



# MCP Sequential Thinking Server – Complete Technical Manual

## Part 1: Technical Foundation

### Architecture & Implementation

The **Sequential Thinking MCP Server** is an official reference tool by Anthropic that provides an “*external scratchpad*” for AI reasoning [1](#) [2](#). Unlike an LLM or planner itself, it performs **no independent reasoning** – it simply manages state and formatting of the AI’s thoughts. Internally, the server defines a `ThoughtData` structure with fields for the thought text and metadata like step numbers, revision markers, branches, etc. [3](#). Each time the AI invokes the tool, the server:

- **Validates Input:** Checks that required fields are present and of correct types (e.g. `thought` must be a string, `thoughtNumber` and `totalThoughts` integers, etc.) [4](#). If validation fails, it returns a JSON error with a message and `isError:true` [5](#) [6](#).
- **Appends to History:** Stores the thought in an in-memory list (`this.thoughtHistory`) for the session [7](#) [8](#). The server maintains this chain of thought across consecutive calls, effectively *persisting the reasoning trail in memory*. The state lives as long as the server process runs (i.e. one MCP session); if the server restarts, history is reset. There is no built-in persistence to disk or cross-session memory – it’s an **ephemeral session state** by design.
- **Branches & Revisions:** If the thought is marked as a **branch** or **revision**, the server updates internal structures accordingly. A branch is indicated by providing `branchFromThought` (the origin step) and a `branchId` – the server then records this thought in a separate branch list indexed by that ID [9](#). Revisions are indicated by `isRevision:true` and `revisesThought` pointing to the original step being revised. These flags are stored in the thought’s data but do not alter past history; they serve as annotations (the original thought remains in history, but now a new thought notes that it revisits that step).
- **Formats Output:** The server logs a nicely formatted **ASCII art box** containing the thought and its metadata (with a `\` prefix for normal thoughts, `\|` for revisions, `\|` for branches) [10](#) [11](#). This appears in the server’s console for developers to inspect the reasoning flow. (Logging can be disabled via an env var `DISABLE_THOUGHT_LOGGING=true` to reduce noise [12](#).) The **tool’s API response** to the AI is a JSON structure embedded in a text block: it echoes back the thought content and metadata like the current step number, total steps, and a list of active branch IDs [13](#). For example, after processing a thought, it returns something like:

```
{  
  "thought": "Check if the user is authenticated",  
  "thoughtNumber": 3,  
  "totalThoughts": 5,  
  "nextThoughtNeeded": true,  
  "branches": [],  
  "thoughtHistoryLength": 3  
}
```

This minimal JSON confirms the input was recorded and provides progress info (how many thoughts so far, etc.). Notably, it does **not** return the full history or any new reasoning – it's up to the AI to remember or utilize prior thoughts. The server's role is essentially a **deterministic state-tracker and mirror** <sup>2</sup> <sub>1</sub>. It's been aptly described as "*receives structured input..., tracks the thought in an in-memory array, and provides a pretty-printed version for inspection*" <sup>2</sup>.

Under the hood, the server is implemented using the MCP SDK (TypeScript in the reference version) and runs as a JSON-RPC service (typically over STDIO for local use) <sup>14</sup> <sub>15</sub>. It registers a single tool named `sequentialthinking` which the AI can call <sup>16</sup>. The server doesn't actively push anything to the AI; it reacts to `CallTool` requests for `sequentialthinking` by invoking the above `processThought` logic and returning the result <sup>17</sup>. **Lifecycle:** Once launched (e.g. via `npx @modelcontextprotocol/server-sequential-thinking` or Docker), it runs persistently in the background waiting for requests <sup>18</sup> <sub>19</sub>. Each connected client (e.g. Claude or VS Code's MCP client) will typically spawn its own instance or maintain its own session. There is no special session-handling beyond the singular history – so if one were to reuse the same server across unrelated tasks or users without resetting, the `thoughtHistory` would accumulate. In practice, tools are usually started/stopped per context or isolated to avoid cross-talk.

Because the state is in-memory, **long thought chains** increase the server's memory usage slightly (storing each thought string and metadata), but this is minimal overhead for reasonable chain lengths. The reference docs note "*sequential-thinking - minimal overhead*" in terms of performance impact <sup>20</sup>. The JSON-RPC processing is lightweight, and logging is the only potentially heavy operation (string formatting). Latency added per thought is negligible (typically a few milliseconds). In other words, the server should not be a bottleneck even in multi-step reasoning – the main cost of using it is in **token overhead** and complexity (discussed later), not CPU time.

## Protocol & Interface Details

**Tool Name:** The server registers the tool under the name `"sequentialthinking"` (or sometimes shown as `"sequential_thinking"` in documentation, but the actual name is without spaces) <sup>21</sup>. The AI invokes it like any other MCP tool call, passing a JSON object with the required fields.

**Input Schema:** The JSON input has the following fields (with their types and meanings) <sup>22</sup> <sub>23</sub>:

- `thought` (string, required): The content of the current thought – essentially the step's reasoning or conclusion in natural language <sup>24</sup>. This is the *substance* of what the AI is thinking at this step.
- `nextThoughtNeeded` (boolean, required): A flag indicating whether the AI anticipates needing another thought after this one <sup>25</sup>. If `true`, it signals the chain should continue; if `false`, it means the AI believes this thought completes the reasoning (i.e. it has reached a final answer or no further analysis is needed).
- `thoughtNumber` (integer, required): The index of this thought in the sequence <sup>26</sup>. Counting starts at 1 for the first thought. The AI is responsible for incrementing this as it generates each step. (The server does not enforce sequential incrementing, but passing non-sequential numbers would be logically inconsistent. If `thoughtNumber` ever exceeds `totalThoughts`, the server will automatically bump up `totalThoughts` to match it <sup>8</sup>.)
- `totalThoughts` (integer, required): The AI's current estimate of how many thoughts will be needed in total <sup>26</sup>. This is an **elastic planning parameter** – it's not a hard limit, just an initial guess or plan length. The AI can adjust it on the fly: for instance, start with 5, but if step 5 ends and the problem isn't solved, it can set a higher total and `nextThoughtNeeded: true` to

continue. Conversely, it might finish early (setting `nextThoughtNeeded:false` before reaching the originally stated total). The server doesn't enforce using all steps; it's mainly used for logging context and guiding the AI's own thinking.

- `isRevision` (boolean, optional): Indicates this thought is a revision of a previous one <sup>27</sup>. When `true`, the AI is effectively saying "*I am reconsidering or correcting an earlier thought.*" It should accompany `revisesThought` to point to which step is being revised.
- `revisesThought` (integer, optional): If revising, this is the number of the thought that is being revisited or corrected <sup>28</sup>. For example, `isRevision:true` and `revisesThought:2` means "*I'm updating my 2nd thought with new understanding.*" This does **not** remove or overwrite the original thought; it just marks the relationship. The server will label the formatted log as "`Revision (revising thought 2)`" <sup>29</sup>.
- `branchFromThought` (integer, optional): If this thought initiates a new branch of reasoning, this field is the number of the thought from which the branch diverged <sup>30</sup>. For example, `branchFromThought:3` with a new `branchId` indicates the AI is exploring an alternative path starting from what was established at step 3. The server logs it as "`Branch (from thought 3, ID: X)`" in the console <sup>31</sup>.
- `branchId` (string, optional): An identifier for the branch <sup>32</sup>. The AI chooses a simple ID (like "B1" or "alternativeA"). All thoughts sharing this `branchId` are tracked together. The server will list active branch IDs in the output JSON (`"branches": [...]`) so the AI knows which branches exist <sup>33</sup>. It's up to the AI to manage how it uses branch IDs (e.g. continuing a branch vs. mainline).
- `needsMoreThoughts` (boolean, optional): A hint that even after reaching the originally planned total thoughts, more steps are needed <sup>34</sup>. In practice, this is similar to setting `nextThoughtNeeded:true` at the final step; it explicitly flags that the initial plan was insufficient. The AI can set `needsMoreThoughts:true` when it realizes late in the process that it must extend the chain (this is more of a semantic marker; functionally, the AI could also just increase `totalThoughts` and keep `nextThoughtNeeded:true`).

**Output Schema:** The server's response (as seen by the AI) is an object: `{ content: [ {type: "text", text: "<JSON string>"}, isError: false }]` (or `isError:true` with an error message JSON if something went wrong) <sup>5</sup>. The `text` is a JSON-formatted string containing at least the fields: `thought`, `thoughtNumber`, `totalThoughts`, `nextThoughtNeeded` echoing back the input, plus `branches` (an array of branch IDs currently in use) and `thoughtHistoryLength` (the count of thoughts stored so far) <sup>13</sup>. This information can be used by the AI to gauge progress or branch context. *Important:* The **actual thought text content is not fed back** beyond echoing the last input. The server does *not* supply the AI with earlier thoughts automatically (the assumption is the AI retains them in its own context or doesn't need them explicitly). The JSON is mostly a confirmation and bookkeeping. In some client implementations, the MCP client might not even expose this JSON to the AI's prompt (it could intercept it just for tool state management). But typically in Claude/VSCode, the AI does see this text (as a tool response), which it may largely ignore or parse minimally.

**Error Handling & Edge Cases:** The server is strict about types: if any required field is missing or of the wrong type, it throws an error with a message like "`Invalid thoughtNumber: must be a number`" <sup>35</sup>. Optional fields can be omitted (they'll just be `undefined` in the internal object). The server doesn't check logical consistency (e.g., you can set `revisesThought:5` even if there was no thought 5 yet, or branch from a future thought) – such usage would be nonsensical, but the server wouldn't know; it's the AI's responsibility to use the fields coherently. There's also no cap on `thoughtNumber` / `totalThoughts` values (other than being integers  $\geq 1$ ), so the AI must avoid infinite loops or runaway counting. If the AI sets `nextThoughtNeeded:true` indefinitely, the server will happily keep accepting thoughts until the universe ends – there is no built-in termination condition other than the AI eventually sending

`nextThoughtNeeded: false`. (The server's own guidance explicitly states to only end when truly done: "Only set `next_thought_needed` to **false** when a satisfactory answer is reached." <sup>36</sup>). Thus, avoiding infinite reasoning loops is entirely up to the AI/agent (discussed more under pitfalls).

The server also doesn't enforce uniqueness of branch IDs – if the AI reuses the same `branchId` for different branch origins, those will co-mingle in one list. By convention, the AI should generate distinct IDs for distinct branches. Also, each branch thought still increments the global `thoughtNumber` sequence (so branches don't have independent numbering; e.g. you might have thoughts 1,2,3 then branch at 3 creating thought 4 with branchId "B", thought 5 with branchId "B", then maybe back to main as thought 6, etc.). The linear history holds *all* thoughts in order of creation, and `branchId` lets the AI filter which belong to an alternate path. This design keeps things simple in implementation but means the AI must manage not to confuse itself when juggling branches.

**Connection & Lifecycle:** Typically, you add the Sequential Thinking server to your MCP config (e.g. in `claude_desktop_config.json` as shown in official docs) and start it. Once running, it appears as an available tool. Notably, the AI usually **won't call it automatically** unless prompted or instructed. In many UIs (Claude Desktop, VS Code etc.), you must enable the server and often hint the AI to use it. For example, after installing you might prompt Claude: "*Use the sequential thinking tool to solve this.*" The documentation notes "*Claude understands which tool you mean regardless [of spacing]*" and that if you don't direct it, it likely won't engage the tool on its own <sup>37</sup>. This is a crucial detail: **the agent decides when to invoke the tool**. The server itself just waits; it doesn't initiate anything.

Finally, since the Sequential Thinking server only defines one tool and a simple protocol, it has **no special resource requirements or external dependencies**. It doesn't access network, files, etc., so its security surface is minimal. It also means the tool is **model-agnostic** – any LLM client using MCP can use it (Claude, GPT-4 with an MCP client, etc.) <sup>38</sup>. It's basically a universal plugin for structuring reasoning.

## Behavioral Analysis

Using the Sequential Thinking server can significantly **change the AI's reasoning style and output structure**. It forces the model into an explicit stepwise format, in contrast to the opaque chain-of-thought that normally happens hidden inside the model. Some key behavioral effects and considerations include:

- **Structured, Transparent Reasoning:** With Sequential Thinking, the AI's problem-solving becomes more like a checklist or journal of thoughts. Each step is recorded and visible (to the developer, if not always to the end-user). This yields a *transparent and auditable trail of reasoning* <sup>39</sup> <sup>40</sup>. For example, an agent might output: "Thought 1: I should clarify the requirements...", "Thought 2: Now I will plan sub-tasks...", etc., if those thoughts are user-visible. Even if hidden, the developer can inspect the logged steps. This transparency makes it easier to debug why an AI reached a conclusion or where it went wrong, addressing the "black box" nature of AI reasoning <sup>39</sup>.
- **Stateful "Memory" of Process:** Because the server keeps a context of all previous thoughts (outside the model's own context window), it can help the AI maintain continuity in long analyses without reiterating everything in the prompt. The AI can implicitly rely on the fact it *has* a chain of thought recorded. However, recall that the server doesn't feed back the content of those thoughts – the benefit is indirect. In practice, the model still needs to remember earlier conclusions (or re-derive them) unless the system or developer surfaces the thought history.

Some advanced setups do parse the `thoughtHistoryLength` or the log and provide summaries back to the model if needed, but the vanilla use doesn't include an automatic memory injection beyond the model's own capacity. Nonetheless, **the mere act of writing thoughts in sequence can improve the model's focus**, much like how humans writing down steps helps avoid losing track. It externalizes some of the cognitive load.

- **Improved Problem Decomposition:** The requirement to specify `totalThoughts` upfront encourages the AI to **plan out a series of steps** rather than jumping straight to an answer <sup>41</sup>. Even though the number can change, starting with "I think I'll need about N steps" sets a mindset of methodical progression. Studies and community reports indicate that this approach (akin to "chain-of-thought prompting") often yields better performance on complex tasks that need lookahead or multi-step reasoning <sup>42</sup> <sup>43</sup>. The AI is effectively prompted to "*think step-by-step*", which can reduce mistakes. For instance, an AI might avoid a logical error because it catches it in a later thought and marks a revision, instead of giving a flawed answer immediately. One analysis notes that with sequential thinking, "*confident-but-wrong answers don't survive*" because the process encourages self-checking and alternate approaches <sup>44</sup>.
- **Ability to Revise and Branch:** In native reasoning, an LLM can of course correct itself mid-response, but the Sequential Thinking schema makes revisions explicit. This can influence Claude (or other models) to be more critical of its prior steps, since it has a formal way to mark a mistake and try again. Similarly, branching allows exploration of multiple possibilities in parallel (though serialized in practice). For example, if uncertain between two approaches, the AI can fork: *Thought 4A*: do path A, *Thought 4B*: do path B (with branches "A" and "B"), then later compare. This *tree-of-thoughts* style reasoning is facilitated by the `branchId` system, and is known to improve problem-solving by not getting stuck in one line of thought <sup>45</sup>. In effect, the server supports a kind of **non-linear thinking** that standard prompt-response does not – the AI can pursue different lines and then converge. This often leads to more robust solutions, at the cost of more computation.
- **Changes in Output Quality:** When used appropriately, Sequential Thinking tends to yield **more thorough and well-justified answers** for complex queries. The final answer typically comes after a series of tool calls (thoughts), meaning the AI has effectively "double-checked" its work. Users have reported that for complicated tasks (like debugging a tricky code issue or designing an architecture), the results were more reliable: "*it gave me better output for more complex problems like building a data pipeline end-to-end*" <sup>46</sup>. The structure helps ensure no important step is skipped – which can prevent subtle failures. On the other hand, if overused on simple tasks, it can be overkill. Several developers noted that for straightforward queries, enabling sequential thinking made no noticeable improvement – or even made things worse by adding verbosity. The consensus is that **the benefit scales with problem complexity**. For trivial Q&A or single-step problems, the model's native reasoning is fine (and faster); for anything with ambiguity, multiple steps, or high risk of error, the structured approach shines <sup>47</sup>.
- **Token & Time Overhead:** Because every thought is communicated via a tool call and the model must generate the JSON arguments and read the JSON response, there is an overhead in tokens and latency. Each thought might consume a handful of extra tokens (for the JSON boilerplate and the reasoning text that might otherwise be internal). Over a long chain, this adds up. In addition, the model might use some of its context window to recall prior thoughts (unless offloaded to the server memory alone). In practice, users have mixed opinions on this overhead. The official perspective is that the overhead is minimal relative to the gains – "*Sequential-thinking [has] minimal overhead*" in performance terms <sup>48</sup>. And indeed, the server itself is lightweight. However, from a cost perspective, one must be mindful that **each thought is essentially an**

**additional prompt-response cycle.** If an agent uses 10 thoughts to solve something it could have done in 2, it may burn five times more tokens. Some advanced models or environments have built-in “extended thinking” that achieves a similar stepwise reasoning internally without extra API calls, which can be more token-efficient. For example, GitHub Copilot’s agent had a built-in think/todo mechanism, leading some to call Sequential Thinking MCP redundant in that context <sup>49</sup>. In Claude’s case, when using Claude Code, some found that **disabling the Sequential Thinking MCP did not degrade the reasoning** on simpler tasks – Claude was able to “reason just fine on its own” <sup>50</sup>. So the value must justify the cost: it’s most justified for difficult tasks where the cost of a mistake is high, and less so for trivial tasks (more on this in Part 4 decision criteria).

- **Comparison to Native CoT:** Modern LLMs like Claude and GPT-4 already use chain-of-thought style reasoning internally (Claude has things like “Windsurf” or extended reasoning modes). The difference with the MCP approach is **portability and standardization**. The Sequential Thinking server externalizes the chain-of-thought in a way that’s consistent across models and clients <sup>51</sup>. This means an AI agent using MCP can have the same structured thinking process whether it’s Claude today or GPT tomorrow, because the logic lives partly in the skill (the way it calls the tool) rather than being an intrinsic model quirk. Also, the external trace is *auditable and reproducible*, which is useful for debugging or compliance (you can save the reasoning steps and verify them). If you rely only on the model’s hidden reasoning, you might not know why it did something if it fails. One community member summed it up: *Use Sequential Thinking MCP when you care about portability, auditability, and reproducible reasoning; otherwise, the host’s built-in planning is usually enough.* <sup>51</sup> In short, the **server doesn’t necessarily make the AI “smarter,” but it makes its reasoning more structured and observable**, which can indirectly improve outcomes and definitely improve traceability.
- **Potential Drawbacks:** If misused, sequential thinking can lead to **overthinking or looping**. An AI might keep generating thoughts without ever deciding it’s done (especially if it’s unsure when to set `nextThoughtNeeded: false`). If not well-prompted, it might also produce very superficial thoughts just to fill out the sequence (e.g., “Thought 1: I will solve the problem. Thought 2: Now I solved it.” – not actually helpful). Some users observed that if the agent goes down a wrong path, the sequential format might actually *increase confidence in the wrong path* because the AI keeps elaborating a flawed plan step by step <sup>52</sup>. Ideally, the revision mechanism mitigates this (the AI should revise if it notices a mistake), but it requires the model to realize its error. There’s also the scenario where the AI spends too many steps “thinking” and not enough time **doing** – especially if it doesn’t integrate other tools when needed. We’ll discuss strategies to avoid these anti-patterns in Part 2 and Part 4.

In summary, the Sequential Thinking server *augments Claude’s reasoning by enforcing a disciplined approach*: break problems into steps, document each step, allow mid-course corrections, and only finalize when satisfied <sup>53</sup> <sup>36</sup>. This tends to produce more **logical, well-vetted solutions** at the cost of some additional interactions. It transforms the AI from a one-shot respondent into a more “*methodical expert*” that can show its work <sup>54</sup> <sup>55</sup>. When used in the right scenarios, this leads to better quality and user trust. In scenarios where the overhead outweighs the complexity, it can be unnecessary. The key is knowing when to invoke this tool – something we explore next.

## Part 2: Best Practices & Patterns

Using Sequential Thinking effectively requires understanding *when* it adds value and how to structure the chain-of-thought for maximum benefit. This section distills proven usage patterns, common pitfalls, and techniques from official guidance and community experience.

### Proven Usage Patterns

**When to Use Sequential Thinking – Decision Criteria:** As a rule of thumb, reserve the Sequential Thinking tool for **non-trivial, multi-step problems**. If the task is “*genuinely hard*” or has many interconnected parts, sequential thinking can provide clarity. Examples of ideal scenarios include:

- **Complex debugging:** e.g. tracing a bug in code that might have multiple causes. The stepwise approach ensures you check one hypothesis at a time and see the effect <sup>47</sup>.
- **Architecture or design decisions:** when you have to evaluate multiple approaches (design patterns, algorithms, etc.) – sequential thinking allows explicitly comparing alternatives (potentially via branching) and gradually converging on the best solution <sup>47</sup>.
- **Strategic planning and multi-phase tasks:** for instance, planning a multi-stage project or deployment. The tool excels at laying out a phased plan and adjusting as needed <sup>56</sup>.
- **Research and analysis:** breaking down a research question into sub-questions, fact-checking each, and synthesizing results. If the problem requires gathering and verifying information across sources, sequential steps can structure that process <sup>47</sup>.
- **Code review and refactoring:** analyzing code with multiple issues or planning a large refactor (as a step-by-step plan of what to change first, next, etc.) <sup>57</sup> <sup>58</sup>.
- **Requirement decomposition:** taking a complex or vague requirement and systematically clarifying and subdividing it into actionable tasks. The thought chain can first interpret requirements, then break them down, then tackle each piece.

In short, **use it for tasks with ambiguity, multiple steps, or when mistakes are costly**. A community guideline puts it clearly: “*For simple tasks ('add two numbers'), it's overkill... Use it when the problem is genuinely hard. Or when getting it wrong is expensive.*” <sup>59</sup> If a task can be answered in one step or is straightforward, forcing a structured chain can waste time and tokens without improving results. Many users have found that enabling sequential thinking on trivial prompts just adds unnecessary verbosity.

**Initial Planning – Choosing Total Thoughts:** When starting a thought chain, the AI should make a reasonable initial guess for `totalThoughts`. This doesn't have to be exact (and can change), but it sets the scope. A good practice is to err on the side of a *moderate number* of steps – not too low (which might pressure the AI to rush or prematurely conclude) and not absurdly high (which might lead to verbose, low-value thoughts). For example, if faced with a moderately complex problem, an initial plan of 3-5 thoughts is common; for a truly complex project, maybe 8-10. The key is to **reflect problem complexity in the plan length**. If during reasoning the AI realizes more steps are needed, it should update `totalThoughts` upward and set `needsMoreThoughts:true` or simply continue beyond the original count (the server will accept it) <sup>8</sup>. Conversely, if it can conclude early, it can stop before using all planned steps (set `nextThoughtNeeded:false`).

Concrete example: “*Design a system architecture for a scalable web app.*” The AI might start with `totalThoughts: 5` (e.g., 1: gather requirements, 2: outline components, 3: consider bottlenecks, 4: pick tech stack, 5: finalize diagram) and `thoughtNumber:1`. If by thought 5, some open issues remain, it could add a 6th thought, increasing `totalThoughts` to 6 or just marking `needsMoreThoughts:true`. **It's okay for the plan to evolve** – in fact, the ability to adjust is a feature

of the tool <sup>60</sup>. Just ensure `thoughtNumber` and `totalThoughts` remain consistent (the server's auto-adjust will make `totalThoughts` at least as large as the highest `thoughtNumber`).

**Step-by-Step Clarity:** Each thought should ideally represent a **clear, distinct step** in reasoning. Avoid overly broad or vague thoughts that encompass too much at once. The power of sequential thinking comes from granularity. For example, if solving a math word problem, a good sequence might be: (1) restate the problem in own words, (2) identify relevant formulas, (3) do the calculation step 1, (4) do calculation step 2, (5) verify result and conclude. By contrast, a poor usage would be to cram multiple logical leaps into one "thought." A practical tip from the official description is to **treat each thought as one unit of reasoning** – whether that's an analysis, an action, a check, or a realization <sup>61</sup>.

**Progressive Refinement:** One effective pattern is **progressive refinement** of a solution. Early thoughts can be high-level outlines, and later thoughts fill in details. The Sequential Thinking tool encourages this by letting the AI express uncertainty or incomplete ideas initially, then refine them. For instance, *Thought 1*: "I recall the problem asks for X; I plan to consider A or B approaches. I'll need more info." (`nextThoughtNeeded = true`). *Thought 2*: "Let me analyze Approach A..." followed by perhaps discovering an issue. *Thought 3 (revision)*: "I realized Approach A has flaw, revising thought 2..." and so on. This **incremental deepening** mirrors how a human might start with a rough plan and iteratively improve it. The AI should not be afraid to use a thought to *pause and reflect*, or even say "I'm not fully sure, I will gather more information or think of alternatives." In fact, explicitly **flagging uncertainty** in a thought (e.g. "I'm not confident in this approach; I'll branch out.") is encouraged – it can then either branch or revise in subsequent steps rather than plowing ahead blindly <sup>61</sup>.

**Integrating Tools Within the Chain:** A powerful usage pattern is to **interleave sequential thinking with other tools**. The sequential thoughts can serve as the "planner" or reasoning layer, and between thoughts the AI can call action tools (filesystem, web search, database, etc.) as needed. For example, in a coding scenario:

1. *Thought 1*: "I need to understand the user's code structure before refactoring. Next step: read the main file." (`nextThoughtNeeded:true`) <sup>62</sup>
2. AI then calls the Filesystem MCP to `readFile` (this is outside the ST tool – an actual action).
3. *Thought 2*: "I have the file content. It looks like there's a performance bottleneck in function X. Next, I should identify specific inefficiencies."
4. Perhaps calls a profiling tool or just continues reasoning...
5. *Thought 3*: "One inefficiency is the nested loop. I think using a hashmap would improve this. I will proceed to implement that. `nextThoughtNeeded:true, thoughtNumber:3/ maybe 5.`"
6. AI calls Filesystem MCP to write changes to the file (action).
7. *Thought 4*: "Changes made. I should run tests to verify nothing broke. `nextThoughtNeeded:true.`"
8. Calls Shell MCP to run tests. Suppose tests fail.
9. *Thought 5*: "Test failed, indicating an issue with null inputs. I'll revise my approach to handle that." (`isRevision:true, revisesThought:3`)
10. ... and so on, eventually concluding with a final thought with `nextThoughtNeeded:false` and the answer or summary.

This pattern – **Planner/Doer** – is highly effective <sup>63</sup> <sup>64</sup>. Sequential Thinking acts as the *planner (brain)* that figures out what needs doing, step by step, while other MCP servers act as the *doers (hands)* executing those steps <sup>65</sup>. The best practice is to use ST to decide *what to do next*, then actually do it with a specialized tool, then feed the result into the next thought. This keeps the reasoning grounded in reality (the AI can verify assumptions through action) and prevents the thought chain from becoming detached or purely hypothetical. It's essentially implementing an observe-think-act loop.

**Example:** A user tasked the agent to **migrate a project from JavaScript to TypeScript**. The agent used Sequential Thinking to outline a plan: *Thought 1*: list all steps (rename files, update imports, add types, update config, run tests) – it set `totalThoughts:10` <sup>66</sup>. Then for each step, the agent called the appropriate tools: renaming files with Git/MV commands, writing new type definitions with Filesystem, running `tsc` with Shell, etc., interleaving with ST thoughts that checked off each step <sup>58</sup> <sup>67</sup>. This orchestrated approach “*saved hours of back-and-forth*” and ensured the agent didn’t skip crucial tasks <sup>68</sup>. The key takeaway: **don’t just think in a vacuum – use the thinking steps to guide real actions as needed**.

**Branching Strategies:** When faced with uncertainty or multiple approaches, use branching deliberately. A common pattern is to branch when you have **distinct options to evaluate**. For example, “I could solve this either with method A or method B.” The AI can then do: *Thought 3*: “I will explore method A first. (`branchFromThought:2, branchId:'A'`) ... analysis of A ... `nextThoughtNeeded:false` (if concluding branch A).” Then it can *backtrack* to the main line or start branch B: *Thought 3 again but with branchId:'B'* exploring method B. After exploring both, it might have a Thought 4 (back on main, no `branchId`) to compare results and choose the best solution. Each branch is essentially a sub-sequence of thoughts. It’s wise to **keep branches relatively short** – they are for parallel exploration, not infinite side quests. Also, clearly name or label the branches if possible (even just “A” / “B”) and keep track of them via the `branches` list returned. The server will list all branch IDs that have been used in `branches` <sup>33</sup>, which can remind the AI if it left any branch unfinished. Once a branch yields a conclusion or is deemed unpromising, the AI should return to either the main line or a higher-level revision.

One caution: **Do not branch frivolously**. Branches add complexity. Only branch when there is a genuine decision point or alternative worth exploring. Otherwise, you can confuse yourself (the agent) with too many threads. A good practice is to branch at most a couple of times in a single session, and not too deeply. E.g., one branch for an alternative approach, or a branch to handle an edge-case scenario separately. The aim is to use branching to *compare solutions or cover different angles* without conflating them in one linear narrative. If branches are used, eventually the AI should converge – possibly via a revision or a final thought that reconciles branch outcomes.

**Revision Patterns:** Use revisions when the AI identifies that a previous thought was incorrect or based on a false assumption. The pattern is: state the realization and mark the correction. For instance, *Thought 5*: “*I realize my earlier assumption in Thought 2 was wrong because X. I will revise that approach.*” (`set isRevision:true, revisesThought:2`). The next thought (or the content of this revision thought itself) should then provide the corrected thought 2. Often the AI may effectively “replace” the logic of thought 2 and then proceed from there. In practice, because the history isn’t altered, the agent must keep track that thought 2 is superseded. A wise approach is after a revision, to branch or continue fresh from that point. For example, after revising thought 2, the agent might treat the next thought as a new branch (since the subsequent original thoughts 3,4 might not apply anymore). Or it can simply proceed as if thought 2 was fixed and do thought 6 as the new step 3 (conceptually). There’s flexibility – the main point is to **acknowledge mistakes explicitly** rather than silently carrying them forward. This builds more reliable solutions.

A pattern that emerged from users is to trigger a revision whenever a test or validation fails. E.g., *Thought 7*: “*The test in step 6 failed. I will revise thought 4 (where I implemented the function) to account for the edge-case.*” (`isRevision:true, revisesThought:4`). This way, the plan itself gets updated mid-flight. The final chain of thoughts will show a clear story: it tried something, found it flawed, corrected it – which is great for traceability.

**Example Prompt Templates:** To get the AI to use sequential thinking, you often have to nudge it. Some effective prompt starters: - *"Let's solve this step by step. Use the sequential thinking tool to outline your plan."* – This explicitly tells the model to invoke the tool and start a chain.

- *"Proceed using the Sequential Thinking method."* – This instruction (as noted in a user forum <sup>69</sup>) can trigger the agent to call the tool, as it recognizes the name.

- *"Think out loud with the sequentialthinking tool and show your reasoning."* – This encourages transparent reasoning via the MCP server.

Once the chain is initiated, the AI will continue calling the tool for each subsequent thought until it concludes. It's important to include in the system or developer prompt something like: *"Only set nextThoughtNeeded to false when you have a final answer."* (This is basically quoting the built-in guideline <sup>36</sup>). That ensures the agent knows not to stop prematurely.

**Optimal Chain Length:** While there's no fixed "optimal" number of thoughts (depends on the task), a pattern from community usage is: - **Small tasks:** 0-1 thoughts (just do it directly, no tool) - **Medium tasks (some complexity):** ~3-5 thoughts often suffices <sup>70</sup>. This covers things like a moderate coding problem or a single-file refactor, where you might plan a few steps (understand problem, implement fix, test, done). - **Large tasks:** Could be 6-10 or more thoughts, broken into phases. For instance, a project plan might have 10 steps (as in the TypeScript migration case) <sup>66</sup>. Going beyond ~10 thoughts in one chain is possible but relatively rare; at that point it might be better to break the problem into sub-problems or use memory to handle the complexity. Very long chains also risk earlier context falling out of the model's window (if the intermediate results aren't being summarized).

Always be willing to adjust: It's perfectly fine if the AI planned 6 steps but solves everything in 4 – then it should cut it short (set `nextThoughtNeeded:false` at thought 4). The **goal is not to exactly hit the initial count, but to have a plan and adapt as needed.**

## Anti-Patterns & Pitfalls

Knowing what *not* to do is as important as knowing best practices. Here are common pitfalls and how to avoid them:

- **Over-engineering Simple Problems:** Perhaps the most frequent anti-pattern is invoking the sequential thinking tool for something that doesn't need it. If the user asks "What is 5+7?" and the agent starts a 3-thought chain ("Thought 1: Understand the problem – it's addition. Thought 2: Calculate 5+7=12. Thought 3: Therefore the answer is 12."), this is overkill. It wastes tokens and time. It's crucial to have a **decision step (or heuristic)** in the agent: *Is this problem complex enough to warrant Sequential Thinking?* If not, skip it. Otherwise, you risk annoying users and consuming budget for zero benefit. In general, **don't use sequential thinking for one-step Q&A or trivial tasks** <sup>59</sup>. A related pitfall is continuing a chain longer than needed – once the solution is clear, the agent should stop rather than inventing steps just because it planned for more. Always remember to finish when done (and set that `nextThoughtNeeded:false` promptly).

- **Shallow or Redundant Thoughts:** Sometimes the AI might generate thoughts that are basically no-ops, just to fill space ("Thought 3: I will continue thinking." or repeating what's already known). This defeats the purpose of the structured approach. Every thought should *progress* the solution in some way – either by adding new insight, making a decision, or revising something. If you find the chain-of-thought includes filler or obvious statements that don't aid progress, that's an anti-pattern. To combat this, the agent can be guided with a principle: *each thought must*

*meaningfully contribute*. It might help to ask after each thought, “*Did this step resolve something or lead to a next actionable point?*” If not, it wasn’t a good thought. Avoid overly granular thoughts that don’t add value (e.g., “Thought 1: I will begin. Thought 2: I have begun.”). Balance is key – not too big steps, but also not absurdly small trivial steps.

- **Infinite Loops / Never Concluding:** Because the tool leaves termination to the AI, a poorly managed chain can loop forever. For example, if the AI is unsure and keeps saying `nextThoughtNeeded:true` without ever deciding, you have a problem. Users have noted that sometimes agents could get stuck, especially if they’re not confident in the answer, constantly revising or branching but never finalizing <sup>71</sup>. To avoid this, implement **guardrails for convergence**: The agent should have a logical stopping criterion (e.g., “if I’ve branched and revised and still have no solution after X thoughts, I should stop or escalate”). Often, simply reminding the AI of the goal each step helps – the content of each thought can reiterate how close it is to solving the original question. Also, the design of the tool itself nudges not to loop indefinitely by requiring that final answer only when ready <sup>36</sup>. If a chain does go unexpectedly long (past a reasonable number of steps), the agent or calling application might consider forcibly halting and summarizing. But ideally, the AI’s policy should prevent infinite loops (perhaps by incrementing a hidden counter and deciding to break after a threshold).
- **Branch Proliferation / Unmanaged Branches:** Branching can become chaotic if not handled carefully. An anti-pattern is branching on a branch on a branch... without ever reconciling them. While the server would technically allow multiple nested branches (just as more entries in its branch map), the cognitive load on the AI skyrockets and it may lose track. It’s recommended to **keep branching breadth and depth limited**. The AI should label branches clearly and merge or abandon them explicitly. For instance, if it explored two branches, it should at some point say “Now I will compare Branch A and Branch B results and pick one.” Leaving branches hanging (e.g., starting branch “B” and never coming back to it) is like leaving a sub-problem unsolved. Also avoid using branching when a simple linear revision would do. If the alternative approach comes as a result of realizing the first is flawed, that’s more of a revision than a branch. Branch when two approaches *both seem viable initially* and you want to see which works better. If one is clearly inferior halfway, you can drop that branch with a thought noting “Branch X turned out to be a dead end.” Marking that explicitly is good practice – it tells anyone reading the chain that the branch was closed.
- **Lack of Integration (Analysis Paralysis):** Another pitfall is the agent staying in “thinking mode” too long without using available tools to actually get answers. For example, imagine an agent tasked with answering a question that requires data from a file. If it keeps making sequential thoughts guessing what might be in the file instead of just reading it via the Filesystem MCP, that’s an anti-pattern. The fix is to always consider, “*Do I have a tool that could answer this question directly?*” If yes, use it rather than purely theorizing. Sequential Thinking should not replace tool use; it should guide tool use. If you see the AI writing a multi-step thought process about something it could simply fetch (like “*Thought: Perhaps the API will return X, I wonder what data it has...*” and it has a Fetch tool available), it should instead call the Fetch tool and then think based on actual data. Don’t let the agent fall into a self-dialogue loop when an external action is needed.
- **Premature Finalization:** The opposite of endless looping is stopping too early, possibly because the AI is eager to produce an answer. If the chain stops before truly resolving key parts, that’s a failure mode. Signs include: the final answer is given but it hasn’t addressed part of the question, or the agent ended the chain because it hit the initial `totalThoughts` count even though unresolved issues remain. To avoid this, emphasize that `nextThoughtNeeded` should be true

until *all aspects of the problem are solved*<sup>36</sup>. The AI should not feel compelled to stop just because it reached the planned number of thoughts – plans can change. Also, if the final answer is derived without using all information or verifying certain things mentioned in earlier thoughts, that's suspect. Encourage the agent to **double-check** before concluding: e.g., one of the guidelines from the tool description is “*verify the hypothesis based on the chain of thought steps and repeat until satisfied*”<sup>72</sup>. That implies the agent should only finalize when it's confident.

- **Ignoring User Input or Deviating from Task:** While engaged in deep sequential thinking, the agent must remain anchored to the user's request. An anti-pattern would be the agent getting so engrossed in its thought process that it drifts off-topic or solves a different problem. Each thought should still be in service of the original question or task. It's good to occasionally restate the goal in a thought (e.g., “Our goal is X, so far I have done Y, next I will do Z”) to ensure alignment. If the user provides new input mid-way (in interactive settings), the agent might need to adjust or even restart the chain – the server has no direct concept of a “restart,” but the agent could simply drop the current context and begin a new chain of thought for the new info (or integrate it via a revision). Be cautious that the existence of a chain doesn't lock the agent into a tunnel vision – remain responsive to the actual problem requirements.

In summary, **use Sequential Thinking judiciously**. It's a powerful tool, but misuse can lead to verbose, roundabout solutions or confusion. Always aim for clarity, relevance, and purposeful steps. When in doubt, ask: “Is this thought helping solve the problem?” If not, refocus or end the chain.

## Domain-Specific Applications

Sequential Thinking is a general tool, but certain domains especially benefit from its structured approach. Here we outline how it can be applied in various development and problem-solving contexts, with domain-specific tips:

- **Complex Debugging & Root Cause Analysis:** When debugging a tough issue (e.g., a bug that only appears under certain conditions), sequential thinking helps systematically narrow down the cause. A typical pattern:
  - Thought 1: Restate the observed problem and initial hypotheses (e.g., “The app crashes on input X. Possible causes: A, B, or C.”).
  - Thought 2: Plan investigations (e.g., “I will first check if component A receives the correct value.”).
  - Use appropriate tools (read logs, debug data) to gather info.
  - Thought 3: Analyze results, form next hypothesis.
  - (Potential branch if multiple hypotheses to test in parallel).
  - Thought N: Identify root cause and propose fix, then conclude.

This stepwise elimination of possibilities ensures you don't jump to conclusions. It's especially useful if the bug is due to a combination of factors – you can branch to explore different subsystems. Also, by logging each reasoning step, if you have to later explain the bug or your fix, you have a clear trail. The revision feature is handy here: if a hypothesis was wrong, mark it as a revision (“I assumed it was the DB, but logs show DB is fine – revising that thought”). This approach mirrors how a human would debug systematically, which is far more reliable than an AI guessing in one shot.

- **Architecture & Design Decisions:** High-level design questions (like choosing an architecture or framework for a project) benefit from a structured pro/con analysis. Using sequential thinking, an AI can break down the decision:
  - Thought 1: Enumerate goals and criteria (e.g., scalability, ease of maintenance, etc.).
  - Thought 2: Option A – describe it, assess it against criteria.

- Thought 3: Option B – do the same.
- Thought 4: Perhaps branch if more options or deeper dive needed (Option C, or variant of A or B).
- Thought 5: Synthesize comparison (maybe revision of Thought 1 now that info is gathered).
- Thought 6: Conclude with recommended choice and reasoning, then  
nextThoughtNeeded: false.

The branch capability is useful if two designs are radically different – treat each as a branch path and then come back to main to pick one. This ensures **all viable approaches get fair consideration**. It also documents why an approach was accepted or rejected, which is valuable for future reference (or team communication). The structured output basically becomes a mini design document. This pattern was explicitly identified by users as helpful: “*Architecture decisions with multiple valid approaches*” is a scenario where sequential thinking shines <sup>47</sup>.

- **Code Review & Refactoring:** When reviewing code or planning a refactor, sequential thinking can organize the process:
  - For code review: The AI can iterate through concerns stepwise – e.g., *Thought 1*: check code style/conventions, *Thought 2*: check for potential bugs, *Thought 3*: check performance issues, *Thought 4*: summarize main issues found. Each thought targets a category of issues. This is more effective than trying to hold the whole code review in one monologue. It also allows inserting tool calls (like using a static analysis tool via MCP if available) at specific steps.
  - For refactoring: The AI can create a plan of attack (similar to the case study in Part 1). For instance, *Thought 1*: “We need to refactor module X for better performance. Plan: 1) Identify slow parts, 2) Optimize those, 3) Update related docs, 4) Test.” Then each thought corresponds to these steps, with tool usage in between (read code, apply changes, run tests). This structured approach ensures the refactor doesn’t miss anything and can handle large-scale changes systematically. In one case, an agent did a multi-step refactor using sequential thinking and multiple tools and successfully migrated code without getting lost <sup>58</sup> <sup>62</sup>. So, use ST as the project manager for your refactors.
- **Requirements Analysis & Decomposition:** Often, as developers, we get a broad requirement that needs to be clarified and broken down. An AI can use sequential thinking to simulate that process. Example:
  - Thought 1: Clarify any ambiguous terms in the requirement (maybe even ask the user if interactive, but if not, at least state assumptions).
  - Thought 2: Break the requirement into sub-tasks or components.
  - Thought 3: For each sub-task, identify what needs to be done or what questions need answering. (This could even spawn branches per sub-task if complex).
  - Thought 4: Possibly identify risks or open questions.
  - Final Thought: Summarize the plan or list of derived requirements.

This is essentially performing a mini requirements workshop in the AI’s mind. It ensures nothing is overlooked and that the AI doesn’t rush to a solution without fully understanding the problem. As a bonus, this structured analysis can be presented to a user (if desired) for validation: “Here is how I understand your requirements and how I plan to tackle them, step by step.” That can increase user confidence that the AI isn’t hallucinating features that weren’t asked for.

- **Data Analysis & Pipeline Planning:** If an agent is tasked with a data-related problem (say, analyze some data or build a data pipeline), sequential thinking can outline the stages:
  - Thought 1: Understand the dataset or problem (maybe branch if multiple data sources).

- Thought 2: Decide what transformations or analyses are needed.
- Thought 3: Execute or simulate those steps (here the AI might call database or analysis tools at each step).
- Thought 4: Verify results at each stage (with potential revisions if results are unexpected).
- ... etc.

Particularly in analytics, where there's a natural sequence (extract, transform, load, then analyze), the tool aligns well. It can also help in scenarios of exploratory data analysis by keeping track of what has been investigated and what remains.

- **Knowledge-Based Reasoning (with Memory Tools):** Some users pair Sequential Thinking with a **Memory or knowledge graph MCP** <sup>73</sup>. The idea is the AI can think through a complex question and at each step store or retrieve facts from a persistent memory. For instance, an agent doing a lengthy research report might at Thought 2 save some intermediate findings to a Memory server (so it doesn't lose them), then later retrieve them to formulate conclusions. The sequential steps ensure the process is orderly (first gather data, then analyze, then conclude), and the memory integration ensures that even if the context window is exceeded, important facts remain accessible. A best practice in such cases is to clearly delineate when the AI is storing info vs. when it's reasoning, so as not to mix up the two. (E.g., *Thought 3*: "I will now store these key points in memory for later reference." followed by the memory tool call, then *Thought 4*: continuing analysis using those key points.)
- **Multi-Agent or Multi-Model Scenarios:** In advanced setups, one might have multiple AI "agents" or specialist models collaborating (as hinted by some community approaches <sup>74</sup> <sup>75</sup>). Sequential Thinking can serve as the coordinator. For example, a "planner" model could use sequential thinking to break the task and then delegate subtasks to other models (like a coder model, a tester model, etc.). Each thought could be like: "Ask coder model to implement X... (then that happens)... next, ask tester to verify Y...". The server doesn't directly support multi-agent communication, but by structuring the plan, it implicitly manages the workflow. This is an emerging pattern, and if you design such a system, treat each external action (calling another model or API) as a tool invocation between thoughts. The sequential plan ensures the overall mission is accomplished in a logical sequence. This is similar to the Planner/Doer pattern but extended to multiple AI specialists rather than multiple tools.

The key thread across domains is **structured breakdown**. Whether it's code, design, data, or research, using sequential thinking means the AI will approach it systematically: clarify -> plan -> execute -> verify -> adjust -> conclude. In each domain, that basic cycle can be adapted to the specifics (debugging cycle, design deliberation, etc.). Developers who have employed Sequential Thinking in these areas report more reliable and insightful outcomes, as well as a clearer record of what the AI did (which is great for reviews or audits) <sup>40</sup>.

One user even noted using Sequential Thinking plus a memory server "*before I start any project... saved me on a number of issues I was having*" <sup>73</sup> – indicating that in complex projects, kicking off with a structured thinking session can preempt problems that would arise from a more haphazard approach.

## MCP Integration Patterns

Sequential Thinking often works best not in isolation, but in tandem with other MCP servers. Here are some common integration patterns and tips for orchestrating multi-tool workflows:

- **Filesystem + Sequential Thinking (Development Workflow):** This combo is useful for tasks like code refactoring, generating project scaffolding, or any development that requires reading/writing files. Pattern:
  - Use **Sequential Thinking** to outline the dev task: for example, “Step 1: open file X, Step 2: make change Y, Step 3: run test Z, Step 4: verify output.” <sup>76</sup>
  - After each thought that involves file ops or code, call the **Filesystem server** to do it. E.g., Thought says “I should create a new module file”, then the agent calls `filesystem.writeFile` with the content, then next thought picks up after creation.
  - This ensures the agent doesn’t lose track of which files to modify and in what order. It’s essentially a safe to-do list execution. It also helps avoid the agent hallucinating file contents – it can actually open and inspect them as needed at each step.
- Best practice: after using filesystem actions, have a thought to interpret the result (e.g., after reading a file, the next thought should summarize what was found and plan accordingly). This keeps the chain cohesive.
- **Web/API Fetch + Sequential Thinking (Research Workflow):** If a question requires data from the web or an API, combining Fetch (or a specialized API server) with Sequential Thinking is powerful <sup>77</sup>:
  - The agent can plan out what needs to be fetched: “Step 1: get data from API endpoint, Step 2: analyze it for trend X, Step 3: compare with threshold Y.”
  - Then use **Fetch server** (or any HTTP tool) at the appropriate step to retrieve the data. The next thought then reasons about that data.
  - This prevents the AI from making assumptions – it explicitly pulls facts then reasons. It’s exactly the research process a human analyst might follow. If multiple calls are needed (maybe gather data from two sources), the plan can include each as separate steps.
  - One thing to watch: network calls can fail or be slow. The agent should be prepared to handle errors (maybe a thought “If data not available, reconsider approach”). But generally, sequential thinking plus fetch means the AI won’t skip the crucial step of actually obtaining information.
- **GitHub/Git + Sequential Thinking (Code Analysis & PR Review):** When doing repository-wide analysis or reviewing a pull request:
  - The agent can enumerate tasks: “Step 1: fetch the diff of PR, Step 2: analyze changes for bugs, Step 3: check coding standards, Step 4: output review comments.”
  - It then uses a **Git MCP** to get the diff or file contents in step 1 <sup>76</sup>.
  - Thoughts 2 and 3 reason about those changes (potentially using branch if there are different areas to consider).
  - The final thought (or a couple of thoughts) formulate the review.
  - The structured approach ensures the agent thoroughly covers different review aspects systematically (no forgetting to check tests, for example). And it can intermix tool usage, like calling GitHub API to see test results or issue references if needed, as part of the plan.
  - This pattern results in very comprehensive code reviews that are traceable (each thought explains the rationale behind each comment).

- **Database (SQL) + Sequential Thinking (Data Q/A or ETL):** For instance, using the SQLite MCP or PostgreSQL MCP in read-only mode:
  - The agent can outline a mini-ETL or query plan: “Step 1: inspect schema, Step 2: form query to get relevant data, Step 3: run query, Step 4: analyze results, Step 5: provide answer.”
  - It then calls **SQLite** (or other DB tool) to get schema info, then next thought decides which tables to query, then calls the query, then next thought interprets the result.
  - This pattern avoids blindly querying without context – by including a schema inspection step, the AI knows what to ask. Each subsequent thought refines the query or analysis. It’s much like how a data analyst would progressively dig into a database.
  - If multiple queries are needed (perhaps to gather different pieces of data), the agent can do them in sequence, each preceded by a thought explaining why.
  - The final output is well-supported by the data it fetched along the way. This reduces error because any misunderstanding of the data can be caught in an intermediate step and revised.
- **Memory + Sequential Thinking (Long Context Reasoning):** As mentioned earlier, coupling a memory store (like a knowledge graph or vector memory) with sequential thinking is valuable for long, complex tasks:
  - The pattern is often: Thought process uses sequential thinking, and at logical checkpoints, the agent **stores intermediate results** to Memory. For example, after a heavy analysis step, it might call a memory tool to save the conclusions so far (with an ID or key).
  - Later in the chain, if needed, it retrieves that info (especially if there’s a risk the context was lost or if branching and needing to recall what was found in another branch).
  - Essentially, memory can act as an extension of the chain-of-thought beyond the immediate session. It’s particularly useful if the reasoning spans user sessions or needs to persist (the Sequential Thinking server itself doesn’t persist once turned off). So memory MCP can be used to checkpoint the thinking so it can be resumed or referenced even if the chain is interrupted.
  - One caution: avoid dumping every thought into memory unnecessarily (that’s redundant since the server already keeps them in its log). Instead, store key outcomes or facts that are needed later or globally. This hybrid approach yields a more **persistent and resilient reasoning process**.
- **Orchestration in Multi-MCP Workflows:** If you have many tools at your disposal, you might wonder how to coordinate them. A strong pattern is to let **Sequential Thinking be the top-level orchestrator** that decides which tool to call at each juncture. Essentially, treat it as the “brain” as mentioned, and tools as extensions. For example, a complex task like “Build me a web app” could involve:
  - Thought: “First, I’ll generate a project structure.” (Then calls a template generator tool or filesystem to create files)
  - Thought: “Now populate the main code.” (Calls code generation or writes files)
  - Thought: “Fetch latest library versions.” (Calls Fetch tool to get data from web)
  - Thought: “Write configuration files.” (Filesystem tool again)
  - Thought: “Run the app to test.” (Shell tool)
  - Thought: “Test failed, analyze logs.” (Maybe uses filesystem to read a log file)
  - Thought: “Found bug, fix code.” (Filesystem to fix)
  - Thought: “Run tests again.” (Shell)
  - Thought: “All tests passed now. Done.” (Final thought, nextThoughtNeeded:false)

This illustrates an **end-to-end orchestration** where ST managed the flow between multiple servers (template, fetch, filesystem, shell). The agent basically did high-level planning and low-level execution seamlessly. The benefit is each tool is used optimally and in the right sequence, guided by reasoning. Without ST, the agent might have tried to do too much in one step or forgotten a step. With ST, it had a checklist and followed it.

A tip for implementing this: in your skill prompt or agent logic, you might intersperse a bit of guidance like "After each thought, if an action is required, perform it with the appropriate tool before continuing." This ensures the agent doesn't just say the plan but actually executes along the way.

- **Interleaving Tool Calls Within Thoughts:** Typically, the model stops its thought when calling the tool (that's how tools are invoked in MCP – the thought content becomes the tool arguments). However, conceptually one might consider if the agent could do partial reasoning, then call a tool, then continue the same thought. In current implementations, that's not how it works: each tool call is atomic and returns to the model with the result, then the model generates a new thought. So practically, thoughts and tool calls alternate (which is fine). A pattern to note is: sometimes the AI might produce a thought that actually includes the desire to call a tool next. For example, Thought: "I should search for X to get more info." It then immediately does a Fetch tool call. Some clients might combine this (where the thought content is directly triggering a tool). Others might require the agent to explicitly separate them (first thought explaining plan to search, then actual search). Either way, it's wise to ensure that when a tool is needed, the plan (thought) explicitly mentions it, so it's clear why the agent is calling that tool. This makes the chain easier to follow and justifies each action.

**In summary**, treat the Sequential Thinking server as the **conductor** in an orchestra of MCP tools. It keeps everything in sync and on track. The best patterns involve using it to manage workflow complexity while delegating specialized tasks to other servers. Many guides suggest combinations like "*fetch + filesystem*", "*puppeteer + memory*", "*github + sequential-thinking*" for various use cases <sup>77</sup> <sup>76</sup>. These combinations have emerged because they naturally complement each other: sequential thinking provides the reasoning, the others provide capability. By following the above patterns, you can maximize the synergy between MCP servers and solve very complex tasks that would be hard to do with a single tool or one-off prompt.

Next, we'll look at what the community has learned from real-world usage of Sequential Thinking, which will further reinforce some of these practices and provide additional insight.

## Part 3: Community Intelligence

To truly understand the Sequential Thinking MCP server's impact and quirks, it helps to learn from the community – the developers and users who have employed it in various scenarios. This section compiles real-world insights, success stories, limitations discovered, and tips that aren't obvious from the docs.

### Real-World Usage Insights

**Improved Reliability for Complex Tasks:** Many users have reported that Sequential Thinking "tends to produce more reliable results for multi-stage technical problems with many dependencies." <sup>78</sup> For example, in one Reddit discussion a user noted that building a data pipeline involving multiple transformations was more successful with ST enabled – the agent didn't skip steps and handled the sequence correctly <sup>46</sup>. Similarly, tasks like configuring a cloud environment or solving intricate puzzles have seen better outcomes because the AI can break them down.

A common theme is that **Claude (and other models) become more “deterministic” or consistent when using ST**. Because the reasoning is explicit, two runs of the same prompt are more likely to follow the same chain-of-thought, especially if the tool’s plan is guided. This reproducibility is valuable. One user mentioned reproducible reasoning traces as a benefit, useful for compliance or review <sup>79</sup>. In contrast, without ST, the model might solve it one way one time and a different way another time, with no trace.

**Auditability and Collaboration:** Community members have praised the audit trail that ST provides. In team settings, having the AI’s chain-of-thought logged means a human can later review how a conclusion was reached. If the AI made an error, it’s easier to pinpoint where the reasoning went wrong. One user in an enterprise context highlighted this auditability for compliance – you can show regulators or stakeholders the exact reasoning the AI used (which is otherwise hidden) <sup>79</sup>. This fosters trust and also allows **collaborative problem-solving**: a developer can intervene mid-chain if they spot the AI going astray, because they see the intermediate thoughts.

**Training the AI to Use ST:** A question that came up frequently in the community is “*Will the AI automatically use sequential thinking, or do I need to prompt it?*” The consensus: usually you need to prompt or configure it. Out of the box, many found that Claude (especially Claude Code) won’t call ST on its own for a given query unless the query or system instruction encourages that style. People have tried phrases like “think step-by-step” or explicitly referencing the tool name. In Claude’s case, even asking it directly about the tool, it responded that it might not need it and can do things internally <sup>80</sup> <sup>81</sup>. Thus, if you want the benefits, it’s recommended to **include an instruction to use the sequentialthinking tool for appropriate tasks**. This can be in your skill’s system prompt, or it can be injected when certain conditions are met (like a complexity threshold).

Interestingly, some advanced agents have logic to decide on the fly. For instance, one user mentioned “*when I explicitly ask it to think or when it understands the complexity of the task, it uses this tool... that makes sense.*” <sup>82</sup> This suggests some agents (like Augment Code AI, as in that context) have heuristics to trigger ST. In practice, you can implement something similar: e.g., if the user prompt is long or contains words like “plan”, “step-by-step”, or multiple questions, have the AI use ST.

**Perception of Overhead vs Benefit:** The community seems split at times on the necessity of ST given new model improvements: - Some argue it’s “*wasted overhead, just as Claude said*”, especially with Claude 2 or GPT-4 which can internally plan quite well <sup>83</sup>. They point out that built-in features like the Todo list or extended reasoning mode cover much of the same ground. For example, GitHub Copilot’s agent includes a `thinkingTool` internally, making ST redundant in that environment <sup>49</sup>. - Others maintain “*I still use it and it’s very efficient with Claude*” <sup>84</sup>, especially Claude Instant or earlier versions that might not plan deeply on their own. And even with advanced models, some users *prefer having the explicit steps shown*, if only for their own understanding. One user said: “*I like that with ST you can clearly see steps that lead to a conclusion.*” <sup>85</sup> They later decided to stop using it because the AI wasn’t auto-invoking it mid-task as they expected, but the sentiment of clarity was noted. - There are also pragmatic users who toggle ST on or off depending on the situation: “*I completely disabled it and got better results. When I don’t achieve what I want, I turn it back on and it sometimes does the trick.*” <sup>86</sup>. This suggests treating ST as a tool in the toolbox – use it when needed, bypass it when not. From a skill design perspective, you might similarly allow a fallback: try without ST for simpler queries, and escalate to ST if the first approach fails or if the user explicitly asks for detailed reasoning.

**Limitations Discovered:** Through usage, several limitations of the official Sequential Thinking server have been highlighted: - **No Built-in Tool Recommendations:** Some early documentation or adaptations implied the server might suggest which tool to use next (e.g., recommending a search tool after a thought) <sup>87</sup> <sup>88</sup>. In reality, the **official server does not have any AI logic or tool integration**

**beyond recording thoughts.** Any “intelligent tool selection” has to come from the AI itself. A Redditor pointed out that claims of the server recommending `search_docs` or such were “dreamed up” or referring to a modified version <sup>88</sup>. Indeed, a community fork by spences10 added such features <sup>89</sup>, but the reference implementation doesn’t. So, if you saw references to ST doing more than tracking thoughts, know that those refer to custom servers. The **official ST is simple:** it won’t decide to call other tools or suggest actions – that’s on the agent. - **No Automatic History Recall:** Another limitation is that, as discussed, the server doesn’t feed back previous thoughts. Some assumed it might return the whole chain or keep the model’s context updated with all steps, but it doesn’t. This means if the model’s memory of earlier steps fades (due to token limits), it might lose context. One has to mitigate this either by the model summarizing as it goes, or by using an external memory. This isn’t a flaw per se, but a design choice – the server isn’t meant to be a memory store (though it internally is one, but not exposed). A few users have wished for a feature where the server could output, say, a summary of all thoughts so far after each call, which the model could use to refresh its memory. That’s not implemented, but nothing stops an agent from calling a separate summary tool on the thought history if needed. - **Session Resetting:** People have asked, “How do I reset or start a new thought chain?” Currently, the only way is to restart the server or, if the MCP client supports it, disconnect and reconnect (some clients might treat each conversation as separate, automatically starting a fresh instance). There’s no specific “reset” command in the protocol. This means if in one conversation you finish a task and want to do another unrelated one, you might not want to reuse the same ST server instance (as it would continue appending to history). In practice, tools like Claude Desktop may handle this by instantiating a fresh server per session or giving each conversation thread its own state. If not, a workaround is to stop and start the ST server. Some community members have noted this inconvenience – ideally one would have a command to clear history, but it’s not there yet. - **No Concurrency / Multi-user Awareness:** The reference server isn’t built to handle multiple concurrent requests or separate user contexts. It’s a single pipeline of thoughts. If, theoretically, two agents pointed at one ST instance at once, their thoughts would intermix in the history, which would be disastrous. So each agent should have its own ST process (or a multiplexing MCP client that keeps contexts separate). This is usually handled by the MCP framework itself, but it’s a noted limitation: **one ST server = one chain-of-thought context**. Community enterprise users discussing scaling have pointed out that for using ST in a production setting with many parallel tasks, you’d need to scale out the server (run many copies) or wait for future MCP protocol features that handle multi-session (the roadmap hints at statelessness improvements to aid scalability <sup>90</sup>). - **Token Overhead Confirmed:** Users have quantitatively observed the token overhead. For example, one person mentioned burning through \$40 of Claude API in a day by forcing very detailed step-by-step prompting (with or without ST) <sup>91</sup>. The chain-of-thought was “horrendous dog shit” in their words – meaning it under-delivered for the cost. This was perhaps more about the prompt strategy than ST itself, but it underscores that forcing an AI to articulate every step can be expensive. Another example: some found that Augment Code’s “extended thinking” mode might achieve a similar outcome at lower cost than using the generic ST MCP, because it’s optimized in the host model. This kind of feedback led some to turn ST off by default and only enable for truly knotty problems.

**Success Stories & Creative Uses:** On the positive side, community members have shared scenarios where ST was a game-changer: - A user on Hacker News (or similar forum) described using ST to have Claude **reason through a novel puzzle** – by breaking down the puzzle’s rules stepwise and systematically exploring solutions, something the raw model struggled with in a single go. The structured approach got to a correct solution whereas before it was giving partial answers. - Some have integrated ST into custom agent frameworks (like one called *TachiBot-MCP* described in a blog <sup>92</sup>). They used ST as part of a larger “team of AIs” approach – e.g., a coordinator AI uses ST to manage specialist AIs. They reported this significantly improved the outcomes for complex tasks (like designing a rate limiter in the example, where ST helped catch a flaw that a single-pass solution would have missed <sup>93</sup>). - In educational or explainer contexts, ST has been used to have an AI *show its work*. One Reddit user mentioned enabling ST in brainstorming mode so Claude would produce a chain-of-thought visible to

them while solving a problem. This can be used as a learning tool or just to increase confidence in the answer (similar to how a math student shows steps). The ability to branch in reasoning was highlighted as uniquely valuable for exploring “*What-if*” scenarios in planning – e.g., “Plan A vs Plan B” in parallel, which few other AI tools allow so explicitly. - **Combining ST with Knowledge Graph Memory** was mentioned (as we saw) to kick off projects. One creative use: before writing any code for a new project, the user had Claude run sequential thinking to outline everything and store the plan in a memory graph, acting as a living design spec. Throughout development, the AI could query that memory (which is basically the saved chain-of-thought) to stay on track. This is like a dynamic project plan that the AI refers to – a novel workflow enabled by MCP tool chaining. - Some have tried ST in non-coding domains, like writing an essay or story. While not its main intent, an AI could sequentially outline an essay (intro, arguments, conclusion) using the tool, ensuring it covers all points. Or for a story, list plot points then flesh them out. The benefit is a more organized output. However, anecdotal feedback suggests that for creative writing, the rigid structure might hinder more than help (since creative writing benefits from some spontaneity). Still, it’s been experimented with – showing the versatility of ST beyond coding tasks.

**Common Questions & Answers:** - Q: “Who actually does the thinking in Sequential Thinking? The server or the model?”

**A:** The model does all the thinking. The server just records it. This was emphasized often [2](#) [94](#). It’s a common misunderstanding that ST itself might add some logic. But no, it’s just a structured notepad that the AI writes to. So you still need a capable model behind it – ST doesn’t compensate for a weak model’s reasoning ability, it only channels a good model’s reasoning. - Q: “Does ST guarantee a better answer?”

**A:** Not necessarily. It provides a better process. A flawed chain-of-thought can still lead to a wrong answer – but you’ll at least be able to trace why. The quality of final answer still depends on the model’s intelligence and the correctness of each step. That said, by encouraging self-checks and structure, it often leads to better answers, just not guaranteed. If the model lacks knowledge about the domain of the question, ST won’t magically give it that knowledge (that’s where retrieval tools or human input come in). - Q: “How do I get the AI to stop using ST once it’s not needed?”

**A:** Since we often explicitly tell it to use ST, you also need to manage turning it off. Usually, the agent stops calling it when it sets `nextThoughtNeeded:false` and produces the final answer. After that, the chain is done. If a new user query comes that’s simple, you wouldn’t instruct ST usage again. Some advanced agents might decide: “the last query was complex, we used ST; this new query is trivial, we won’t.” Implementing that logic can be part of the skill: detect complexity. Alternatively, one can always run ST in a passive way – the AI might just do one thought with `nextThoughtNeeded false` (effectively a single-step chain) for simple things, which is not harmful aside from a tiny overhead. - Q: “Is ST MCP outdated by newer features (e.g., Claude’s internal planning) or new models like GPT-4.5?”

**A:** It’s true that many advancements in LLMs incorporate better internal planning. Some Reddit threads discussed that “*Sonnet 4.5 extended thinking*” (a Claude model version) and the like made ST less obviously useful [95](#). However, ST remains relevant in contexts where you want standardization across models or you want that external trace. One expert said: “*It’s now redundant in Copilot*” [49](#) but others pointed out they still use it in Claude and find value. So the answer is: it depends on the environment. If you are exclusively using a single model that has its own chain-of-thought mechanism, ST might be optional. But if you are building a general agent that might use different models or want consistent behavior, ST is still very useful. Also, enterprise adoption of MCP suggests keeping logic out of the black-box model and in a transparent layer is preferred for control reasons [51](#). - Q: “What’s the difference between the Anthropic ST server and forks or others I’ve seen (Zalab, Arben, etc.)?”

**A:** The official one (Anthropic’s reference) is minimal and demonstration-focused [96](#). Others, like Zalab Inc.’s version or Arben Ademi’s, often add features: some add “cognitive stages”, structured formats, persistence of thoughts to disk, or integration with certain workflows [97](#). For instance, one fork may automatically categorize thoughts (like goal, question, solution) or store them. These are community experiments. When using the *public Claude Code skill*, you’ll likely be dealing with the official server

unless specified. So the behaviors described here apply to that. It's good to be aware, though, that if you see references to ST doing something beyond what's documented, it might be referring to those enhanced versions. For our manual's purpose, we focus on the official implementation (which is also what's on NPM and Docker as [@modelcontextprotocol/server-sequential-thinking](#)).

**Frustrations and Workarounds:** A few frustrations have been voiced and corresponding workarounds found: - Frustration: "*The AI sometimes doesn't call ST even after I installed it.*" - Workaround: provide a more explicit nudge in the prompt (mention the tool by name, or use a trigger word like "step-by-step"). Also ensure in the MCP config it's actually enabled (some forgot to flip the switch in UI) <sup>98</sup>. - Frustration: "*Claude's built-in Chain-of-Thought seems fine; ST just duplicates it.*" - Response from community: The benefit is in the explicit tracking and cross-model ability. Also, some found that while Claude can chain-of-thought internally, it can also get distracted or compress its reasoning too much. ST forces a certain rigor. However, it's valid that for straightforward coding tasks, you might not see a difference except slower output. The workaround is basically using ST only for certain types of prompts and not for others (which we've covered how to choose). - Frustration: "*ST made my agent too slow, it waits for each thought to return.*" - Indeed, each thought is a round-trip. One user in a realtime setting (maybe a live coding assistant) found it sluggish. They solved this by limiting the number of thoughts (maybe set an expectation to finish in <=3 thoughts) or by upgrading to a faster model for the reasoning part. Another idea is parallelizing branches if possible (though not trivial; would require multiple instances and is not how MCP normally works). Generally, if latency is an issue, consider whether you truly need the intermediate confirmations for every step – in some cases, you might trust the model to do 2-3 steps internally and only call ST at milestones. But that requires custom prompt hacking (like the model writing multiple "thoughts" in one tool call – which it could do by setting thoughtNumber and totalThoughts accordingly, though the server would treat it as one giant thought string). This is an advanced tweak not commonly done. - Frustration: "*It sometimes reinforces wrong assumptions*" - As mentioned, if the AI doesn't realize it's wrong, it might produce a very convincing but wrong multi-step answer. The fix is to try to instill self-checking. One user wished for a more "open-minded thinking" where the AI would doubt itself more rather than increasing confidence in a wrong path <sup>52</sup>. A tip here: include a prompt hint for the AI to consider alternatives if things aren't adding up. E.g., "If at any point you are unsure or the solution doesn't look correct, branch or revise rather than continuing the same line." This encourages second-guessing when appropriate, which can catch errors. It leverages the branching feature more proactively.

**Community Wishlist:** A few features people said they'd love to see (which indicate limitations now): - A way to **persist thought history** between sessions or to file (some want to log it for record-keeping beyond just console output). - **Better integration with UI:** e.g., showing the thought boxes in the chat interface nicely. (Claude's UI doesn't natively display them, though VS Code might show MCP server logs in output panel; some users probably copy them from console to see the boxes). Having the chain-of-thought shown to the user optionally was floated. - **Automated branch evaluation:** one idea was if the server could help choose the best branch, but that again is a bit beyond its scope as a reference. That might be left to the model or a higher-level orchestrator. - **Easier resetting/parallel sessions:** possibly coming with future MCP updates (stateless mode) <sup>99</sup>. - **More semantic structure:** e.g., tag thoughts as "hypothesis", "analysis", "counterexample" etc. This is more of a UI/interpretation layer wish; the data is there to do it if someone builds a tool on top.

Overall, the community experience suggests that when used in the right context, Sequential Thinking MCP is a beloved tool for adding rigor to AI reasoning. When misapplied, it can be frustrating or superfluous. The next part will synthesize these lessons and the technical details into concrete guidance for developing a skill (like a Claude Code skill) that uses Sequential Thinking optimally.

## Part 4: Synthesis for Skill Development

Bringing together the technical knowledge and practical insights, this section provides actionable guidance to develop a Claude Code skill (or any AI agent logic) that effectively leverages the Sequential Thinking server. We'll outline when to invoke the tool, how to formulate calls, how to manage the chain, and how to integrate it into complex workflows, ensuring the agent makes intelligent decisions about its usage.

### Decision Framework for Using Sequential Thinking

The first step is deciding **whether and when** the Sequential Thinking tool should be invoked. Hard-coding the agent to use it every time is not ideal; instead, design a decision framework based on the task at hand. Key criteria to consider:

- **Problem Complexity:** As repeatedly noted, only engage ST for problems that require multi-step reasoning or planning. A simple check can be the estimated difficulty or number of steps. For instance, if you can identify sub-tasks in the user's query (e.g. "Do X, then Y, then Z"), that's a sign to use ST. If the query is just one clear question, handle it directly. Some skill implementations use prompt length or keywords as proxies for complexity. For example, if the user's prompt is long and contains multiple sentences or asks, it might imply a more complex task – good trigger for ST. On the other hand, a short question likely doesn't need ST (unless it's a tricky one that hides complexity).
- **Need for Auditability:** If the use case demands that the reasoning process be documented (for debugging, compliance, learning, etc.), lean towards using ST. This might be a user preference (some users might explicitly ask the assistant to "show your thought process" – in which case ST is perfect) or a system requirement. If auditability is critical, even slightly complex tasks could use ST just to have that trace.
- **Cross-Tool Planning:** If the solution will obviously require multiple tools or data sources, ST can serve as the glue. For example, if the user asks for something that involves reading a file and then doing a web search and then writing output, using ST to orchestrate those multiple steps ensures nothing is missed. In contrast, if only one tool is needed one time, ST might be overkill.
- **Portability & Consistency:** If your skill is intended to work across different models or contexts, and you want consistent behavior, you might favor ST. As Patiencing on Reddit pointed out, "*Do you need to reuse the same planner across multiple clients or models?*" If yes, ST provides that uniform planner <sup>100</sup>. For example, a skill that might run on Claude one day and GPT-4 the next – using ST ensures both follow the same structured approach rather than relying on their internal differences.
- **Long-Term Planning / Sessions:** Is the task something that will persist or evolve over a long session or multiple sessions (like an ongoing project)? If the plan or reasoning needs to be carried forward, ST can help maintain continuity. If it's a one-off short query, internal reasoning might suffice. Patiencing also mentioned: "*Is the plan long-lived across sessions?*" If yes, ST might be beneficial <sup>101</sup> – one could even save the thought log and resume later.
- **User or Developer Instruction:** If the user explicitly requests a step-by-step solution or if the developer debugging the agent decides they want to see the reasoning, that's a clear flag to use ST. For instance, some advanced users might know about ST and say, "Use the sequential thinking tool to explain your answer." Your skill should honor that.
- **Time/Quality Trade-off:** Consider the context in which the agent operates. If it's an interactive setting with a user waiting for an answer, too many tool calls might slow things down or seem verbose. If it's a backend automation where quality is more important than speed, the overhead is fine. Essentially, gauge the tolerance for a slower, thorough approach versus a faster, heuristic one. If correctness and completeness are paramount (e.g., code generation where a bug could

cause significant issues), lean toward ST; if a quick guess is acceptable (like a casual question where stakes are low), maybe skip ST for brevity.

One way to implement this decision process in a skill is to have a preliminary analysis of the user prompt using the model itself or some heuristics, then set a flag to enable ST. For example, a pseudo-code:

```
if is_complex(user_prompt):
    use_sequential_thinking = True
else:
    use_sequential_thinking = False
```

Where `is_complex` could check length, presence of conjunctions ("and", "then"), question complexity, etc., or even call the LLM with "Do you need multi-step reasoning for this? (yes/no)".

**Checklist (from community, by patiencing):** The Reddit comment provided a nice checklist which can be turned into yes/no questions for the agent: - *Portability needed?* *Auditability needed?* *Multi-tool workflow?* *Team policy constraints?* *Long-lived plan?* <sup>102</sup>. If majority yes, use ST. This encapsulates a lot of the above.

Finally, if uncertain, a conservative approach is: try solving without ST first, and if the model's initial attempt fails or is incomplete, then switch to ST. Some have done this manually – they get an answer, find it's not good, then say "Okay, let's break this down step by step." You can automate that: check the model's direct answer, if it's unsatisfactory (maybe using evaluation heuristics or user feedback), then invoke ST for a second attempt.

## Invocation Templates & Examples

When it's time to use Sequential Thinking, how to start the chain matters a lot. The first thought often sets the tone for the whole sequence. Here are some patterns and example templates for initiating and guiding thought chains in various scenarios:

### 1. Problem Restatement & Plan Outline (General Template):

For many problems, a strong opening is to have the AI restate or clarify the problem and outline a high-level plan. This ensures understanding and provides a roadmap for subsequent thoughts. For example:

```
Tool: sequentialthinking
Input: {
    "thought": "The question is asking for a strategy to optimize the database query. I'll break this down:
    1) Understand the current query and its bottlenecks,
    2) Consider indexing or query refactoring options,
    3) Test the performance after changes,
    4) Choose the best optimization.
    Thus, I estimate 4 thoughts for this.",
    "nextThoughtNeeded": true,
    "thoughtNumber": 1,
```

```
        "totalThoughts": 4  
    }
```

In this **first thought**, the AI clarifies the task ("optimize the database query") and enumerates steps. This provides a clear plan. The `totalThoughts:4` corresponds to the number of steps outlined. This template works for planning tasks, troubleshooting steps, or any scenario where an outline is beneficial. The actual text in `thought` can be less explicitly numbered (it could be a paragraph), but often numbering or bullet-style in the thought content helps structure the plan (even though it's in a JSON string, it's for the AI's own benefit and any developer reading logs).

## 2. Asking Questions / Identifying Unknowns:

If the prompt is ambiguous or missing info, the sequential approach can start by raising those questions (which either the AI can research or ask back to the user). For example:

```
Tool: sequentialthinking  
Input: {  
    "thought": "First, I need to clarify some details that are not specified:  
    What is the size of the data? What is the performance goal?  
    I will assume a large data set and aim to reduce query time by at least 50%.  
    Given that, I'll proceed to identify the query bottleneck as the next step.",  
    "nextThoughtNeeded": true,  
    "thoughtNumber": 1,  
    "totalThoughts": 5  
}
```

Here, the AI uses Thought 1 to surface assumptions or questions. In an interactive setting, the user could answer. If not, the AI made assumptions and moved on. This template is useful for **requirements analysis or disambiguation** tasks, where step 1 is always "clarify." It sets up context for subsequent thoughts. Notice it still set `totalThoughts:5` (assuming after clarification, maybe 4 more steps to solve problem).

## 3. Branch Initiation Example:

Suppose by Thought 3, the AI encounters two possible paths. How to invoke a branch? Example:

```
Tool: sequentialthinking  
Input: {  
    "thought": "Approach A: Use a composite index on the columns. Approach B:  
    Denormalize the table.  
    I will explore Approach A first.",  
    "nextThoughtNeeded": true,  
    "thoughtNumber": 3,  
    "totalThoughts": 6,  
    "branchFromThought": 2,  
    "branchId": "A"  
}
```

This indicates that at Thought 3, the AI has branched from Thought 2 along path "A". It plans to also explore "B" presumably. After one or two thoughts in branch A, it might then do:

```

Tool: sequentiaalthinking
Input: {
  "thought": "Now I will explore Approach B for comparison.",
  "nextThoughtNeeded": true,
  "thoughtNumber": 5,
  "totalThoughts": 6,
  "branchFromThought": 2,
  "branchId": "B"
}

```

Here Thought 5 is the start of branch B (both branch A and B originate from thought 2 in this scenario). The numbering continues globally (so we had thought 3 & 4 in branch A, and now thought 5 starts branch B). The server's `branches` list would contain ["A","B"] by now <sup>33</sup>. The agent should keep track of which thought numbers correspond to which branch to not mix them up.

#### 4. Revision Invocation Example:

Let's say at Thought 6, after exploring A and B, the AI realizes its initial plan (maybe thought 1's assumption) was flawed. To revise:

```

Tool: sequentiaalthinking
Input: {
  "thought": "I realize the earlier assumption about the data size was wrong.
Revising thought 1: The dataset is actually small, so some of the optimization steps might be overkill.
I'll adjust the plan to focus more on query logic than on indexing overhead.",
  "nextThoughtNeeded": true,
  "thoughtNumber": 6,
  "totalThoughts": 6,
  "isRevision": true,
  "revisesThought": 1
}

```

The content explicitly says it's revising thought 1, and `isRevision:true, revisesThought:1` mark it. The prefix "I" is optional (the server will add one in its log formatting automatically <sup>103</sup>, but including it in the `thought` text can be a stylistic choice by the AI). After this, the AI might continue with thought 7 as a new step reflecting the adjusted plan.

#### 5. Concluding Thought Example:

The final thought should present the answer or conclusion, with `nextThoughtNeeded:false`. For instance:

```

Tool: sequentiaalthinking
Input: {
  "thought": "In conclusion, the best optimization is to add a composite index on (user_id, created_at), which reduced query time from 10s to 1s in testing."
}

```

```

    The query was the main bottleneck, and this index addresses it effectively.",
    "nextThoughtNeeded": false,
    "thoughtNumber": 6,
    "totalThoughts": 6
}

```

This would be the termination of the chain. The agent might then output this conclusion to the user (depending on how you design the skill – possibly it might strip away the JSON and just give a nice answer, since the user doesn't need the JSON part).

**Important:** The actual format you'll use in code might not literally be `Tool: sequentialthinking ...` as above; in implementation, you might call a function to send the JSON via MCP client. The examples are illustrative for how the JSON payload should look for each scenario. Ensure the keys match exactly (`thoughtNumber` not `thought_number`, etc., per schema <sup>104</sup>).

Also, note that `thought` content can include newlines, lists, etc., since it's just a string. The server doesn't mind, as long as it's valid JSON. However, the AI should be careful to properly escape any special JSON characters if it includes them (like quotes inside the thought string). Typically, the model will handle that if it's following a format template. Some SDKs might let you pass a structured object rather than raw JSON string, which is safer.

**Guiding the First Thought:** In the skill prompt, you might provide a template or few-shot example to Claude of how to format the first thought. For instance: - System message could include: "When using `sequentialthinking`, the first call should outline your plan in JSON with `thoughtNumber` 1 and an estimated `totalThoughts`." Possibly show a dummy example in the prompt. This can prime the model to produce good first call. - Or after the model decides to use ST, you intercept and format it appropriately. Some devs find it easier to have the model output a plan in natural language and then have some code convert it to the JSON call. But if we assume the model can do it, giving it a direct example is best.

## Chain Management Guidelines

Once the chain is underway, the agent must manage it: deciding what each next thought should be, when to branch or revise, and how to integrate external info. Here are guidelines to maintain an effective chain:

**Thought Progression Strategy:** Generally, follow the plan or outline from the first thought, but be ready to deviate if needed. Each thought should logically follow from the last. A good mental model is: - At the end of each thought, ask: **Did I reach a sub-goal? What's the next sub-goal?** Then formulate the next thought accordingly. - If you planned N thoughts but realize at thought k that you need more detail, don't hesitate to expand. You can either refine what "`totalThoughts`" means (increase it) or just let `thoughtNumber` exceed it and update total as the server does automatically <sup>8</sup>. - Keep thoughts **focused**. If a thought starts becoming too long or tackling multiple things, consider splitting that into multiple thoughts or using branching for the separate concerns. - Use the `needsMoreThoughts` flag optionally as a signal: for example, at what was originally the final thought, if the AI isn't fully satisfied, set `needsMoreThoughts:true`. This is somewhat redundant with just increasing `totalThoughts` and continuing, but it's an explicit way to say "I'm extending the plan". It might be useful for clarity (for a developer reading the JSON). The server doesn't do anything special with it except include it in output.

**When and How to Branch:** Branching should be used deliberately at decision points: - Trigger a branch **only after you have enough information to justify exploring an alternative**. Usually, that means by

the time you finish a thought, you either have uncertainty or a clear divergence in options. That is the moment to branch. - Use meaningful branch IDs. They can be simple ("A", "B") or descriptive ("index\_approach", "denormalize\_approach"). The server will accept any string <sup>32</sup>. Descriptive IDs can help if the AI refers to them later, but they also might consume tokens. Often single letters or numbers are fine, given the agent itself keeps track contextually ("Approach A/B"). - Manage branch depth: Ideally branch from a mainline thought, rather than branching off a branch (branching off a branch is effectively depth 2). Depth 2 might be manageable but beyond that it's really hard to handle. If you find yourself wanting to branch inside a branch, perhaps reframe it as a new mainline approach or revise something instead. - After exploring branches, **merge results**. The AI should come back to a single line of reasoning. Typically, you do branch A (with one or several thoughts), then branch B, then a thought that compares or chooses. This final comparison thought would have no branchId (meaning it's back on mainline), or you could consider it implicitly revising the earlier decision point. - If one branch clearly fails or is inferior, you can terminate it early by setting `nextThoughtNeeded: false` on that branch's last thought, *but* that ends the entire chain if you do it (since the server doesn't handle branches independently in terms of stopping – a `false` will likely signal the agent to wrap up overall). A better approach is to keep `nextThoughtNeeded: true` but note that you'll switch branch. E.g., "Conclusion of Branch A: it doesn't solve the problem. I will now try Branch B" – and then proceed to B in the next thought. Essentially, treat the branch exploration as part of the single overall chain, not separate processes that can individually halt. Only the final answer to the user should produce `false`.

**When and How to Revise:** Revisions are for correcting or improving prior reasoning: - Use them when you catch an error or oversight in a previous thought (not just to restate – if nothing was wrong, better to continue forward). - Mark `isRevision: true` and point `revisesThought` to the exact number of the thought you're fixing <sup>27</sup>. In the revision thought text, explicitly mention what you're changing. This clarity is helpful for anyone reading the chain and also for the AI itself to stay consistent. - After a revision, you have two paths: - **Overwrite approach:** Essentially ignore the content of the original thought in future reasoning and use the revised one instead. For the AI, this is just a mental shift – the server won't remove the old thought from history. The AI can treat the old thought as "replaced". If the thought being revised had subsequent dependent thoughts, you might need to branch or adjust those too. - **Branch as alternative:** Alternatively, you could treat the revision as a branch scenario: the original timeline remains but you spin off a new approach from the revised point. However, since a revision semantically means the original had a mistake, it's usually best to consider the original path deprecated. So typically, you won't continue the original sequence after a revision; you'll proceed along the corrected line. - Don't overuse revisions. If every thought becomes a revision of the previous, then something's off – the AI might be second-guessing too much or the plan was poor. Revisions should be relatively infrequent, and ideally, significant (fixing a real issue, not minor tweaks). For minor adjustments, sometimes just explaining in the next thought without formal revision is fine. Save revision markers for notable course corrections. - Example scenario: In coding, you wrote an algorithm in Thought 4, and in Thought 5 testing it found a bug. Thought 6 revision of 4 is warranted. Then Thought 7 can be re-testing the new algorithm. That's a clean use of revision.

**Using External Info & Tools in Chain:** When interleaving other tools: - Always let a thought **inform the next action**. For instance, don't call a tool arbitrarily; have a thought concluding "I should use X tool now to get Y." Then do it. This way the chain-of-thought documents why an action was taken. This also often prevents wasted tool calls because the AI reasoned about the necessity first. - After getting a result from a tool, the very next thought should interpret that result. The pattern should look like: *Thought -> tool call -> Thought analyzing tool output*. This ensures the info fetched is integrated into the reasoning. If a tool's output is large (e.g., reading a big file), the AI might consider summarizing it in the thought to keep context size down (unless it needs details). - If a tool result triggers a change in plan (like discovering an error, or new data that suggest a different approach), that is a good time for a revision or branch. For example, if you fetch from an API and find the data is different than assumed,

maybe revise an earlier assumption. - The chain-of-thought should be flexible to incorporate new information. This is a big advantage: rather than generating a static answer that might ignore new info, ST allows the agent to dynamically adjust when tools return surprising results.

**Monitoring Chain Health (Quality Heuristics):** As the chain progresses, the agent (and/or the skill developer) should monitor whether it's productive: - Signs of a *productive chain*: each thought brings the solution closer, uncertainties are resolved over time, no repetitive loops, the branches (if any) are converging, and the final answer is forming with confidence. The content of thoughts should become more concrete as you go (starting maybe abstract, ending very specific). - Signs of an *unproductive chain*: thoughts start to become circular (repeating earlier points without progress), the AI seems stuck ("I'm still not sure... maybe do this, maybe that" back-and-forth), or the chain is growing very long without reaching a conclusion. If 10+ thoughts in, and it's still as confused as at start, something's wrong. - **Self-correction indicators:** The model should use revisions or branch switching as indicators of self-correction. If it never uses them at all, either it never encountered any issue (which is unlikely in complex tasks) or it isn't recognizing mistakes. Encourage some reflection at certain intervals. For instance, after a few thoughts, it might be good to have a thought where the AI pauses to recap: "So far, I have done X and Y. Is my approach working? If not, I should revise." This meta-reasoning thought can be built-in as a step (like after a branch or before final). - If the agent gets stuck (detected via looping or repeated statements), a strategy is to introduce a **new angle**: maybe branch to a radically different approach, or call a different tool for insight (like ask a knowledge base if available, or attempt a simpler heuristic to see if it yields clues). Essentially, break the loop by injecting new information or perspective. The agent could even escalate: e.g., if it can't solve a coding issue after many thoughts, perhaps it should call a search tool to see if anyone has solved something similar. That can be integrated: *Thought: "I'm stuck, let's search documentation."* -> *use fetch* -> *next thought with new info*. This can often snap it out of a dead-end. - As a developer, if you see your skill's output chains often going awry, you might refine the system instructions or provide more examples. For instance, show an example of a proper short chain vs a runaway chain to the model, and instruct it to avoid the latter.

## Termination Criteria and Final Output

Knowing when and how to conclude the sequential thinking process is critical. Termination should be graceful and accompanied by a synthesis of the solution.

**When to Conclude the Chain:** - The obvious case: when the problem posed is solved or the question answered to satisfaction. The model should set `nextThoughtNeeded:false` on that final thought. A guideline to embed (and is embedded in the server's description) is: "*Only set next\_thought\_needed to false when truly done and a satisfactory answer is reached.*"<sup>36</sup>. This means the agent should double-check: *Is this answer correct and complete?* If yes, conclude. If not, continue reasoning. - If the chain reaches the initially estimated `totalThoughts` and the AI thinks it's done, conclude. If it reaches `totalThoughts` but is not done, it must extend (increase total or just keep going with `needsMoreThoughts:true`). So hitting the planned number is not a reason to stop by itself – only stop if done or out of ideas/paths. - In some rare cases, the AI might determine the problem is unsolvable or out of scope. In those cases, concluding the chain early with a statement of inability might be appropriate (e.g., "I cannot solve this with the information given." and `nextThoughtNeeded:false`). However, usually there's something the AI can provide, so this should be rare. - Avoid stopping just because of diminishing returns if the goal isn't met. Instead, as mentioned, if it's looping without progress, try a different tactic or tool. But if truly nothing works, a termination with best-effort answer may be necessary (perhaps accompanied by an apology or note of uncertainty in user-facing output).

**Synthesizing Final Output from the Chain:** - The final thought should ideally contain the answer or result in a clear form, as if you were directly answering the user. Often, you will design the agent to take the content of the final thought and present it (maybe with some formatting) as the answer. Alternatively, the agent might generate a separate final answer after concluding the chain, using the chain-of-thought as internal context. Different frameworks handle it differently: - Some might treat the chain-of-thought entirely hidden, and only after finishing does the model produce the final answer outside the tool. In such a case, the model might set `nextThoughtNeeded:false` on the last tool call and then the next model output is the answer. If doing that, ensure the final thought encapsulates what should be in the answer, or at least all pieces are known by then. - Other times, the final thought content itself can be the answer to display. For instance, if the final thought says "Therefore, the answer is 42 because ...", one could present that. You might trim the "Therefore" and just give the answer. Up to your skill design. - It's good practice for the final thought to **summarize the reasoning briefly while giving the answer**. This way, the user gets context if needed. But be mindful: if you don't want to show the entire reasoning to the user (which usually you do not, unless user asked), you might have the AI give a concise answer. The chain-of-thought logs are for you/the agent, not necessarily the end user. Often the final answer can be more polished than the raw thoughts. For instance, an agent might think in steps messy, but then present a clean solution. That's fine as long as the information is consistent. - Ensure that if the chain uncovered any caveats or conditions (like "this approach works only if X is true"), the final answer mentions them. The worst outcome would be the chain finds an issue but the final answer ignores it. This could happen if, say, the AI branched into two solutions, found one had a flaw, so concluded the other is better – the final answer should state the solution and possibly note "(the alternative was considered but had flaw Y)." This provides completeness and also demonstrates the thorough reasoning (which can impress users). - If the sequence was used for something like code generation, the final output might be the code or the result of the plan's execution. In such cases, the final thought might literally include the final code or a statement like "All steps completed successfully." Then the skill might retrieve the artifacts (like the modified files) to present or save. Consider how the outcomes of the sequential process are delivered.

**Graceful Exit Strategies:** - If aborted early or if something went wrong (e.g., a tool failed in middle and the agent can't continue), the chain should handle it gracefully. Perhaps produce a final thought explaining the issue ("I couldn't complete the plan because the data was unavailable.") with `nextThoughtNeeded:false`. The skill can then propagate that to the user suitably. The ST server itself won't know if a tool fails (unless the AI informs it by altering its plan), so it's on the agent to decide to conclude or try alternate steps. - In interactive use, sometimes the user might interrupt or ask for the result when the chain is not finished. In that case, the agent should either quickly wrap up in the next thought or explain it needs more steps. As a developer, you might set a max number of thoughts for safety so it doesn't annoy users by going on and on. If that max is hit without conclusion, you can either force a termination or ask the user if they want the reasoning to continue. Usually, though, if using ST, the user expects some wait for a thorough answer, so it's okay if it takes a few iterations.

**Final Check Before Presenting Answer:** - After concluding, the agent can do a quick self-review: Did the chain fully address the original query? Are all parts answered? If anything is missing, ideally it would have caught that before finalizing. But if in doubt, maybe add a concluding note like "If more information is needed, I can elaborate further." (However, that might invite more conversation, which could be fine.) - Since ST encourages correctness, by the final step the agent should have verified its solution (especially in math or code: maybe it ran tests, etc.). If not, that's a missed opportunity: the chain-of-thought allows verification steps. A quality heuristic: never finalize without at least mentally verifying consistency. For instance, if the final step in reasoning was to "Provide the answer," maybe the second-to-last could be "Double-check all previous steps are consistent." – This is optional but can catch last-minute slip-ups. - **Token cleanup:** If earlier thoughts included a lot of detail that isn't needed in the final answer, the final answer should not regurgitate all that unless asked. Summarize. The user likely

doesn't want the step log, just the outcome (unless they explicitly want the reasoning too). So instruct the model to keep the final answer user-friendly.

In the skill design, you might literally have:

```
if sequential_thinking_used:  
    answer = final_thought_text  
    (maybe strip any structure)  
    return answer
```

So ensuring that `final_thought_text` is well-formed as an answer is key.

## Integration Playbook for Multi-Tool Workflows

To wrap up, here's a concise **playbook** for using Sequential Thinking in tandem with other MCP servers in common development workflows (as might be implemented in a Claude Code skill):

### 1. Complex Code Refactoring (Filesystem + Git + SequentialThinking):

Use ST to plan the refactor, step by step: - Thought 1: Outline refactor plan (files to change, steps to verify). - Thought 2: (If using Git) perhaps create a new branch via Git tool for safety. - Thought 3: For each file, call Filesystem to open it, next thought decide modifications, then call Filesystem to apply changes. - Thought ...: After changes, call Shell or appropriate tool to run tests. - Final Thought: Summarize outcome (e.g., "Refactor complete and tests passed.") 105 106.

### 2. Debugging a Failing Test (Filesystem/Runtime + SequentialThinking):

- Thought 1: Read the test error message or failing test file (call Filesystem to open log or test). - Thought 2: Analyze the failure (in thought). - Thought 3: Form hypothesis and decide to inspect a specific part of code (Filesystem open that file). - Thought 4: Identify the bug, decide on fix. - Thought 5: Apply fix (Filesystem write). - Thought 6: Run tests again (Shell exec). - If tests still fail, branch or revise the hypothesis. - Final: either "All tests passed, bug fixed" or report inability.

### 3. Answering a Question with Documentation (Fetch/Browser + SequentialThinking):

- Thought 1: Plan: figure out what needs to be looked up (e.g., "I will search the API docs for X"). - Thought 2: Use Fetch tool to retrieve relevant doc page. - Thought 3: Read result (it will be in the tool output), summarize key info. - Thought 4: If more info needed, perhaps branch out to another search. - Final Thought: Provide answer citing the found info. Possibly include references if required (the chain-of-thought could store the URL or content needed).

### 4. Multi-Step Data Processing (DB + Filesystem + SequentialThinking):

- Thought 1: Outline pipeline (e.g., "1) Query database for data, 2) Process data, 3) Save results to file"). - Thought 2: Formulate SQL query criteria. - Thought 3: Execute query via DB tool. - Thought 4: Process data in memory (the AI can do this if small, or call a transformation tool if available). - Thought 5: Save results via Filesystem tool. - Thought 6: Verify file saved or perhaps open it to confirm. - Final Thought: Confirm pipeline done, maybe output summary of results (like number of records processed, etc.).

### 5. Research & Write Report (Memory + Fetch + SequentialThinking):

- Thought 1: Plan sections of the report and info needed. - Thought 2: For section 1, decide to fetch some data (use Fetch). - Thought 3: Summarize fetched data (maybe store in Memory for later reference). - Thought 4: For section 2, fetch or recall from Memory if already have info. - ... - Thought N:

Use Memory to retrieve any earlier points needed for the conclusion. - Final Thought: Outline the final report structure or even draft it (maybe as bullet points or paragraphs). - (Then outside ST, the agent might actually format the final document from these pieces).

Each workflow should intermix the sequential thinking and action tools fluidly. The **Sequential Thinking tool provides the coherent narrative**, while the other tools provide the capabilities.

**Template for Interleaving in Code** (conceptual):

```
# Pseudocode for agent loop
while True:
    if current_task_requires_reasoning and not done_reasoning:
        thought = model.generate_next_thought(context)
        call_tool(sequentialthinking, thought)
        response = get_tool_response()
        update_context_with_tool_output(response)
        if response["nextThoughtNeeded"] == False:
            done_reasoning = True
            final_answer = response["thought"]
            break # exit loop
        else:
            # possibly decide next action (maybe another tool or next
            thought)
            continue
    else:
        # If current step in plan is an action:
        action = decide_action_from_current_thought()
        call_tool(action.tool, action.params)
        action_result = get_tool_response()
        update_context_with_action_result(action_result)
        # then loop will have model generate next thought incorporating this.
```

This pseudo-loop indicates how your skill might alternate between reasoning steps and action steps until the reasoning is done. Ensure that after each action, the next sequential thinking call incorporates what happened.

**Monitoring and Adjusting:** As the skill developer, you should test these patterns with representative scenarios and adjust prompt phrasing or logic as needed. Look at the chain-of-thought outputs during testing to see if the agent is following the guidelines. If you notice, say, it's not branching when it should, you might add a line in the system prompt like "If you identify multiple viable approaches, use branching to explore them." Similarly, if it's verbose or looping, add "Don't repeat thoughts and be concise in reasoning, focusing on progress."

By implementing the above frameworks, templates, and patterns, your Claude Code skill should be well-equipped to leverage the Sequential Thinking MCP server optimally – calling it when appropriate, guiding the AI to use it effectively, and integrating it with other tools for complex, multi-step tasks. The result will be an AI agent that can tackle difficult problems in a methodical, transparent way, much like a human engineer, and ultimately deliver more reliable and explainable outcomes.

# Quick Reference

For quick consultation, here's a summary of key points and best practices when using the Sequential Thinking MCP server in your AI development workflows:

- **What It Is:** A state-tracking tool for AI reasoning. It doesn't think for the model; it logs and structures the model's own thoughts <sup>2</sup>. Use it to break down complex tasks into steps, with the AI explicitly writing each step.
- **When to Use (Decision Checklist):** Use ST if the task is complex, multi-step, or requires careful planning and verification. Specifically consider:
  - Does the problem have multiple parts or unclear scope? (Yes → use ST)
  - Need the reasoning trace for audit/debug? (Yes → use ST)
  - Will solving involve multiple tool calls or interactions? (Yes → ST can orchestrate)
  - Is the cost of error high (code bugs, critical decision)? (Yes → ST for thoroughness)
  - Is it a trivial question or one-step task? (Yes → skip ST, avoid overhead) <sup>59</sup>.
- **Starting a Thought Chain (Invocation):** Begin with a clear plan or analysis:
  - **Thought 1:** Restate the problem, list a plan of attack or questions to answer <sup>41</sup>. Set `thoughtNumber: 1`, and `totalThoughts` to the planned number of steps (adjust later if needed). Example initial input:

```
{ "thought": "Let's break this problem down: (1)... (2)... (3)...", "nextThoughtNeeded": true, "thoughtNumber": 1, "totalThoughts": 3 }
```
  - Prompt the model to use the tool by saying something like “Proceed step-by-step using the sequential-thinking tool.” <sup>107</sup>.
- **Tool Input Parameters:** Always provide:
  - `thought` (string): The content of the reasoning step (what you're thinking/deciding in this step).
  - `nextThoughtNeeded` (bool): True if more steps will follow, False if this is the final step/answer <sup>25</sup>.
  - `thoughtNumber` (int): 1 for first thought, increment each time <sup>26</sup>.
  - `totalThoughts` (int): current planned total steps (can change as you go) <sup>108</sup>.
- Optional flags when appropriate:
  - `isRevision: true` if correcting a previous thought, with `revisesThought: <number>` indicating which one <sup>27</sup>.
  - `branchFromThought: <number>` and `branchId: "ID"` if branching from a particular prior step <sup>30</sup>.
  - `needsMoreThoughts: true` if you reached an end but realize you need to extend (or set a new larger `totalThoughts`) <sup>34</sup>.
- **Managing Sequence:**

- Keep each **thought focused** on one aspect or sub-problem. Short paragraphs or bullet points in the thought can help clarity (the server just treats it as text).
- After each thought, **decide the next action**:
  - If it's another reasoning step, prepare the next thought input.
  - If it's time to use a tool (e.g., fetch data, read file), call that tool *outside* the ST chain, then continue the chain incorporating the result <sup>109</sup> <sub>62</sub>.
- **Integrate results:** whenever you use another tool, the next thought should reflect on that output (summarize what was found, how it affects the plan).
- Adjust `totalThoughts` upward if you realize more steps are needed than initially thought <sup>8</sup>. It's okay for `thoughtNumber` to exceed the original `totalThoughts` (the server will update it in output).
- **Branching Best Practices:** Use branches to explore alternatives:
  - Introduce a branch when you hit a decision point with multiple options worth exploring (e.g., "Option A vs Option B") <sup>47</sup>.
  - Assign a simple `branchId` (like "A", "B"). Start branch with `branchFromThought` = the point of divergence <sup>9</sup>.
  - Complete one branch path (a few thoughts) then go back and explore the other. Finally, converge: have a thought that compares outcomes and picks one (no `branchId` on that converging thought, returning to main sequence).
  - **Limit branches:** avoid more than 2-3 parallel paths and don't branch recursively without resolution.
  - Example:
    - Thought 3A (`branchId": "A"`): exploring first approach.
    - Thought 3B (`branchId": "B"`): later, exploring second approach <sup>109</sup>.
    - Thought 4 (no branch, main): "Comparing A and B, I choose B as it's better."

- **Revision Best Practices:** Revise if you find a mistake in earlier reasoning:
  - Set `isRevision:true` and `revisesThought:X` for the thought you're correcting <sup>27</sup>.
  - In the revision thought, explain the mistake and the correction ("I revise thought X: ...").
  - After revising, continue the chain with the new understanding. Do not treat the original thought as valid anymore.
  - Use revisions sparingly – only for significant changes in approach or fixing errors, not for minor tweaks.

- **Signs of a Good Chain (Quality Heuristics):**
  - Each thought yields progress or a new insight (no repetitive stalling).
  - The chain converges to a solution within a reasonable number of steps for the problem complexity (usually <10 for most cases; more only if truly needed).
  - The AI uses branching and revision where appropriate *but not excessively*. It questions earlier steps when warranted and explores alternatives if stuck, rather than plowing ahead incorrectly <sup>93</sup>.
  - The final answer addresses the original query fully, and any uncertainties have been resolved or noted.

- **Avoiding Pitfalls:**

- **Don't overthink trivial tasks:** Skip ST for simple queries – it adds overhead with no benefit [59](#).
- **Don't loop indefinitely:** If you're not making progress after several thoughts, try a new angle or tool, or conclude with the best you have. Always ensure to eventually set `nextThoughtNeeded: false` – do not leave it true forever [36](#).
- **No filler thoughts:** Every thought should have substance. Avoid generic statements like "I will continue thinking" with no content.
- **Use tools when needed:** Don't try to reason out something that you can directly obtain via a tool (e.g. don't guess a file's content – open it with Filesystem).
- **Keep context in mind:** The server does not remind you of past thoughts (except via the JSON you get back). You must remember important details or re-summarize if needed. Use `thoughtHistoryLength` or branch list in the output if it helps track where you are [33](#).

- **Integration Pattern:** Treat ST as the **planner** and other MCP servers as executors:

- Insert tool calls in between thoughts as actions and feed results into next thoughts. E.g., Thought: "I will now fetch user data from API" -> use Fetch tool -> next Thought analyzes that data.
- This interleaving ensures the reasoning stays grounded in real data and effects [64](#) [110](#).
- Common combos:
  - *sequential-thinking + filesystem* (for stepwise code edits, refactoring) [76](#),
  - *sequential-thinking + fetch/web* (for research tasks) [77](#),
  - *sequential-thinking + memory* (for long projects, to store intermediate findings) [73](#),
  - *sequential-thinking + git/github* (for multi-file or PR analysis) [76](#),
  - etc., with ST orchestrating multi-tool usage.

- **Ending the Chain:**

- Set `nextThoughtNeeded: false` on the final thought when you have a satisfactory answer/solution [36](#).
- The final thought's content should ideally provide the solution clearly. This can be a direct answer, a summary of findings, or the completed code, etc., depending on the task.
- If appropriate, have the final thought synthesize the reasoning ("Therefore, I conclude that ...") so the answer is well-justified (though typically, you'll present just the conclusion to the user, the justification having been in the hidden chain).
- Ensure any remaining branches are resolved: don't leave alternate solutions dangling – close them out by stating why the chosen answer is best.
- Example final tool call:

```
{ "thought": "After this analysis, the best solution is to implement caching. This will reduce load times significantly.",  
  "nextThoughtNeeded": false, "thoughtNumber": 6, "totalThoughts": 6 }
```

- The agent should then deliver this conclusion as the final output to the user (in a user-friendly format, stripping JSON).

By following this quick-reference and the detailed guidance above, you can confidently integrate the Sequential Thinking server into your AI agent, enabling it to tackle complex tasks systematically and transparently – improving the quality of solutions and providing insight into the AI's reasoning process.

40 39

---

1 38 39 57 58 62 63 64 65 66 67 68 90 94 97 99 105 106 109 110 Mastering Structured AI

Reasoning: A Deep Dive into Zalab Inc.'s Sequential Thinking MCP Server

<https://skywork.ai/skypage/en/mastering-structured-ai-reasoning/1978348626220806144>

2 37 40 55 96 98 Using the Sequential Thinking MCP Server to go from Generative to Agentic AI - Phase 2

<https://phase2online.com/2025/05/23/sequential-thinking-mcp-server-anthropic/>

3 4 5 6 7 8 9 10 11 13 14 15 16 17 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36  
41 53 60 61 72 103 104 108 sequentialthinking.ts

<https://github.com/CherryHQ/cherry-studio/blob/c53d6f3fd0df3bbf77f18300510f0f5d4ea37ed3/src/main/mcpServers/sequentialthinking.ts>

12 18 19 21 README.md

<https://github.com/modelcontextprotocol/servers/blob/e6b0b0f5d355530ebdd09bc547240e556fb6c19/src/sequentialthinking/README.md>

20 48 56 70 76 77 server-guide.md

<https://github.com/LETHALDOSE1300/johns-private-mcp-setup/blob/c4ba37169a977abf3d0d62eef67456af75d9e7bf/servers/server-guide.md>

42 43 44 45 47 59 74 75 92 93 How Sequential Thinking Changes Everything | byPawel

<https://www.bypawel.com/how-sequential-thinking-changes-everything>

46 49 84 85 Is Sequential Thinking still relevant? : r/GithubCopilot

[https://www.reddit.com/r/GithubCopilot/comments/1neflwm/is\\_sequential\\_thinking\\_still\\_relevant/](https://www.reddit.com/r/GithubCopilot/comments/1neflwm/is_sequential_thinking_still_relevant/)

50 80 81 83 91 Cluae Code's take on Sequential Thinking MCP : r/ClaudeAI

[https://www.reddit.com/r/ClaudeAI/comments/1l2zcf/ecluae\\_codes\\_take\\_on\\_sequential\\_thinking\\_mcp/](https://www.reddit.com/r/ClaudeAI/comments/1l2zcf/ecluae_codes_take_on_sequential_thinking_mcp/)

51 79 87 88 89 100 101 102 How does the Sequential Thinking MCP work? : r/mcp

[https://www.reddit.com/r/mcp/comments/1jwjagw/how\\_does\\_the\\_sequential\\_thinking\\_mcp\\_work/](https://www.reddit.com/r/mcp/comments/1jwjagw/how_does_the_sequential_thinking_mcp_work/)

52 71 73 82 86 Better results without sequential thinking, MCP? : r/AugmentCodeAI

[https://www.reddit.com/r/AugmentCodeAI/comments/1neojmj/better\\_results\\_without\\_sequential\\_thinking\\_mcp/](https://www.reddit.com/r/AugmentCodeAI/comments/1neojmj/better_results_without_sequential_thinking_mcp/)

54 Building Smarter AI Agents: How Sequential Thinking MCP Transforms Complex Problem-Solving | by Micheal Lanham | Medium

<https://medium.com/@Micheal-Lanham/building-smarter-ai-agents-how-sequential-thinking-mcp-transforms-complex-problem-solving-443e68b4d487>

69 107 How to Use Sequential Thinking - Help - Cursor - Community Forum

<https://forum.cursor.com/t/how-to-use-sequential-thinking/50374>

78 95 Sequential thinking MCP vs Claude 3.7 Extended Thinking - Reddit

[https://www.reddit.com/r/ClaudeAI/comments/1j6zi68/sequential\\_thinking\\_mcp\\_vs\\_claude\\_37\\_extended/](https://www.reddit.com/r/ClaudeAI/comments/1j6zi68/sequential_thinking_mcp_vs_claude_37_extended/)