Microsoft

**Generative AI** with JavaScript

# Streaming generative AI output with the AI Chat Protocol

# Agenda

- Streaming GenAI output

- Introducing the AI Chat Protocol

- Usage Sample

# Streaming GenAI output

# Why use streaming?

# Why use streaming?

Streaming is the expectation for GenAI apps

- Reduced latency
- Enhanced user experience

# Why use streaming?

Streaming is the expectation for GenAI apps

- Reduced latency

- Enhanced user experience

2 options for implementing

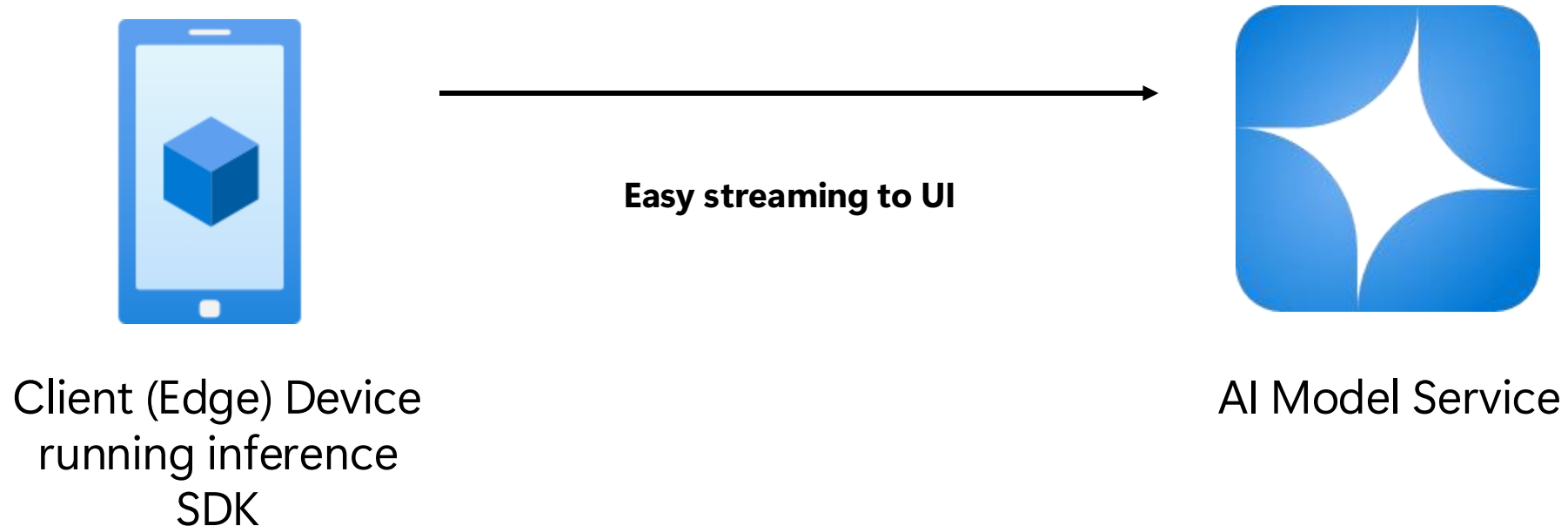- Use a GenAI Inference SDK directly in the browser

- Use an AI inference server to stream to your client

# Streaming - Inference in Browser

· The simplest approach to streaming is by using the SDK in the edge device
· For example, using Azure OpenAI SDK for JS in your frontend browser code

**Easy streaming to UI**

Client (Edge) Device
running inference
SDK

AI Model Service

# Streaming - Inference in Browser

- The simplest approach to streaming is by using the SDK in the edge device
- For example, using Azure OpenAI SDK for JS in your frontend browser code

**Easy streaming to UI**

Client (Edge) Device
running inference
SDK

AI Model Service

# Streaming - Inference in Browser

- The simplest approach to streaming is by using the SDK in the edge device
- For example, using Azure OpenAI SDK for JS in your frontend browser code

**Easy streaming to UI**

Client (Edge) Device
running inference
SDK

AI Model Service

# Browser-side inference is <u>unsafe</u> and not best practice

# Browser-side inference is <u>unsafe</u> and not best practice

## Security risks

- API key exposure
- No data sanitization of input
- No integration with data compliance

# Browser-side inference is <u>unsafe</u> and not best practice

## Security risks

- API key exposure
- No data sanitization of input
- No integration with data compliance

## Application limitations

- No rate limit/quota handling
- No caching for performance
- No integration with business logic
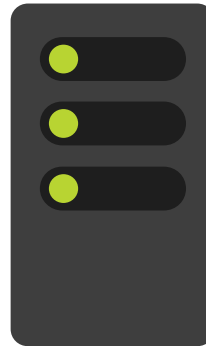- No handling of excess data to AI Model service ($$)

# Streaming – AI Inference Service

A server sits in between the user's device and the AI Model Service

· Clean separation of concerns with AI inference and backend logic
· **Added challenge**: streaming to the client/UI is now tricky

Client Browser

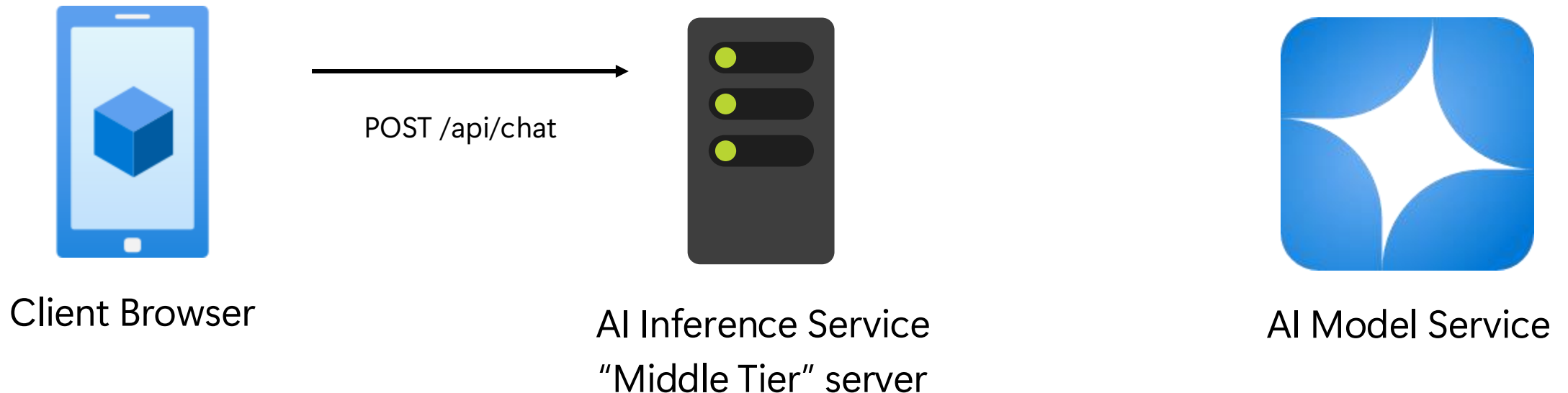AI Inference Service
"Middle Tier" server

AI Model Service

# Streaming – AI Inference Service

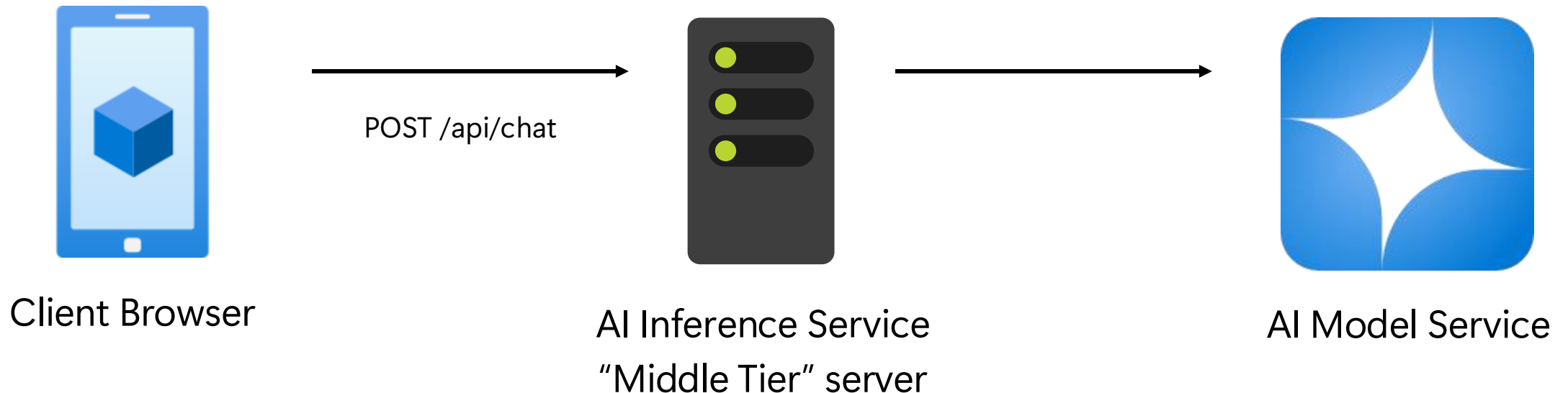A server sits in between the user's device and the AI Model Service

· Clean separation of concerns with AI inference and backend logic
· **Added challenge**: streaming to the client/UI is now tricky

POST /api/chat

Client Browser

AI Inference Service
"Middle Tier" server

AI Model Service

# Streaming – AI Inference Service

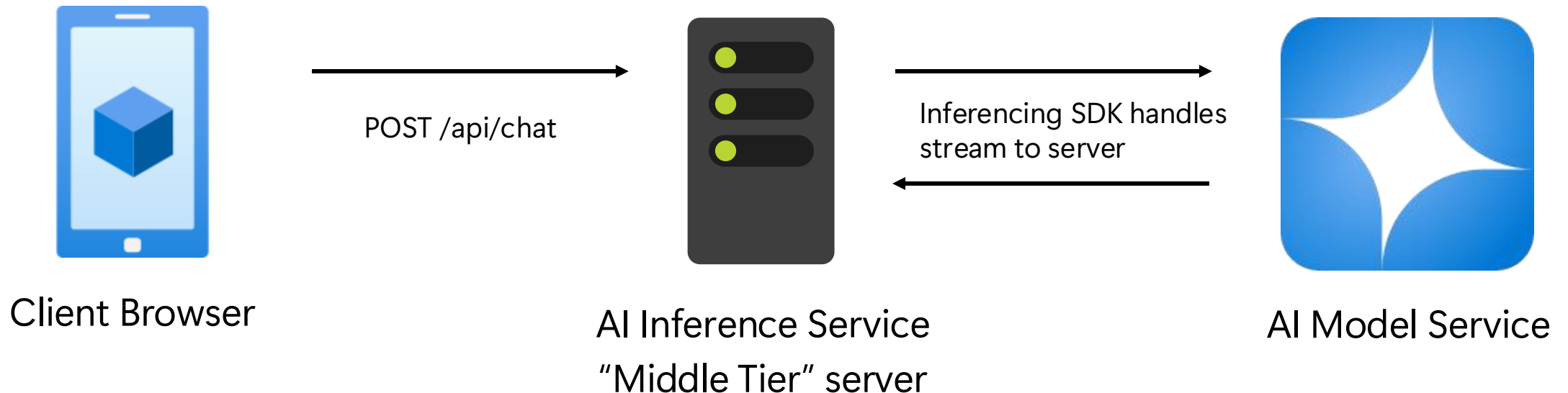A server sits in between the user's device and the AI Model Service

· Clean separation of concerns with AI inference and backend logic
· **Added challenge**: streaming to the client/UI is now tricky

POST /api/chat

Client Browser

AI Inference Service
"Middle Tier" server

AI Model Service

# Streaming – AI Inference Service

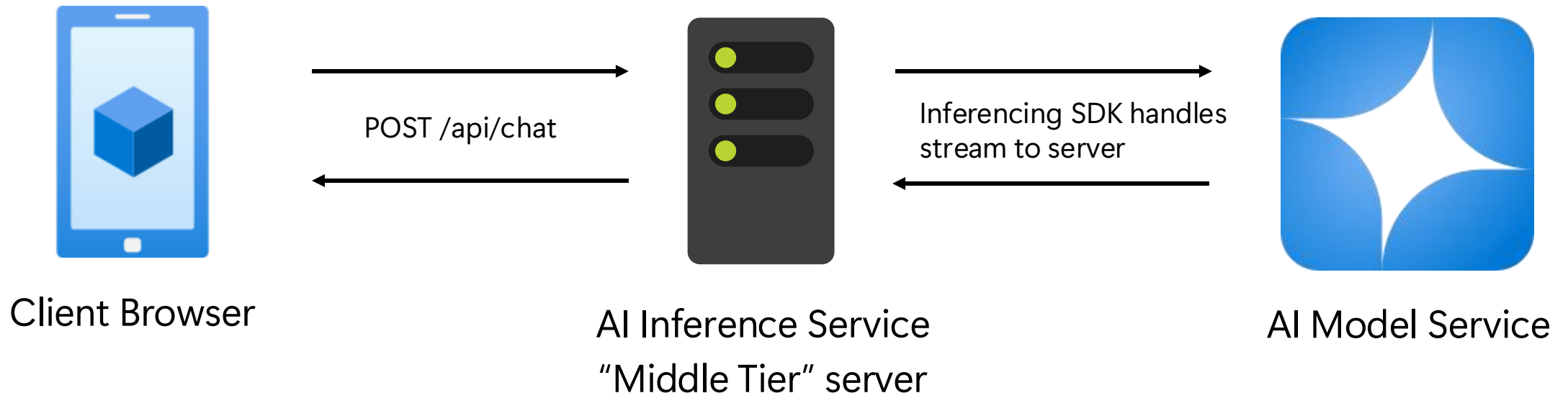A server sits in between the user's device and the AI Model Service

· Clean separation of concerns with AI inference and backend logic

· **Added challenge**: streaming to the client/UI is now tricky

POST /api/chat

Inferencing SDK handles stream to server

Client Browser

AI Inference Service
"Middle Tier" server

AI Model Service

# Streaming – AI Inference Service

A server sits in between the user's device and the AI Model Service
· Clean separation of concerns with AI inference and backend logic
· **Added challenge**: streaming to the client/UI is now tricky

POST /api/chat

Inferencing SDK handles stream to server

Client Browser

AI Inference Service
"Middle Tier" server

AI Model Service

# Streaming is hard!

# Streaming is hard!

1. Select a network protocol for streaming (HTTP, WebSockets, SSE...)
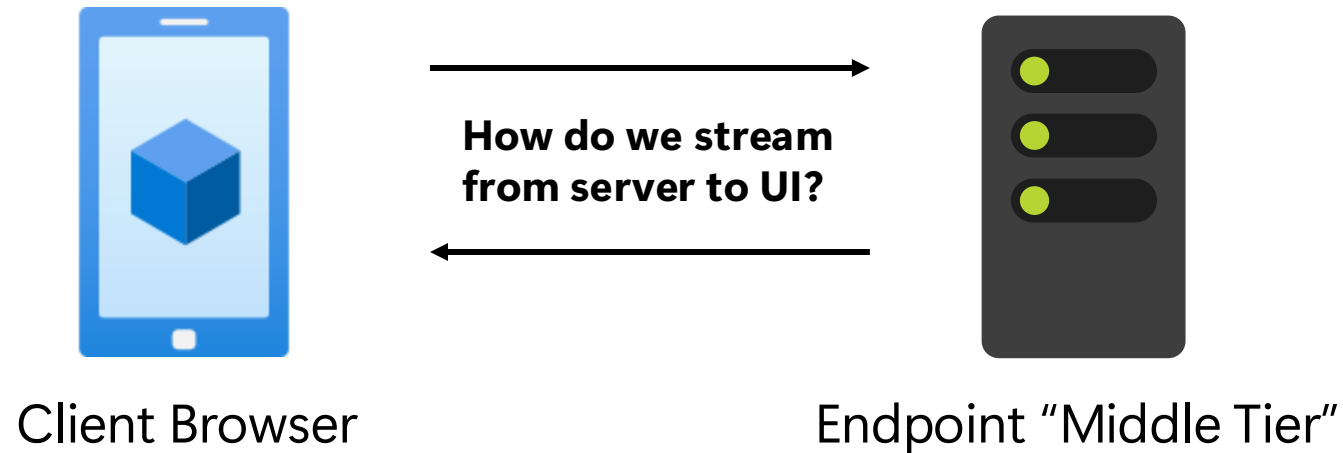
# Streaming is hard!

1. Select a network protocol for streaming (HTTP, WebSockets, SSE…)
2. Define an exchange format

# Streaming is hard!

1. Select a network protocol for streaming (HTTP, WebSockets, SSE…)
2. Define an exchange format
3. Parse the output (hundreds of lines of browser-side JS to write)

# Streaming is hard!

1. Select a network protocol for streaming (HTTP, WebSockets, SSE...)
2. Define an exchange format
3. Parse the output (hundreds of lines of browser-side JS to write)

**How do we stream from server to UI?**

Client Browser

Endpoint "Middle Tier"

# Introducing the AI Chat Protocol

# Streaming to UI with AI Chat Protocol

The **AI Chat Protocol** (JS) is designed to make streaming easy

# Streaming to UI with AI Chat Protocol

The **AI Chat Protocol** (JS) is designed to make streaming easy

- 2 Requirements:
  - Server uses API Specification
  - Frontend uses lightweight parsing SDK

# Streaming to UI with AI Chat Protocol

The **AI Chat Protocol** (JS) is designed to make streaming easy

- 2 Requirements:
  - Server uses API Specification
  - Frontend uses lightweight parsing SDK
- One client - `AIChatProtocolClient`

# Streaming to UI with AI Chat Protocol

The **AI Chat Protocol** (JS) is designed to make streaming easy

- 2 Requirements:
  - Server uses API Specification
  - Frontend uses lightweight parsing SDK
- One client - `AIChatProtocolClient`
- Two methods: `getCompletion` + `getStreamedCompletion`
  - Now you have easy streaming along with TypeScript models!
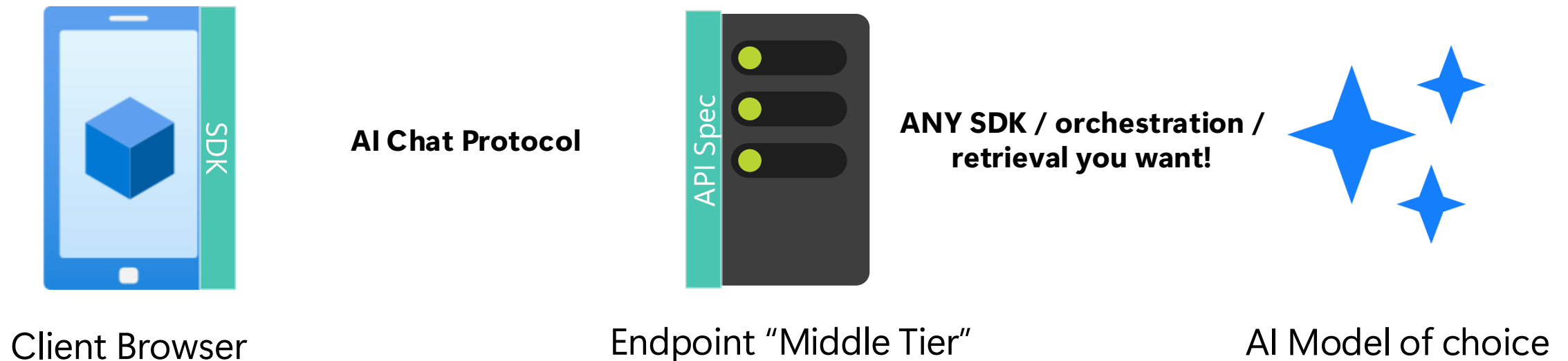  - Also includes a flexible 'context' options bag for any other data needed

> **For our sample, I would roughly estimate it to remove something between 200-400 lines of code, maybe a bit more because of the parser"**

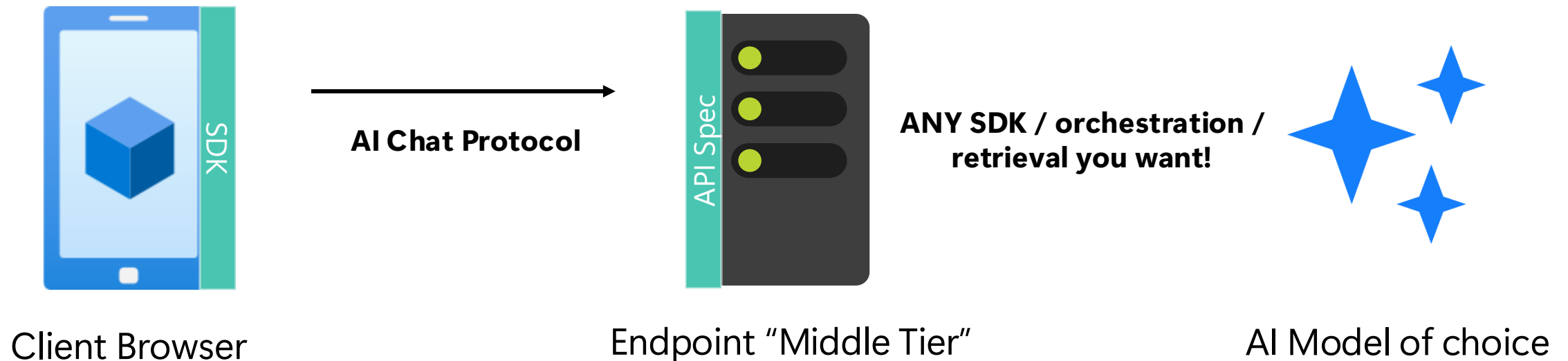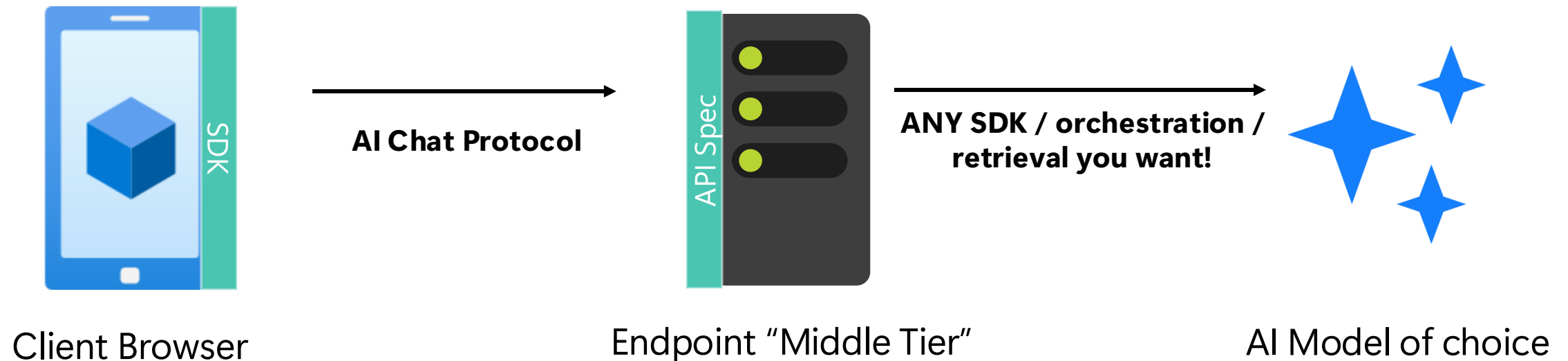Yohan Lasorsa, Cloud Advocate, Microsoft

# Building with the AI Chat Protocol

- **Server-to-client only** design pattern, not server-to-AI model
- AI backend choices are flexible
    - Any orchestration, RAG option, programming language, and language model



**AI Chat Protocol**

**ANY SDK / orchestration / retrieval you want!**

Client Browser

Endpoint "Middle Tier"
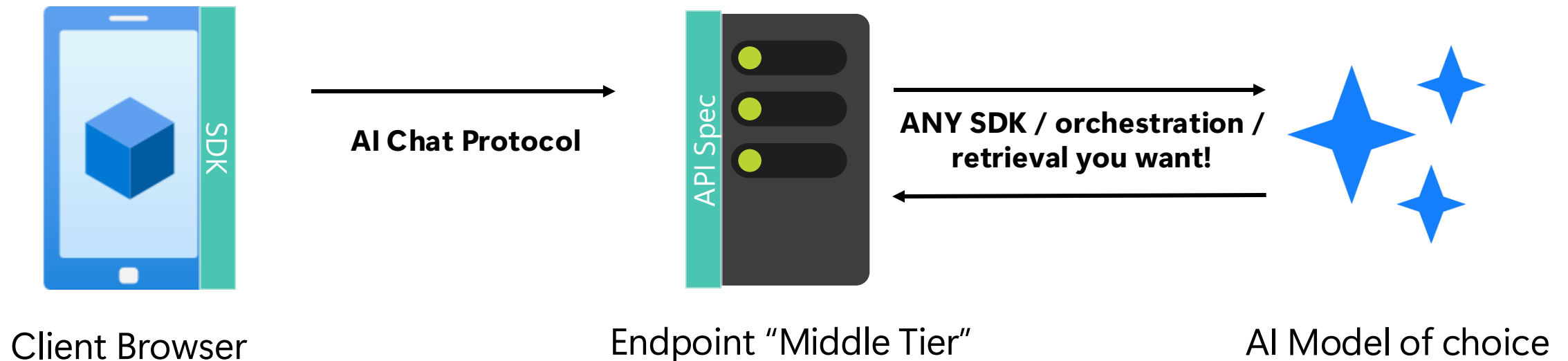
AI Model of choice

# Building with the AI Chat Protocol

- **Server-to-client only** design pattern, not server-to-AI model
- AI backend choices are flexible
  - Any orchestration, RAG option, programming language, and language model



Client Browser        Endpoint "Middle Tier"        AI Model of choice

AI Chat Protocol

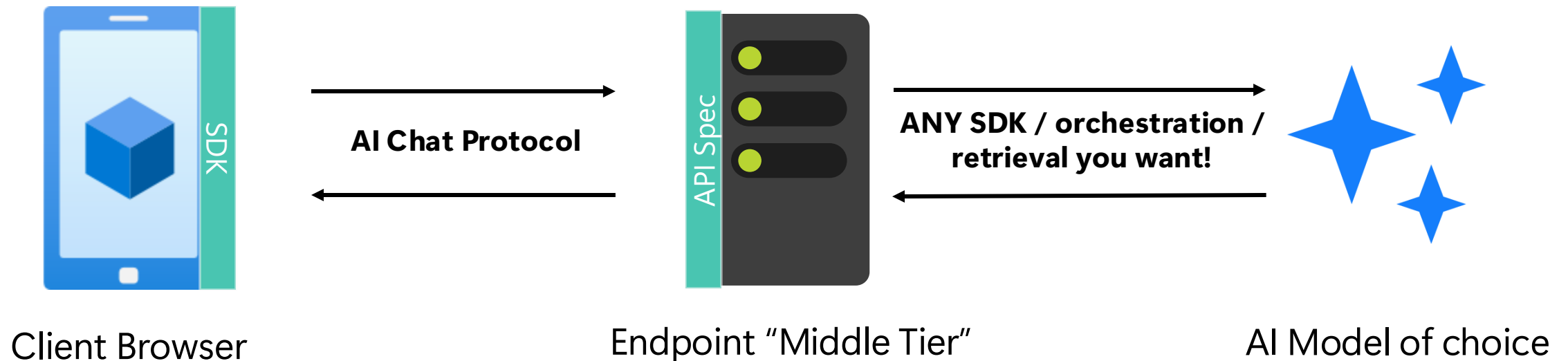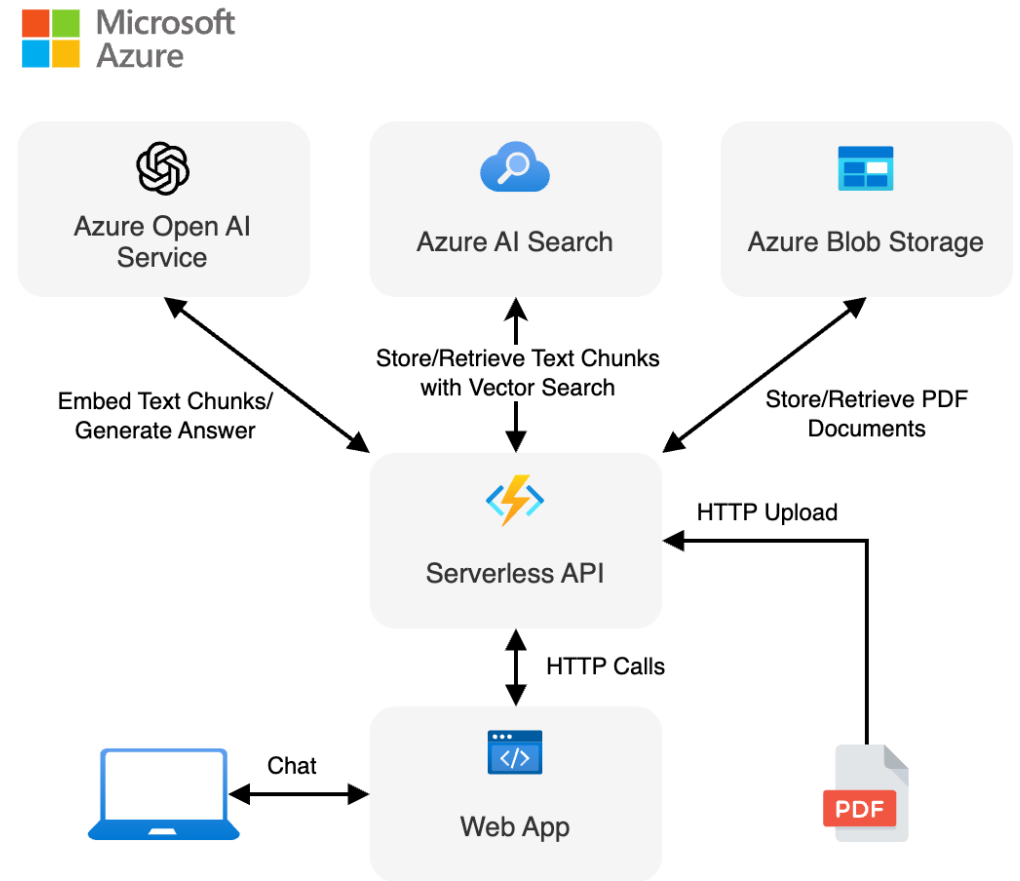ANY SDK / orchestration / retrieval you want!

# Building with the AI Chat Protocol

- **Server-to-client only** design pattern, not server-to-AI model
- AI backend choices are flexible
  - Any orchestration, RAG option, programming language, and language model



Client Browser        Endpoint "Middle Tier"        AI Model of choice

**AI Chat Protocol**

**ANY SDK / orchestration / retrieval you want!**

# Building with the AI Chat Protocol

- **Server-to-client only** design pattern, not server-to-AI model
- AI backend choices are flexible
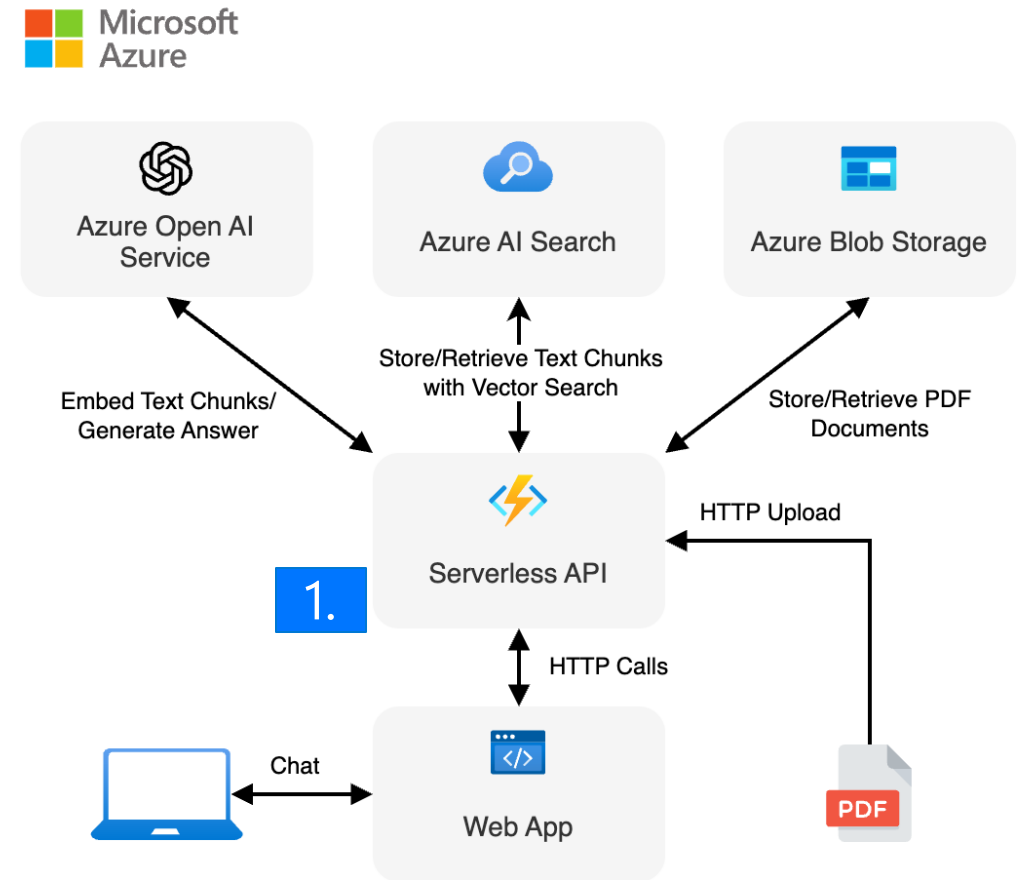  - Any orchestration, RAG option, programming language, and language model



**AI Chat Protocol**

**ANY SDK / orchestration / retrieval you want!**

Client Browser                    Endpoint "Middle Tier"                    AI Model of choice

# Building with the AI Chat Protocol

- **Server-to-client only** design pattern, not server-to-AI model
- AI backend choices are flexible
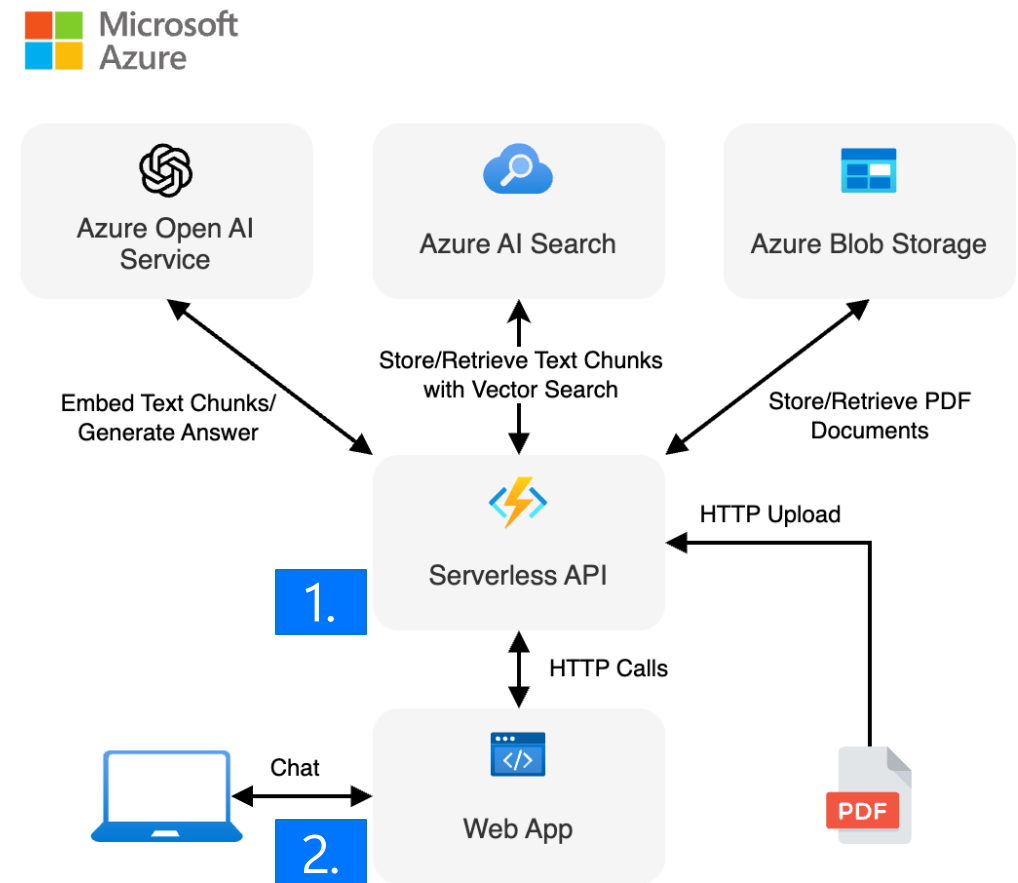  - Any orchestration, RAG option, programming language, and language model



AI Chat Protocol

ANY SDK / orchestration / retrieval you want!

SDK

API Spec

Client Browser

Endpoint "Middle Tier"

AI Model of choice

# Adding the AI Chat Protocol

# Adding the AI Chat Protocol

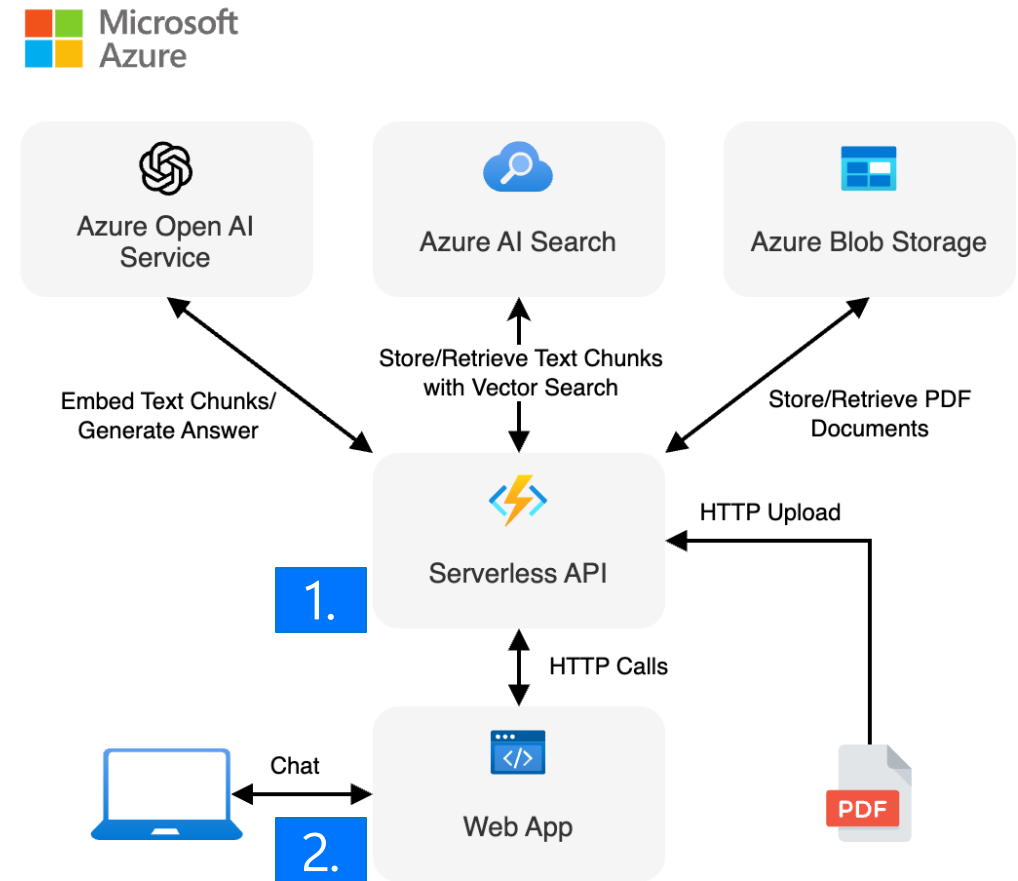1. Ensure service API conforms to Chat Protocol

# Adding the AI Chat Protocol

1. Ensure service API conforms to Chat Protocol

2. Incorporate SDK into frontend code

# Adding the AI Chat Protocol

1. Ensure service API conforms to Chat Protocol

2. Incorporate SDK into frontend code

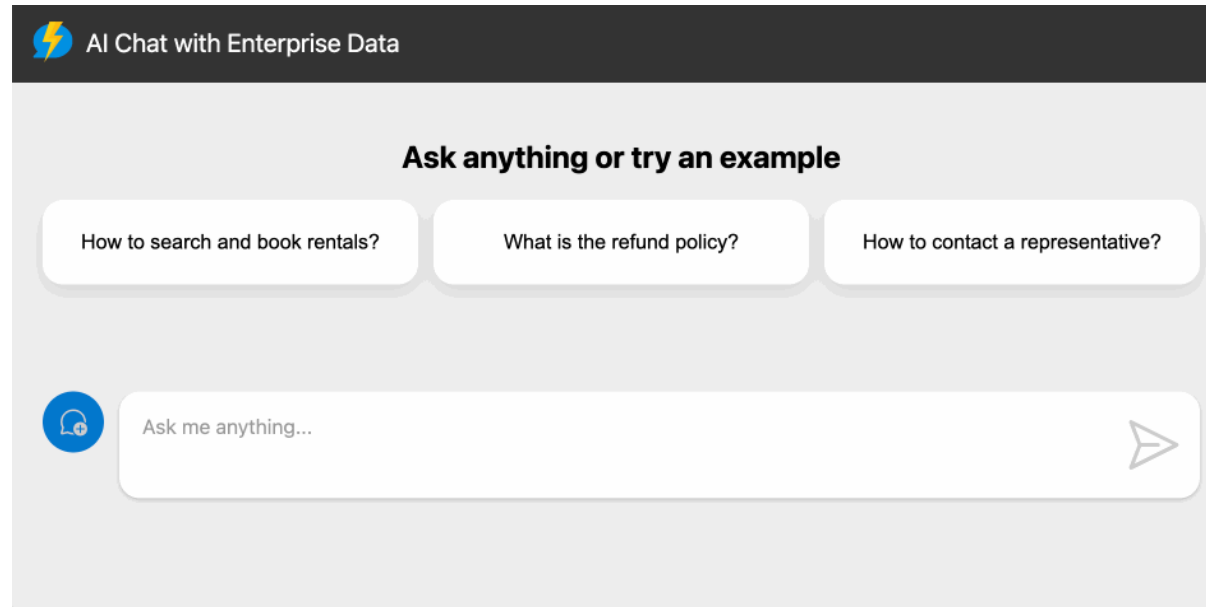3. Enjoy streamed output to the client 🙂

# Example integration

Sample Code

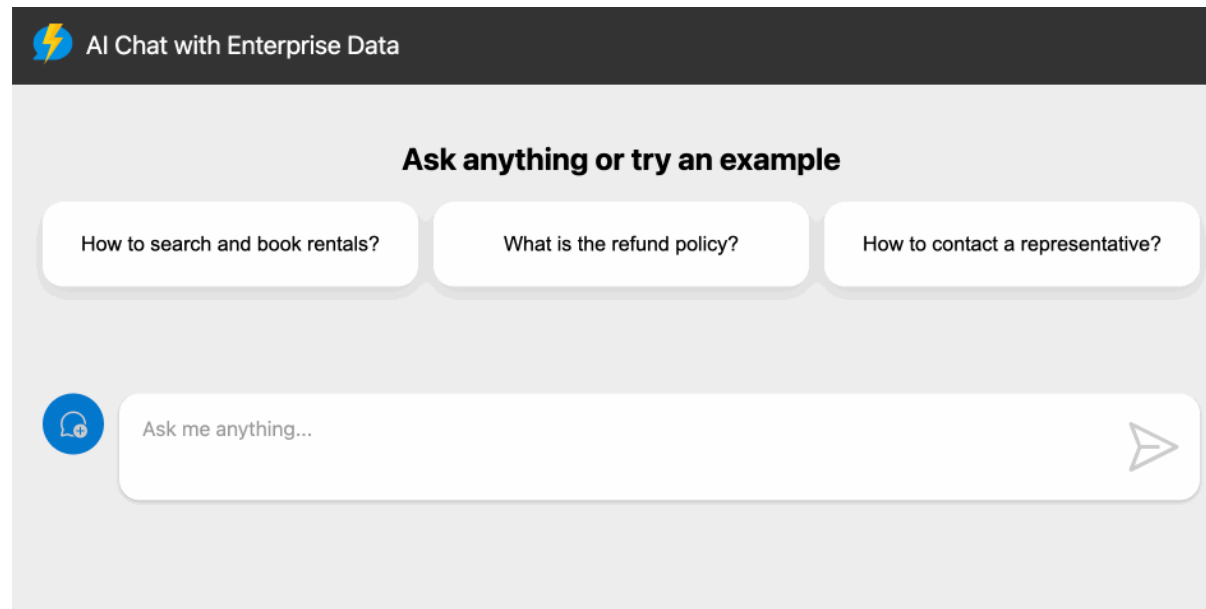# Serverless RAG with Langchain.js Sample

Sample uses the AI Chat Protocol for streaming to UI

# Serverless RAG with Langchain.js Sample

Sample uses the AI Chat Protocol for streaming to UI

- 2 core pieces:
  - RAG Service API – Conforms to AI Chat Protocol API Spec
  - Web App – Uses AI Chat Protocol library

# RAG Service API – ndjson Stream

```typescript
104   // Transform the response chunks into a JSON stream
105   async function* createJsonStream(chunks: AsyncIterable<string>) {
106     for await (const chunk of chunks) {
107       if (!chunk) continue;
108
109       const responseChunk: AIChatCompletionDelta = {
110         delta: {
111           content: chunk,
112           role: 'assistant',
113         },
114       };
115
116       // Format response chunks in Newline delimited JSON
117       // see https://github.com/ndjson/ndjson-spec
118       yield JSON.stringify(responseChunk) + '\n';
119     }
120   }
```

# RAG Service API – AI Chat Protocol Spec

```
86     const responseStream = await ragChain.stream({
87       input: question,
88       context: await retriever.invoke(question),
89     });
90     const jsonStream = Readable.from(createJsonStream(responseStream));
91
92     return data(jsonStream, {
93       'Content-Type': 'application/x-ndjson',
94       'Transfer-Encoding': 'chunked',
95     });
96   } catch (_error: unknown) {
97     const error = _error as Error;
98     context.error(`Error when processing chat-post request: ${error.message}`);
99
100    return serviceUnavailable('Service temporarily unavailable. Please try again later.');
101  }
102 }
```

# Web App – AI Chat Protocol Library

```
1    import { AIChatMessage, AIChatCompletionDelta, AIChatProtocolClient } from '@microsoft/ai-chat-protocol';
2
3    export const apiBaseUrl: string = import.meta.env.VITE_API_URL || '';
4
5    export type ChatRequestOptions = {
6      messages: AIChatMessage[];
7      chunkIntervalMs: number;
8      apiUrl: string;
9    };
10
```

# Web App – AI Chat Protocol Library

```typescript
11    export async function* getCompletion(options: ChatRequestOptions) {
12        const apiUrl = options.apiUrl || apiBaseUrl;
13        const client = new AIChatProtocolClient(`${apiUrl}/api/chat`);
14        const result = await client.getStreamedCompletion(options.messages);
15
16        for await (const response of result) {
17          if (!response.delta) {
18            continue;
19          }
20
21          yield new Promise<AIChatCompletionDelta>((resolve) => {
22            setTimeout(() => {
23              resolve(response);
24            }, options.chunkIntervalMs);
25          });
26        }
27    }
```

# Easy streaming with context!

Microsoft

# Resources

- AI Chat Protocol GitHub     aka.ms/aichat

- AI Chat Protocol API Spec     aka.ms/chatprotocol

- Serverless AI Chat sample     aka.ms/ai/js/chat