

Veo 3 is now available in the Gemini API!

Learn more (<https://developers.googleblog.com/en/veo-3-now-available-gemini-api/>)

Session management with Live API

In the Live API, a session refers to a persistent connection where input and output are streamed continuously over the same connection (read more about [how it works](/gemini-api/docs/live) (/gemini-api/docs/live)). This unique session design enables low latency and supports unique features, but can also introduce challenges, like session time limits, and early termination. This guide covers strategies for overcoming the session management challenges that can arise when using the Live API.

Session lifetime

Without compression, audio-only sessions are limited to 15 minutes, and audio-video sessions are limited to 2 minutes. Exceeding these limits will terminate the session (and therefore, the connection), but you can use [context window compression](#) (#context-window-compression) to extend sessions to an unlimited amount of time.

The lifetime of a connection is limited as well, to around 10 minutes. When the connection terminates, the session terminates as well. In this case, you can configure a single session to stay active over multiple connections using [session resumption](#) (#session-resumption). You'll also receive a [GoAway message](#) (#goaway-message) before the connection ends, allowing you to take further actions.

Context window compression

To enable longer sessions, and avoid abrupt connection termination, you can enable context window compression by setting the [contextWindowCompression](#) (/api/live#BidiGenerateContentSetup.FIELDS.ContextWindowCompressionConfig.BidiGenerateContentSetup.context_window_compression) field as part of the session configuration.

In the [ContextWindowCompressionConfig](/api/live#contextwindowcompressionconfig) (/api/live#contextwindowcompressionconfig), you can configure a [sliding-window mechanism](/api/live#ContextWindowCompressionConfig.FIELDS.ContextWindowCompressionConfig.SlidingWindow.ContextWindowCompressionConfig.sliding_window)

(/api/live#ContextWindowCompressionConfig.FIELDS.ContextWindowCompressionConfig.SlidingWindow.ContextWindowCompressionConfig.sliding_window)

and the [number of tokens](/api/live#ContextWindowCompressionConfig.FIELDS.int64.ContextWindowCompressionConfig.trigger_tokens)

(/api/live#ContextWindowCompressionConfig.FIELDS.int64.ContextWindowCompressionConfig.trigger_tokens)

that triggers compression.

PythonJavaScript (#javascript)
(#python)

```
from google.genai import types

config = types.LiveConnectConfig(
    response_modalities=["AUDIO"],
    context_window_compression=(
        # Configures compression with default parameters.
        types.ContextWindowCompressionConfig(
            sliding_window=types.SlidingWindow(),
        )
    ),
)
```

Session resumption

To prevent session termination when the server periodically resets the WebSocket connection, configure the [sessionResumption](/api/live#BidiGenerateContentSetup.FIELDS.SessionResumptionConfig.BidiGenerateContentSetup.session_resumption)

(/api/live#BidiGenerateContentSetup.FIELDS.SessionResumptionConfig.BidiGenerateContentSetup.session_resumption)

field within the [setup configuration](/api/live#BidiGenerateContentSetup) (/api/live#BidiGenerateContentSetup).

Passing this configuration causes the server to send [SessionResumptionUpdate](/api/live#SessionResumptionUpdate)

(/api/live#SessionResumptionUpdate) messages, which can be used to resume the session by

passing the last resumption token as the [SessionResumptionConfig.handle](/api/live#SessionResumptionConfig.FIELDS.string.SessionResumptionConfig.handle)

(/api/live#SessionResumptionConfig.FIELDS.string.SessionResumptionConfig.handle) of the subsequent connection.

Python LiveScript (#livescript)

```
import asyncio
from google import genai
from google.genai import types

client = genai.Client()
model = "gemini-live-2.5-flash-preview"

async def main():
    print(f"Connecting to the service with handle {previous_session_handle}")
    async with client.aio.live.connect(
        model=model,
        config=types.LiveConnectConfig(
            response_modalities=["AUDIO"],
            session_resumption=types.SessionResumptionConfig(
                # The handle of the session to resume is passed here,
                # or else None to start a new session.
                handle=previous_session_handle
            ),
        ),
    ) as session:
        while True:
            await session.send_client_content(
                turns=types.Content(
                    role="user", parts=[types.Part(text="Hello world!")]
                )
            )
            async for message in session.receive():
                # Periodically, the server will send update messages that may
                # contain a handle for the current state of the session.
                if message.session_resumption_update:
                    update = message.session_resumption_update
                    if update.resumable and update.new_handle:
                        # The handle should be retained and linked to the s
                        return update.new_handle

                # For the purposes of this example, placeholder input is co
                # to the model. In non-sample code, the model inputs would
                # the user.
                if message.server_content and message.server_content.turn_c
                    break

if __name__ == "__main__":
```

```
asyncio.run(main())
```

Receiving a message before the session disconnects

The server sends a GoAway (/api/live#GoAway) message that signals that the current connection will soon be terminated. This message includes the timeLeft (/api/live#GoAway.FIELDS.google.protobuf.Duration.GoAway.time_left), indicating the remaining time and lets you take further action before the connection will be terminated as ABORTED.

PythonJavaScript (#javascript)
(#python)

```
async for response in session.receive():
    if response.go_away is not None:
        # The connection will soon be terminated
        print(response.go_away.time_left)
```

Receiving a message when the generation is complete

The server sends a generationComplete (/api/live#BidiGenerateContentServerContent.FIELDS.bool.BidiGenerateContentServerContent.generation_complete) message that signals that the model finished generating the response.

PythonJavaScript (#javascript)
(#python)

```
async for response in session.receive():
    if response.server_content.generation_complete is True:
        # The generation is complete
```

What's next

Explore more ways to work with the Live API in the full [Capabilities](/gemini-api/docs/live) (/gemini-api/docs/live) guide, the [Tool use](/gemini-api/docs/live-tools) (/gemini-api/docs/live-tools) page, or the [Live API cookbook](https://colab.research.google.com/github/google-gemini/cookbook/blob/main/quickstarts/Get_started_LiveAPI.ipynb) (https://colab.research.google.com/github/google-gemini/cookbook/blob/main/quickstarts/Get_started_LiveAPI.ipynb)

.

Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/) (https://creativecommons.org/licenses/by/4.0/), and code samples are licensed under the [Apache 2.0 License](https://www.apache.org/licenses/LICENSE-2.0) (https://www.apache.org/licenses/LICENSE-2.0). For details, see the [Google Developers Site Policies](https://developers.google.com/site-policies) (https://developers.google.com/site-policies). Java is a registered trademark of Oracle and/or its affiliates.

Last updated 2025-07-08 UTC.