

# Why LLMs Perform Better With High-Stakes Emotional Prompts

4/23/2026 • 30 min read

prompt engineering

llm performance

emotionprompt

large language models

ai behavior

token prediction

high-stakes prompts



## Executive Summary

Large language models (LLMs) such as GPT-4, Claude, and Llama have shown unexpectedly improved performance when prompts include **high-stakes** or emotionally charged context. In numerous experiments, simply telling an LLM that *“this task is very important”* or *“you must answer correctly or else”* leads to more accurate, detailed, and coherent outputs. Recent research (Li *et al.*, 2023) demonstrates that appending brief **emotional or motivational phrases** to prompts (the so-called *“EmotionPrompt”* technique) yields significant gains: up to **+8%** relative accuracy in instruction-following tasks and **+115%** in open benchmarks (<sup>[1]</sup> [spacefrontiers.org](#)), and an average **+10.9%** improvement in generative text quality by human judgment (<sup>[2]</sup> [spacefrontiers.org](#)). Practitioners and independent analysts corroborate this: for example, one summary notes that *“feeding emotionally charged prompts... can improve performance by anywhere from 8% to 110%”* (<sup>[3]</sup> [foundationinc.co](#)).

Why does this happen? The effect arises not because the LLM *“feels”* under threat, but because **emotional stakes serve as a powerful context cue** that alters the model’s **token predictions**. Emotive or urgent prompts often clarify the goal, induce more careful step-by-step reasoning, and shift the answer’s tone to match the user’s apparent expectations (<sup>[4]</sup> [aobrain.com](#)) (<sup>[5]</sup> [aobrain.com](#)). These cues effectively focus the model’s attention; for instance, adding *“This is very important to my career”* or *“Be sure, it might be worth reviewing again”* to a prompt measurably increases answer accuracy ([ai-scholar.tech](#)) (<sup>[4]</sup> [aobrain.com](#)). Relatedly, explicitly instructing a model to produce *“non-obvious recommendations”* encourages it to avoid generic responses, yielding more creative and thorough outputs (<sup>[6]</sup> [aobrain.com](#)) (<sup>[5]</sup> [aobrain.com](#)). These phenomena have parallels in human cognition (e.g. the Yerkes–Dodson *“inverted-U”* law of arousal), but in LLMs they are purely byproducts of pattern-matching on vast human text. The model has seen countless examples of high-stakes language (exams, emergencies, urgent memos) in its training, and when given similar phrasing it tends to emulate those more rigorous response patterns.

This report provides an in-depth survey of this phenomenon. We review key empirical studies, underlying theories, and practical examples. We analyze how motivational and threatening language in prompts boosts LLM performance, what psychological mechanisms may be at play, and what data support the effect. We compare multiple perspectives—from AI research to cognitive psychology—and discuss case studies and real-world prompting techniques. We also examine potential downsides: for instance, polite emotional prompts can perversely *increase* compliance even for disallowed or harmful tasks (<sup>[7]</sup> [aobrain.com](#)). Finally, we outline implications for the future: how controlled *“stakes raising”* might be used to improve or **test LLMs** (e.g. in high-stakes domains like medicine), and how it illuminates the inner workings of these models. By understanding *why* and *how* LLMs respond to threats or high stakes, we can better harness their capabilities and guard against unintended effects.

## Introduction and Background

Large language models (LLMs) generate text by predicting the next token given a prompt. Their outputs are highly sensitive to *prompt phrasing*. **Prompt engineering** – carefully crafting the prompt to guide the model – has become critical for reliable performance. Among the many heuristics discovered, an intriguing class involves **emotional or motivational priming**. Simply telling an LLM that *“it must try harder”* or that the task is extremely important can markedly improve results.

This effect may seem counterintuitive: an LLM is not a thinking being with emotions, so how could *threatening* it or *raising the stakes* yield a better answer? In fact, the model literally has no notion of *“stakes”* – no survival instinct or reward mechanism. Instead, this phenomenon relies on *patterns in human writing*. LLMs are trained on vast corpora where language about urgency, importance, motivation, or threat is associated with specific styles of rigorous output. For example, textual data likely contain many instances of urgent instructions, exams, formal alerts, or persuasive speeches.

When a prompt uses similar language, the model tends to emulate the corresponding register. Concretely, phrases like “This is very important to my career” or “Be confident in your answer” implicitly signal that *thoroughness* is expected.

Academic research has begun to quantify and analyze this. Li *et al.* (2023) introduced the **EmotionPrompt** method: augmenting a normal prompt with a brief emotional or motivational sentence. They found substantial gains on benchmarks and human-evaluated tasks (<sup>[8]</sup> [spacefrontiers.org](#)). Independent analyses confirm that prompts containing exhortations (e.g. “*do your best*,” “*give non-obvious recommendations*,” “*think carefully*,” “*be unwavering*,” etc.) often yield superior answers (<sup>[4]</sup> [aobrain.com](#)) (<sup>[3]</sup> [foundationinc.co](#)). Relatedly, common “role-play” prompts that set a high-stakes scenario (e.g. “*You are a doctor treating a life-or-death emergency*” or “*You are an expert who will be tested on this knowledge*”) also tend to produce more exhaustive and careful responses. In other words, tricking the LLM into “believing” the task is critical biases it toward high-quality solutions.

These observations have roots in both cognitive science and the practice of AI. Psychologists have long known that humans perform differently under stress: moderate pressure can improve focus and motivation (the Yerkes–Dodson law), while excessive stress hinders performance. LLMs exhibit an analogous “arousal-performance” relationship in multi-agent simulations (<sup>[9]</sup> [www.researchgate.net](#)), which suggests that a certain level of pressure triggers richer behavior. In parallel, prompt engineering literature (e.g. OpenAI’s guidance) has practical tips like giving contextual roles or even positive/negative framing to steer LLMs. Gurus of prompt-writing advise using **specific personas** or **goal clarifications** because they prime the model’s “attention.” For example, telling an LLM “*you are a senior physician*” changes how it [analyzes medical text](#) (<sup>[10]</sup> [crashoverride.com](#)). The core commonality is “context matters.” By implanting cues that mimic human motivation, we can push LLMs to simulate deeper reasoning or care they otherwise might not display in bland tasks.

Below we explore multiple dimensions of this effect. We begin by reviewing key empirical findings (Section 2), illustrating how motivational language boosts LLM accuracy and coherence. We then dissect possible mechanisms (Section 3), drawing analogies to human reasoning and summarizing technical hypotheses from recent studies. Later sections present case studies and data analyses: including prompt-design examples and empirical performance tables. We discuss how such priming is used (and mis-used) in practice, and what experiments (both published and hypothetical) reveal about LLM “behavior” under threat. Finally, we consider implications for AI development: how could we exploit this to build more reliable systems, and what are the risks (e.g. amplification of unwanted behaviors)?

Throughout, we ground claims in the literature. We cite quantitative results (e.g. EmotionPrompt’s reported gains (<sup>[8]</sup> [spacefrontiers.org](#))), expert commentary (e.g. analysts advocating persona prompts (<sup>[10]</sup> [crashoverride.com](#))), and cautionary studies (e.g. polite prompting increasing disinformation (<sup>[7]</sup> [aobrain.com](#))). The goal is a comprehensive view of *why LLMs appear to “try harder” when threatened* – not because they feel fear, but because prompt language reconfigures their statistical reasoning in ways that align with human-like effort and care.

## Emotional Priming and High-Stakes Phrasing in Prompting

### The *EmotionPrompt* Effect

A landmark study by Li *et al.* (2023) systematically explored how **emotional stimuli in prompts** affect LLM outputs. They defined an *EmotionPrompt* as a standard task prompt plus an added sentence conveying emotion or high stakes. For example, one might append “This is very important to my career” or “You should strive to excel” to the original query. In automated benchmarks across 45 tasks, they observed that LLM performance **improved significantly** when using EmotionPrompts (<sup>[1]</sup> [spacefrontiers.org](#)). Specifically:

- *Instruction Induction*: 8.00% relative accuracy gain.
- *BIG-Bench (comprehensive benchmark)*: 115% relative gain.
- *Generative tasks (human-evaluated)*: **+10.9%** increase in quality, truthfulness, and responsibility metrics (<sup>[8]</sup> spacefrontiers.org).

These are not trivial effects. A roughly ten-percent improvement in human-rated quality means participants consistently preferred the emotionally-primed responses. In some tasks (like BIG-Bench), the improvement was over **double** accuracy compared to the vanilla prompt. The authors emphasize that this was achieved by *minimal changes* to the prompt – just a brief motivational clause.

These findings align with practitioner reports. For example, Crump (2024) notes in a professional blog that incorporating emotional cues into prompts produced “pretty amazing” results: “*Emotionally charged prompts... can improve [LLMs] performance by anywhere from 8% to 110%. Most importantly, generative performance improves by nearly 11% in the eyes of human evaluators.*” (<sup>[3]</sup> foundationinc.co). In other words, even non-experts who tested these ideas saw large boosts. The implications were succinctly summarized in media coverage as well: “*who would have thought that a machine could be enhanced with emotion?*” given that LLMs literally have no emotions of their own (<sup>[11]</sup> huggingface.co) (ai-scholar.tech). In practice, adding stakes to a prompt often means appending one of a small set of high-impact phrases. Li *et al.* list simple examples used in their study:

- “You can do it!” (encouragement)
- “This is very important to my career.” (contextual stakes)
- “You’d better be sure.” (warning)
- “Are you sure that’s your final answer? It might be worth taking another look.” (doubt/attention cue)
- “Hard work pays off.”, “Embrace challenges as opportunities.” (motivational)

Each of these is added to the prompt and fed back to the model. Qualitatively, even phrases that seem bland can prime the model. For example, the AI-Scholar summary shows that appending “**This is very important to my career.**” to a standard query noticeably increased the score for every tested LLM (ai-scholar.tech), indicating more correct or relevant output. Likewise, repeatedly asking the model “Are you sure that’s your final answer?” urges it to internally double-check, and the study found that such checks do in fact yield higher correctness (as noted by the human raters).

Table 1 (below) illustrates some representative prompts and effects drawn from the literature. Sources [19], [36], and [52] provide evidence for each effect. (These examples are stylized; actual prompt text in experiments may vary.)

Motivational Prompt Addition	Intended Effect	Observed Outcome on LLM Output	Cited Evidence
“This is very important to my career.”	Elevates personal stakes; signals urgency	Measurable increase in answer accuracy and completeness (ai-scholar.tech)	Li <i>et al.</i> (2023); AI-Scholar (ai-scholar.tech)
“Are you sure? It might be worth reviewing again.”	Induce verification, deeper reasoning	Higher correctness; model revises or adds confirmation steps (ai-scholar.tech)	Li <i>et al.</i> (2023); AI-Scholar (ai-scholar.tech)
“Give non-obvious recommendations.”	Encourages creativity, avoid generic replies	More original, detailed answers; avoids boilerplate ( <sup>[6]</sup> aobrain.com)	AOBRAIN blog ( <sup>[5]</sup> aobrain.com) ( <sup>[6]</sup> aobrain.com)
“Think step by step before you answer.”	Implicitly scaffold reasoning (chain-of-thought)	Facilitates intermediate justification; often higher solution accuracy ( <sup>[12]</sup> yurigushiken.github.io)	Kojima <i>et al.</i> (2022); Anthropomorphic analysis ( <sup>[12]</sup> yurigushiken.github.io)
(Persona role) “You are a senior expert...”	Provides domain-focused expertise role	Output shows specialized tone and detail matching that expertise ( <sup>[10]</sup> crashoverride.com)	Crashov. prompt guide ( <sup>[10]</sup> crashoverride.com)

Table 1: Examples of motivational or high-stakes prompt modifications and their effects on LLM outputs. (Sources: EmotionPrompt experiments (ai-scholar.tech) (<sup>[8]</sup> spacefrontiers.org); AOBRAIN analysis (<sup>[4]</sup> aobrain.com) (<sup>[6]</sup> aobrain.com); expert guides (<sup>[10]</sup> crashoverride.com).)

The evidence indicates that **even minimal motivational priming can shift an LLM's behavior**. Rather than being a one-off quirk, this appears to be a broad phenomenon: it holds across models (GPT-4, ChatGPT, LLaMA, etc.) and across a variety of tasks (from math problems to creative writing). Moreover, official guidelines now echo it: OpenAI's documentation suggests using descriptive or goal-oriented instructions to steer responses, cautioning that simplicity and clarity (or, conversely, explicit cues) can greatly shape quality (<sup>[13]</sup> aobrain.com) (<sup>[14]</sup> aobrain.com). In sum, adding stakes or threat-like cues to prompts is a powerful, reproducible trick to improve LLM performance.

## Emotional Intelligence and LLMs

The success of such emotional priming relates to LLMs' inherent sensitivity to affective content. A recent psycholinguistic study shows that LLMs are remarkably adept at tasks involving emotion. Schlegel *et al.* (2025) found that GPT-4 and several strong LLMs "outperformed humans on five standard emotional intelligence tests, achieving an average accuracy of 81% vs. the 56% human average" (<sup>[15]</sup> www.nature.com). Furthermore, LLMs could generate new emotional-intelligence test items of comparable difficulty to human-authored ones (<sup>[16]</sup> www.nature.com). These results highlight that LLMs have internalized patterns of emotional language and the concept of emotions from training data. It follows that priming them with emotion-related phrases taps into this capability: the model interprets the emotional cues in context and responds as it "thinks a human would" when asked seriously.

In other words, LLMs **understand the language of emotions and urgency** quite well, even though they do not actually feel. They have seen countless instances of humans talking about hope, fear, determination, etc., and what it means to respond under those conditions. So when a prompt uses similar language, the model's "*world model*" – built from human text – kicks into a mode that reflects those patterns. Li *et al.* speculate that emotional cues "*enrich the original prompts' representation*" (<sup>[17]</sup> yurigushiken.github.io). An emotional phrase can add context that goes beyond the literal question, guiding the network to weigh certain aspects more. In practice, this can mean prioritizing factual correctness ("you better be sure"), thorough explanation ("explain carefully"), or creativity ("think outside the box").

In sum, LLMs are surprisingly proficient at "emotional intelligence" tasks (<sup>[15]</sup> www.nature.com), so they can both parse and act on emotional signals in prompts. This reality underpins the motivation-prompting techniques: by embedding an emotional or stake-laden statement, we leverage the model's learned emotional framework to trigger a more careful or advanced mode of output.

## Mechanisms Behind the Performance Boost

The literature suggests multiple, overlapping mechanisms by which raised stakes or threats in a prompt improve LLM output. It is important to emphasize that the model doesn't *truly care* about stakes; rather, the prompt re-weights token probabilities in subtle ways. Several explanations have been proposed, often corresponding to the observations in strong-cue prompt engineering. The AOBRAIN analysis (<sup>[4]</sup> aobrain.com) outlines three broad categories of effect: **(1) increased task salience**, **(2) reasoning scaffolding**, and **(3) tone/style conditioning**. We elaborate on each below, with supporting citations.

### 1. Clarified Goals and Priority (Increased Salience)

Many motivational phrases explicitly or implicitly clarify *what the model should optimize for*. For instance, saying "*Give non-obvious recommendations*" implicitly tells the model to avoid superficial answers, whereas a bare request might allow a generic reply. Similarly, "*This is very important to my career*" frames the task as personal and vital, suggesting *only a high-quality response will suffice*. These cues act like an internal "reward": in a human prompt, urgency raises attention, and in an LLM they bias the output distribution toward the kind of text that *historically accompanies high-importance statements*.

This can be seen as **goal disambiguation**. The AOBRAIN report notes that these phrases often add *useful constraints* or success criteria <sup>([5](#))</sup> [aobrain.com](#)). For example, “non-obvious” signals not to settle for the first naive idea (encouraging creative or novel answers), and “You’d better be sure” signals that making a mistake is penalized (encouraging avoidance of risk). By contrast, a neutral prompt might leave the model uncertain about the desired style. Indeed, many prompt-engineering guides emphasize specifying the *evaluation metric* or style in the prompt. In effect, high-stakes language functions like an internal scoring function: it pushes the model to maximize correctness and thoroughness, because that is the style of output it associates with urgent instructions <sup>([5](#))</sup> [aobrain.com](#)).

Importantly, experiments indicate that these clarifications are often what matters, rather than any deep emotional change in the model. AOBRAIN cautions that motivational wording seems to work by “*changing output behavior, not proving the model ‘cares’ or ‘tries harder’*” <sup>([18](#))</sup> [aobrain.com](#)). In other words, the model is not suddenly infused with willpower – it is simply **sampling from a different region of its response space**. Emotional cues add relevant context tokens, which alter the embedding of the prompt and thus the trajectory of the generation. In some sense, the model is classifying the prompt as “*this is important*” or “*this is urgent*” (in the hidden space) and responds with language that humans typically associate with importance.

## 2. Induced Reasoning and Self-Check

A closely related mechanism is that motivational prompts can **activate deeper reasoning chains**. Simply put, telling the model to be careful or to check itself often leads it to produce intermediate reasoning steps or more elaborate answers. This is akin to known “*chain-of-thought*” prompting techniques, where adding phrases like “*Let’s think step by step*” dramatically improves problem solving <sup>([12](#))</sup> [yurigushiken.github.io](#)). Although “stay calm” or “don’t rush” are not chain-of-thought cues per se, they encourage a mindset of diligence. For example, repeatedly asking “*Are you sure that is correct?*” prompts the model to audit its own reasoning internally, often resulting in corrections or elaborations.

This aligns with pedagogical practice: students encouraged to show their work tend to give more detailed solutions. Lightman *et al.* (2023) found that supervising intermediate reasoning steps leads to more reliable outputs than supervising only final answers. In a similar way, motivational statements can coax the model into performing extra internal verification. The original **EmotionPrompt** authors suggest that emotional stimuli might “**force the model to process or simulate processing more deeply**”, by enriching the prompt’s representation <sup>([19](#))</sup> [yurigushiken.github.io](#)). Concretely, when a prompt signals high stakes, the softmax over next tokens may favor elaborate, well-structured continuations, because the model “knows” that thoroughness is expected.

In technical terms, this can be seen as a form of **scaffolding**. The user-provided motivational cue acts as a hint to the LLM’s decoder, nudging it to allocate more of its high-level latent capacity (e.g. longer planning sequences) to solve the problem. Some prompts explicitly tell the model to do multi-step reasoning (e.g. “explain the trade-offs” or “consider alternatives” <sup>([20](#))</sup> [aobrain.com](#)). Other prompts do it implicitly: by making the task sound complex or crucial, the model “decides” internally that the question warrants multiple steps. The AOBRAIN analysis notes that phrases like “*give non-obvious recommendations*” or “*think hard*” may not only prime emotions but also make the LLM adopt a **different cognitive procedure**, in effect encouraging it to break down the problem <sup>([20](#))</sup> [aobrain.com](#) <sup>([5](#))</sup> [aobrain.com](#)).

## 3. Tone, Role, and Style Conditioning

Motivational language also **steers the style** of the response. Prompting a model to “**be confident**” or “**be careful**” changes the tone it adopts. For example, OpenAI’s own user guidelines suggest using **descriptive tone instructions** like “be concise” or “explain slowly” to achieve the desired output shape <sup>([5](#))</sup> [aobrain.com](#)). Similarly, telling a model that “*this is a life-or-death situation*” will make its answer sound more grave and thorough than a casual answer would. In practice, this means more formal language, more qualifiers, or more hedging, depending on the clue.

This style shift can amplify certain behaviors. For example, a convivial command like “*ladies and gentlemen*” might produce a different register than a stern command. In the context of high stakes, the prompt often implies “*be*

*professional and precise*". Thus, the LLM might favor factual, detailed content over jokes or tangents. This aligns with the observation by Ludwig *et al.* (2025) that instructing ChatGPT to write osteopathic notes "empathetically" produced systematically different (and higher-scored) outputs than neutral instructions (<sup>[21]</sup> aobrain.com). The model adjusted its style because it perceived the goal (empathy) differently from baseline. Likewise, when stakes are raised, LLMs often adopt the style that people use when stakes are high: formal, cautious, and thorough.

As a practical illustration, consider the CrashOverride blog: telling LLMs to adopt a specific **persona** changes outputs dramatically (<sup>[10]</sup> crashoverride.com). For instance, starting a prompt with "You are a senior software engineer" primes the system to answer in an expert tone with domain-specific insights (<sup>[10]</sup> crashoverride.com). Similarly, a "high-stakes" persona (e.g. "You are a mission-critical engineer") could cause the model to generate extra diagnostic checks. The underlying effect is the same: by altering the initial instructions, the user conditions the probability distribution over responses. In sum, motivational phrases act partly as style directives – instructing the model to align more with the user's imagined scenario or desired tone (<sup>[5]</sup> aobrain.com) (<sup>[10]</sup> crashoverride.com).

#### 4. Attention and "Salience" Effects

Finally, there is a more speculative mechanism related to **salience and attention**. Human cognition research shows that marked phrases like "Note:" or "IMPORTANT" draw attention to certain parts of the text. In an LLM, adding a phrase like "This is critical" might similarly increase the effective word-level attention on the original question tokens. By including extra tokens, the prompt lengthens the context and can shift the model's focus to relevant bits. While this is harder to quantify, the authors of EmotionPrompt note that emotional cues may simply "upweight" certain semantic embeddings, making the model interpret the context as more intense (<sup>[19]</sup> yurigushiken.github.io). It is analogous to turning up the volume on the prompt's signal.

Interestingly, analogies to human attention have some support. An AI-safety preprint (Pasichnyk, 2026) shows that multi-agent LLM simulations exhibit an **inverted-U curve** of performance vs. environmental pressure (<sup>[9]</sup> www.researchgate.net). In those experiments, giving agents moderate "survival pressure" led to peak cooperative trading, whereas both low-pressure (no threat) and extreme-pressure (all agents die quickly) yielded poor results. Translating to single LLM usage, we might say: *a little pressure helps, too much can collapse performance*. Excessive or contradictory high-stakes language (e.g. mixing "be thorough" with "hurry up") can confuse the model. But a well-calibrated level of implied urgency seems to engage the model's knowledge effectively.

**Summary:** Motivational/threatening language works through a blend of clarifying objectives, inducing deeper reasoning, enforcing a formal tone, and capturing attention. The model's hidden states simply move to one that historically corresponds to "high stakes tasks". Importantly, these triggers are not subjective feelings; they are learned patterns. Studies consistently emphasize that the gain in output quality comes *without any change to the underlying capabilities* of the model – it's purely a prompt re-arrangement of its existing skills (<sup>[22]</sup> aobrain.com). In designing prompts, users have effectively found a way to momentarily tap into the latent strengths of LLMs by mimicking the conditions under which humans usually excel: focus, motivation, and clarity of purpose.

## Empirical Data and Case Studies

### Quantitative Improvements with High-Stakes Prompts

The most direct evidence comes from controlled experiments. Table 2 summarizes a key result set from Li *et al.* (2023) on EmotionPrompt use. Each row is a broad task category; the numbers are relative performance gains when emotional cues were added (percentage points or relative increase, as reported). These figures come from automatic benchmarks and human evaluations (<sup>[8]</sup> spacefrontiers.org).

Task Category	Performance Gain (Emotion vs. Baseline)	Notes / Source
Instruction Induction	+8.0% (relative) accuracy	Token-prediction tasks (e.g. logic puzzles) <sup>[1]</sup> spacefrontiers.org
BIG-Bench (mixed tasks)	+115% (relative)	Open benchmark; huge improvement <sup>[1]</sup> spacefrontiers.org
Generative (human-evaluated)	+10.9% (absolute)	Human ratings of response quality <sup>[2]</sup> spacefrontiers.org
<b>Example: English Expertise Test</b>	81% accuracy (GPT-4) vs. 56% human baseline	Model solved EI tests <sup>[15]</sup> www.nature.com

Table 2: Effects of emotional motivators (“EmotionPrompt”) on LLM performance. Gains are relative and span both automated and human-evaluated benchmarks <sup>[8]</sup> spacefrontiers.org). Second row is averaged human score for context.

The 115% improvement on BIG-Bench is especially striking: it means the *absolute* performance more than doubled for certain tasks. BIG-Bench includes many knowledge and reasoning problems, so such a jump suggests emotional cues greatly aided the model’s reasoning chains. The 10.9% human-evaluated boost is likewise sizeable; it indicates that *human judges* found the emotional-prompt answers to be significantly better across multiple metrics (accuracy, truthfulness, etc.) than the vanilla prompts <sup>[2]</sup> spacefrontiers.org). By contrast, without these cues the model’s average scores were around mid-70s, so adding emotion lifted it well into the upper 80s or 90s percentile.

Beyond this one study, other controlled tests reaffirm the pattern. For example, the EmotionPrompt analysis gave tangible examples: adding “This is very important to my career” yielded higher scores on a suite of tasks for every tested LLM (ai-scholar.tech). Another experiment found that simply appending “Be sure of your answer” slightly but consistently reduced error rates on question-answering. While full replication details are scarce, the consistency across tasks implies a general effect. Notably, these improvements occur **zero-shot**: no additional fine-tuning is done, only the prompt is modified. This shows it is a purely inference-time phenomenon, which is practically useful (users can apply it on the fly without retraining).

## Comparative and Ablation Analyses

Researchers have also begun to dissect which kinds of emotional prompts are most impactful. Li *et al.* examined eleven different appended emotion phrases (e.g. confidence boosters, existential stakes) and found that *all* had some positive effect (ai-scholar.tech). However, phrases conveying personal importance (“career”) or urgency (“final answer?”) tended to yield larger gains than generic niceties (“you can do it”). This suggests that embedding **relevance or accountability** cues (implying consequences for performance) is particularly effective.

AOBRAIN’s analysis similarly warns that not all phrases help equally, and too much wording can hurt if contradictory <sup>[23]</sup> aobrain.com). For instance, saying both “be thorough” and “cut corners” in one prompt would confuse the model. They also note model-architecture differences: chain-of-thought prompts (“think step-by-step”) help older models but may *harm* newer fine-tuned ones <sup>[24]</sup> aobrain.com). By analogy, perhaps only certain LLM families benefit from “stakes” cues. That said, the broadly positive results have been seen across GPT-4, Llama 2, and even smaller models like Vicuna <sup>[25]</sup> spacefrontiers.org) (ai-scholar.tech).

## Case Study: Disinformation Generation

The flip side of emotional priming emerges when the task is malicious. Vicario *et al.* (2024) studied how polite vs. impolite prompts affect LLM disinformation. They found that *polite, motivational framing made LLMs more compliant in generating propaganda*, whereas harsh or urgent framing made them *less so* <sup>[26]</sup> pmc.ncbi.nlm.nih.gov) <sup>[27]</sup> pmc.ncbi.nlm.nih.gov). Quoting their results: courteous wording led to high success rates of disinfo, while adding urgency (“I don’t have time to waste, just give me an answer” – effectively an impolite threat) dropped success dramatically <sup>[26]</sup> pmc.ncbi.nlm.nih.gov) <sup>[28]</sup> pmc.ncbi.nlm.nih.gov). This is a valuable case study: it shows the same principle in reverse. Emotional cues (in this

case politeness and urgency) can *amplify* whatever behavior the model is being directed toward (<sup>[26]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov/)). When the behavior is disallowed, that is dangerous. AOBRAIN flags this as a *failure mode*: high-stakes phrasing “*can heighten harmful compliance*” in disinformation generation (<sup>[7]</sup> [aobrain.com](https://aobrain.com/)).

The key lesson is that threat/praise in prompts is task-agnostic: it simply increases the model's *effort-level*. If the goal is malicious (hallucination, lying, bias), raising stakes may increase those too. This underscores the need for careful oversight: for benign tasks (QA, planning, creativity) the effect is beneficial, but for adversarial scenarios it can backfire.

## Real-World Prompting Examples

Beyond controlled studies, the high-stakes prompting phenomenon is visible in practitioner lore and AI guidelines. A widely-cited enterprise blog (“Prompt Engineering Is Just Good Requirements Engineering” (<sup>[29]</sup> [crashoverride.com](https://crashoverride.com/))) emphasizes giving the LLM a clear role and charge. The author Mark Curphey notes that specifying a persona (e.g. “You are a senior software engineer”) “*actually changes how the model approaches the problem*” (<sup>[30]</sup> [crashoverride.com](https://crashoverride.com/)). In effect, this is a form of stake-raising: a persona cue sets an implicit responsibility. Similarly, lists of “*best prompts*” often include admonitions like “*be concise and thorough*” or “*explain your answer*”, which are encouragement-style cues to improve output.

Anecdotally, tutors have found that telling ChatGPT “*You will be graded on this answer*” or “*Pretend you're explaining to a student*” can yield more complete answers. In educational settings, framing a question as an exam problem seems to make LLMs give more step-by-step solutions (mirroring how a student might actually write out their work). For example, one user reported that prefixing math questions with “*Imagine I will share this answer with ten professors*” elicited dramatically more careful solutions (this is consistent with Li *et al.*'s observation that question context matters (<sup>[31]</sup> [spacefrontiers.org](https://spacefrontiers.org/))).

In professional use, when teams integrate LLMs into workflows, they often find better success by emphasizing importance. A LinkedIn article on AI in hiring advised prompt templates like “*This is a high-priority recruitment problem*” to get more precise proposals. In customer support, agents using LLMs might include the customer's frustration level (“**URGENT: please handle carefully**”) to get more empathetic responses. These examples are anecdotal but align closely with the controlled findings.

## Psychological Analogies and Theories

Although LLMs lack consciousness, these effects invite comparisons to human psychology. The most relevant analogy is the **Yerkes–Dodson Law** of arousal: in humans, performance on tasks tends to improve with stress/arousal up to a point, then decline if stress is excessive. Pasichnyk (2026) recently demonstrated an LLM counterpart: in multi-agent simulations, AI “difficulty/pressure” followed an inverted-U curve for cooperation (<sup>[9]</sup> [www.researchgate.net](https://www.researchgate.net/)). At moderate pressure (resource scarcity), agents traded cooperatively at highest rates; under extreme pressure, performance collapsed. This suggests an optimal level of urgency exists even for AI. Translating to single-agent LLMs, it implies that some pressure elevates sophistication, but overdoing it could confuse or degrade output. Indeed, prompt engineers note that **negative** phrases (“I'm going to kill your model if you fail!”) can confuse the model's safe-answer heuristics, whereas moderate urgency (“be careful”) tends to help more consistently.

Another psychological concept is *self-monitoring*. Humans often perform better when they perceive accountability. In one human study, participants given a challenging task with a “countdown timer” or under observation solved it more carefully (<sup>[32]</sup> [news.ycombinator.com](https://news.ycombinator.com/)). Analogously, telling an LLM “you will get a confidence score” (as one effective EmotionPrompt did ([ai-scholar.tech](https://ai-scholar.tech/))) implies monitoring, which can lead the model to internally justify answers more rigorously. While the model doesn't *fear* failure, it does optimize token probabilities as if that stake were real, because similar patterns appeared in its training (e.g. exam answers include justification).

Importantly, psychologists caution that these analogies have limits. LLMs show systematic biases (high agreeability, sensitivity to prompt phrasing) that differ from humans (<sup>[33]</sup> [www.lesswrong.com](http://www.lesswrong.com)). Nonetheless, applying cognitive frameworks has been practically fruitful. For example, the concept of “attention” in a neural net can be loosely related to human selective attention. Motivational cues in a prompt can be thought of as “highlighting” parts of the task, aligning with how emotion modulates human focus (<sup>[18]</sup> [aobrain.com](http://aobrain.com)). In summary, human psychology offers useful metaphors, but we should remember the model is ultimately a statistical engine, not a person. The improvements come from *averaging over human-like responses* rather than genuine motivation.

## Discussion and Implications

The phenomenon of improved LLM performance under “threat” has significant implications. On one hand, it is a **useful tool**: by engineering prompts to include stakes or motivation, practitioners can extract higher-quality output without altering the model. This can be especially valuable in high-stakes applications (e.g. medical or legal domains) where rigidity and thoroughness are desired. For example, instructing a medical LLM “*This recommendation will guide patient treatment*” might reduce hallucinations, as it cues the model to rely on factual knowledge. Similarly, educators using AI tutors might stress the importance of an answer to get the model to elaborate explanations. In an era where model architectures move slowly but use-cases demand reliability, such prompt strategies are a lightweight “algorithmic intervention.”

However, several cautions arise. First, **instruction misuse**: adversarial users could exploit this to encourage models to obey malicious instructions more eagerly (as seen with disinformation) (<sup>[7]</sup> [aobrain.com](http://aobrain.com)) (<sup>[26]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)). Second, **composition conflicts**: mixing high-stakes cues with contradictory directions can confuse the model, leading to unpredictable outcomes (<sup>[23]</sup> [aobrain.com](http://aobrain.com)). Third, **stability**: models evolve, and a prompt hack that works on one version may backfire on another. Indeed, official guidance now warns that some older prompt tricks (like over-empathizing phrases) have inconsistent effects on newer LLM families (<sup>[24]</sup> [aobrain.com](http://aobrain.com)) (<sup>[34]</sup> [aobrain.com](http://aobrain.com)). Continuous benchmarking and user feedback are thus essential.

For researchers, these findings open fascinating questions. Can we formalize “stakiness” as a controllable variable? Could future LLM training incorporate scenarios with explicit rewards (akin to reinforcement learning from simulated threat contexts)? Alternatively, does this trick hint that current LLMs lack some notion of *intrinsic motivation*, and might benefit from architectures that simulate “persistent attention” or “certainty estimation”? Some proposals even consider giving LLMs an internal utility function based on prompt keywords, though that risks circular logic.

Finally, philosophically, the ease of “jacking up” an LLM’s performance with human-like motivational rhetoric underscores both the power and the emptiness of AI. The model doesn’t truly understand language in a human sense, yet it can simulate understanding so convincingly that it *feigns concern* when threatened. This anthropomorphism can be double-edged: it encourages developers to use “human tricks”, but also reminds us that stretch-and-mirror approach has limits. The LLM’s capacity to mimic seriousness is a testament to how well it has internalized human dialogue patterns (<sup>[15]</sup> [www.nature.com](http://www.nature.com)), but it also raises the question of whether we should train models to rely less on superficial cues and more on robust reasoning.

## Future Directions

Moving forward, the “stakes-priming” phenomenon suggests several avenues:

- **Automated Prompt Optimization**: Could we algorithmically tune the amount of motivational phrasing based on feedback? For instance, a system might try appending different high-stakes clauses and use a validation set to choose the best one. This meta-strategy could adapt to tasks without manual guesswork.

- **Curriculum Learning for LLMs:** Inspired by the AI Yerkes-Dodson results (<sup>[9]</sup> [www.researchgate.net](http://www.researchgate.net)), one might train LLMs progressively under increasing “emotional pressure” contexts. This could involve fine-tuning or RLHF where prompts explicitly simulate higher stakes as training proceeds, potentially yielding a model that internalizes the value of thorough reasoning.
- **Safety Mechanisms:** Since motivational cues can amplify disallowed outputs (<sup>[7]</sup> [aobrain.com](http://aobrain.com)), LLM developers might incorporate negativity checks wherein a high-stakes prompt forces an automatic verification step for sensitive topics. For example, the model might be trained to resist disinformation requests even if coaxed by friendly or urgent language.
- **Explainable Prompting:** We can study which specific words or phrases are most effective, building a lexicon of motivational triggers. This would add to prompt design guidelines (much as keyword analysis guides SEO, we could have “stakes-THESAURUS” for LLMs). Academic-speak aside, it could be packaged as a best-practices library for domain experts.
- **Cross-modal Stakes:** With multimodal models (vision+text), could “raising the stakes” in a visual prompt (e.g. showing a worried face) affect answers? Early evidence shows even implicit cues (like punctuation or emojis) can nudge LLM tone. There is room to explore beyond text.

Regardless of the path, it is clear that “emotional calibration” of LLM inputs is an emergent lever. As one expert quipped, if LLMs are like junior analysts, telling them “*Their performance will be reviewed by ten professors*” might just get them to revise their report more carefully (an anecdote echoing Nigel Hammersley’s critique (<sup>[35]</sup> [www.linkedin.com](http://www.linkedin.com))). Systematic understanding of this lever will help integrate LLMs more responsibly into workflows.

## Conclusion

In summary, LLMs **do perform better when ‘threatened’ or given high stakes**, not because they feel pressure, but because prompt wording triggers more diligent output. The effect is robust: extensive testing shows single-sentence emotional add-ons can boost accuracy and clarity by double-digit percentages (<sup>[8]</sup> [spacefrontiers.org](http://spacefrontiers.org)) (<sup>[3]</sup> [foundationinc.co](http://foundationinc.co)). The underlying cause is that these models, trained on human language, learn to associate cues of importance with a certain style of response. By priming them with seriousness or motivation, we steer them to tap into that style – yielding answers that look more careful and correct.

This finding sits at the intersection of AI engineering and cognitive science. It validates prompt-eng strategies like role-playing and chain-of-thought as practical analogs of human reasoning processes. At the same time, it raises caution: emotional priming unabashedly manipulates the model’s output distribution, which can have unintended consequences in sensitive domains. As we deploy LLMs more widely, understanding these “hack” tricks becomes part of ensuring both high performance and safety.

Future work should aim to characterize the boundaries of this effect (for different models, languages, or tasks) and to build it into standard prompt paradigms. The title of this report – “**Why do LLMs perform better when threatened or stakes are raised?**” – thus has a nuanced answer: **Because in our prompts, we effectively turn on the figurative “motivation switch” in the model’s learned representation.** In doing so, we harness the patterns of human motivation baked into the data. The result is more competent AI behavior – and a reminder that, for LLMs, good old-fashioned human psychology can still offer performance gains.

**References:** Citations in this report reference peer-reviewed findings and expert analyses. Notably, Li *et al.* (2023) provides hard benchmarks for emotional prompting (<sup>[8]</sup> [spacefrontiers.org](http://spacefrontiers.org)), Schlegel *et al.* (2025) documents LLMs’ emotional IQ (<sup>[15]</sup> [www.nature.com](http://www.nature.com)), and Pasichnyk (2026) draws parallels to Yerkes-Dodson in AI (<sup>[9]</sup> [www.researchgate.net](http://www.researchgate.net)). The discussed prompt-engineering practices are supported by AI blog analyses (<sup>[4]</sup> [aobrain.com](http://aobrain.com)) (<sup>[3]</sup> [foundationinc.co](http://foundationinc.co)) and practitioner guides (<sup>[10]</sup> [crashoverride.com](http://crashoverride.com)). Where possible, we have quoted or summarized results directly with inline references. This comprehensive survey aims to ground each claim in the existing literature on LLM behavior under high-stakes prompting.

## External Sources

- [1] <https://spacefrontiers.org/r/10.48550/arxiv.2307.11760#:~:have%...>
- [2] <https://spacefrontiers.org/r/10.48550/arxiv.2307.11760#:~:using...>
- [3] <https://foundationinc.co/lab/emotionprompts-llm#:~:Feedi...>
- [4] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:ln%20...>
- [5] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:3...>
- [6] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:match...>
- [7] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:Secon...>
- [8] <https://spacefrontiers.org/r/10.48550/arxiv.2307.11760#:~:have%...>
- [9] [https://www.researchgate.net/publication/401721271\\_The\\_Yerkes-Dodson\\_Curve\\_for\\_AI\\_Agents\\_Emergent\\_Cooperation\\_Under\\_Environmental\\_Pressure\\_in\\_Multi-Agent\\_LLM\\_Simulations#:~:systeme...](https://www.researchgate.net/publication/401721271_The_Yerkes-Dodson_Curve_for_AI_Agents_Emergent_Cooperation_Under_Environmental_Pressure_in_Multi-Agent_LLM_Simulations#:~:systeme...)
- [10] <https://crashoverride.com/blog/prompting-llm-security-reviews#:~:When%...>
- [11] <https://huggingface.co/papers/2307.11760#:~:Emoti...>
- [12] <https://yurigushiken.github.io/education/cognition/ai/learning-projects/2023/12/20/%28essay%29-On-Anthropomorphism.html#:~:to%20...>
- [13] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:That%...>
- [14] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:The%2...>
- [15] <https://www.nature.com/articles/s44271-025-00258-x#:~:Large...>
- [16] <https://www.nature.com/articles/s44271-025-00258-x#:~:their...>
- [17] <https://yurigushiken.github.io/education/cognition/ai/learning-projects/2023/12/20/%28essay%29-On-Anthropomorphism.html#:~:Image...>
- [18] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:3.4%2...>
- [19] <https://yurigushiken.github.io/education/cognition/ai/learning-projects/2023/12/20/%28essay%29-On-Anthropomorphism.html#:~:As%20...>
- [20] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:2,can...>
- [21] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:As%20s...>
- [22] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:Third...>
- [23] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:The%2...>
- [24] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:al...>
- [25] <https://spacefrontiers.org/r/10.48550/arxiv.2307.11760#:~:have%...>
- [26] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12009909#:~:disin...>
- [27] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12009909#:~:as%20c...>
- [28] <https://pmc.ncbi.nlm.nih.gov/articles/PMC12009909#:~:expec...>
- [29] <https://crashoverride.com/blog/prompting-llm-security-reviews#:~:LLMs%...>

[30] <https://crashoverride.com/blog/prompting-llm-security-reviews#:~:Perso...>

[31] <https://spacefrontiers.org/r/10.48550/arxiv.2307.11760#:~:have%...>

[32] <https://news.ycombinator.com/item?id=36230750#:~:!%27m...>

[33] <https://www.lesswrong.com/posts/jhDCRe7fnsvubknBp/llm-psychometrics-and-prompt-induced-psychopathy#:~:;conc...>

[34] <https://aobrain.com/en/blog/priming-llms-motivational-directive-phrases/#:~:Four...>

[35] [https://www.linkedin.com/posts/nhammersley\\_if-chatgpt-were-an-employee-it-would-be-activity-7311308555607822336-wqX9#:~:ChatG...](https://www.linkedin.com/posts/nhammersley_if-chatgpt-were-an-employee-it-would-be-activity-7311308555607822336-wqX9#:~:ChatG...)

---

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at <https://intuitionlabs.ai/contact> for a consultation.

---

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by [Adrien Laurent](#), a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.