

Bi-Directional Audio Streaming Integration Document

This document outlines the Web Socket (WS) event structure for bi-directional audio streaming between our system and an endpoint. It enables to send us the voice data along with the information of the caller to the endpoint (webSocket) which can return the voice data back and it would be played out to the caller. The integration involves sending and receiving audio data and metadata in real-time using predefined event types.

1. Overview

The integration stream audio data and metadata between our system and the webSocket endpoint. The following events are exchanged to manage the audio stream lifecycle:

- **Events Sent to the Vendor (endpoint):** Used to initiate, manage, and terminate the audio stream.
 - **Events Received from the Vendor (endpoint):** Used to receive audio data and metadata from the vendor.
-

2. Events Sent to the Vendor

2.1. Connected

The `connected` event acts as a handshake response and sets expectations between the client and server. It is the first message sent after establishing the Web Socket connection.

We can send the `connected` event once a Web Socket connection is established. This is the first message the Web Socket server receives.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case, <code>connected</code> .

Payload Structure:

```
1 {  
2   "event": "connected"  
3 }
```

2.2. Start Message

The `start` message contains metadata about the Stream and is sent immediately after the `connected` message. It is only sent once at the start of the Stream.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case, <code>start</code> .
<code>sequenceNumber</code>	Number used to keep track of message sending order. The first message has a value

	of <code>1</code> and then is incremented for each subsequent message.
<code>start</code>	An object containing Stream metadata
<code>start.streamSid</code>	The unique identifier of the Stream
<code>start.accountSid</code>	The unique identifier of the Account for which the stream was created
<code>start.callSid</code>	The unique identifier of the Call for which the Stream was started
<code>from</code>	The number from which the call was originated to the above-mentioned account.
<code>to</code>	The number of the account to which the call originated to.
<code>start.mediaFormat</code>	An object containing the format of the payload in the <code>media</code> messages.
<code>start.mediaFormat.encoding</code>	The encoding of the data in the upcoming payload. Value is always <code>audio/x-mulaw</code> . (also known as G.711 μ -law (PCMU))
<code>start.mediaFormat.sampleRate</code>	The sample rate in hertz of the upcoming audio data. Value is always <code>8000</code>
<code>start.mediaFormat.bitRate</code>	The number of bits used to represent one second of audio in the input audio data. Value is always 64 kbps .
<code>start.mediaFormat.bitDepth</code>	It refers to the number of bits used to represent each sample. (<code>8-bit</code>)
<code>start.customParameters</code>	An object containing the custom parameters that were set when defining the Stream
<code>streamSid</code>	The unique identifier of the Stream

Payload Structure:

```

1  {
2    "event": "start",
3    "sequenceNumber": "1",
4    "start": {
5      "accountSid": "ACXXXXXXXXXXXXXXXXXXXXXX",
6      "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX",
7      "callSid": "CAXXXXXXXXXXXXXXXXXXXXXXX",
8      "from": "XXXXXXXXXX",
9      "to": "XXXXXXXXXX",
10     "mediaFormat": {
11       "encoding": "audio/x-mulaw",
12       "sampleRate": 8000,
13       "bitRate": 64,
14       "bitDepth": 8 },
15     "customParameters": {
16       "FirstName": "Jane",

```

```

17     "LastName": "Doe",
18     "RemoteParty": "Bob",
19   },
20 },
21   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX"
22 }

```

2.3. Media message

This message type encapsulates the raw audio data.

Property	Description
event	Describes the type of Web Socket message. In this case, "media".
sequenceNumber	Number used to keep track of message sending order. The first message has a value of 1 and then is incremented for each subsequent message.
media	An object containing media metadata and payload
media.chunk	The chunk for the message. The first message will begin with 1 and increment with each subsequent message.
media.timestamp	Presentation Timestamp in Milliseconds from the start of the stream.
media.payload	Raw audio encoded packets in base64
streamSid	The unique identifier of the Stream

Payload Structure:

```

1 {
2   "event": "media",
3   "sequenceNumber": "3",
4   "media": {
5     "chunk": "1",
6     "timestamp": "5",
7     "payload": "no+JhoaJjpz..."
8   },
9   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX"
10 }

```

2.4. Stop Message

This message indicates when the Stream has stopped, or the call has ended.

Property	Description
----------	-------------

<code>event</code>	Describes the type of Web Socket message. In this case, <code>stop</code> .
<code>sequenceNumber</code>	Number used to keep track of message sending order. The first message has a value of <code>1</code> and then is incremented for each subsequent message.
<code>stop</code>	An object containing Stream metadata
<code>stop.accountSid</code>	The Account identifier that created the Stream
<code>stop.callSid</code>	The Call identifier that started the Stream
<code>stop.reason</code>	The reason for ending the Stream.
<code>streamSid</code>	The unique identifier of the Stream

Payload Structure: ⚡

```

1 {
2   "event": "stop",
3   "sequenceNumber": "5",
4   "stop": {
5     "accountSid": "ACXXXXXXXXXXXXXXXXXXXXXX",
6     "callSid": "CAXXXXXXXXXXXXXXXXXXXXXXX",
7     "reason": "The caller disconnected the call"
8   },
9   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX"
10 }
```

2.5 DTMF message ⚡

A `dtmf` message is sent when someone presses a touch-tone number key in the inbound stream, typically in response to a prompt in the outbound stream.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case, <code>dtmf</code> .
<code>streamSid</code>	The unique identifier of the Stream
<code>sequenceNumber</code>	Number used to keep track of message sending order. The first message has a value of <code>1</code> and then is incremented for each subsequent message.
<code>dtmf.digit</code>	the number-key tone detected

An example `dtmf` message is shown below. The `dtmf.digit` value is `1`, indicating that someone pressed the `1` key on their handset.

Payload:

```
1 {
```

```

2   "event": "dtmf",
3   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXXXXXXXXX",
4   "sequenceNumber": "5",
5   "dtmf": {
6     "digit": "1"
7   }
8 }
```

2.6. Mark message ↗

When endpoint sends a media message, it could then send a mark message with a label; When that `media` message's playback is complete, we send the `mark` message to the endpoint using the same label `mark.name` indicating that the media has been played.

If the endpoint (WebSocket server) sends a clear message, we will empty the audio buffer and send back the `mark` messages matching any remaining `mark` messages from the server.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case, <code>"mark"</code> .
<code>sequenceNumber</code>	Number used to keep track of message sending order. The first message has a value of <code>1</code> and then is incremented for each subsequent message.
<code>streamSid</code>	The unique identifier of the Stream
<code>mark</code>	An object containing the mark metadata
<code>mark.name</code>	A custom value. We send back the <code>mark.name</code> you specify when it receives a <code>mark</code> message

Payload Structure: ↗

```

1  {
2   "event": "mark",
3   "sequenceNumber": "4",
4   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXXXXXXXXX",
5   "mark": {
6     "name": "mark label"
7   }
8 }
```

3. Events Received from the Vendor ↗

3.1. Media ↗

The payload must be encoded `audio/x-mulaw` with a sample rate of `8000` encoded with base64 PCM mono audio. The audio can be of any size.

Property	Description
----------	-------------

<code>event</code>	Describes the type of Web Socket message. In this case, "media".
<code>streamSid</code>	The SID of the Stream that should play the audio
<code>media</code>	An object containing the media payload
<code>media.payload</code>	Raw <code>mulaw/8000</code> audio in encoded in base64
<code>media.chunk</code>	The chunk for the message. The first message will begin with <code>1</code> and increment with each subsequent message.

Payload Structure: ↴

```

1 {
2   "event": "media",
3   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX",
4   "media": {
5     "payload": "a3242sa...",
6     "chunk" : 1
7   }
8 }
```

3.2. Mark ↴

Sends a `mark` event message after sending a `media` event message to be notified when the audio that they have sent has been completed. We send back a `mark` event with a matching `name` when the audio ends (or if there is no audio buffered).

The Web Socket Server also receives an incoming `mark` event message if the buffer was cleared using the `clear` event message.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case "mark".
<code>streamSid</code>	The SID of the Stream that should receive the mark
<code>mark</code>	An object containing mark metadata and payload
<code>mark.name</code>	A name specific to your needs that will assist in recognizing future received <code>mark</code> event

Payload Structure: ↴

```

1 {
2   "event": "mark",
3   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX",
4   "mark": {
5     "name": "my label"
```

```
6 }
7 }
```

3.3. Clear

Sends a `clear` message if the server want to interrupt the audio that has been sent in various `media` messages. This empties all buffered audio and causes any `mark` messages to be sent back to the Web Socket server.

Property	Description
<code>event</code>	Describes the type of Web Socket message. In this case, <code>"clear"</code> .
<code>streamSid</code>	The SID of the Stream in which you wish to interrupt the audio.

Payload Structure: 

```
1 {
2   "event": "clear",
3   "streamSid": "MZXXXXXXXXXXXXXXXXXXXXXX",
4 }
```

4. Stream Lifecycle

1. Connection Establishment:

- Establish a Web Socket connection with the vendor.
- Send the `connected` event to initiate the handshake.

2. Stream Initialization:

- Send the `start` event with stream metadata.

3. Audio Streaming:

- Send `media` events with base64 encoded audio data.
- Receive `media` events from the vendor with base64 encoded audio data.

4. Stream Termination:

- Send the `stop` event to terminate the stream.
- Handle the `clear` event from the vendor to reset the stream.

5. End of Input:

- Send the `mark` event when all media from the bot has been played.
- Handle the `mark` event from the vendor to denote the end of input.

5. Example Workflow

1. Client to Vendor:

- Send `connected` → `start` → `media` → `stop`.

2. Vendor to Client:

- Receive `media` → `mark` → `clear`.

6. Notes

- Ensures that the `streamSid` is unique for each stream and consistent across all events for a given stream.
 - Base64 encoding is used for audio data to ensure compatibility and ease of transmission.
 - The `mark` event is used to synchronize the end of input between the client and vendor.
-

Version	Date	Comment
Current Version (v. 8)	Feb 25, 2025 16:56	Kartikay Saini
v. 7	Feb 24, 2025 16:58	Kartikay Saini
v. 6	Feb 21, 2025 14:09	Kartikay Saini
v. 5	Jan 28, 2025 15:45	Ansh Tyagi
v. 4	Jan 28, 2025 15:43	Ansh Tyagi
v. 3	Jan 28, 2025 15:43	Ansh Tyagi
v. 2	Jan 28, 2025 14:24	Ansh Tyagi
v. 1	Jan 28, 2025 13:36	Ansh Tyagi