

# API Calls to Large Language Models

Seth J. Chandler, with help from AI

August 21, 2025



# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Basic LLM API Call</b>                       | <b>1</b> |
| <b>2</b> | <b>Advanced LLM API Usage</b>                   | <b>3</b> |
| <b>3</b> | <b>Basic LLM API Call with Web Search</b>       | <b>7</b> |
| <b>4</b> | <b>Basic LLM API Call in Multiple Languages</b> | <b>9</b> |
| 4.1      | JavaScript . . . . .                            | 9        |
| 4.2      | Node.js . . . . .                               | 10       |
| 4.3      | C# . . . . .                                    | 10       |
| 4.4      | Go . . . . .                                    | 11       |
| 4.5      | TypeScript . . . . .                            | 13       |
| 4.6      | Wolfram Language . . . . .                      | 13       |



# Chapter 1

## Basic LLM API Call

This chapter provides a simple Python script demonstrating a basic API call to a Large Language Model (LLM). The script sends a user prompt to the model and retrieves the response, using minimal configuration options to illustrate the core functionality of an LLM API for educational purposes.

```
1 import requests
2
3 # API endpoint and credentials (hypothetical)
4 API_URL = "https://api.x.ai/v1/chat/completions"
5 API_KEY = "your-api-key-here" # Replace with actual API key
6
7 def basic_llm_api_call(user_prompt: str) -> str:
8     """
9     Make a basic API call to an LLM with a user prompt.
10
11     Args:
12         user_prompt: The user's input prompt
13
14     Returns:
15         The model's response as a string or an error message
16     """
17     headers = {
18         "Authorization": f"Bearer {API_KEY}",
19         "Content-Type": "application/json"
20     }
21
22     payload = {
23         "model": "grok-3",
24         "messages": [
25             {"role": "user", "content": user_prompt}
26         ]
27     }
28
29     try:
30         response = requests.post(API_URL, headers=headers, json=payload)
31         response.raise_for_status()
32         return response.json()[0]["choices"][0]["message"]["content"]
33     except requests.exceptions.RequestException as e:
34         return f"Error: {str(e)}"
35
36 # Example usage
37 if __name__ == "__main__":
38     prompt = "Explain what a transformer model is in one sentence."
39     result = basic_llm_api_call(prompt)
40     print("Model Response:", result)
```

Listing 1.1: Simple Python script for an LLM API call



## Chapter 2

# Advanced LLM API Usage

This chapter presents a Python script that demonstrates how to interact with a Large Language Model (LLM) API using advanced configuration options. The code showcases parameters such as temperature, system prompts, output formats, and tool integration, providing insight into the inner workings of LLMs for educational purposes.

```
1 import requests
2 import json
3 from typing import List, Dict, Any
4
5 # API endpoint and credentials (hypothetical)
6 API_URL = "https://api.x.ai/v1/chat/completions"
7 API_KEY = "your-api-key-here" # Replace with actual API key
8
9 # Define a tool for function calling (e.g., a calculator tool)
10 tools = [
11     {
12         "type": "function",
13         "function": {
14             "name": "calculate",
15             "description": "Perform mathematical calculations",
16             "parameters": {
17                 "type": "object",
18                 "properties": {
19                     "expression": {
20                         "type": "string",
21                         "description": "Mathematical expression to evaluate (e.g., '2 + 3 * 4')",
22                     }
23                 },
24                 "required": ["expression"]
25             }
26         }
27     }
28 ]
29
30 def make_llm_api_call(
31     user_prompt: str,
32     system_prompt: str = "You are a helpful assistant with expertise in explaining technical concepts.",
33     temperature: float = 0.7,
34     max_tokens: int = 512,
35     top_p: float = 0.9,
36     presence_penalty: float = 0.1,
37     frequency_penalty: float = 0.2,
38     response_format: Dict[str, str] = {"type": "json_object"},
39     tools: List[Dict[str, Any]] = None,
40     tool_choice: str = "auto"
41 ) -> Dict[str, Any]:
42     """
```

```

43  Make an API call to the LLM with advanced configuration options.
44
45  Args:
46      user_prompt: The user's input prompt
47      system_prompt: System prompt to guide the model's behavior
48      temperature: Controls randomness (0.0 to 2.0)
49      max_tokens: Maximum number of tokens in the response
50      top_p: Nucleus sampling parameter
51      presence_penalty: Penalty for new tokens
52      frequency_penalty: Penalty for frequent tokens
53      response_format: Desired output format (e.g., JSON)
54      tools: List of available tools for function calling
55      tool_choice: Strategy for tool usage ('auto', 'none', or specific tool)
56
57  Returns:
58      API response as a dictionary
59  """
60  headers = {
61      "Authorization": f"Bearer {API_KEY}",
62      "Content-Type": "application/json"
63  }
64
65  payload = {
66      "model": "grok-3",
67      "messages": [
68          {"role": "system", "content": system_prompt},
69          {"role": "user", "content": user_prompt}
70      ],
71      "temperature": temperature,
72      "max_tokens": max_tokens,
73      "top_p": top_p,
74      "presence_penalty": presence_penalty,
75      "frequency_penalty": frequency_penalty,
76      "response_format": response_format,
77      "tools": tools,
78      "tool_choice": tool_choice
79  }
80
81  try:
82      response = requests.post(API_URL, headers=headers, json=payload)
83      response.raise_for_status()
84      return response.json()
85  except requests.exceptions.RequestException as e:
86      return {"error": str(e)}
87
88  # Example usage
89  def demonstrate_llm_api():
90      user_prompt = """
91      Explain how a transformer model processes input for a college audience.
92      Include a simple calculation example: what is 5 + 3 * 4?
93      Return the response in JSON format with sections for 'explanation' and 'calculation'.
94      """
95
96      # Make the API call with custom parameters
97      response = make_llm_api_call(
98          user_prompt=user_prompt,
99          system_prompt="""
100      You are a patient, clear instructor teaching college students.
101      Break down complex concepts into simple terms and provide structured JSON output.
102      Use the calculator tool for any mathematical operations.
103      """,
104          temperature=0.8, # Moderate randomness for clear but varied responses
105          max_tokens=1000, # Allow longer responses for detailed explanations
106          top_p=0.95, # Consider 95% of probability mass for token sampling
107          presence_penalty=0.3, # Encourage new topics
108          frequency_penalty=0.4, # Reduce repetition
109          response_format={
110              "type": "json_object",

```



```

111         "schema": {
112             "type": "object",
113             "properties": {
114                 "explanation": {"type": "string"},
115                 "calculation": {"type": "string"}
116             },
117             "required": ["explanation", "calculation"]
118         }
119     },
120     tools=tools,
121     tool_choice="auto" # Let the model decide when to use the calculator tool
122 )
123
124 # Process and display the response
125 if "error" in response:
126     print(f"Error: {response['error']}")
127 else:
128     print("API Response:")
129     print(json.dumps(response, indent=2))
130
131 # Handle tool call if present
132 if "choices" in response and response["choices"][0].get("tool_calls"):
133     tool_call = response["choices"][0]["tool_calls"][0]
134     if tool_call["function"]["name"] == "calculate":
135         expression = json.loads(tool_call["function"]["arguments"])["expression"]
136         print(f"Tool called: calculate({expression})")
137         # Simulate tool execution (in reality, you'd evaluate the expression)
138         result = eval(expression, {"__builtins__": {}}) # Safe eval for demo
139         print(f"Tool result: {result}")
140
141 if __name__ == "__main__":
142     demonstrate_llm_api()

```

Listing 2.1: Python script for advanced LLM API call



## Chapter 3

# Basic LLM API Call with Web Search

This chapter presents a simple Python script that demonstrates an API call to a Large Language Model (LLM) endpoint with integrated web search capabilities. The script sends a user query to retrieve real-time web information and returns a summarized response, illustrating the core functionality of a web search-enabled LLM API for educational purposes.

The API endpoint used in this example is a hypothetical `https://api.x.ai/v1/web-search/completions`, designed to mimic real-world APIs like Tavily or SerpAPI, which provide web search results as context for LLMs. The script uses minimal parameters: only the user query and a limit on the number of web results (set to 5 for brevity), making it accessible for students. The response is extracted from a JSON structure, assuming the LLM summarizes the web search results. For real-world usage, the API URL and key would need to be replaced with values from a service like Tavily (`https://tavily.com`) or SerpAPI (`https://serpapi.com`). For xAI-specific APIs, students should refer to `https://x.ai/api`, as no public xAI web search API is available as of August 21, 2025.

```
1 import requests
2 import json
3
4 # API endpoint and credentials (hypothetical)
5 API_URL = "https://api.x.ai/v1/web-search/completions"
6 API_KEY = "your-api-key-here" # Replace with actual API key
7
8 def web_search_llm_api_call(user_query: str) -> str:
9     """
10     Make a basic API call to an LLM with web search capabilities.
11
12     Args:
13         user_query: The user's search query
14
15     Returns:
16         The model's summarized response based on web search results or an error message
17     """
18     headers = {
19         "Authorization": f"Bearer {API_KEY}",
20         "Content-Type": "application/json"
21     }
22
23     payload = {
24         "model": "grok-3-web",
25         "query": user_query,
26         "max_results": 5 # Limit to 5 web search results for context
27     }
28
29     try:
30         response = requests.post(API_URL, headers=headers, json=payload)
31         response.raise_for_status()
32         return response.json()["response"]["content"]
33     except requests.exceptions.RequestException as e:
34         return f"Error: {str(e)}"
```

```
35
36 # Example usage
37 if __name__ == "__main__":
38     query = "What are the latest advancements in transformer models?"
39     result = web_search_llm_api_call(query)
40     print("Model Response:", result)
```

Listing 3.1: Python script for basic LLM API call with web search

## Chapter 4

# Basic LLM API Call in Multiple Languages

This chapter presents a simple API call to a Large Language Model (LLM) implemented in six programming languages: JavaScript, Node.js, C#, Go, TypeScript, and Wolfram Language. Each example sends a user prompt to a hypothetical LLM endpoint (<https://api.x.ai/v1/chat/completions>) and retrieves the response, using minimal configuration to illustrate the core functionality for educational purposes. The Wolfram Language example leverages the high-level `ChatSubmit` function for a terser implementation. The code blocks are labeled by language for clarity.

### 4.1 JavaScript

This example uses the `fetch` API in JavaScript to make an HTTP POST request, suitable for browser-based environments.

```
1 // API endpoint and credentials (hypothetical)
2 const API_URL = "https://api.x.ai/v1/chat/completions";
3 const API_KEY = "your-api-key-here"; // Replace with actual API key
4
5 async function basicLLMApiCall(userPrompt) {
6     // Make a basic API call to an LLM
7     const headers = {
8         "Authorization": `Bearer ${API_KEY}`,
9         "Content-Type": "application/json"
10    };
11
12    const payload = {
13        model: "grok-3",
14        messages: [{ role: "user", content: userPrompt }]
15    };
16
17    try {
18        const response = await fetch(API_URL, {
19            method: "POST",
20            headers: headers,
21            body: JSON.stringify(payload)
22        });
23        if (!response.ok) throw new Error(`HTTP error: ${response.status}`);
24        const data = await response.json();
25        return data.choices[0].message.content;
26    } catch (error) {
27        return `Error: ${error.message}`;
28    }
29 }
30
31 // Example usage
```

```

32 (async () => {
33     const prompt = "Explain what a transformer model is in one sentence.";
34     const result = await basicLLMApiCall(prompt);
35     console.log("Model Response:", result);
36 }) ();

```

Listing 4.1: JavaScript implementation of basic LLM API call

## 4.2 Node.js

This example uses the `axios` library in Node.js for the HTTP request, suitable for server-side environments.

```

1  const axios = require('axios');
2
3  // API endpoint and credentials (hypothetical)
4  const API_URL = 'https://api.x.ai/v1/chat/completions';
5  const API_KEY = 'your-api-key-here'; // Replace with actual API key
6
7  async function basicLLMApiCall(userPrompt) {
8      // Make a basic API call to an LLM
9      const headers = {
10         'Authorization': `Bearer ${API_KEY}`,
11         'Content-Type': 'application/json'
12     };
13
14     const payload = {
15         model: 'grok-3',
16         messages: [{ role: 'user', content: userPrompt }]
17     };
18
19     try {
20         const response = await axios.post(API_URL, payload, { headers });
21         return response.data.choices[0].message.content;
22     } catch (error) {
23         return `Error: ${error.message}`;
24     }
25 }
26
27 // Example usage
28 (async () => {
29     const prompt = 'Explain what a transformer model is in one sentence.';
30     const result = await basicLLMApiCall(prompt);
31     console.log('Model Response:', result);
32 }) ();

```

Listing 4.2: Node.js implementation of basic LLM API call

## 4.3 C#

This example uses `HttpClient` in C# to make the API call, suitable for .NET applications.

```

1  using System;
2  using System.Net.Http;
3  using System.Text;
4  using System.Threading.Tasks;
5  using System.Text.Json;
6
7  class Program {
8      // API endpoint and credentials (hypothetical)
9      private static readonly string API_URL = "https://api.x.ai/v1/chat/completions";
10     private static readonly string API_KEY = "your-api-key-here"; // Replace with actual API key

```

```

11
12     static async Task<string> BasicLLMApiCall(string userPrompt) {
13         // Make a basic API call to an LLM
14         using (var client = new HttpClient()) {
15             client.DefaultRequestHeaders.Add("Authorization", $"Bearer {API_KEY}");
16
17             var payload = new {
18                 model = "grok-3",
19                 messages = new[] { new { role = "user", content = userPrompt } }
20             };
21
22             var content = new StringContent(JsonSerializer.Serialize(payload), Encoding.UTF8,
23                 "application/json");
24
25             try {
26                 var response = await client.PostAsync(API_URL, content);
27                 response.EnsureSuccessStatusCode();
28                 var jsonResponse = await response.Content.ReadAsStringAsync();
29                 var data = JsonSerializer.Deserialize<JsonElement>(jsonResponse);
30                 return data.GetProperty("choices")[0].GetProperty("message").GetProperty("content").GetString();
31             } catch (Exception ex) {
32                 return $"Error: {ex.Message}";
33             }
34         }
35
36         // Example usage
37         static async Task Main() {
38             string prompt = "Explain what a transformer model is in one sentence.";
39             string result = await BasicLLMApiCall(prompt);
40             Console.WriteLine("Model Response: " + result);
41         }
42     }

```

Listing 4.3: C# implementation of basic LLM API call

## 4.4 Go

This example uses the `net/http` package in Go to make the API call, suitable for lightweight server applications.

```

1 package main
2
3 import (
4     "bytes"
5     "encoding/json"
6     "fmt"
7     "net/http"
8 )
9
10 // API endpoint and credentials (hypothetical)
11 const apiURL = "https://api.x.ai/v1/chat/completions"
12 const apiKey = "your-api-key-here" // Replace with actual API key
13
14 type Message struct {
15     Role    string `json:"role"`
16     Content string `json:"content"`
17 }
18
19 type Payload struct {
20     Model    string    `json:"model"`
21     Messages []Message `json:"messages"`
22 }
23

```

```

24 type Response struct {
25     Choices []struct {
26         Message struct {
27             Content string `json:"content"`
28         } `json:"message"`
29     } `json:"choices"`
30 }
31
32 func basicLLMApiCall(userPrompt string) (string, error) {
33     // Make a basic API call to an LLM
34     payload := Payload{
35         Model: "grok-3",
36         Messages: []Message{{Role: "user", Content: userPrompt}},
37     }
38
39     payloadBytes, err := json.Marshal(payload)
40     if err != nil {
41         return "", err
42     }
43
44     req, err := http.NewRequest("POST", apiURL, bytes.NewBuffer(payloadBytes))
45     if err != nil {
46         return "", err
47     }
48
49     req.Header.Set("Authorization", "Bearer "+apiKey)
50     req.Header.Set("Content-Type", "application/json")
51
52     client := &http.Client{}
53     resp, err := client.Do(req)
54     if err != nil {
55         return "", err
56     }
57     defer resp.Body.Close()
58
59     if resp.StatusCode != http.StatusOK {
60         return "", fmt.Errorf("HTTP error: %d", resp.StatusCode)
61     }
62
63     var response Response
64     err = json.NewDecoder(resp.Body).Decode(&response)
65     if err != nil {
66         return "", err
67     }
68
69     return response.Choices[0].Message.Content, nil
70 }
71
72 func main() {
73     // Example usage
74     prompt := "Explain what a transformer model is in one sentence."
75     result, err := basicLLMApiCall(prompt)
76     if err != nil {
77         fmt.Println("Error:", err)
78     } else {
79         fmt.Println("Model Response:", result)
80     }
81 }

```

Listing 4.4: Go implementation of basic LLM API call



## 4.5 TypeScript

This example uses TypeScript with the `fetch` API, adding type definitions for better type safety, suitable for browser-based applications.

```

1 // API endpoint and credentials (hypothetical)
2 const API_URL: string = "https://api.x.ai/v1/chat/completions";
3 const API_KEY: string = "your-api-key-here"; // Replace with actual API key
4
5 interface Message {
6   role: string;
7   content: string;
8 }
9
10 interface Payload {
11   model: string;
12   messages: Message[];
13 }
14
15 interface ApiResponse {
16   choices: { message: { content: string } }[];
17 }
18
19 async function basicLLMApiCall(userPrompt: string): Promise<string> {
20   // Make a basic API call to an LLM
21   const headers: HeadersInit = {
22     "Authorization": `Bearer ${API_KEY}`,
23     "Content-Type": "application/json"
24   };
25
26   const payload: Payload = {
27     model: "grok-3",
28     messages: [{ role: "user", content: userPrompt }]
29   };
30
31   try {
32     const response = await fetch(API_URL, {
33       method: "POST",
34       headers,
35       body: JSON.stringify(payload)
36     });
37     if (!response.ok) throw new Error(`HTTP error: ${response.status}`);
38     const data: ApiResponse = await response.json();
39     return data.choices[0].message.content;
40   } catch (error: any) {
41     return `Error: ${error.message}`;
42   }
43 }
44
45 // Example usage
46 (async () => {
47   const prompt: string = "Explain what a transformer model is in one sentence.";
48   const result: string = await basicLLMApiCall(prompt);
49   console.log("Model Response:", result);
50 })();

```

Listing 4.5: TypeScript implementation of basic LLM API call

## 4.6 Wolfram Language

This example uses the high-level `ChatSubmit` function in Wolfram Language, specifying the model and API key via `LLMConfiguration`, to make a concise API call, suitable for Mathematica environments.

```
1  (* API credentials and configuration (hypothetical) *)
2  apiKey = "your-api-key-here"; (* Replace with actual API key *)
3
4  basicLLMApiCall[userPrompt_] := Module[{response},
5      (* Make a basic API call to an LLM *)
6      response = ChatSubmit[
7          userPrompt,
8          LLMConfiguration[
9              <|
10                 "Model" -> "grok-3",
11                 "Service" -> "Custom",
12                 "BaseURL" -> "https://api.x.ai/v1/chat/completions",
13                 "APIKey" -> apiKey
14             |>
15          ]
16      ];
17
18      If[FailureQ[response],
19          "Error: " <> ToString[response["Message"]],
20          response["Content"]
21      ]
22  ]
23
24  (* Example usage *)
25  prompt = "Explain what a transformer model is in one sentence.";
26  result = basicLLMApiCall[prompt];
27  Print["Model Response: ", result]
```

Listing 4.6: Wolfram Language implementation of basic LLM API call