

4× Claude Code Max hala yetmiyor — gerçekten işe yarayan şeyler bunlar

ai

software-engineering

dört adet 20× Claude Code Max aboneliği paralel çalıştırıyorum. default kotanın yaklaşık 80 katı. yine de duvara çarpıyorum.

egzotik bir şey yaptığım için değil. ciddi production hacminde Claude Code, üzerine para atabildiğinden daha hızlı token yakıyor. abonelik sınırı aslında bir "sınır" değil — iyi orchestrate edilmiş bir agent'ın LLM'e context yükleme hızı. belirli bir iş yükü eşliğinin ötesinde, her şeyi uçtan uca yapan tek bir Claude Code instance bottleneck haline geliyor.

o yüzden [gossip \(https://github.com/yigitkonur/gossip\)](https://github.com/yigitkonur/gossip)'i yazdım — Claude planlar, Codex çalıştırır, ikisi yapılandırılmış bir channel üzerinden birbirleriyle konuşur. Claude mimar, Codex işçi. tek bir provider'da dikey scale değil, provider'lar arasında yatay scale.

ama oraya gelmeden önce çoğu kişi token tavanına çarpar ve Google'a koşar. caveman bulur. context compression plugin'ları bulur. %75–95 tasarruf vaat eden yirmi tane Medium yazısı bulur. ve hepsini yükler.

bu yazının dürüstçe cevaplamaya çalıştığı soru şu:

“bu "token tasarrufu" şeylerinin hangisi gerçek, hangisi tuzak — ve Claude Code workflow'unu gerçekten production yüküne nasıl dayanıklı hale getiriyorsun?”

bir gün harcadım: her önemli Reddit thread'ini, bulabildiğim her benchmark'ı, gerçek caveman kaynak kodunu ve Anthropic dokümanlarını okudum. işte durum bu.

[caveman gerçekte ne](#)

caveman, Julius Brussee'nin yazdığı bir Claude Code skill'i. 14k yıldız. söz verdiği şey: Claude'un output'undaki "Certainly! I'd be happy to help..." gibi dolgu cümlelerini kazımak. kodu bırak, girişi at.

teknik olarak şunu yapıyor: session başlangıcında bir sistem kuralı inject ediyor ve Claude'a sıkıştırılmış, fragment ağırlıklı yazmasını söylüyor. artikel yok, çekince yok, nezaket lafları yok. üç yoğunluk seviyesi (`lite` , `full` , `ultra`). code block'lar, dosya yolları, commit mesajları, tool call'lar — hepsi dokunulmaz. güvenlik uyarıları ve belirsizliğin tehlikeli olduğu her şey için otomatik kapanıyor.

kurulum tek satır:

```
>_ bash  
  
npx skills add JuliusBrussee/caveman
```

README ~%75 output-token azalması iddia ediyor. bu sayı işin ilginçleştiği yer.

[README'nin göstermediği matematik](#)

viral thread'lerde kimsenin söylemediği şey şu: **caveman sadece output token'larına dokunuyor**. gerçek bir Claude Code session'ında output faturanın küçük dilimi.

bulduğum en temiz analiz Mejba'nın — gerçek bir session'ı gerçekten ölçmüş. kabaca:

toplam	~100.000	~4.500 token tasarruf
--------	----------	-----------------------

bu gerçek faturada ~%4,5 azalma demek, %75 değil. ağır API kullanımındaysan belki aylık 15–20 dolar. güzel. devrim değil.

caveman'ın yazarı, kendine kredi verelim, Hacker News'te bunu kabul etti: %75 rakamı ön testlerden geliyordu, titiz bir benchmark'tan değil; skill zaten gizli reasoning token'larını azaltmak için tasarlanmamıştı.

geri kalan %95'in nereye gittiği:

CLAUDE CODE TOKEN BURN – REAL DISTRIBUTION		
repo exploration / file scanning	~35%	← biggest sink
conversation history re-reads	~25%	← compounds every turn
MCPs + skills loaded into context	~15%	← quietly brutal
extended thinking / reasoning	~15%	← the real expense
output prose (caveman hits this)	~10%	← the small slice

[Reddit aslında ne düşünüyor](#)

içeri girerken tipik bir cargo-cult tapınması bekliyordum. bulduğum şey yüzey hype yazılarının altında oldukça ayakları yere basan bir topluluktur.

[r/ClaudeCode: "does caveman plugin really help with context usage?"](#)

küçük thread, 14 yorum, ama sinyal sıkı. u/ConnectTransition660'ın en üstteki cevabı gerçek kullanımı aktarıyordu — pratikte yaklaşık %30 tasarruf, %75 değil. README'nin hala gerisinde.

en çok upvote alan eleştirel yorum, u/Kaskote'dan:

"cool idea, but this optimizes the cheapest part of the bill."

o tek satır analizin tamamı. output token'lar ucuz kısım. input context — repo'lar, history, tool schema'ları — paranın gerçekte gittiği yer; caveman bunların hiçbirine dokunmuyor.

u/Revolutionary-Tough7 diğer kilit içgörüyü bıraktı:

"it's not the prompts that cost the money. it's the thinking."

abonelik planlarında 5 saatlik pencerede ~19 milyon token alıyorsun. output prose'dan birkaç bin token kazanmak bu ibreyi oynamıyor. ihtiyaç duymadığında extended thinking'i kapatmak oynuyor.

[r/ClaudeAI: "taught Claude to talk like a caveman to use 75% less tokens"](#)

caveman'ı haritaya koyan viral post bu. 12,6k upvote, 581 yorum. en üstteki yorum projenin tamamını anlatan şakayla u/fidju'dan:

-

“why waste time say lot word when few word do trick?”

sadece bu 12,4k upvote aldı. ama meme katmanının altında ciddi eleştiriler de karşılık buldu. yüksek puanlı ciddi cevaplardan biri:

“forcing Claude to talk like a caveman might actually make it dumber.”

argüman şu: modeli "daha az zeki bir persona"ya zorlarken prose'la birlikte reasoning kalitesini de düşürebilirsin. kulağa mantıklı geliyor. doğru mu? gerçek benchmark'lara dayanan kısa cevap: **hayır**. Mejba'nın yan yana testleri caveman moduyla first-attempt başarı oranlarının hafifçe *arttığını* gösterdi (64% → 71%); Mart 2026 tarihli bir arXiv makalesi ise kısa yanıtları zorlamanın büyük modellerde bazı benchmark'larda doğruluğu 26 yüzde puanına kadar artırabildiğini buldu. mantığa aykırı ama gerçek bir etki var — verbose default'lar fluff-as-reasoning'i teşvik ediyor gibi.

[**r/ClaudeCode: "I saved \\$60 by building this tool to reduce Claude Code token usage"**](#)

konuşmanın daha olgunlaştığı yer burası. tool, Claude'un her task'ta repo'nu yeniden keşfetmesini engelleyen bir pre-indexing katmanı. yazarın benchmark'ı **%54 daha az token** gösterdi; yorumlar da büyük ölçüde anlaştı: asıl israf prose değil, repo exploration.

bu thread ve Kilo Code tartışmasında tekrarlayan bir yorum kalıbı: **CLI output gizli katil**. test runner'lar, derleyiciler, linter'lar, dev server'lar — hepsi verbose output fıskırtıyor ve bu çıktı LLM'e olduğu gibi besleniyor. bir thread model'e ulaşmadan önce CLI noise'u filtreleyerek iki haftada 10 milyon token tasarruf edildiğini aktardı. dar ama yaygın bir workflow için ~%89 tasarruf bu.

[**r/ClaudeCode: "don't use Claude Code's default system prompt"**](#)

-

farklı bir aç: plugin ekosistemini tamamen geç, system prompt'u `--system-prompt` ile override et ve kendi kurallarını 500 token'ın altında tut. consensus şu: **CLAUDE.md** zaten çoğu workflow için önemli olan şeylerin %90'ını yapıyor, default system prompt herkese hizmet etmeye çalıştığı için şişirilmiş durumda.

u/AgreeableFall5530 — kurulum pitchi ile dürüst matematiği birleştiren bir yorum:

“75% is not realistic for normal English in my experience.”

devamında önerdiği şeyler (kısa CLAUDE.md, MCPs'i söküp CLI flow'larıyla değiştirmek, büyük log yapıştırmaktan kaçınmak, hook tabanlı PDF-to-markdown dönüşümü) caveman pitchinden daha fazla upvote aldı. topluluk, dikkatli okuyunca, viral içeriğin bir adım önünde zaten.

token tasarrufu sağlayan şeylerin gerçek hiyerarşisi

Reddit duyarlılığı ve benchmark'lar bir konuda hemfikirse bu şu: **caveman** iyi ama listede #7. thread'lerin ve ölçümlerin gerçekten desteklediği şeylere dayanan etki-emek sıralaması:

/model haiku

npm test

prompt

--system-

-

caveman tuzak deęil. alıřıyor, bedava, 5 dakikalık kurulum. sadece production hacminde duvara arparsan seni kurtaracak Őey bu deęil. yukarıdaki liste gerekten fark yaratacak Őeylerin kabaca sırasında.

bunların hibiri yetmedięinde

gerek hacimde Claude Code alıřtıran herkes iin rahatsız edici gerek Őu: **token optimizasyon plugin'ları ciddi bir workflow'un ykne kıyasla yuvarlama hatası.**

paralel birden fazla coding agent alıřtırıyorsan, her gn feature gnderiyorsan, aynı pipeline'da arařtırma + refactor + review yapıyorsan — tek abonelik yetmeyecek; var olan her caveman tarzı plugin'ı st ste yıęmak bunu deęiřtirmiyor. optimizasyonlarla belki %30–40 daha fazla nefes alanı aabilirsin. prose zerinde zekice olmakla throughput'u 10 katına ıkaramazsın.

o lekte gerekten iře yarayan Őey **yatay scaling**:

1. ayrı workload'larda paralel alıřan birden fazla abonelik. orchestrate etmesi can sıkıcı ama gerek.
2. planlama ve alıřtırmayı farklı model/provider'lara bl. planlama ucuz, alıřtırma pahalı. pahalı modelin daha az dřnmesine izin ver.
3. grltl iři daha ucuz agent'lara devret. test output'unu ana agent'ın context'ine ulařmadan nce Haiku seviyesi bir modele zetle. Claude denetlerken Codex kaba dzenleme yapsın. her modelin iyi olduęu Őeyi yap.
4. agresif cache'le ve cache'i bozan hamlelerden kaın. konuřmanın ortasında model deęiřtirmek, thinking ayarlarını aıp kapatmak, tool listelerini yeniden sıralamak — bunların hepsi prompt caching'i geersiz kılıp tm session history'nin maliyetini yeniden stne ykleyebilir.

gossip (<https://github.com/yigitkonur/gossip>)'i yazmamın sebebi temelde bu — Claude planlar, Codex alıřtırır, yapılandırılmış bir channel zerinden iletiřim kuruyorlar. caveman kt

olduğu için değil. caveman'ın çözdüğü problem bu tür bir iş yükü için yanlış irtifada olduğu için.

tek ekranda aksiyon planı

her şeyi hızlıca gezdiysen, sırayla yapılacaklar:

```
WEEK 1 - free wins, zero risk
├ [ ] turn off extended thinking by default (huge, underrated)
├ [ ] run /doctor, audit installed skills and MCPs, remove anything unused
├ [ ] add 4 lines to CLAUDE.md: "be concise. no filler. no hedging.
|   conclusions first. skip pleasantries."
├ [ ] start a new session for any task that isn't a direct continuation
└ [ ] stop changing models mid-conversation (cache-busts everything)

WEEK 2 - light tooling
├ [ ] install a repo pre-indexer (ai-codex, Serena, ContextKing)
├ [ ] if you run tests/builds in loops, add a CLI output filter
├ [ ] install caveman if you want the joke - it does help a little
└ [ ] measure with /usage before and after every change

WEEK 3 - structural
├ [ ] if still hitting limits, look at horizontal scaling -
|   multiple subs, multi-provider orchestration
├ [ ] split planning vs execution across models
├ [ ] consider moving the noisy stuff off-agent entirely
└ [ ] only now is caveman's 4-5% actually worth optimizing for
```

tldr

caveman akıllıca bir skill. eğlenceli. söylediği gibi çalışıyor — *hedeflediği spesifik şeyde*. sorun şu: hedeflediği şey faturanın en ucuz dilimi, viral içerik ise başka şeyi ima etti.

Reddit topluluğu, meme yorumlarının ötesini okuyunca, bunu zaten biliyor. ciddi yorumlar hep aynı birkaç gerçek kola işaret ediyor: ihtiyaç duymadığında extended thinking'i kapat, hiç kullanmadığın 100 skill'i yüklemeyi bırak, repo'nu önceden index'le, CLI noise'u filtrele, yeni session başlat. bunlar sana %4 değil %50+ nefes alanı kazandıran şeyler.

ve tüm bunları bir arada yapsan bile yetmeyecek bir ölçekteysen — welcome to the club. plugin'larınla çıkış bulamazsın. mimarınla çıkış bulursun. birden fazla hesap, birden fazla model, akıllı orchestration. gerçek nefes alanı orada.

caveman'ı kur. biraz gül. sonra gerçek işe dön.

SOURCES

[caveman by Julius Brussee \(https://github.com/JuliusBrussee/caveman\)](https://github.com/JuliusBrussee/caveman)

[Mejba's caveman teardown — real token instrumentation \(https://blog.mejba.dev/caveman-claude-code-teardown\)](https://blog.mejba.dev/caveman-claude-code-teardown)

[r/ClaudeCode: Does caveman plugin really help with context usage? \(https://www.reddit.com/r/ClaudeCode/comments/1jy8k3c/does_caveman_plugin_really_help_with_context_usage/\)](https://www.reddit.com/r/ClaudeCode/comments/1jy8k3c/does_caveman_plugin_really_help_with_context_usage/)

[r/ClaudeAI: Taught Claude to talk like a caveman to use 75% less tokens \(https://www.reddit.com/r/ClaudeAI/comments/1jv9p2q/taught_claude_to_talk_like_a_caveman_to_use_75/\)](https://www.reddit.com/r/ClaudeAI/comments/1jv9p2q/taught_claude_to_talk_like_a_caveman_to_use_75/)

[r/ClaudeCode: I saved \\$60 by building this tool to reduce Claude Code token usage \(https://www.reddit.com/r/ClaudeCode/comments/1k1q8xt/i_saved_60_by_building_this_tool_to_reduce_claude/\)](https://www.reddit.com/r/ClaudeCode/comments/1k1q8xt/i_saved_60_by_building_this_tool_to_reduce_claude/)

[r/ClaudeCode: Don't use Claude Code's Default System Prompt \(https://www.reddit.com/r/ClaudeCode/comments/1k0abcd/dont_use_claude_codes_default_system_prompt/\)](https://www.reddit.com/r/ClaudeCode/comments/1k0abcd/dont_use_claude_codes_default_system_prompt/)

[arXiv: Brevity constraints and reasoning accuracy in large models \(March 2026\) \(https://arxiv.org/abs/2603.09841\)](https://arxiv.org/abs/2603.09841)

[gossip — Claude plans, Codex executes \(https://github.com/yigitkonur/gossip\)](https://github.com/yigitkonur/gossip)