building with one of xAI's latest models ³⁷.

Tools Overview

The **Tools Overview** page explains xAI's agentic server-side tool calling. Unlike traditional tool calling where clients handle each tool invocation, xAI's agentic API manages the entire reasoning and tool-execution loop on the server ³⁸. Using version 1.3.1 or higher of the xAI SDK is required to access these features ³⁹.

Tools Pricing

Agentic requests incur costs based on **token usage** and **tool invocations** ⁴⁰. Because the agent autonomously decides how many tools to call, costs scale with query complexity. See the [pricing page](#) for more details ⁴¹.

Agentic Tool Calling

When you provide server-side tools in a request, the xAI server runs an autonomous reasoning loop. Instead of returning tool calls for the client to execute, the agent researches, analyzes and responds automatically ⁴². Behind the scenes, the model follows an iterative reasoning process ⁴³:

1. **Analyzes the query and current context** to determine what information is needed.
2. **Decides what to do next** – either make a tool call to gather more data or provide a final answer ⁴⁴.
3. **If making a tool call** – selects the appropriate tool and parameters ⁴⁵.
4. **Executes the tool** in real time on the server and receives the results ⁴⁶.

5. **Processes the tool response** and integrates it with previous context ⁴⁷ .
6. **Repeats** the reasoning loop as needed until enough information is gathered ⁴⁸ .
7. **Returns the final response** once sufficient information is available ⁴⁹ .

Core Capabilities

xAI's agentic tool calling supports several core tools ⁵⁰ :

- **Web Search:** Real-time search across the internet and web browsing ⁵¹ .
- **X Search:** Semantic and keyword search across X posts, users and threads ⁵² .
- **Code Execution:** Write and execute Python code for calculations, data analysis and complex computations ⁵³ .
- **Image/Video Understanding:** Optional analysis of visual content in search results ⁵⁴ .

Quick Start Example

It's recommended to use the xAI Python SDK in **streaming** mode to get real-time observability and immediate feedback ⁵⁵ . A quick start example to get the latest updates from xAI might look like this ⁵⁶ :

```
import os

from xai_sdk import Client
from xai_sdk.chat import user
from xai_sdk.tools import web_search, x_search, code_execution

client = Client(api_key=os.getenv("XAI_API_KEY"))
chat = client.chat.create(
    model="grok-4-fast", # reasoning model
    tools=[
        web_search(),
        x_search(),
        code_execution(),
    ],
)

# Ask a question
chat.append(user("What are the latest updates from xAI?"))

is_thinking = True
for response, chunk in chat.stream():
    # Display tool calls as they occur
    for tool_call in chunk.tool_calls:
        print(f"\nCalling tool: {tool_call.function.name} with arguments: {tool_call.function.arguments}")
        if response.usage.reasoning_tokens and is_thinking:
            print(f"\rThinking... ({response.usage.reasoning_tokens} tokens)",
end="", flush=True)
```



```

    if chunk.content and is_thinking:
        print("\n\nFinal Response:")
        is_thinking = False
    if chunk.content and not is_thinking:
        print(chunk.content, end="", flush=True)

print("\n\nCitations:")
print(response.citations)
print("\n\nUsage:")
print(response.usage)
print(response.server_side_tool_usage)

```

This script streams the response, prints each tool call as it happens, and finally prints citations and usage statistics. During streaming you will see `Thinking...` messages and tool call notifications ⁵⁶.

A synchronous (non-streaming) version waits for the full response before printing anything ⁵⁷:

```

import os

from xai_sdk import Client
from xai_sdk.chat import user
from xai_sdk.tools import web_search, x_search, code_execution

client = Client(api_key=os.getenv("XAI_API_KEY"))
chat = client.chat.create(
    model="grok-4-fast",
    tools=[
        web_search(),
        x_search(),
        code_execution(),
    ],
)
chat.append(user("What is the latest update from xAI?"))

# Wait for the entire agentic process to finish
response = chat.sample()

print("\n\nFinal Response:")
print(response.content)

print("\n\nCitations:")
print(response.citations)

print("\n\nUsage:")
print(response.usage)
print(response.server_side_tool_usage)

```

```
print("\n\nServer Side Tool Calls:")
print(response.tool_calls)
```

Understanding the Agentic Tool Calling Response

The agentic tool calling API provides rich observability into the research process. You can view real-time tool call decisions via the `tool_calls` attribute on each `chunk` during streaming ⁵⁸. Only tool call invocations are shown; the outputs of those calls are not returned directly ⁵⁹. When using the xAI SDK in streaming mode, the SDK accumulates `tool_calls` in the final `response` object for later inspection ⁶⁰.

Real-time Server-side Tool Calls

To print each tool call as it happens, iterate over `chunk.tool_calls` during streaming and inspect each `function.name` and `function.arguments` ⁵⁸:

```
for tool_call in chunk.tool_calls:
    print(f"\nCalling tool: {tool_call.function.name} with arguments:
    {tool_call.function.arguments}")
```

Citations

The `citations` attribute on the `response` object lists all URLs encountered during the agent's search process ⁶¹. Citations are returned only when the agentic request completes and are not available in real time ⁶². Not every URL may be relevant to the final answer because the agent may discard sources that are not useful ⁶³.

Server-side Tool Calls vs Tool Usage

The API distinguishes between two metrics ⁶⁴:

- `tool_calls` – **All Attempted Calls**: A list of every attempted tool call, including those that fail ⁶⁵. Each entry is a `ToolCall` object containing an `id`, `function.name` and `function.arguments` ⁶⁶.
- `server_side_tool_usage` – **Successful Calls (Billable)**: A map of tools that executed successfully and how many times each was called ⁶⁷. Only these calls are billable ⁶⁸.

Tool Call Function Names vs Usage Categories

The `tool_calls` list uses precise function names, while `server_side_tool_usage` uses higher-level categories. The mapping is ⁶⁹:

Usage Category	Function Name(s)
SERVER_SIDE_TOOL_WEB_SEARCH	<code>web_search</code> , <code>web_search_with_snippets</code> , <code>browse_page</code>
SERVER_SIDE_TOOL_X_SEARCH	<code>x_user_search</code> , <code>x_keyword_search</code> , <code>x_semantic_search</code> , <code>x_thread_fetch</code>
SERVER_SIDE_TOOL_CODE_EXECUTION	<code>code_execution</code>
SERVER_SIDE_TOOL_VIEW_X_VIDEO	<code>view_x_video</code>
SERVER_SIDE_TOOL_VIEW_IMAGE	<code>view_image</code>

When Tool Calls and Usage Differ

`tool_calls` and `server_side_tool_usage` will differ when ⁷⁰ :

- **Failed tool executions:** The agent browses a non-existent web page, fetches a deleted X post or encounters other execution errors.
- **Invalid parameters:** Tool calls with malformed arguments.
- **Network or service issues:** Temporary failures in the tool execution pipeline.

The agent is robust and continues searching even after failures ⁷¹ . Only successful calls are billed ⁶⁸ .

Understanding Token Usage

Agentic requests have unique token usage patterns ⁷² :

- `completion_tokens` – tokens for the final output text ⁷³ .
- `prompt_tokens` – cumulative input tokens across all inference requests ⁷⁴ . Agentic workflows involve multiple reasoning steps; however prompt caching allows shared context to be reused, improving efficiency ⁷⁵ .
- `reasoning_tokens` – tokens used for the model's internal reasoning, excluding final outputs ⁷⁶ .
- `cached_prompt_text_tokens` – number of prompt tokens served from cache ⁷⁷ .
- `prompt_image_tokens` – tokens from visual content (images or videos) processed ⁷⁸ .
- `prompt_text_tokens` and `total_tokens` – text tokens in prompts (excluding special tokens) and the sum of all token types ⁷⁹ .

Synchronous Agentic Requests (Non-streaming)

You can make non-streaming requests for simpler use cases, waiting for the entire agentic process to finish before receiving the response ⁸⁰ . The synchronous example shown earlier demonstrates how to access final content, citations, usage and tool calls in one go ⁵⁷ .

Using Tools with OpenAI Responses API

Agentic tool calling also works via the OpenAI Responses API. For example ⁸¹ :

```

import os

from openai import OpenAI

api_key = os.getenv("XAI_API_KEY")
client = OpenAI(
    api_key=api_key,
    base_url="https://api.x.ai/v1",
)

response = client.responses.create(
    model="grok-4-fast",
    input=[
        {
            "role": "User",
            "content": "What is the latest update from xAI?",
        },
    ],
    tools=[
        {"type": "web_search"},
        {"type": "x_search"},
    ],
)

print(response)

```

Agentic Tool Calling Requirements and Limitations

- **Model Compatibility:** Supported models are `grok-4`, `grok-4-fast` and `grok-4-fast-non-reasoning` ⁸². It is strongly recommended to use `grok-4-fast` for agentic tool calling ⁸³.
- **Tool Configuration:** Only server-side tools are allowed; you cannot mix server-side and client-side tools in the same request ⁸⁴.
- **Request Constraints:** Batch requests (`n > 1`) are not supported. Structured output formats are not available yet. Only `temperature` and `top_p` sampling parameters are respected ⁸⁵. These constraints may be relaxed in future updates ⁸⁶.

FAQ and Troubleshooting

If you see empty or incorrect content when using agentic tool calling with the xAI SDK, ensure you are using SDK version `1.3.1` or higher ⁸⁷.

This combined document brings together key points from multiple xAI guides, preserving the original instructions, code examples and context. Use it as a single reference for reasoning models, chat completions, getting started with Grok and understanding the agentic tools ecosystem.

1 2 3 4 5 6 7 8 9 Reasoning

<https://docs.x.ai/docs/guides/reasoning>

10 11 12 13 14 15 16 17 18 19 20 21 22 Chat

<https://docs.x.ai/docs/guides/chat>

23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 The Hitchhiker's Guide to Grok | xAI Docs

<https://docs.x.ai/docs/tutorial>

38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67

68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 Overview

<https://docs.x.ai/docs/guides/tools/overview>