

Dec 4, 2025

Codex Prompting Guide



Noah MacCallum (OpenAI)



Open in GitHub



View as Markdown

Codex models advance the frontier of intelligence and efficiency and our recommended agentic coding model. Follow this guide closely to ensure you're getting the best performance possible from this model. This guide is for anyone using the model directly via the API for maximum customizability; we also have the [Codex SDK](#) for simpler integrations.

In the API, the Codex-tuned model is `gpt-5.2-codex` (see the [model page](#)).

Recent improvements to Codex models

- Faster and more token efficient: Uses fewer thinking tokens to accomplish a task. We recommend “medium” reasoning effort as a good all-around interactive coding model that balances intelligence and speed.
- Higher intelligence and long-running autonomy: Codex is very capable and will work autonomously for hours to complete your hardest tasks. You can use `high` or `xhigh` reasoning effort for your hardest tasks.
- First-class compaction support: Compaction enables multi-hour reasoning without hitting context limits and longer continuous user conversations without needing to start new chat sessions.
- Codex is also much better in PowerShell and Windows environments.

Getting Started

If you already have a working Codex implementation, this model should work well with relatively minimal updates, but if you’re starting with a prompt and set of tools that’s optimized for GPT-5-series models, or a third-party model, we recommend making more significant changes. The best reference implementation is our fully open-source codex-cli agent, available on [GitHub](#). Clone this repo and use Codex (or any coding agent) to ask questions about how things are implemented. From working with customers, we’ve also learned how to customize agent harnesses beyond this particular implementation.

Key steps to migrate your harness to codex-cli:

1. Update your prompt: If you can, start with our standard Codex-Max prompt as your base and make tactical additions from there.
 - a) The most critical snippets are those covering autonomy and persistence, codebase exploration, tool use, and frontend quality.
 - b) You should also remove all prompting for the model to communicate an upfront plan, preambles, or other status updates during the rollout, as this can cause the model to stop abruptly before the rollout is complete.
2. Update your tools, including our apply_patch implementation and other best practices below. This is a major lever for getting the most performance.

Prompting

Recommended Starter Prompt

This prompt began as the default [GPT-5.1-Codex-Max prompt](#) and was further optimized against internal evals for answer correctness, completeness, quality, correct tool usage and parallelism, and bias for action. If you’re running evals with this model, we recommend turning up the autonomy or prompting for a “non-interactive” mode, though in actual usage more clarification may be desirable.

```
You are Codex, based on GPT-5. You are running as a coding agent in the
# General
- When searching for text or files, prefer using `rg` or `rg --files` re
- If a tool exists for an action, prefer to use the tool instead of shel
```

- When multiple tool calls can be parallelized (e.g., todo updates with
 - Code chunks that you receive (via tool calls or from user) may include
 - Default expectation: deliver working code, not just a plan. If some de
- # Autonomy and Persistence
- You are autonomous senior engineer: once the user gives a direction, p
 - Persist until the task is fully handled end-to-end within the current
 - Bias to action: default to implementing with reasonable assumptions; c
 - Avoid excessive looping or repetition; if you find yourself re-reading
- # Code Implementation
- Act as a discerning engineer: optimize for correctness, clarity, and i
 - Conform to the codebase conventions: follow existing patterns, helpers
 - Comprehensiveness and completeness: Investigate and ensure you cover a
 - Behavior-safe defaults: Preserve intended behavior and UX; gate or fla
 - Tight error handling: No broad catches or silent defaults: do not add
 - No silent failures: do not early-return on invalid input without log
 - Efficient, coherent edits: Avoid repeated micro-edits: read enough cor
 - Keep type safety: Changes should always pass build and type-check; avo
 - Reuse: DRY/search first: before adding new helpers or logic, search fo
 - Bias to action: default to implementing with reasonable assumptions; c
- # Editing constraints
- Default to ASCII when editing or creating files. Only introduce non-AS
 - Add succinct code comments that explain what is going on if code is no
 - Try to use apply_patch for single file edits, but it is fine to explor
 - You may be in a dirty git worktree.
 - * NEVER revert existing changes you did not make unless explicitly r
 - * If asked to make a commit or code edits and there are unrelated ch
 - * If the changes are in files you've touched recently, you should re
 - * If the changes are in unrelated files, just ignore them and don't
 - Do not amend a commit unless explicitly requested to do so.
 - While you are working, you might notice unexpected changes that you d
 - **NEVER** use destructive commands like `git reset --hard` or `git che
- # Exploration and reading files
- **Think first.** Before any tool call, decide ALL files/resources you
 - **Batch everything.** If you need multiple files (even from different
 - **multi_tool_use.parallel** Use `multi_tool_use.parallel` to parallel:
 - **Only make sequential calls if you truly cannot know the next file w
 - **Workflow:** (a) plan all needed reads → (b) issue one parallel batch
 - Additional notes:
 - Always maximize parallelism. Never read files one-by-one unless lo
 - This concerns every read/list/search operations including, but not
 - Do not try to parallelize using scripting or anything else than `m
- # Plan tool
- When using the planning tool:
- Skip using the planning tool for straightforward tasks (roughly the ea
 - Do not make single-step plans.
 - When you made a plan, update it after having performed one of the sub-
 - Unless asked for a plan, never end the interaction with only a plan. F
 - Plan closure: Before finishing, reconcile every previously stated int
 - Promise discipline: Avoid committing to tests/broad refactors unless y

- For any presentation of any initial or updated plans, only update the # Special user requests

- If the user makes a simple request (such as asking for the time) which

- If the user asks for a "review", default to a code review mindset: pr:

Frontend tasks

When doing frontend design tasks, avoid collapsing into "AI slop" or sa

Aim for interfaces that feel intentional, bold, and a bit surprising.

- Typography: Use expressive, purposeful fonts and avoid default stacks

- Color & Look: Choose a clear visual direction; define CSS variables; c

- Motion: Use a few meaningful animations (page-load, staggered reveals)

- Background: Don't rely on flat, single-color backgrounds; use gradient

- Overall: Avoid boilerplate layouts and interchangeable UI patterns. Va

- Ensure the page loads properly on both desktop and mobile

- Finish the website or app to completion, within the scope of what's po

Exception: If working within an existing website or design system, pres

Presenting your work and final message

You are producing plain text that will later be styled by the CLI. Follo

- Default: be very concise; friendly coding teammate tone.

- Format: Use natural language with high-level headings.

- Ask only when needed; suggest ideas; mirror the user's style.

- For substantial work, summarize clearly; follow final-answer formattin

- Skip heavy formatting for simple confirmations.

- Don't dump large files you've written; reference paths only.

- No "save/copy this file" – User is on the same machine.

- Offer logical next steps (tests, commits, build) briefly; add verify s

- For code changes:

* Lead with a quick explanation of the change, and then give more deta

* If there are natural next steps the user may want to take, suggest t

* When suggesting multiple options, use numeric lists for the suggesti

- The user does not command execution outputs. When asked to show the ou

Final answer structure and style guidelines

- Plain text; CLI handles styling. Use structure only when it helps scar

- Headers: optional; short Title Case (1-3 words) wrapped in **...**; no h

- Bullets: use - ; merge related points; keep to one line when possible;

- Monospace: backticks for commands/paths/env vars/code ids and inline e

- Code samples or multi-line snippets should be wrapped in fenced code b

- Structure: group related bullets; order sections general → specific →

- Tone: collaborative, concise, factual; present tense, active voice; se

- Don'ts: no nested bullets/hierarchies; no ANSI codes; don't cram unrelated

- Adaptation: code explanations → precise, structured with code refs; si

- File References: When referencing files in your response follow the best

* Use inline code to make file paths clickable.

* Each reference should have a stand alone path. Even if it's the same

* Accepted: absolute, workspace-relative, a/ or b/ diff prefixes, or t

* Optionally include line/column (1-based): :line[:column] or #Lline[C

* Do not use URIs like file://, vscode://, or https://.

* Do not provide range of lines

* Examples: src/app.ts, src/app.ts:42, b/server/index.js#L10, C:\repo\

Mid-Rollout User Updates

The Codex model family uses reasoning summaries to communicate user updates as it's working. This can be in the form of one-liner headings (which updates the ephemeral text in Codex-CLI), or both heading and a short body. This is done by a separate model and therefore is **not promptable**, and we advise against adding any instructions to the prompt related to intermediate plans or messages to the user. We've improved these summaries for Codex-Max to be more communicative and provide more critical information about what's happening and why; some of our users are updating their UX to promote these summaries more prominently in their UI, similar to how intermediate messages are displayed for GPT-5 series models.

Using agents.md

Codex-cli automatically enumerates these files and injects them into the conversation; the model has been trained to closely adhere to these instructions.

1. Files are pulled from `~/codex` plus each directory from repo root to CWD (with optional fallback names and a size cap).
2. They're merged in order, later directories overriding earlier ones.
3. Each merged chunk shows up to the model as its own user-role message like so:

```
# AGENTS.md instructions for <directory>
<INSTRUCTIONS>
...file contents...
</INSTRUCTIONS>
```



Additional details

- Each discovered file becomes its own user-role message that starts with `# AGENTS.md instructions for <directory>`, where `<directory>` is the path (relative to the repo root) of the folder that provided that file.
- Messages are injected near the top of the conversation history, before the user prompt, in root-to-leaf order: global instructions first, then repo root, then each deeper directory. If an `AGENTS.override.md` was

used, its directory name still appears in the header (e.g., #AGENTS.md instructions for backend/api), so the context is obvious in the transcript.

Compaction

Compaction unlocks significantly longer effective context windows, where user conversations can persist for many turns without hitting context window limits or long context performance degradation, and agents can perform very long trajectories that exceed a typical context window for long-running, complex tasks. A weaker version of this was previously possible with ad-hoc scaffolding and conversation summarization, but our first-class implementation, available via the Responses API, is integrated with the model and is highly performant.

How it works:

1. You use the Responses API as today, sending input items that include tool calls, user inputs, and assistant messages.
2. When your context window grows large, you can invoke /compact to generate a new, compacted context window. Two things to note:
 1. The context window that you send to /compact should fit within your model's context window.
 2. The endpoint is ZDR compatible and will return an "encrypted_content" item that you can pass into future requests.
3. For subsequent calls to the /responses endpoint, you can pass your updated, compacted list of conversation items (including the added compaction item). The model retains key prior state with fewer conversation tokens.

For endpoint details see our [/responses/compact docs](#).

Tools

1. We strongly recommend using our exact `apply_patch` implementation as the model has been trained to excel at this diff format. For terminal commands we recommend our `shell` tool, and

for plan/TODO items our `update_plan` tool should be most performant.

2. If you prefer your agent to use more “terminal-like tools” (like `file_read()` instead of calling `sed` in the terminal), this model can reliably call them instead of terminal (following the instructions below)
3. For other tools, including semantic search, MCPs, or other custom tools, they can work but it requires more tuning and experimentation.

Apply_patch

The easiest way to implement `apply_patch` is with our first-class implementation in the Responses API, but you can also use our freeform tool implementation with [context-free grammar](#). Both are demonstrated below.

```
# Sample script to demonstrate the server-defined apply_patch tool 
```

```
import json
from pprint import pprint
from typing import cast

from openai import OpenAI
from openai.types.responses import ResponseInputParam, ToolParam

client = OpenAI()

## Shared tools and prompt
user_request = """Add a cancel button that logs when clicked"""
file_excerpt = """
export default function Page() {
  return (
    <div>
      <p>Page component not implemented</p>
      <button onClick={() => console.log("clicked")}>Click me</button>
    </div>
  );
}
"""

input_items: ResponseInputParam = [
  {"role": "user", "content": user_request},
  {
    "type": "function_call",
    "call_id": "call_read_file_1",
```

```
"name": "read_file",
"arguments": json.dumps({"path": ("/app/page.tsx")}),
},
{
  "type": "function_call_output",
  "call_id": "call_read_file_1",
  "output": file_excerpt,
},
]

read_file_tool: ToolParam = cast(
  ToolParam,
  {
    "type": "function",
    "name": "read_file",
    "description": "Reads a file from disk",
    "parameters": {
      "type": "object",
      "properties": {"path": {"type": "string"}},
      "required": ["path"],
    },
  },
)

### Get patch with built-in responses tool
tools: list[ToolParam] = [
  read_file_tool,
  cast(ToolParam, {"type": "apply_patch"}),
]

response = client.responses.create(
  model="gpt-5.1-Codex-Max",
  input=input_items,
  tools=tools,
  parallel_tool_calls=False,
)

for item in response.output:
  if item.type == "apply_patch_call":
    print("Responses API apply_patch patch:")
    pprint(item.operation)
    # output:
    # {'diff': '@@\n'
    #      '  return (\n'
    #      '      <div>\n'
    #      '          <p>Page component not implemented</p>\n'
    #      '          <button onClick={() => console.log("clicked")}\n'
    #      '              <button onClick={() => console.log("cancel c1\n'
    #      '                  </div>\n'
```

```

#      );\n'
#      ' }\n',
#  'path': '/app/page.tsx',
#  'type': 'update_file'}

### Get patch with custom tool implementation, including freeform tool
apply_patch_grammar = """
start: begin_patch hunk+ end_patch
begin_patch: "*** Begin Patch" LF
end_patch: "*** End Patch" LF?

hunk: add_hunk | delete_hunk | update_hunk
add_hunk: "*** Add File: " filename LF add_line+
delete_hunk: "*** Delete File: " filename LF
update_hunk: "*** Update File: " filename LF change_move? change?
change_move: "*** Move to: " filename LF
change: (change_context | change_line)+ eof_line?
change_context: ("@@" | "@@ " /(.+)/) LF
change_line: ("+" | "-" | " ") /(.+)/ LF
eof_line: "*** End of File" LF

%import common.LF
"""

tools_with_cfg: list[ToolParam] = [
    read_file_tool,
    cast(
        ToolParam,
        {
            "type": "custom",
            "name": "apply_patch_grammar",
            "description": "Use the `apply_patch` tool to edit files. This tool uses a Lark grammar to parse and apply patches to files. It supports adding, deleting, and updating files, as well as moving them to different locations. The patch format is defined by the 'apply_patch_grammar' configuration below.",
            "format": {
                "type": "grammar",
                "syntax": "lark",
                "definition": apply_patch_grammar,
            },
        },
    ),
],
]

response_cfg = client.responses.create(
    model="gpt-5.1-Codex-Max",
    input=input_items,
    tools=tools_with_cfg,
)

```

```

        parallel_tool_calls=False,
    )

    for item in response_cfg.output:
        if item.type == "custom_tool_call":
            print("\n\nContext-free grammar apply_patch patch:")
            print(item.input)
            # Output
            # *** Begin Patch
            # *** Update File: /app/page.tsx
            # @@
            #     <div>
            #         <p>Page component not implemented</p>
            #         <button onClick={() => console.log("clicked")}>Click me</button>
            #         <button onClick={() => console.log("cancel clicked")}>Cancel</button>
            #     </div>
            # );
            #
            # ***
            # *** End Patch

```

Patches objects the Responses API tool can be implemented by following this [example](#) and patches from the freeform tool can be applied with the logic in our canonical GPT-5 [apply_patch.py](#) implementation.

Shell_command

This is our default shell tool. Note that we have seen better performance with a command type “string” rather than a list of commands.

```
{
  "type": "function",
  "function": {
    "name": "shell_command",
    "description": "Runs a shell command and returns its output.\nAlways run as root.",
    "strict": false,
    "parameters": {
      "type": "object",
      "properties": {
        "command": {
          "type": "string",
          "description": "The shell script to execute in the user's default shell."
        },
        "workdir": {
          "type": "string",
          "description": "The working directory to execute the command in."
        }
      }
    }
  }
}
```

```

    "timeout_ms": {
        "type": "number",
        "description": "The timeout for the command in milliseconds"
    },
    "with_escalated_permissions": {
        "type": "boolean",
        "description": "Whether to request escalated permissions. Set
    },
    "justification": {
        "type": "string",
        "description": "Only set if with_escalated_permissions is true
    }
},
"required": ["command"],
"additionalProperties": false
}
}
}

```

If you're using Windows PowerShell, update to this tool description.

Runs a shell command and returns its output. The arguments you pass will

You can check out codex-cli for the implementation for `exec_command`, which launches a long-lived PTY when you need streaming output, REPLs, or interactive sessions; and `write_stdin`, to feed extra keystrokes (or just poll output) for an existing `exec_command` session.

Update Plan

This is our default TODO tool; feel free to customize as you'd prefer. See

instructions to maintain hygiene and tweak behavior.

```

{
    "type": "function",
    "function": {
        "name": "update_plan",
        "description": "Updates the task plan.\nProvide an optional explanation
        "strict": false,
        "parameters": {
            "type": "object",
            "properties": {

```

```

    "explanation": {
        "type": "string"
    },
    "plan": {
        "type": "array",
        "items": {
            "type": "object",
            "properties": {
                "step": {
                    "type": "string"
                },
                "status": {
                    "type": "string",
                    "description": "One of: pending, in_progress, completed"
                }
            },
            "additionalProperties": false,
            "required": [
                "step",
                "status"
            ]
        },
        "description": "The list of steps"
    }
},
"additionalProperties": false,
"required": [
    "plan"
]
}
}
}

```

View_image

This is a basic function used in codex-cli for the model to view images.

```

{
    "type": "function",
    "function": {
        "name": "view_image",
        "description": "Attach a local image (by filesystem path) to the cor
        "strict": false,
        "parameters": {
            "type": "object",
            "properties": {
                "path": {

```

```

        "type": "string",
        "description": "Local filesystem path to an image file"
    },
},
"additionalProperties": false,
"required": [
    "path"
]
}
}
}

```

Dedicated terminal-wrapping tools

If you would prefer your codex agent to use terminal-wrapping tools (like a dedicated `list_dir('..')` tool instead of `terminal('ls .')`, this generally works well. We see the best results when the name of the tool, the arguments, and the output are as close as possible to those from the underlying command, so it's as in-distribution as possible for the model (which was primarily trained using a dedicated terminal tool). For example, if you notice the model using git via the terminal and would prefer it to use a dedicated tool, we found that creating a related tool, and adding a directive in the prompt to only use that tool for git commands, fully mitigated the model's terminal usage for git commands.

```

GIT_TOOL = {
    "type": "function",
    "name": "git",
    "description": (
        "Execute a git command in the repository root. Behaves like running "
        "the terminal; supports any subcommand and flags. The command can be "
        "a full git invocation (e.g., `git status -sb`) or just the argument "
        "(e.g., `status -sb`)."
    ),
    "parameters": {
        "type": "object",
        "properties": {
            "command": {
                "type": "string",
                "description": (
                    "The git command to execute. Accepts either a full command "
                    "or only the subcommand/args."
                ),
            }
        }
    }
}
```

```

    "timeout_sec": {
        "type": "integer",
        "minimum": 1,
        "maximum": 1800,
        "description": "Optional timeout in seconds for the git
    },
},
"required": ["command"],
},
}

...
PROMPT_TOOL_USE_DIRECTIVE = "- Strictly avoid raw `cmd`/terminal when a

```

Other Custom Tools (web search, semantic search, memory, etc.)

The model hasn't necessarily been post-trained to excel at these tools, but we have seen success here as well. To get the most out of these tools, we recommend:

1. Making the tool names and arguments as semantically “correct” as possible, for example “search” is ambiguous but “semantic_search” clearly indicates what the tool does, relative to other potential search-related tools you might have. “Query” would be a good param name for this tool.
2. Be explicit in your prompt about when, why, and how to use these tools. including good and bad examples.
3. It could also be helpful to make the results look different from outputs the model is accustomed to seeing from other tools, for example ripgrep results should look different from semantic search results to avoid the model collapsing into old habits.

Parallel Tool Calling

In codex-cli, when parallel tool calling is enabled, the responses API request sets `parallel_tool_calls: true` and the following snippet is added to the system instructions:

```
## Exploration and reading files
```



- ****Think first.**** Before any tool call, decide ALL files/resources you need.
 - ****Batch everything.**** If you need multiple files (even from different sources), batch them together.
 - ****multi_tool_use.parallel**** Use `multi_tool_use.parallel` to parallelize tool calls.
 - ****Only make sequential calls if you truly cannot know the next file we want to read.****
 - ****Workflow:**** (a) plan all needed reads → (b) issue one parallel batch of tool calls.
- **Additional notes:****
- Always maximize parallelism. Never read files one-by-one unless logically required.
 - This concerns every read/list/search operations including, but not only limited to, `read_file`.
 - Do not try to parallelize using scripting or anything else than `multi_tool_use.parallel`.



items and responses are ordered in the following way:

```
function_call
function_call
function_call_output
function_call_output
```



Tool Response Truncation

We recommend doing tool call response truncation as follows to be as in-distribution for the model as possible:

- Limit to 10k tokens. You can cheaply approximate this by computing `num_bytes/4`.
- If you hit the truncation limit, you should use half of the budget for the beginning, half for the end, and truncate in the middle with `...3 tokens truncated...`