Generative AI · Follow publication

# Your AI Coding Assistant Is Wasting 40% of Your Tokens. Here's How to Fix It.

DeepContext brings semantic search to large codebases, slashing token costs and making Claude Code actually useful for real projects

7 min read · 21 hours ago

Adham Khaled  ( Follow )

▶ Listen      ⬆ Share      ••• More

You're five hours into debugging, and your AI coding assistant has burned through $50 worth of API tokens. Again.

The problem isn't the AI model. It's what you're feeding it. When you ask Claude Code or Cursor to "find authentication logic" in your 10,000-file codebase, traditional grep dumps hundreds of irrelevant matches into the context window. The AI wades through test files, config snippets, and comments before finding what you actually need.

Every wasted token costs money. Every irrelevant code chunk increases latency. And every hallucination sends you down another debugging rabbit hole.

DeepContext changes this equation completely.

## What DeepContext Actually Does

DeepContext is an open-source Model Context Protocol (MCP) server that brings symbol-aware semantic search to AI coding agents. Instead of flooding your AI with keyword matches, it understands what your code means and returns only the five most relevant snippets.

Think of it as the difference between asking a librarian "show me books with the word 'dragon'" versus "find fantasy novels about dragon trainers." One gives you cookbooks and biology textbooks. The other gives you exactly what you need.

Here's what makes it different:

- Combines vector embeddings with BM25 lexical search and Jina reranking for hybrid accuracy

- Uses Tree-sitter AST parsing to understand code structure at the compiler level

- Reduces token usage by approximately 40% compared to grep-based search

- Operates 50% faster while returning more relevant results

- Works seamlessly with Claude Code, Cursor, and any MCP-compatible AI client

The best part? **It's completely free and open-source.**

## Why Traditional Search Fails AI Agents

When you grep for "PaymentProcessor" in your codebase, you get everything. Class definitions, test mocks, type annotations, comments, documentation, and that TODO note from three years ago.

Your AI agent doesn't know which matters. So it processes everything, burning tokens to analyze code that's completely irrelevant to your actual query.

I learned this the hard way working on a V2X collision avoidance system. Our embedded C codebase had 15,000+ files spanning CAN protocols, RTOS implementations, and wireless communication stacks. When I asked my AI assistant to "find where we handle authentication timeouts," grep returned 847 matches.

The AI dutifully processed them all. My token budget didn't survive.

## How Semantic Search Actually Works

DeepContext doesn't just match text patterns. It understands relationships between code elements.

When you search for "authentication logic," it recognizes that involves:

- Login handlers and credential validation

- Session management and token generation

- User database interactions

- API authentication middleware

- Error handling for auth failures

It retrieves functions that perform these operations, even if they never use the word "authentication." This is possible because of three key technologies working together.

Tree-sitter AST parsing breaks your code into a structured tree that captures functions, classes, imports, and their relationships. It knows the difference between a function definition and a function call.

Vector embeddings capture semantic meaning. Code that does similar things generates similar embeddings, regardless of naming conventions.

BM25 lexical search handles cases where you want exact identifier matches — useful when you know the specific class or function name.

A Jina reranking layer synthesizes these approaches, ordering results by actual relevance rather than just similarity scores.

## Getting Started in 15 Minutes

DeepContext installation is straightforward if you're comfortable with Python and have an MCP-compatible AI client.

Prerequisites:

- Python 3.9 or higher

- Node.js and npm

- Claude Code, Cursor, or another MCP-enabled development tool

- Git for repository management

Installation Steps:

```
# Clone the repository
git clone https://github.com/Wildcard-Official/deepcontext-mcp.git
cd deepcontext-mcp

# Set up Python virtual environment
python -m venv venv
source venv/bin/activate  # On Windows: venv\Scripts\activate
# Install dependencies
pip install -r requirements.txt
```

Configure Your MCP Server:

Create a configuration file for your AI client. For Claude Code:

```
# Add DeepContext to your MCP servers
claude mcp add deepcontext /path/to/deepcontext-mcp/server.sh

# Verify installation
claude mcp list
```

Index Your Codebase:

The beauty of incremental indexing is that you only need to do this once.
DeepContext automatically updates when files change.

Navigate to your project directory and trigger the initial index through your AI
client. The first query you make will start the indexing process automatically. For a
10,000-file repository, expect this to take 10–15 minutes. Subsequent updates process
only modified files.

## Real-World Usage Examples

Let me show you the difference between traditional grep and DeepContext with
actual scenarios I've encountered.

Finding Cross-Cutting Concerns:

Instead of: `grep -r "logging" . | wc -l` (returns 2,847 matches)

Use: "Show me all places where we log user authentication events"

DeepContext returns the five functions that actually log auth events, not every file that imports the logging library.

Understanding Dependencies:

Query: "Find all code that depends on the UserRepository class"

DeepContext traces actual usage through imports, method calls, and inheritance. It distinguishes between type annotations and actual instantiation. This is critical when refactoring — you need to know where code actually runs, not just where types are mentioned.

Discovering API Patterns:

Query: "Show me examples of how the EmailService is used"

Instead of searching for class names, DeepContext returns actual usage contexts with surrounding code. Your AI can learn your team's error handling patterns, retry logic, and best practices directly from real implementations.

## Current Limitations You Should Know

DeepContext focuses on Python and TypeScript. If your stack is Java, Go, Rust, or Ruby, you'll need to wait for expanded language support or contribute parsers yourself.

The tool requires MCP-compatible AI clients. Standalone LLMs without agent frameworks won't work without additional integration.

Initial indexing of massive monorepos takes time and computational resources. Once set up, incremental updates are fast, but that first index can be a barrier for 50,000+ file repositories.

Semantic search quality depends on embedding model training data. Niche frameworks or domain-specific code patterns might not be well-represented.

## Why This Matters for Your Team

Token costs are just the visible expense. The real cost is developer time.

Every query that returns irrelevant results interrupts your flow. Every hallucination from insufficient context requires debugging. Every manual search through your

codebase to find "that function that does the thing" burns cognitive energy you could spend solving actual problems.

DeepContext optimizes for the right metric: getting developers the right information immediately.

At scale, this compounds. A team of ten developers each saving 30 minutes daily adds up to 100+ hours monthly. That's time spent shipping features instead of fighting tools.

## Getting the Most Out of DeepContext

Start with exploratory queries when learning a new codebase. Ask "how does the checkout flow work" instead of trying to guess file names.

Combine natural language with specific identifiers when you know them. "Authentication using JWT tokens" works better than just "auth."

Use it for impact analysis before refactoring. "Where is PaymentProcessor instantiated" reveals every place you'll need to update.

Monitor your token usage before and after adoption. The 40% reduction is an average — your results may vary based on codebase size and query patterns.

## The Open Source Advantage

DeepContext's MIT license means you own your installation completely. No usage limits, no pricing tiers, no vendor lock-in.

The codebase is auditable for security-conscious teams. You can customize search algorithms, add language support, or integrate with internal tools.

Wildcard offers a hosted version with free trial tokens if you want to test before self-hosting. But the power of open source means you're never dependent on their infrastructure.

## What Comes Next

The MCP ecosystem is exploding with new tools and integrations. DeepContext benefits from this standardization — as new AI coding agents adopt MCP, they automatically work with DeepContext.

Community requests include Ruby/Rails support, additional language parsers, and multi-repository indexing for microservices architectures.

The broader trend is clear: AI-assisted development is moving from "better autocomplete" to "intelligent context-aware assistants." Tools like DeepContext represent the infrastructure layer that makes this possible.

## Making the Switch

If you're spending more than $20 monthly on AI coding assistant API calls, DeepContext will pay for itself in saved tokens within weeks.

If your team wastes time searching large codebases, the productivity gains justify the setup investment immediately.

If you're building AI-powered development tools yourself, DeepContext demonstrates how hybrid search architectures outperform pure semantic or pure lexical approaches.

The 30-minute setup is worth it. Your token budget will thank you.

This story is published on Generative AI. Connect with us on LinkedIn and follow Zeniteq to stay in the loop with the latest AI stories.

Subscribe to our newsletter and YouTube channel to stay updated with the latest news and updates on generative AI. Let's shape the future of AI together!

Software Development    Programming    Coding    Artificial Intelligence

Technology

# Published in Generative AI

62K followers · Last published 20 hours ago

Stay updated with the latest news, research, and developments in the world of generative AI. We cover everything from AI model updates, comprehensive tutorials, and real-world applications to the broader impact of AI on society. Work with us: jimclydegm@gmail.com

# Written by Adham Khaled

87 followers · 88 following

Embedded Systems Engineer || AI & Tech enthusiast || https://linktr.ee/adhamhidawy

## No responses yet

Bgerby

What are your thoughts?

## More from Adham Khaled and Generative AI

In AI Mind by Adham Khaled

## The Developer's Co-pilot is Dead: How Factory AI's Droid Ushers in the Age of Agent-Native...

Forget code completion—Droid autonomously handles entire software development workflows while beating Claude Code and GPT-5 on industry...

Sep 29 · 👏 66

In Generative AI by Dr. GenAI

## Hierarchical Reasoning Model (HRM): a tiny brain that embarrasses giant LLMs

HRM is a 27M-parameter model trained in hours that beats giant LLMs on reasoning puzzles, proving architecture can trump scale.

✦ Sep 7 · ✋ 306 · 💬 6

In Generative AI by Joe Njenga

## 15 Ways I'm Using Google's Nano Banana for UI/UX Design (Like a Pro)

I am not a UI/UX designer, but Google Nano Banana is filling the gap, and that doesn't mean I'm becoming one overnight.

✦ Oct 9 · ✋ 500 · 💬 6

See all from Adham Khaled

See all from Generative AI

## Recommended from Medium

Joe Njenga

### 8 Little-Known Books Every AI Founder Should Read First (Before Building a Unicorn)

We all have dreams, but few reach the finish line. Why?

In AI Advances by Pooja Kashyap

## Speech-to-Meaning: Why the Future of Voice AI Isn't About Better Microphones

"Tea, Earl Grey, Hot", remember this line by Picard from Star Trek?

In Coding Beauty by Tari Ibaba

## Yet Another Claude Model Just Shocked The World — Faster Than Sonnet 4.5 😮

Anthropic just released another incredible coding model — and the speed is unbelievable.

In CodeX by AI Rabbit

## Anthropic Combined 3 Years of AI Lessons Into One Feature

Anthropic just released a new feature that might change the way we interact with AI entirely. After years of experience with generative AI...

Bytefer

## A 0.9B Open-Source Model for SOTA Document Parsing: Outperforming GPT-4o and Gemini 2.5 Pro

Beyond OCR: Parsing text, tables, formulas & charts in 109 languages. Get SOTA document parsing that runs locally.

✦  5d ago  👋 36  💬 3

In **MeetCyber** by **Jules May**

## Vibe coding: the antidote

LLMs have shown themselves to be a terrible way to write programs. But the problem they address is real, and we aren't solving it any other...

✦  5d ago  👋 438  💬 29

See more recommendations