

Veo 3 is now available in the Gemini API!

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

Tool use with Live API

Tool use allows Live API to go beyond just conversation by enabling it to perform actions in the real-world and pull in external context while maintaining a real time connection. You can define tools such as [Function calling](/gemini-api/docs/function-calling) (/gemini-api/docs/function-calling), [Code execution](/gemini-api/docs/code-execution) (/gemini-api/docs/code-execution), and [Google Search](/gemini-api/docs/grounding) (/gemini-api/docs/grounding) with the Live API.

Overview of supported tools

Here's a brief overview of the available tools for each model:

Tool	Cascaded models gemini-live-2.5- flash-preview gemini-2.0-flash- live-001	gemini-2.5-flash-preview- native-audio-dialog	gemini-2.5-flash-exp-native- audio-thinking-dialog
Search	Yes	Yes	Yes
Function calling	Yes	Yes	No
Code execution	Yes	No	No
Url context	Yes	No	No

Function calling

Live API supports function calling, just like regular content generation requests. Function calling lets the Live API interact with external data and programs, greatly increasing what your applications can accomplish.

You can define function declarations as part of the session configuration. After receiving tool calls, the client should respond with a list of `FunctionResponse` objects using the `session.send_tool_response` method.

See the [Function calling tutorial](/gemini-api/docs/function-calling) (/gemini-api/docs/function-calling) to learn more.

Note: Unlike the `generateContent` API, the Live API doesn't support automatic tool response handling. You must handle tool responses manually in your client code.

PythonJavaScript (#javascript)
(#python)

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

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

# Simple function definitions
turn_on_the_lights = {"name": "turn_on_the_lights"}
turn_off_the_lights = {"name": "turn_off_the_lights"}

tools = [{"function_declarations": [turn_on_the_lights, turn_off_the_lights]}]
config = {"response_modalities": ["TEXT"], "tools": tools}

async def main():
    async with client.aio.live.connect(model=model, config=config) as session:
        prompt = "Turn on the lights please"
        await session.send_client_content(turns={"parts": [{"text": prompt}]}

        async for chunk in session.receive():
            if chunk.server_content:
                if chunk.text is not None:
                    print(chunk.text)
```

```

elif chunk.tool_call:
    function_responses = []
    for fc in chunk.tool_call.function_calls:
        function_response = types.FunctionResponse(
            id=fc.id,
            name=fc.name,
            response={ "result": "ok" } # simple, hard-coded fu
        )
        function_responses.append(function_response)

    await session.send_tool_response(function_responses=function_responses)

if __name__ == "__main__":
    asyncio.run(main())

```

From a single prompt, the model can generate multiple function calls and the code necessary to chain their outputs. This code executes in a sandbox environment, generating subsequent [BidiGenerateContentToolCall](/api/live#bidigeneratecontenttoolcall) (/api/live#bidigeneratecontenttoolcall) messages.

Asynchronous function calling

Note: Asynchronous function calling is only supported in [half-cascade](/gemini-api/docs/live#half-cascade) (/gemini-api/docs/live#audio-generation) audio generation.

Function calling executes sequentially by default, meaning execution pauses until the results of each function call are available. This ensures sequential processing, which means you won't be able to continue interacting with the model while the functions are being run.

If you don't want to block the conversation, you can tell the model to run the functions asynchronously. To do so, you first need to add a **behavior** to the function definitions:

PythonJavaScript (#javascript)
(#python)

```

# Non-blocking function definitions
turn_on_the_lights = {"name": "turn_on_the_lights", "behavior": "NON_BLOCKING"}

```

```
turn_off_the_lights = {"name": "turn_off_the_lights"} # turn_off_the_ligh
```

NON-BLOCKING ensures the function runs asynchronously while you can continue interacting with the model.

Then you need to tell the model how to behave when it receives the `FunctionResponse` using the `scheduling` parameter. It can either:

- Interrupt what it's doing and tell you about the response it got right away (`scheduling="INTERRUPT"`),
- Wait until it's finished with what it's currently doing (`scheduling="WHEN_IDLE"`),
- Or do nothing and use that knowledge later on in the discussion (`scheduling="SILENT"`)

```
PythonJavaScript (#javascript)
(#python)
```

```
# for a non-blocking function definition, apply scheduling in the function
function_response = types.FunctionResponse(
    id=fc.id,
    name=fc.name,
    response={
        "result": "ok",
        "scheduling": "INTERRUPT" # Can also be WHEN_IDLE or SILENT
    }
)
```

Code execution

You can define code execution as part of the session configuration. This lets the Live API generate and execute Python code and dynamically perform computations to benefit your results. See the [Code execution tutorial](/gemini-api/docs/code-execution) (/gemini-api/docs/code-execution) to learn more.

```
PythonJavaScript (#javascript)
(#python)
```

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

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

tools = [{'code_execution': {}}]
config = {"response_modalities": ["TEXT"], "tools": tools}

async def main():
    async with client.aio.live.connect(model=model, config=config) as session:
        prompt = "Compute the largest prime palindrome under 100000."
        await session.send_client_content(turns={"parts": [{"text": prompt}]}

        async for chunk in session.receive():
            if chunk.server_content:
                if chunk.text is not None:
                    print(chunk.text)

                model_turn = chunk.server_content.model_turn
                if model_turn:
                    for part in model_turn.parts:
                        if part.executable_code is not None:
                            print(part.executable_code.code)

                        if part.code_execution_result is not None:
                            print(part.code_execution_result.output)

if __name__ == "__main__":
    asyncio.run(main())
```

Grounding with Google Search

You can enable Grounding with Google Search as part of the session configuration. This increases the Live API's accuracy and prevents hallucinations. See the [Grounding tutorial \(/gemini-api/docs/grounding\)](https://ai.google.dev/gemini-api/docs/grounding) to learn more.

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

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

tools = [{'google_search': {}}]
config = {"response_modalities": ["TEXT"], "tools": tools}

async def main():
    async with client.aio.live.connect(model=model, config=config) as session:
        prompt = "When did the last Brazil vs. Argentina soccer match happen?"
        await session.send_client_content(turns={"parts": [{"text": prompt}]}

        async for chunk in session.receive():
            if chunk.server_content:
                if chunk.text is not None:
                    print(chunk.text)

                # The model might generate and execute Python code to use S
                model_turn = chunk.server_content.model_turn
                if model_turn:
                    for part in model_turn.parts:
                        if part.executable_code is not None:
                            print(part.executable_code.code)

                        if part.code_execution_result is not None:
                            print(part.code_execution_result.output)

if __name__ == "__main__":
    asyncio.run(main())
```

Combining multiple tools

You can combine multiple tools within the Live API, increasing your application's capabilities even more:

```
prompt = """
Hey, I need you to do three things for me.

1. Compute the largest prime palindrome under 100000.
2. Then use Google Search to look up information about the largest earthqua
3. Turn on the lights

Thanks!
"""

tools = [
    {"google_search": {}},
    {"code_execution": {}},
    {"function_declarations": [turn_on_the_lights, turn_off_the_lights]},
]

config = {"response_modalities": ["TEXT"], "tools": tools}

# ... remaining model call
```

What's next

- Check out more examples of using tools with the Live API in the [Tool use cookbook](https://colab.research.google.com/github/google-gemini/cookbook/blob/main/quickstarts/Get_started_LiveAPI_tools.ipynb) (https://colab.research.google.com/github/google-gemini/cookbook/blob/main/quickstarts/Get_started_LiveAPI_tools.ipynb).
- Get the full story on features and configurations from the [Live API Capabilities guide](#) ([/gemini-api/docs/live-guide](#)).

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.