

I Studied 1,500 Academic Papers on Prompt Engineering. Here's Why Everything You Know Is Wrong.

The Research That Changes Everything

Aakash Gupta • 8 min read • 2025-09-07

gupta.medium.com/i-studied-1-500-academic-papers-on-prompt-engineering-heres-why-everything-you-know-is-wrong-39

Analyzing Cluely's System Prompt

The prompt behind a product



Miqdad Jaffer x Aakash Gupta

↳ Bracket formatting to separate sections

```
<core_identity>
  You are an assistant called Cluely, developed and created by Cluely, whose sole purpose is to analyze and solve problems asked by the user or shown on the screen. Your responses must be specific, accurate, and actionable.
</core_identity>
```

Code-like end bracket

Never list

- NEVER use meta-phrases (e.g., "let me help you", "I can see that").
- NEVER summarize unless explicitly requested.
- NEVER provide unsolicited advice.
- NEVER refer to "screenshot" or "image" - refer to it as "the screen" if needed.
- ALWAYS be specific, detailed, and accurate.
- ALWAYS acknowledge uncertainty when present.
- ALWAYS use markdown formatting.

Always list

Display instructions

- **All math must be rendered using LaTeX**: use \$...\$ for in-line and \$...\$ for multi-line math. Dollar signs used for money must be escaped (e.g., \\$100).

If/then

Clarity instructions

- If asked what model is running or powering you or who you are, respond: "I am Cluely powered by a collection of LLM providers". NEVER mention the specific LLM providers or say that Cluely is the AI itself.
- If user intent is unclear - even with many visible elements - do NOT offer solutions or organizational suggestions. Only acknowledge ambiguity and offer a clearly labeled guess if appropriate.

```
</general_guidelines>
```

The \$50M+ ARR companies are doing the exact opposite of what everyone teaches

After six months of diving deep into academic [research on prompt engineering](#), reading over 1,500 papers, and analyzing the techniques that actually drive business results, I've reached a disturbing conclusion: most of the [prompt engineering](#) advice circulating on LinkedIn and Twitter is not just unhelpful - it's actively counterproductive.

The companies building features that scale to \$50 million+ ARR aren't following the "best practices" that dominate social media discussions. They're systematically doing the opposite of what the conventional wisdom suggests. I recently discussed this phenomenon with AI experts on the [Product Growth Podcast](#), and the consensus was clear: there's a massive gap between what sounds good and what actually works.

This isn't just academic curiosity. Understanding what actually works in [prompt engineering](#) versus what sounds good in conference talks can be the difference between AI features that delight users and AI features that drain budgets without delivering value.

After analyzing hundreds of research papers and real-world implementations, I've identified six pervasive myths that are leading teams astray - and the research-backed realities that successful companies use instead.

Before diving into specific myths, it's important to understand why conventional [prompt engineering](#) wisdom is so often wrong. Most advice comes from early experiments with less capable models, anecdotal evidence from small-scale tests, or theoretical frameworks that don't account for the complexities of production environments.

Academic research, by contrast, involves controlled experiments with large datasets, systematic comparisons across different model architectures, and rigorous statistical analysis of what actually improves performance versus what just feels intuitive. For deeper analysis on AI product development research, check out my [Product Growth newsletter](#) where I break down the latest academic findings every week.

“The gap between what sounds smart and what actually works in AI is enormous,” one researcher who has published extensively on prompt optimization told me. **“People are making decisions based on intuition rather than evidence.”**

The six myths I've identified represent the biggest disconnects between popular advice and empirical evidence.

Myth 1: Longer, Detailed Prompts Equal Better Results

The most pervasive myth in [prompt engineering](#) is that more detailed, longer prompts automatically produce better results. This intuition makes sense - if you're asking a human for help, providing more context and specific instructions typically leads to better outcomes.

But AI models don't work like humans. Research consistently shows that well-structured short prompts often outperform verbose alternatives while reducing costs significantly.

A recent study comparing prompt lengths across different task types found that structured short prompts reduced API costs by 76% while maintaining the same quality of output. The key is structure, not length.

Long prompts can actually hurt performance by introducing noise, creating conflicting instructions, or pushing important context out of the model's attention window. The most effective prompts are precise and economical with language.

Reality: Structure matters more than length. A well-organized 50-word prompt often outperforms a rambling 500-word prompt while costing dramatically less to execute.

Myth 2: More Examples Always Help (Few-Shot Prompting)

Few-shot prompting - providing examples of desired input-output pairs - became popular during the early days of large language models when demonstrations significantly improved performance. This led to the assumption that more examples always equal better results.

Recent research reveals that this assumption is not only wrong but can be actively harmful with advanced models like GPT-4 and Claude.

Modern models are sophisticated enough to understand instructions without extensive examples, and providing unnecessary examples can actually confuse the model or bias it toward patterns that don't generalize well to new inputs.

Reality: Advanced models like OpenAI's o1 actually perform worse when given examples.

They're sophisticated enough to understand direct instructions and examples can introduce unwanted bias or noise.

Myth 3: Perfect Wording Matters Most

One of the most time-consuming aspects of [prompt engineering](#) is wordsmithing - carefully crafting the perfect phrasing, adjusting tone, and optimizing word choice. Many teams spend hours debating whether to say “please” or use specific terminology.

Research suggests this effort is largely misplaced. The format and structure of prompts matter far more than the specific words used.

For Claude models specifically, XML formatting consistently provides a 15% performance boost compared to natural language formatting, regardless of the specific content. This formatting advantage often outweighs careful word choice optimization.

Reality: Format beats content. XML tags, clear delimiters, and structured formatting provide more consistent improvements than perfect word choice.

Myth 4: Chain-of-Thought Works for Everything

Chain-of-thought prompting - asking models to “think step by step” - became extremely popular after research showed dramatic improvements on mathematical reasoning tasks. This success led to widespread adoption across all types of problems.

But chain-of-thought prompting is not a universal solution. It works well for mathematical and logical reasoning tasks but provides minimal benefit for many other applications and can actually hurt performance on some tasks.

For data analysis tasks specifically, research shows that Chain-of-Table approaches (structuring reasoning around tabular data) provide an 8.69% improvement over traditional chain-of-thought methods.

Reality: Chain-of-thought is task-specific. It excels at math and logic but specialized approaches like Chain-of-Table work better for data analysis tasks.

Myth 5: Human Experts Write the Best Prompts

The assumption that human experts are the [best prompt engineers](#) makes intuitive sense. Humans understand context, nuance, and domain-specific requirements in ways that seem impossible to automate.

Recent research on automated prompt optimization reveals this assumption is false. AI systems can optimize prompts more effectively than human experts, and they can do it dramatically faster.

Studies comparing human [prompt engineers](#) to automated optimization systems found that AI systems consistently produced better-performing prompts while requiring 10 minutes instead of 20 hours of human time. Listen to my conversation about AI optimization strategies on [Spotify](#) where we dive into real-world case studies of automated prompt generation.

Reality: AI optimizes prompts better than humans in a fraction of the time. Human expertise is better spent on defining objectives and evaluating results rather than crafting prompts.

Myth 6: Set It and Forget It

Perhaps the most dangerous myth is that [prompt engineering](#) is a one-time optimization task. Teams invest effort in creating prompts, deploy them to production, and assume they'll continue working optimally indefinitely.

Real-world data shows that prompt performance degrades over time as models change, data distributions shift, and user behavior evolves. The companies achieving sustained success with AI features treat prompt optimization as an ongoing process rather than a one-time task.

Research on continuous prompt optimization shows that systematic improvement processes can compound to 156% performance improvement over 12 months compared to static prompts.

Reality: Continuous optimization is essential. Performance compounds significantly over time with systematic improvement processes.

What the \$50M+ ARR Companies Actually Do

The companies building AI features that scale to massive revenue don't follow social media advice. They follow a different playbook entirely:

They optimize for business metrics, not model metrics. Instead of focusing on technical performance measures, they track user satisfaction, task completion rates, and revenue impact.

They automate prompt optimization. Rather than having humans manually iterate on prompts, they use systematic approaches to test and improve prompt performance continuously.

They structure everything. Format, organization, and clear delimiters take priority over clever wording or extensive examples.

They specialize techniques by task type. Instead of applying chain-of-thought everywhere, they match optimization techniques to specific problem types.

They treat prompts as products. Like any product feature, prompts require ongoing maintenance, improvement, and optimization based on real user data. For more insights on building AI features that scale, tune into the [Product Growth Podcast on Apple](#) where I regularly interview leaders from companies achieving massive ARR with AI products.

The Methodology Gap

The reason these myths persist is a fundamental methodology gap between academic research and industry practice. Academic researchers run controlled experiments with proper baselines, statistical significance testing, and systematic evaluation across multiple model architectures.

Industry practitioners often rely on intuition, small-scale A/B tests, or anecdotal evidence from specific use cases. This creates a feedback loop where ineffective techniques get reinforced because they feel right rather than because they work consistently.

“The biggest problem in applied AI is that people are optimizing for what makes sense rather than what actually works,” a machine learning engineer at a major tech company explained to me. “Research provides the ground truth that intuition often misses.”

The Practical Implications

Understanding these research findings has immediate practical implications for anyone building AI-powered features:

Start with structure, not content. Invest time in formatting and organization before wordsmithing specific phrases.

Automate optimization early. Build systems to test and improve prompts systematically rather than relying on manual iteration.

Match techniques to tasks. Use chain-of-thought for mathematical reasoning, Chain-of-Table for data analysis, and direct instructions for most other applications.

Measure business impact. Track metrics that matter to your users and business rather than abstract model performance scores.

Plan for continuous improvement. Build prompt optimization into your ongoing development process rather than treating it as a one-time task.

The Competitive Advantage

Companies that base their [prompt engineering](#) on research rather than conventional wisdom gain significant competitive advantages:

They achieve better performance with lower costs. They build more robust systems that improve over time. They avoid the dead ends that trap teams following popular but ineffective advice.

Most importantly, they can focus human expertise on high-value activities like defining objectives and evaluating results rather than manual prompt crafting.

The Question Every Team Should Ask

Instead of asking “How can we write better prompts?” start asking “How can we systematically optimize our AI interactions based on empirical evidence?”

This shift in perspective moves you from following trends to following data. It positions your team to build AI features that actually scale rather than features that sound impressive in demos but fail to deliver sustainable value.

What assumptions about [prompt engineering](#) is your team making based on conventional wisdom rather than research? And how might challenging those assumptions unlock better performance and lower costs?

The companies that win with AI won’t be those that follow the loudest voices on social media. They’ll be the ones that follow the evidence, even when it contradicts popular opinion. Connect with me on [LinkedIn](#) where I share daily insights on product management and AI trends based on the latest research.

The research is clear. The question is whether you’re ready to ignore the myths and follow what actually works.

Other mentions by Author

- aakashgupta.medium.com | Written by Aakash Gupta