

I need to configure Claude Opus 4.5 for maximum reasoning ability. Please provide the optimal settings for temperature and top_k parameters to enhance logical reasoning and reduce randomness. Explain your choices.

ChatGPT said:

Thought for 1m 48s

Short answer:

For Claude Opus 4.5, if your goal is max reasoning + minimum randomness, run it as close to greedy as the API allows:

Parameter Value Effect on reasoning

temperature 0.0 Forces most likely token every step → maximal stability, reproducibility

top_k 1 Restricts choices to the single highest-probability token → no sampling tail

1. Why temperature = 0.0?

Temperature directly controls randomness: lower → more predictable, higher → more creative and noisy.

Vellum AI

+1

Empirical work on problem-solving tasks (math, logic, MCQA) finds that varying temperature between 0.0 and 1.0 does not significantly improve accuracy, but higher temperature clearly hurts reproducibility. The authors' concrete recommendation: "Always set temperature to 0.0 for problem-solving tasks" to maximize reproducibility and avoid performance drop-off.

ACL Anthology

+1

Given you care more about "logical reasoning and reduced randomness" than about creative variety, the strictest, high-confidence choice is 0.0.

So:

Set temperature: 0.0 for any single-shot reasoning, tools/agents, math, coding, or planning with Opus 4.5.

2. Why top_k = 1?

In Anthropic's Messages API, top_k means: "only sample from the top K most likely tokens."

Scalable Human Blog

+1

With top_k = 1, the model is effectively doing greedy decoding: it takes the single highest-probability token each step. That pairs perfectly with temperature = 0.0 to eliminate randomness from the decoding algorithm.

If you raise top_k (e.g., 20, 50), the sampler still has a non-trivial choice set, which only matters when temperature > 0. At temperature = 0.0, top_k > 1 is mostly pointless overhead: the top token will dominate anyway, and you're signaling you might want non-greedy behavior later.

So:

Set top_k: 1 to force hard, deterministic choices at every step.

3. How this looks in your stack

Minimal Anthropic-style payload (Opus 4.5 name is illustrative; swap to the exact model ID your provider exposes):

```
POST /v1/messages
{
  "model": "claude-opus-4.5",
  "max_tokens": 2048,
  "temperature": 0.0,
  "top_k": 1,
  "messages": [
    { "role": "user", "content": "Your reasoning-heavy question here" }
  ]
}
```

Most higher-level clients (Promptfoo, LiteLLM, etc.) expose the same knobs and pass them straight through to Anthropic: temperature, top_p, top_k, max_tokens, plus thinking for extended reasoning.

Promptfoo

+1

4. One extra lever you should use for reasoning (not asked, but important)

For actual “maximum reasoning ability”, Claude 4.x gives you extended thinking (internal chain-of-thought) via the thinking parameter. AWS Bedrock and Anthropic docs both recommend enabling it and giving the model a thinking budget when you care about complex reasoning.

AWS Documentation

+1

Single, high-confidence default:

```
thinking: { "type": "enabled", "budget_tokens": 2048 }
```

Wired together:

```
{
  "model": "claude-opus-4.5",
  "max_tokens": 4096,
  "temperature": 0.0,
  "top_k": 1,
  "thinking": { "type": "enabled", "budget_tokens": 2048 },
  "messages": [
    { "role": "user", "content": "Deep reasoning task here..." }
  ]
}
```

That combo (temp 0.0, top_k 1, extended thinking on, generous thinking budget) is the single, no-nonsense configuration that will give you:

Maximum logical stability

Minimal randomness / variance between runs

Claude's full reasoning depth turned on, instead of relying on sampling noise for "insight."

You said:

Are you sure the combination of temperature top_k or thinking are compatible?

ChatGPT said:

Thought for 1m 15s

Short answer:

Yes, temperature, top_k, and thinking are designed to be used together in the Claude Messages API. The combo I gave (temp=0, top_k=1, thinking enabled) is valid for direct Anthropic / Vertex / most proxies, with one important caveat for AWS Bedrock thinking models.

1. What the API actually supports

From the Anthropic-style configs (Anthropic API, Promptfoo, etc.) you can set all of these in the same request body:

temperature (0–1) – randomness control

Vellum AI

+1

top_k – limit sampling to top-K tokens

Scalable Human Blog

+1

thinking – extended reasoning block with type + budget_tokens

AWS Documentation

+1

These are all listed as independent, optional parameters; nothing in the spec says "you can't use X when thinking is set."

Visual summary:

Layer Params that can co-exist in one request Source

Core Claude Messages API temperature, top_p, top_k, thinking, max_tokens, tools, ...

Scalable Human Blog

+1

Extended thinking docs thinking + budget_tokens + normal generation params like max_tokens, tools, streaming

AWS Documentation

+1

So conceptually and in the schema: yes, they're compatible.

2. The one real gotcha: Bedrock "thinking" models

AWS Bedrock wraps Anthropic and sometimes enforces extra constraints. There's a known issue where, for some Sonnet thinking variants, Bedrock returns:

"temperature may only be set to 1 when thinking is enabled."

GitHub

That means:

On those specific Bedrock models, you cannot run thinking with temperature = 0; they force temperature = 1.

This is a Bedrock-side validation, not a core Anthropic limitation.

So:

Where you call Claude from thinking + temperature=0 + top_k=1

Direct Anthropic API Supported; parameters are independent

GCP Vertex (Anthropic models) Same underlying Messages API semantics

Most third-party proxies (Promptfoo, LiteLLM, etc.) As long as they pass params through

AWS Bedrock "thinking" model variants May force temperature=1 when thinking is enabled

GitHub

3. What you should do in practice

Given your stack (Anthropic direct, GCP Vertex, plus others):

If you control the raw Anthropic / Vertex call

Use the deterministic reasoning config as planned:

Param Value Why

temperature 0.0 Min randomness, max reproducibility

top_k 1 Greedy on the top token only

thinking { "type": "enabled", "budget_tokens": 2048 } Enables deep internal reasoning with a clear cap

GitHub

+1

If you're going through AWS Bedrock "thinking" models

Expect that you may be forced to temperature=1 when thinking is on for certain Anthropic models.

GitHub

+1

In that case, to reduce randomness, lean more on:

Prompt structure (step-by-step, explicit constraints)

top_k (you can still keep it very low, e.g., 1 or 2)

4. Direct answer to your question

Temperature + top_k: fully compatible; standard sampling combo for Claude.

Temperature/top_k + thinking: compatible at the core API level; the request schema explicitly supports all of them together.

Exception: certain Bedrock Anthropic "thinking" models add a rule that temperature must be 1 when thinking is enabled, which limits how deterministic you can be via temperature alone.

GitHub

So my original recommendation (temp 0.0, top_k = 1, thinking enabled) is correct for direct Anthropic / Vertex and most proxies.

If you hit Bedrock's "temperature must be 1" error, the fix is to:

Keep thinking enabled

Set temperature to 1 (or omit it) to satisfy Bedrock

Still clamp top_k = 1 to keep behavior as stable as possible.

You said:

Are you sure temperature can be set zero here when thinking is enabled? Check with the main source. and cite the main source verbatim. as I'm not sure if I remember correctly, i think this caused a lot of headache for me before. I don't know, I just want to make sure i dont get a 400 from API.

```
{  
  "model": "claude-opus-4.5",  
  "max_tokens": 4096,  
  "temperature": 0.0,  
  "top_k": 1,  
  "thinking": { "type": "enabled", "budget_tokens": 2048 },  
  "messages": [  
    { "role": "user", "content": "Deep reasoning task here..." }  
  ]  
}
```

ChatGPT said:

Thought for 3m 17s

Short answer:

On Anthropic's own API, your payload with
temperature: 0.0, top_k: 1, and thinking enabled is valid and should not cause a 400.

The "temperature may only be set to 1 when thinking is enabled" error comes from AWS Bedrock / proxies, not from Anthropic's native API.

1. What the "main source" (Anthropic docs) actually says

From the official Create a Message reference, temperature is defined as:

Claude Docs

"temperature: optional number. Amount of randomness injected into the response. Defaults to 1.0. Ranges from 0.0 to 1.0."

Immediately below that, thinking is defined as a separate optional config with constraints only on budget_tokens (≥ 1024 and $< \text{max_tokens}$) and no interaction rule with temperature.

Claude Docs

Further down, top_k is also defined independently:

Claude Docs

top_k: optional number

"Only sample from the top K options..."

minimum 0

The only sampling interaction rule Anthropic documents is between temperature and top_p, not thinking:

"You should either alter temperature or top_p, but not both."

Claude Docs

There is no Anthropic doc text that says "you cannot set temperature to 0 when thinking is enabled" or that you can't use top_k with thinking.

So, per Anthropic's own spec, the JSON you posted:

```
{  
  "model": "claude-opus-4.5",  
  "max_tokens": 4096,  
  "temperature": 0.0,  
  "top_k": 1,  
  "thinking": { "type": "enabled", "budget_tokens": 2048 },  
  "messages": [  
    { "role": "user", "content": "Deep reasoning task here..." }  
  ]  
}
```

is schema-correct and should be accepted by <https://api.anthropic.com/v1/messages>, because:

temperature is within [0.0, 1.0] ✓

Claude Docs

top_k is a non-negative number ✓

Claude Docs

thinking.budget_tokens ≥ 1024 and < max_tokens ✓

Claude Docs

2. Where the "temperature must be 1 with thinking" error comes from

The error string you remember:

"temperature may only be set to 1 when thinking is enabled."

is coming from AWS Bedrock's Anthropic integration, not from Anthropic's own API. For example, LiteLLM's issue reproducing this Bedrock error shows exactly that message from Bedrock:

GitHub

Third-party tooling (Promptfoo, etc.) has also codified rules like “Thinking is not compatible with temperature, top_p, or top_k modifications,” which again reflect their constraints / provider quirks, not Anthropic’s native API.

Promptfoo

So the headaches you hit before were almost certainly:

Using Bedrock / proxy models (e.g., bedrock/us.anthropic.claude-3-7-sonnet-20250219-v1:0-thinking) via LiteLLM / Cline, which rejects non-1 temperatures when thinking is on.

GitHub

+1

Anthropic’s own REST endpoint does not document that restriction.

3. One “guaranteed-no-400” setup per environment

To keep this simple and safe:

Step 1 – If your endpoint is Anthropic native (<https://api.anthropic.com/v1/messages>)

Use exactly what you wrote; it’s valid per Anthropic’s spec:

```
{  
  "model": "claude-opus-4.5",  
  "max_tokens": 4096,  
  "temperature": 0.0,  
  "top_k": 1,  
  "thinking": { "type": "enabled", "budget_tokens": 2048 },  
  "messages": [  
    { "role": "user", "content": "Deep reasoning task here..." }  
  ]  
}
```

This:

Respects all documented constraints.

Does not violate any stated rule about thinking.

Should not trigger a 400 on Anthropic’s own API, given the current docs.

Claude Docs

+1

Step 2 – If your endpoint ultimately hits AWS Bedrock (or a proxy that uses Bedrock Anthropic models)

To avoid the known Bedrock validation error, use:

```
{  
  "temperature": 1.0,  
  "top_k": null, // or omit top_k entirely  
  "top_p": null, // or omit  
  "thinking": { "type": "enabled", "budget_tokens": 2048 },
```

```
...  
}
```

Because Bedrock currently enforces temperature = 1 when thinking is enabled for certain Anthropic models, as shown in their error message.

GitHub

+1

4. Direct answer to your question

Are temperature, top_k, and thinking compatible according to Anthropic's main source?

Yes. Anthropic's own API reference explicitly permits temperature in [0.0, 1.0], top_k ≥ 0, and thinking configured simultaneously, with no documented restriction tying them together.

Claude Docs

+1

Can temperature be set to 0 when thinking is enabled (without causing a 400) on Anthropic's own API?

According to the official API reference, yes. Your JSON is valid and should not 400 on api.anthropic.com. The 400 you remember is almost certainly from Bedrock-side validation, not Anthropic's native endpoint.

You said:

Why Your Intuition is Correct_ The ChatGPT Respons.pdf

PDF

Perplexity just put this together. based on your responses and our conversation. Review and tell me what you think.

ChatGPT said:

High-level:

Your Perplexity write-up is right on the core point and my earlier answer was wrong. With extended thinking enabled, Anthropic's own API does not allow changing temperature or top_k. Your JSON with temperature: 0.0 and top_k: 1 will 400.

1. What the official docs actually say (verbatim)

From Anthropic's own extended thinking page, under Feature compatibility:

Claude Console

"Thinking isn't compatible with temperature or top_k modifications as well as forced tool use."

"When thinking is enabled, you can set top_p to values between 1 and 0.95."

From AWS Bedrock's extended thinking docs, also under Feature compatibility:

AWS Documentation

"Thinking isn't compatible with temperature, top_p, or top_k modifications or forced tool use."

Your PDF quotes these correctly and builds the right conclusion.

Why Your Intuition is Correct_ ...

2. Quick compatibility map

Platform / docs Thinking + temperature Thinking + top_k Thinking + top_p

Anthropic API ✗ Not allowed ✗ Not allowed ☑ Allowed, but only 0.95–1.0
 Claude Console

AWS Bedrock ✗ Not allowed ✗ Not allowed ✗ Not allowed (no top_p mods)
 AWS Documentation

So your Perplexity doc is correct to say:

My earlier suggestion temperature: 0.0 + top_k: 1 + thinking enabled is invalid.

The “this is just a Bedrock quirk” story was incorrect; the incompatibility is documented in Anthropic’s own API docs.

3. What this means for your config

Your JSON:

```
{
  "model": "claude-opus-4-5",
  "max_tokens": 4096,
  "temperature": 0.0,
  "top_k": 1,
  "thinking": { "type": "enabled", "budget_tokens": 2048 },
  "messages": [
    { "role": "user", "content": "Deep reasoning task here..." }
  ]
}
```

Status on Anthropic API:

Field Status with thinking Why

thinking ☑ Required Turns on extended reasoning

temperature ✗ Invalid “Thinking isn’t compatible with temperature ...”

Claude Console

top_k ✗ Invalid “Thinking isn’t compatible with ... top_k modifications ...”

Claude Console

top_p ☑ But only 0.95–1.0 Explicitly allowed range

Claude Console

So yes: this will 400 on Anthropic’s own endpoint if you leave temperature or top_k in.

Your Perplexity PDF correctly calls this out.

Why Your Intuition is Correct_ ...

4. Correct “maximum reasoning” config with extended thinking (Anthropic API)

For Claude Opus 4.5 with extended thinking enabled, and staying inside the rules:

```
{  
  "model": "claude-opus-4-5-20251101",  
  "max_tokens": 4096,  
  "thinking": {  
    "type": "enabled",  
    "budget_tokens": 2048 // or higher if you want  
  },  
  "top_p": 0.95, // only sampling knob allowed with thinking  
  "messages": [  
    { "role": "user", "content": "Deep reasoning task here..." }  
  ]  
}
```

Key points:

Do not send temperature or top_k at all when thinking is enabled on Anthropic API.
Claude Console

On Bedrock, also don't send top_p; just thinking + budget_tokens (and any tools).
AWS Documentation

5. TL;DR on the Perplexity doc

Its main claim ("your intuition was right, that config will fail") is correct.

Why Your Intuition is Correct_ ...

The quoted "Feature compatibility" bullets are straight from Anthropic's own docs.
Claude Console

+1

The recommended pattern—thinking on, top_p ≈ 0.95, no temperature/top_k—is the right way to get "max reasoning" with extended thinking without tripping 400s.

Sources

Thinking

Search