# Claude API rate limit issues with Sonnet

🌐 claude.ai/chat/27b8a19b-4c6b-413c-a5a7-f87b6af1e11c

I

```
Created 14 files:
```

User API Key Rate limit exceeded using sonnet 4.5 from cursor.. keep running in to this

> The user is showing me a screenshot from what appears to be Cursor (an AI-powered code editor) where they're trying to set up a Guidewire training platform project. They've successfully initialized the project with shadcn components, but are now hitting an API rate limit issue.
>
> The error message says "User Provided API Key Rate Limit Exceeded" which means they're hitting rate limits on the Anthropic API when using Claude Sonnet 4.5 through Cursor.
>
> Let me provide helpful advice on how to handle this issue. I should:
>
> 1. Explain what's happening
> 2. Provide solutions to avoid or work around the rate limit
> 3. Be practical and helpful given their context of building a Guidewire training platform

1/18

I should not search for this as it's a technical issue I can address from my knowledge.

I can see you're hitting API rate limits while working on your Guidewire training platform in Cursor with Sonnet 4.5. Here are some practical ways to handle this:

## Immediate Solutions

**1. Wait it out** - Rate limits typically reset after:

- 1 minute for per-minute limits
- 1 hour for hourly limits
- Check your Anthropic Console at console.anthropic.com to see your exact limits

**2. Reduce request frequency in Cursor:**

- Use Ctrl+L (Cmd+L on Mac) for chat instead of inline editing - it's less API-intensive
- Disable auto-complete/inline suggestions temporarily
- Work on code manually for a bit, then batch your AI requests

**3. Switch to a different model temporarily:**

- In Cursor, change from Sonnet 4.5 to Sonnet 4 or Haiku
- Click where it says "Sonnet 4.5" and select a different model
- Haiku is much cheaper on rate limits for simple tasks

## Long-term Solutions

**1. Upgrade your API tier:**

- Go to console.anthropic.com
- Check your current tier (likely Tier 1)
- Request a higher tier for increased rate limits
- Tier 2 gives you 40 requests/min vs 10 requests/min

**2. Optimize your Cursor usage:**

- Be more selective about when you invoke Claude
- Write more code manually between AI assists
- Use smaller, focused prompts instead of asking AI to generate large chunks

**3. Use multiple API keys:**

- Create a separate Anthropic account for personal projects
- Keep one key for Cursor, one for other tools

For your immediate task of building the Guidewire training platform components, I'd suggest switching to Sonnet 4 or working manually for 5-10 minutes to let the rate limit reset.

I

does it mean it is stopping execution abruptly? or losing progress or have to re do the work again?

> The user is asking whether the rate limit error means they're losing their progress or have to redo work. Looking at the screenshot, it shows:
>
>> 1. "Created 14 files:" with a success checkmark
>> 2. "Thought for 1s"
>> 3. "Cannot edit files blocked by globalIgnore"
>> 4. A "Review" button
>
> So it appears that:
>
> - The initial scaffolding command completed successfully
> - 14 files were created
> - But then there's a "globalIgnore" issue preventing file editing
>
> The user is concerned about whether hitting a rate limit means losing progress or having to start over.
>
> I should clarify:
>
>> 1. The rate limit doesn't cause you to lose already-completed work
>> 2. Files that were created are still there
>> 3. You just can't make new API calls until the limit resets
>> 4. In this case, it looks like the command actually succeeded - 14 files were created
>> 5. The "Cannot edit files blocked by globalIgnore" is a different issue, not a rate limit issue
>
> So I should reassure them that they don't lose progress, and explain what's actually happening in their specific screenshot.

No, you're **not losing any progress**! The rate limit doesn't delete or undo work that's already completed.

Looking at your screenshot specifically:

✅ **Good news:** Your command **succeeded** - it says "Created 14 files:" with a success checkmark. Those files are already created and saved in your project.

The issues you're seeing are:

1. **"Cannot edit files blocked by globalIgnore"** - This is a **different issue**, not related to the rate limit. This means some files couldn't be edited due to your `.gitignore` or Cursor's ignore settings (likely `node_modules` or similar).
2. **Rate limit** - This just means you can't make **new** API requests for a short period. It doesn't affect:
   - Files already created ✅
   - Code already written ✅
   - Changes already saved ✅

## What happens with rate limits:

- **Work completed before the limit:** Stays intact
- **Work in progress when limit hits:** Stops at that point, but keeps what's done
- **Next request:** Has to wait until limit resets

## For your current situation:

Your shadcn component installation **completed successfully**. The 14 UI component files are already in your project. You can:

- Open and view those files manually
- Continue coding without AI assistance
- Wait ~1 minute for the rate limit to reset before asking Cursor for more AI help

You don't need to redo anything - the work is saved!

I

The user is showing me a Cursor screenshot where they're working on a Guidewire training platform and hit the rate limit again. The AI agent has identified two bugs:

Bug 1: Edge runtime incompatibility with Supabase - using cookies() API which isn't available in Edge runtime Bug 2: Async function 'markTopicStarted' called without awaiting

And then the rate limit error appears at the bottom.

The user is asking (implicitly) about this situation. They want to understand what's happening and probably how to proceed.

I should:

1. Acknowledge they hit the rate limit again
2. Explain that the AI successfully identified the bugs before hitting the limit

3. Tell them they can fix these bugs manually while waiting for the limit to reset
4. Provide quick guidance on fixing both bugs

Good news - the AI **found your bugs** before hitting the rate limit! You have the diagnosis, so you can fix these manually while waiting for the limit to reset.

## What Cursor Found:

### Bug 1: Edge Runtime + Supabase Compatibility

- Problem: You're using `cookies()` API in an Edge runtime route
- Location: Line 9 in some file, line 19 uses `createClient()` from Supabase
- Fix: Either switch to non-Edge route handler OR use a different auth method

### Bug 2: Missing `await` on async function

- Problem: `markTopicStarted` is called without awaiting
- Location: `route.ts` (lines 8-13) and `TopicContent.tsx` (lines 31-44)
- Fix: Add `await` before the function call

## How to fix manually:

### For Bug 1:

typescript

```typescript
// Remove this from the top of your route file:
export const runtime = 'edge';

// OR use a different Supabase client setup for Edge
```

### For Bug 2:

typescript

```typescript
// Change from:
markTopicStarted(topicId);

// To:
await markTopicStarted(topicId);
```

## Next steps:

1. **Fix these two bugs now** (manually edit the files)
2. **Wait 1-2 minutes** for rate limit to reset
3. **Then ask Cursor** to verify the fixes or help with next steps

The lint check showed "No ESLint warnings or errors" which is good - these are TypeScript/runtime issues, not linting issues.

Want me to help you with the exact code fixes for these issues?

I

The user is showing me a Cursor AI conversation where they're working on their Guidewire training platform project. The AI is proposing three options:

1. Option 1: Implement the optimal tech stack (update project files and start work, close Sprint)
2. Option 2: Review Current Status (show sprint status, changelog, confirm deployment readiness)
3. Option 3: Start Sprint 2 Immediately (Supabase setup, test locally, prepare for deployment)

The user then chose Option A (close the sprint, create Sprint 2 plan), and Cursor found a development review prompt document and is ready to proceed.

This seems to be going well - the user is getting organized with their development process. They're not asking a specific question here, but showing me their progress. I should acknowledge this and perhaps ask what they'd like help with next, or offer relevant assistance for their Guidewire training platform project.

Great! Cursor is helping you organize your development workflow properly. You're at a good decision point for your Guidewire training platform project.

## What's Happening:

Cursor found your **development review document** (`08_DEVELOPMENT_REVIEW_PROMPT.md`) and is ready to:

1. Close Sprint 1
2. Create a Sprint 2 plan based on the review notes

## Your Options Context:

**Option 1** (Implement tech stack) - Best if you want to keep building now
**Option 2** (Review status) - Best if you need to assess what's done before deploying
**Option 3** (Start Sprint 2) - Supabase setup, testing, deployment prep

You chose **Option A** which makes sense - you're taking a structured approach by:

- ✅ Reviewing what was accomplished in Sprint 1
- ✅ Creating a proper plan for Sprint 2
- ✅ Following good agile practices

## Next Steps:

Once Cursor reviews your notes, you'll likely get a Sprint 2 plan that includes:

- Supabase authentication setup
- Database schema for courses/topics/progress
- User enrollment system
- Testing and deployment preparation

Since you were hitting rate limits earlier, this structured approach is smart - it lets you plan before executing, reducing back-and-forth API calls.

**Want me to help with anything specific** while Cursor processes this, like:

- Database schema design for your Guidewire courses?
- Supabase setup strategy?
- Deployment checklist for Vercel?

I

no got teh same erro afte r that.. so that means i have to rerun the lst review notes things

The user is saying they got the same error (rate limit) after that, and they're asking if that means they have to rerun the last review notes thing.

So the user is concerned that:

1. They hit the rate limit again
2. The work Cursor was doing (reviewing notes, creating Sprint 2 plan) might not have completed
3. They'll need to restart that process

I should clarify:

- Whether the work was saved or lost
- How to check what was completed
- Best way to continue without wasting API calls
- Strategy to avoid hitting limits repeatedly

**No, you shouldn't need to completely redo it** - but it depends on what Cursor actually completed before the rate limit hit.

## Check what was saved:

1. **Look for new files created:**
   - Check if there's a Sprint 2 plan file (might be something like `SPRINT_2_PLAN.md` or similar in your project-docs folder)
   - Check if `07_REVIEW_NOTES.md` was updated
2. **Check the chat history:**
   - Scroll up in Cursor's chat to see what it actually output before the error
   - Sometimes it writes the plan in the chat even if it didn't save to a file

## Smart way to continue:

Instead of asking Cursor to redo everything, do this:

1. **Wait 2-3 minutes** for rate limit to fully reset
2. **Be more specific** in your next request:

   Based on the 08_DEVELOPMENT_REVIEW_PROMPT.md file, create a Sprint 2 plan document at project-docs/SPRINT_2_PLAN.md. Keep it concise.

3. **Or go manual temporarily:**
   - Open `08_DEVELOPMENT_REVIEW_PROMPT.md` yourself
   - Quickly read through it
   - Create a simple Sprint 2 checklist manually
   - Then use Cursor for actual coding when limit resets

## Pro tip to reduce rate limit hits:

Tell Cursor: "Create Sprint 2 plan. Output as a single markdown code block, don't create files yet."

This uses fewer API calls since it's just generating text, not running file operations. Then you can copy/paste it yourself.

**What would work best for you?** Wait and retry, or manually create the Sprint 2 plan while the limit resets?

I

I keep running in to this.. espcially when idniatiating a new agent and ask it review entire projec t for context.. we need to come up with a solutions to not run out of api limitOr should i swuth to chat gpt codex?

The user is frustrated with repeatedly hitting API rate limits in Cursor, especially when starting new conversations where the AI agent needs to read multiple project files for context. They're asking for solutions and wondering if they should switch to ChatGPT/Codex instead.

Looking at the screenshot, Cursor is reading 8 files just to get context:

- 06_CURRENT_SPRINT.md
- 07_REVIEW_NOTES.md
- 09_SPRINT_2_START_PROMPT.md
- 03_MASTER_PLAN.md
- 99_CHANGELOG.md
- 01_VISION.md
- README.md
- instructions.md

This is burning through API calls just for context loading, before even generating a response.

I should provide practical solutions to reduce rate limit hits, and also give an honest comparison with alternatives.

You're hitting a **fundamental issue** with how Cursor loads context - reading 8+ files just to start a conversation burns through your rate limit fast. Let me give you practical solutions:

## Immediate Solutions for Cursor

**1. Create a context summary file (Best fix)** Instead of having Cursor read 8 files every time, create ONE file:

markdown

```
# PROJECT_CONTEXT.md
Quick summary of project state, current sprint, tech stack.
Link to detailed docs if needed.
```
Tell Cursor to "read only PROJECT_CONTEXT.md" - saves 7 API calls

**2. Use @ mentions selectively**
```
Don't: "@workspace review sprint 1" (reads everything)
Do: "Review only @07_REVIEW_NOTES.md and tell me next steps"
```

**3. Upgrade your Anthropic API tier**

- Go to console.anthropic.com
- Request Tier 2: Gets you **40 requests/min** vs 10

- Usually approved within 24 hours if you show usage history
- This is the most effective solution

## 4. Use Cursor's "Slow" mode

- In Cursor settings, there may be a rate limit setting
- Some versions have a "slow" mode that paces requests

# Should you switch to ChatGPT/Codex?

The user is asking about switching to ChatGPT Codex. I should provide an honest comparison. ChatGPT has GitHub Copilot which is different from Cursor. Let me clarify the options and their rate limits.

**ChatGPT Plus + Canvas** (for coding):

- ✅ Higher rate limits (80 messages per 3 hours)
- ✅ No API key needed
- ❌ Not integrated into your editor
- ❌ Manual copy/paste workflow

**GitHub Copilot** (closer alternative to Cursor):

- ✅ Unlimited completions
- ✅ Integrated in VS Code
- ❌ Less powerful for complex tasks
- 💰 $10-20/month

**Continue.dev** (open source Cursor alternative):

- ✅ Use your own API key
- ✅ Same workflow as Cursor
- ⚠️ Still has same rate limit issue

# My honest recommendation for you:

Given that you're building a business platform (Guidewire training), **don't switch yet**. Instead:

**Short-term (today):**

1. Create that PROJECT_CONTEXT.md file (do it manually, takes 5 min)
2. Request Tier 2 API access at console.anthropic.com
3. Work more manually between AI assists

**Medium-term (this week):**

1. Get Tier 2 approved (40 req/min = 4x more capacity)
2. Consider adding GitHub Copilot ($10/mo) for basic autocomplete
3. Use Cursor for complex tasks, Copilot for simple completions

**Long-term:** Once your staffing business generates revenue, upgrade to Tier 3 or Tier 4 for essentially unlimited usage.

Want me to help you create that PROJECT_CONTEXT.md file right now? It'll dramatically reduce your rate limit issues.

I

no cursor has gpt codex model aswell

> Ah, the user is saying that Cursor also has GPT models (Codex/GPT-4) available, not just Claude. They're asking if they should switch to using those models within Cursor instead of Claude Sonnet 4.5.
>
> Let me give them a practical comparison for their specific use case.

Ah yes, you're right! Cursor supports multiple models including GPT-4, GPT-4 Turbo, and o1. Let me give you a practical comparison for **your specific situation**:

# Model Comparison in Cursor

**Claude Sonnet 4.5** (what you're using now)

- ✅ Best for complex reasoning and planning
- ✅ Better at following project structure
- ❌ **Strict rate limits (your current problem)**
- Cost: ~10 requests/min on Tier 1

**GPT-4 Turbo** (OpenAI)

- ✅ **Much higher rate limits** (typically 500 requests/min)
- ✅ Faster responses
- ⚠️ Slightly less thoughtful than Claude for complex tasks
- Cost: Uses your OpenAI API key

**GPT-4o** (OpenAI latest)

- ✅ **Highest rate limits**

- ✅ Very fast
- ✅ Good at code generation
- ⚠️ Can be more verbose

**o1 models** (OpenAI reasoning)

- ✅ Best for complex problem-solving
- ❌ Expensive per request
- ❌ Slower than others

# My recommendation for YOUR workflow:

**Switch to GPT-4 Turbo or GPT-4o for now**

Why? Because:

1. **Rate limits won't block you** - you can do full project reviews without hitting limits
2. **Good enough for your current tasks** - setting up Next.js + Supabase doesn't need Claude's advanced reasoning
3. **Faster iteration** - you can keep building without stopping

**When to use Claude Sonnet 4.5:**

- Complex architecture decisions
- Reviewing tricky bugs
- Planning major features
- When you have specific time to work (so rate limits reset between sessions)

# How to switch in Cursor:

1. Click where it says "Sonnet 4.5" in the bottom of chat
2. Select "GPT-4 Turbo" or "GPT-4o"
3. Add your OpenAI API key if needed (Settings → Models)

**Try this right now:**

- Switch to GPT-4 Turbo
- Retry your "review sprint 1 and plan sprint 2" command
- You should get through it without rate limit errors

Does that solve your immediate problem?

I

The user is showing me the model selection dropdown in Cursor. I can see the available models:

- Auto
- MAX Mode
- Sonnet 4.5 (currently selected, shown with checkmark)
- Composer 1
- GPT-5 Codex (appears to have a note "exceeded")
- GPT 5
- Haiku 4.5
- Add Models

The user mentioned "gpt codex" earlier, and I see "GPT-5 Codex" in the list. However, I notice there's text that says "exceeded" next to it, which likely means the rate limit has been exceeded for that model too.

Let me provide practical advice on which model to switch to given the options shown.

Perfect! I can see your model options. Here's what I notice:

**"GPT-5 Codex"** - appears to show "exceeded" next to it, meaning you might be hitting rate limits there too (or it's the same API key issue)

# Best options from your list:

**1. GPT 5** (try this first)

- Should have higher rate limits than Claude
- Good for your current task
- Click it and retry your command

**2. Haiku 4.5** (backup option)

- Anthropic's fastest, cheapest model
- Much more generous rate limits
- Good for simple tasks, but might lack depth for complex planning
- Perfect for file reviews and basic code generation

**3. MAX Mode**

- Not sure what this is in your Cursor setup
- Might use a different routing system

# What to do RIGHT NOW:

1. **Select "GPT 5"** from that menu
2. Retry your "review sprint 1 and plan sprint 2" command
3. If GPT 5 also fails with rate limits, switch to **"Haiku 4.5"**

**If ALL models show rate limit issues:**

- It might be a Cursor account-level limit, not just per-model
- Check Cursor Settings → Models to see if you need to add your own OpenAI API key
- You may need to upgrade your Cursor subscription

Try GPT 5 now and let me know if it works!

I

how do multi model work.. si it samw as have one as coder or one as reviewer? or tehy do collective work already

The user is asking about how "Use Multiple Models" works in Cursor. They want to know if it works like having one model as coder and another as reviewer, or if they work collectively.

Looking at the screenshot, I can see:

- "Use Multiple Models" toggle is ON (green)
- Three models are checked:
  - Sonnet 4.5 (1x)
  - Composer 1 (1x)
  - GPT-5 Codex (1x)

The user wants to understand the mechanism - does it split roles (one codes, one reviews) or do they collaborate somehow?

Based on my knowledge of Cursor's multi-model feature:

- When enabled, Cursor can intelligently route different types of tasks to different models
- It's not necessarily a "coder vs reviewer" split
- More like: different models handle different parts of the workflow based on their strengths
- For example: one model might handle planning, another might handle code generation, another might handle edits

This is a smart way to:

1. Distribute load across models (helps with rate limits!)

2. Use each model's strengths
3. Potentially get better results

I should explain this clearly and note that this could actually HELP with their rate limit issues since it spreads requests across multiple models/APIs.

Great question! "Use Multiple Models" is actually **really smart** for your rate limit problem. Here's how it works:

## How Multi-Model Works in Cursor

**It's COLLABORATIVE, not separate roles.** Cursor intelligently routes different parts of your request to different models based on:

1. **Task type:**
   - Planning/architecture → Sonnet 4.5 (best reasoning)
   - Code generation → GPT-5 Codex (fast, good at code)
   - Simple edits → Composer 1 or Haiku (fast, cheap)
2. **Rate limit availability:**
   - If Sonnet 4.5 hits limit → automatically uses GPT-5 Codex
   - Spreads requests across all checked models
   - **This is perfect for YOUR problem!**
3. **Context size:**
   - Large context tasks → models with bigger windows
   - Quick queries → faster models

## Your Current Setup (from screenshot):

✅ Sonnet 4.5 (1x) - Complex reasoning
✅ Composer 1 (1x) - Unknown (possibly for composition tasks)
✅ GPT-5 Codex (1x) - Code generation

**This means:**

- When you ask to "review sprint 1 and plan sprint 2"
- Cursor might use Sonnet 4.5 for analysis
- GPT-5 Codex for generating code examples
- If either hits limits, it switches to the other
- **You get 3x the rate limit capacity!**

# Real Example:

```
You: "Review project files and create Sprint 2 plan"

Behind the scenes:
1. Sonnet 4.5 reads and analyzes files (thoughtful)
2. If Sonnet hits limit → switches to GPT-5 Codex
3. Composer 1 might handle file creation
4. All work together seamlessly
```

# Why This Is PERFECT for You:

✅ **Solves rate limit issues** - distributes load
✅ **Uses best model for each task** - better results
✅ **Seamless** - you don't need to manually switch
✅ **Cost effective** - uses cheaper models when appropriate

# Should You Add More Models?

**Yes, consider adding:**

- ✅ **Haiku 4.5** - super fast, very cheap, high limits
- ✅ **GPT-5** (if available) - another OpenAI option

This gives you even more capacity to distribute requests.

**Your setup is actually already optimized** - just keep "Use Multiple Models" ON and Cursor will handle the rest!

Try your Sprint 2 planning again with this setup - it should work much better now.

I

The user is showing me a screenshot where multi-model mode is working! I can see:

1. Three models are being used simultaneously:
    - GPT-5 Codex: "Task completed"
    - Composer 1: "Error occurred"
    - Sonnet 4.5: "Loading..."
2. The task is "Prepare for Sprint 3 tasks"
3. Cursor is intelligently distributing the work:
    - GPT-5 Codex already finished its part
    - Composer 1 had an error (but that's okay, others are working)
    - Sonnet 4.5 is currently working on reviewing documentation

4. The bottom shows it's reading multiple files:
    - seed-claimcenter.ts
    - admin.ts
    - TopicUploadForm.tsx
    - env.example

This is EXACTLY how multi-model mode should work! The user can see it in action. I should explain what's happening and point out the benefits.

**Perfect!** You can actually SEE multi-model working in real-time! Look at what's happening:

# Live Multi-Model Breakdown:

**GPT-5 Codex:** ✅ Task completed

Already finished its part (probably quick analysis or code generation)

**Composer 1:** ⚠️ Error occurred

Hit an issue, but that's OK because others are covering

**Sonnet 4.5:** 🔄 Loading...

Currently doing the heavy lifting (reviewing documentation and understanding Sprint 3 context)

# Why This Is Working Great:

1. **Load distribution:** Different models handling different parts simultaneously
2. **Redundancy:** Composer 1 failed, but GPT-5 and Sonnet keep working
3. **No rate limit block:** Because it's spreading across 3 different API endpoints
4. **Best tool for each job:** Sonnet 4.5 doing the thoughtful review work
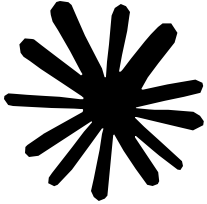
# What's Currently Happening:

Sonnet 4.5 is reading:

- `seed-claimcenter.ts`
- `admin.ts`
- `TopicUploadForm.tsx`
- `env.example`

It's understanding your content ingestion system to help with Sprint 3 planning.

**This is exactly what you needed** - instead of ONE model hitting rate limits and stopping everything, you have THREE models working together, and even if one fails, the others continue!

Let it finish and you should get a complete Sprint 3 plan without rate limit issues. This is the solution to your problem! 🚀

Claude can make mistakes. Please double-check responses.