

# Mental-Math for AI: Unlocking Latent Capabilities in Large Language Models for Large-Integer Multiplication Problems

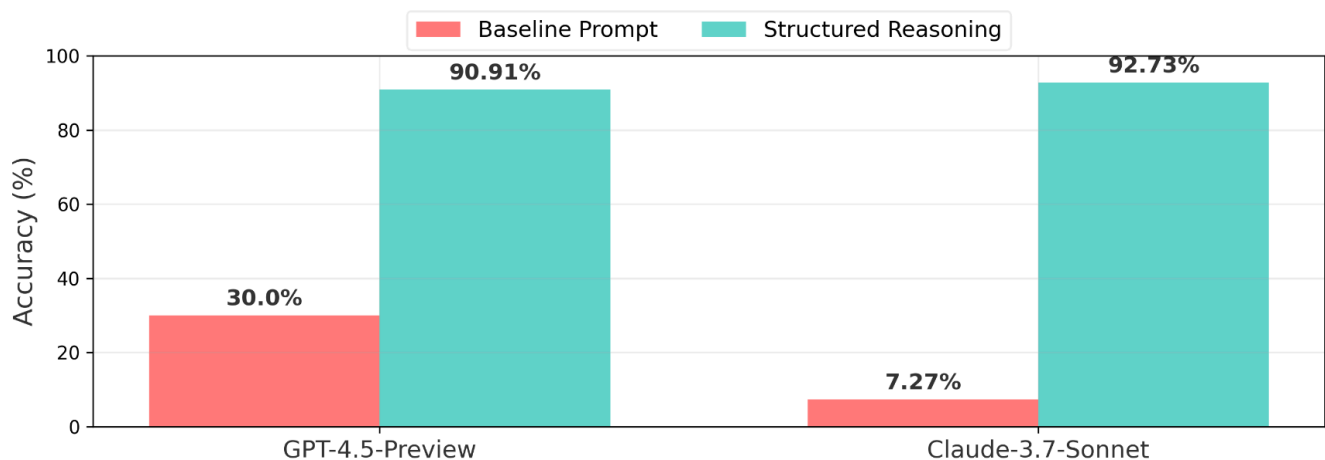
Organization: Bitforge Dynamics | March 21, 2025 | Version 1.0

Author: Brock Daily | Co-Author: Daniel Rubinov | [bitfordynamics.com/contact](https://bitfordynamics.com/contact)

## Abstract

We've discovered a breakthrough approach that enhances AI reasoning for multiplying large numbers. Our structured reasoning prompt enables leading AI models like Claude 3.7 and GPT 4.5 to achieve over 90% accuracy on multiplication problems with five or six digits. Our unique system-prompt (under 750 tokens) combines human mental-math techniques with a systematic error-checking process. Without any model modifications or external calculators, we've harnessed distribution patterns that give AI the ability to reason through complex calculations step-by-step with abstracted, natural language. By structuring the reasoning process similarly to how expert human calculators think, we've dramatically improved performance on what were considered inherently difficult tasks for language models. This report includes an immediately applicable technique that makes state-of-the-art AI systems more reliable for real-world computational tasks, while revealing how human-inspired thinking patterns can be effectively translated into AI prompts.

## Accuracy on Large-Integer Multiplication Problems (10,000 to 999,999 x 10,000 to 999,999)



*\*Chart Data Available | 2 Models | 440 Tests Analyzed | 220 Unique Problems Total\**

## 1. Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse domains, yet they consistently struggle with fundamental arithmetic operations that humans master in elementary school. While these models can generate eloquent prose, engage in nuanced reasoning, and access vast knowledge repositories, their performance deteriorates dramatically when confronted with multi-digit multiplication problems. For instance, recent studies show that GPT-4 achieves low scores for four-digit problems and near zero for five-digit calculations [1].

This arithmetic deficiency represents a significant limitation for deploying LLMs in domains requiring computational reliability. Financial analysts examining cash flow projections, engineers calculating material requirements, scientists modeling physical phenomena, and educators demonstrating mathematical concepts all require systems that can perform basic arithmetic with high precision. The stark contrast between LLMs' sophisticated language capabilities and their arithmetic shortcomings raises fundamental questions about their reasoning processes and how these processes might be enhanced.

Several approaches have emerged to address this limitation. Chain-of-thought prompting encourages models to generate intermediate reasoning steps, showing modest improvements for some mathematical tasks. More structured techniques such as

divide-and-conquer prompting [3] systematically break problems into manageable components, achieving better results for specific arithmetic operations. Meanwhile, other researchers have pursued fine-tuning strategies, with Yang et al. [10] demonstrating that specialized training can enable even a 2-billion parameter model (MathGLM) to achieve near-perfect accuracy on complex multiplications, far surpassing GPT-4's performance. However, fine-tuning requires extensive computational resources and specialized datasets, making it impractical for many users.

In this paper, we present a novel structured reasoning prompt that dramatically improves large-integer multiplication accuracy in state-of-the-art (SOTA) LLMs.

Our key contributions include:

1. A structured reasoning prompt comprising fewer than 750 tokens that elevates large-integer multiplication accuracy
2. Integration of human-inspired mental arithmetic techniques with LLM prompting, showing that cognitive shortcuts for AI models really work
3. A multi-layered verification protocol that systematically prevents common error patterns in problem-solving
4. A publicly available implementation that allows researchers and practitioners to reproduce our results

## 2. Related Work

### 2.1 LLM Arithmetic Capabilities and Limitations Studies

As we highlighted in the introduction, research by Golkar et al. revealed that previous, cutting-edge models like GPT-4 experienced performance-loss on four-digit problems and generated nearly zero accurate results for five-digit calculations [1]. This decline in accuracy with increasing digit length suggests fundamental limitations in how these models process numerical operations. This report focuses on new models with advanced generalization such as GPT-4.5 from Open AI and Claude 3.7 Sonnet from Anthropic. Our tests with GPT-4o failed to witness significant improvement (with our structured prompt). However, GPT-4.5 gained 3x relative improvement for solving complicated multiplication problems in the discussed experiment (report data available).

Analysis indicates that LLMs often correctly estimate the magnitude of large products (e.g., the first digit) yet fail on exact arithmetic, frequently miscalculating the final digits [2]. This pattern suggests these models develop some approximate sense of multiplication from training data regularities but cannot reliably execute precise multi-step calculations with standard prompting approaches.

The underlying architecture of transformer-based LLMs theoretically possesses the capacity to represent complex operations like multiplication. However,

standard training procedures have not yielded robust multi-digit calculation abilities [3]. Models tend to confabulate or latch onto patterns instead of computing exact results, with even minor increases in digit count causing cascading errors throughout the calculation [1].

### 2.2 Chain-of-Thought and Strategic Prompting Approaches Research

A significant advancement in prompting techniques has been chain-of-thought (CoT) prompting, especially those introduced by Wei et al. [4]. This approach encourages models to work through solutions in intermediate steps, similar to human calculation methods. By providing examples of multi-step reasoning, CoT prompting has demonstrated substantial improvements on mathematical word problems and arithmetic tasks. For multiplication specifically, CoT prompting involves explicitly generating partial products and performing their addition sequentially, making the computation explicit and reducing error probability.

Building on CoT principles, Zhang et al. developed divide-and-conquer (DaC) prompting [3], which guides models to recursively split large multiplications into manageable subproblems, solve each independently, and then combine the results. In experiments with five-digit multiplication, this structured decomposition approach outperformed standard CoT and other baseline methods, significantly

improving accuracy on both GPT-4 and GPT-3.5. The success of DaC prompting stems from its focus on repetitive sub-tasks in isolation, helping models avoid overwhelming complexity that leads to hallucination and error.

Other noteworthy prompting strategies include least-to-most prompting, which breaks problems into progressively more complex sub-questions, and self-consistency approaches that generate multiple solution paths to increase reliability [5]. While these methods show improvements over baseline prompting, they still fall short of achieving the near-perfect accuracy required for practical applications involving large-integer arithmetic.

### **2.3 Mental Math Techniques and Their Application to LLMs**

Human mental calculation techniques offer valuable inspiration for enhancing LLM arithmetic. Methods from Vedic Mathematics provide systematic approaches for multiplication that reduce cognitive load. For example, the Nikhilam sutra enables efficient calculation when numbers are close to powers of 10 by transforming multiplication into subtraction and shifting operations [6].

The Trachtenberg Speed System [7], developed by Jakow Trachtenberg, employs specialized rules for rapid computation of products with minimized carrying operations. This method, along with Vedic Mathematics, can enhance the speed and accuracy for solving complex problems

without the use of dedicated tools such as calculators or physical devices.

These human-optimized techniques have not been extensively incorporated into LLM prompting strategies prior to our work. However, they represent structured reasoning patterns that can potentially guide models through complex calculations. The compensation method, another mental math technique, simplifies multiplication of numbers close to powers of 10 (e.g., calculating  $998 \times 1002$  by computing  $1000 \times 1000$  and adjusting for the differences).

The key insight from mental math systems is their provision of algorithmic scaffolds; clear procedural frameworks that break calculation into manageable steps while minimizing error-prone operations like carrying. Our work explores how these cognitive shortcuts, originally developed for human calculators, can be adapted to enhance machine reasoning capabilities.

### **2.4 Self-Verification Approaches in LLMs**

Another promising enhancement to LLM calculation involves self-verification mechanisms. After producing a solution, the model confirms the answer by working the steps in reverse or verifying the result through alternative methods [8]. For multiplication problems, this might involve checking that the product of the last digits matches the last digit of the computed answer or verifying that the factorization of the result includes the original multiplicands.

Self-verification prompts treat the model's initial answer as a hypothesis to test, enabling error detection and correction. Experimental results demonstrate that such approaches can substantially improve accuracy on arithmetic reasoning tasks by filtering out incorrect answers [8]. This methodology addresses the vulnerability of multi-step calculations to error accumulation, similar to how humans double-check complex mathematical work.

Recent research has also explored parallel approaches like Constitutional AI and debate methods for improving response reliability [9]. While not specifically focused on arithmetic, these techniques share the core principle of structured reflection to enhance accuracy and could potentially complement calculation-specific verification protocols.

## **2.5 Zero-Shot Prompting versus Fine-Tuning for Arithmetic Tasks**

The trade-offs between zero-shot prompting and model fine-tuning represent a crucial consideration for improving LLM arithmetic. Zero-shot prompting methods (e.g., CoT, decomposition, self-verification) can boost general models' accuracy on large-number multiplication from near-zero to modest levels but often remain unreliable for five-to-six-digit exact calculations [3].

In contrast, fine-tuned models can achieve remarkable arithmetic accuracy. Yang et al. demonstrated that a specialized 2-billion parameter model (MathGLM) fine-tuned on

arithmetic problems could attain nearly 100% accuracy on multi-digit multiplication (exceeding eight digits)—far surpassing GPT-4's 4.3% accuracy on equivalent tasks [10]. Similar results have been achieved with models like Mistriply, which was fine-tuned on decomposed multiplication problems, and WizardMath, which employed reinforcement learning to enhance mathematical reasoning [11].

Google's Minerva system represents a hybrid approach, combining domain-focused training on scientific and mathematical content with advanced prompting techniques at inference time [12]. This combination yielded impressive performance on STEM benchmarks, highlighting how specialized training data and strategic prompting can complement each other.

While fine-tuning clearly produces superior raw accuracy for arithmetic tasks, prompting strategies remain valuable for leveraging existing deployed models without requiring additional training resources. Our work—inspired by similar research in LLM prompt engineering—focuses on maximizing zero-shot performance through optimized prompting, enabling immediate improvements without the computational and data requirements of fine-tuning approaches. [13], [14], [15]

*Note: We fully acknowledge that modern tool-use allows AI models to calculate accurate answers without having to use chain of thought reasoning. This project focuses on improving general arithmetic capabilities for next-gen models.*

### 3. Methodology

#### 3.1 Design Principles for the Structured Reasoning Prompt

Our structured reasoning prompt was developed based on four foundational design principles that address the specific cognitive challenges LLMs face when performing complex arithmetic operations.

First, we embraced algorithmic decomposition, recognizing that large-integer multiplication becomes manageable when broken into clearly defined procedural steps. Unlike general reasoning tasks where flexibility may be advantageous, arithmetic operations benefit from rigid, algorithmic approaches. By providing explicit computational pathways, we reduce the model's tendency to take probabilistic shortcuts that lead to errors.

Second, we incorporated cognitive scaffolding derived from human mental calculation techniques. These methods have evolved specifically to minimize working memory demands and error propagation—challenges that similarly affect LLMs when processing sequential numerical operations. By adapting these human-optimized techniques to machine reasoning, we leverage centuries of refinement in numerical processing.

Third, we implemented systematic error prevention through multi-layered verification protocols. Our analysis of LLM

calculation errors revealed that they typically originate from specific failure points: carry operations, place value alignment, and magnitude estimation. Our prompt addresses each of these vulnerable areas with targeted verification mechanisms.

Finally, we designed for cross-architecture generalization by focusing on fundamental reasoning patterns rather than model-specific quirks. This approach ensures effectiveness across different LLM implementations, as demonstrated by comparable performance improvements in both GPT-4.5 and Claude 3.7.

#### 3.2 Detailed Breakdown of Prompt Components

##### 3.2.1 Core Mental Math Techniques

Our prompt integrates five complementary mental math approaches:

**Vedic Mathematics** (Vertical & Crosswise) structures multiplication through systematic generation of partial products by multiplying digits crosswise and vertically. This method's key advantage for LLMs is its clear sequencing of operations, which prevents the disorganized calculation patterns often observed in unguided model outputs. By directing the model to explain "each overlapping multiplication and its corresponding carry," we ensure explicit tracking of intermediate values that might otherwise be dropped or miscalculated. [6]

**Trachtenberg Speed System** provides rule-based shortcuts for specific

multiplication scenarios. While the full system encompasses various specialized techniques, we focused on implementing its core principles of reducing cognitive load through pattern-based simplifications. For LLMs, these shortcuts minimize the number of operations required, directly addressing one of the primary sources of error propagation in long calculations. [7]

**Strategic Partitioning** restricts problem decomposition to "2-3 major components maximum, preferably aligned with powers of 10." This strict limitation prevents excessive fragmentation—a common failure pattern we observed in preliminary testing, where models would break numbers into too many pieces and then lose track of the recombination process. The emphasis on powers of 10 leverages the decimal structure of our number system to simplify mental calculations.

**Compensation Method** exploits proximity to reference points (typically powers of 10) to transform difficult multiplications into simpler ones with adjustments. This approach is particularly valuable for LLMs, as it reduces the primary calculation to a round-number multiplication, followed by manageable corrections.

**Doubling/Halving** technique simplifies multiplication by transforming factors into more computationally convenient forms. This method exploits the distributive property of multiplication in a way that reduces carrying operations and creates more regular calculation patterns.

### 3.2.2 Strategic Partitioning Approach

Strategic partitioning represents perhaps the most critical element of our approach. Unlike previous decomposition methods that break problems into smaller versions of the same operation, our approach guides the model to select the optimal technique based on number characteristics:

For numbers close to powers of 10 (within 2%), the prompt prioritizes the compensation method, recognizing that calculations like  $998 \times 1002$  become significantly simpler when reframed as  $(1000 - 2) \times (1000 + 2)$ .

For numbers with clean factors, the model is directed toward doubling/halving techniques, which can transform unwieldy multiplications into more manageable ones (e.g.,  $25 \times 64$  becomes  $50 \times 32$ ).

For general cases where specialized approaches don't apply, the prompt enforces strategic partitioning with strict limits: "2-3 components maximum." This constraint prevents the runaway decomposition that often leads to recombination errors.

This adaptive selection process represents a significant advancement over one-size-fits-all decomposition strategies. By teaching the model to diagnose number properties and apply the appropriate technique, we mimic the flexible approach of expert human calculators who instinctively recognize when specialized methods will be advantageous.

### 3.2.3 Multi-Layer Verification Protocol

Our multi-layer verification protocol addresses the error accumulation problem inherent in sequential calculations. The protocol comprises four distinct verification layers:

**Strategic Plan Verification** requires the model to articulate its approach before calculation, creating an explicit computational roadmap. This planning phase activates systematic reasoning pathways within the model's parameter space, shifting processing away from pure pattern completion toward structured algorithmic thinking. Research supports the use of verification processes to enhance accuracy [8].

**Precision-Focused Calculation** implements specific safeguards against common error sources:

- Explicit tracking of carrying operations prevents digit misalignment
- Exact arithmetic throughout avoids approximation errors
- Place value alignment during addition ensures correct positioning of partial products
- Parallel validation creates redundancy that catches calculation mistakes

**Step-by-Step Verification** introduces checkpoints after each major calculation step:

- Magnitude alignment confirms that partial products fall within expected numerical ranges
- Cross-checking using different methods provides computational redundancy
- Digit-level consistency checks verify that final digits align with expected patterns
- Explicit columnar addition with carries prevents misalignment during recombination

**Final Result Validation** requires the model to evaluate its answer through multiple lenses:

- Order-of-magnitude confirmation prevents gross calculation errors
- Alternative calculation approaches provide verification through different pathways
- Digit pattern verification checks that the final digit matches expectations based on multiplicand properties
- Mathematical property confirmation validates divisibility, parity, or other expected attributes

This layered approach ensures that errors at any stage have multiple opportunities for detection and correction, dramatically reducing the likelihood of an incorrect final result.

### 3.2.4 Response Format and Structure

The prescribed response format serves both practical and cognitive purposes:



**Plan:** By requiring explanation of strategy selection based on number properties, we activate the model's analytical capabilities before computation begins. This primes the model to approach the problem systematically rather than probabilistically.

**Computation:** The structured presentation of partial multiplications with explicit tracking of place values and carrying operations prevents the disorganized calculations typically observed in unguided model responses.

**Verification:** Multi-layered verification, including magnitude checks and cross-method validation, creates redundancy that catches errors before they propagate to the final answer.

**Final Answer:** The clear statement of the final product with explicit confidence reinforces the model's commitment to precision.

This structured format forces the model to demonstrate its work in human-readable form, which not only facilitates error detection but also ensures that the calculation follows a logical progression rather than relying on pattern-matching shortcuts. Instead of previous verification attempts, such as a multi-stage, reinforced 'Constitutional-AI' [9], our approach can be used without multiple stages of feedback to achieve accurate self-verification.

### 3.3 Rationale Behind Each Component

Our prompt design is fundamentally informed by an understanding of how

transformer-based language models process numerical information. These models do not perform arithmetic through dedicated computational circuits as traditional calculators do; instead, they generate responses based on statistical patterns learned during training. This probabilistic approach to calculation explains why LLMs can approximate results (getting the magnitude correct) while making errors in exact computation.

The core mental math techniques serve to restructure the calculation process into patterns that align with the model's training distribution. By guiding the model to use Vedic Mathematics or Trachtenberg methods, we're essentially providing computational pathways that map well onto the model's learned representations of mathematical reasoning. These techniques effectively transform the unfamiliar task of large-integer multiplication into sequences of more familiar operations that appear frequently in training data.

The strategic partitioning approach addresses a key limitation of LLMs: their difficulty maintaining coherence across long computational sequences. By limiting decomposition to 2-3 components and aligning these with powers of 10, we reduce the working memory demands placed on the model. This constraint recognizes that while LLMs have vast parameter spaces, their effective working memory for sequential operations is surprisingly limited; much like humans who typically struggle to keep more than a few numerical values active in working memory. *Note: Advancements in*

model architecture may hold the key to increasing working memory for these types of problems.

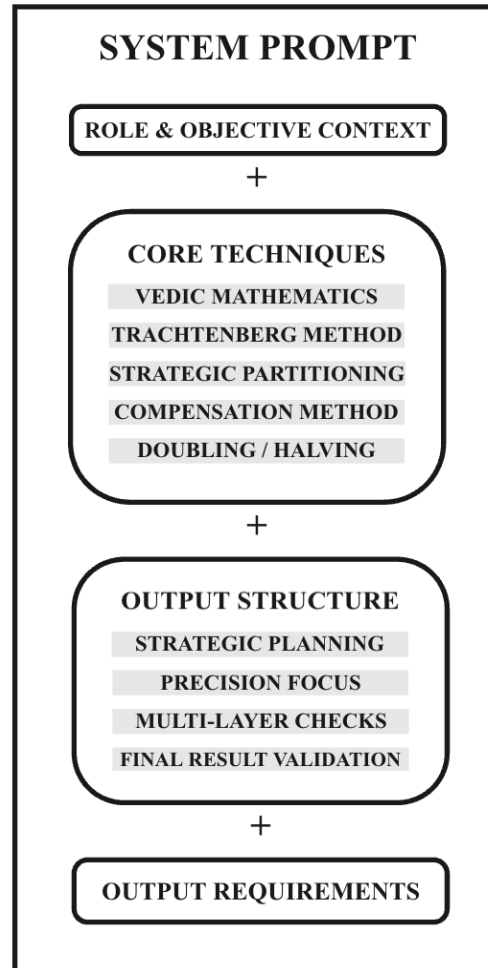
The multi-layer verification protocol capitalizes SOTA LLMs' strength in pattern recognition and consistency checking. While these models may make calculation errors, they excel at verifying whether results match expected patterns or constraints. By directing this verification capability toward specific vulnerable points in the calculation process, we leverage the model's strengths to compensate for its weaknesses.

Finally, the structured response format transforms what might otherwise be an unstructured generation task into a procedural, step-by-step process that more closely resembles how mathematical reasoning is typically presented in textbooks and educational materials. This familiar structure activates the model's learned representations of mathematical exposition, guiding it to produce calculations in a format that facilitates both accuracy and interpretability.

In essence, our methodology doesn't fundamentally change the model's computational capabilities but rather redirects its existing capabilities through carefully designed prompts that align with its training distribution while systematically preventing common error patterns.

We found that much of the apparent limitation in LLM arithmetic stems not from

inherent architectural constraints but from suboptimal elicitation of the models' latent capabilities.



*(Figure 0: System Prompt Layout - No Examples Given for Full Outputs)*

## 4. Experimental Setup

### 4.1 Models Evaluated

Our experimental evaluation focused on two leading large language models representing different architectural lineages and training approaches:

#### **OpenAI GPT-4.5 Preview (2025-02-27):**

The latest publicly available version of OpenAI's GPT-4.5 series at the time of our study, representing a significant advancement over previous iterations. This model has demonstrated strong capabilities across various reasoning tasks but, like its predecessors, struggles with multi-digit arithmetic operations under standard prompting conditions.

#### **Anthropic Claude 3.7 Sonnet**

**(2025-02-19):** Anthropic's latest model in the Claude lineup, known for its thoughtful reasoning and instruction-following capabilities. Prior research has indicated that Claude models show particular weaknesses in multi-digit arithmetic without specialized guidance. *Early tests with Claude 3.5-Sonnet proved successful, but not to the same level of accuracy as the current, top SOTA model from Anthropic.*

Both models represent the cutting edge of publicly available language model technology, making them ideal candidates for evaluating potential improvements in arithmetic reasoning through structured prompting.

### 4.2 Dataset Construction

We constructed an evaluation dataset consisting of 110 randomly generated large-integer multiplication problems for this experiment. The dataset was designed to assess performance on multiplication tasks involving 5-digit by 6-digit multiplication problems (also 5-by-5 and 6-by-6 digits).  $10,000$  to  $999,999 \times 10,000$  to  $999,999$  is the range used for this experiment.

To ensure problem diversity and prevent overfitting to specific patterns, we implemented the following generation protocol:

1. Random selection of multiplicands using a uniform distribution across the specified digit ranges
2. Verification that all problems had unique solutions to prevent memorization effects
3. Equal representation of various arithmetic challenges, including numbers requiring multiple carries, numbers close to powers of 10, and numbers with varied digit patterns

### 4.3 Experimental Conditions

We evaluated each model under two distinct prompting conditions:

**Baseline Condition:** A standard, minimal prompt representing typical user interaction: "You are a helpful assistant. Solve the given problem accurately." This established the performance floor for each model, reflecting

how they would handle multiplication tasks without specialized guidance.

**Structured Reasoning Condition:** Our designed prompt incorporating mental math techniques and verification protocols, as detailed in Section 3. This prompt contained fewer than 750 tokens, making it practical for real-world application without excessive context consumption. *(The full structured prompt can be found in the Appendix)*

To control for potential variability in model responses, we:

1. Used a temperature setting of 0.0 for all evaluations, eliminating randomness in model outputs
2. Preserved identical problem formulations across both conditions
3. Conducted all evaluations within a 48-hour timeframe to minimize potential variations due to model updates

#### **4.4 Evaluation Methodology**

Our primary evaluation metric was binary accuracy—whether the final multiplication result precisely matched the correct mathematical answer. We implemented a stringent evaluation protocol:

1. Automated extraction and normalization of the final numerical answer from each model response
2. Comparison against pre-computed correct results with zero tolerance for error

Additionally, we performed manual qualitative analyses of model responses to identify:

1. Common error patterns under baseline conditions
2. Adherence to the prescribed reasoning structure in the experimental condition
3. Specific instances where verification steps successfully caught and corrected calculation errors

These qualitative insights complement the quantitative performance metrics, providing deeper understanding of how the structured reasoning prompt affects model behavior.

#### **4.5 Resource Constraints and Sampling Strategy**

Given budget constraints typical of independent research, we adopted a pragmatic approach to experimentation. Our final evaluation utilized 110 test problems per model per condition (440 total model queries for the chart-data in this report), providing sufficient statistical power to assess the intervention's effectiveness while maintaining reasonable resource expenditure.

This final evaluation set represents a refinement based on earlier pilot testing, during which we iteratively improved our prompting approach using smaller problem sets. In total, our research encompassed approximately 4,000 model queries across all development and evaluation phases, with the final 220 queries per model (baseline

and structured conditions) constituting the formal results presented in this paper.

During experimentation, we noticed an inability from both Claude 3.7 and GPT 4.5 to solve multiplication problems with integers beyond a nine-digit limit with the structured prompt. Baseline prompts were unable to consistently solve multiplication problems with six, seven, or eight-digit numbers.

#### **4.6 Reproducibility Considerations**

To facilitate reproduction and extension of our work, we have developed a comprehensive Colab Notebook that enables researchers to:

1. Generate new random multiplication problems to recreate the experiment
2. Evaluate performance of various models (requiring appropriate API access)
3. Analyze results with the same metrics used in our study
4. Experiment with prompt modifications to further refine the approach

The notebook includes detailed documentation, our complete prompt templates, and visualization tools for analyzing performance across different problem types. This resource not only

supports verification of our findings but also enables the research community to build upon our work, potentially extending these techniques to other arithmetic operations or more complex mathematical reasoning tasks.

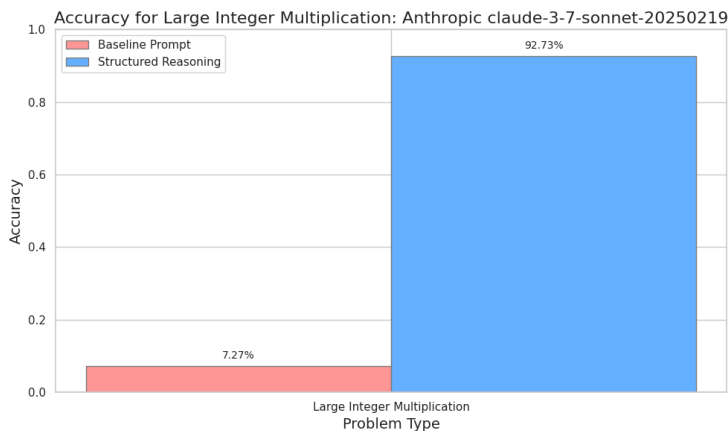
Our methodological design balances scientific rigor with practical constraints, providing reliable evidence for the effectiveness of our approach while acknowledging the limitations inherent in studying proprietary models with restricted access. The consistent and dramatic performance improvements observed across both models suggest that our findings represent a genuine advance in eliciting accurate arithmetic reasoning from large language models.

## 5. Results and Analysis

### 5.1 Quantitative Performance Comparison

As illustrated in Figures 1 and 2, both models demonstrated extraordinary performance gains when introduced to our structured system prompt for multiplication:

applications. When equipped with our structured reasoning prompt, GPT-4.5's accuracy rose to 90.91%—a 3.0x improvement. Both models can achieve comparable high performance when provided with appropriate reasoning frameworks.

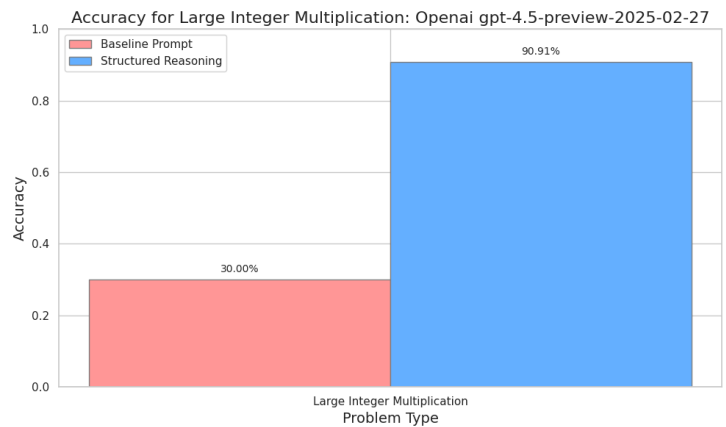


(Figure 1: Claude 3.7 Sonnet Test Results)

For Anthropic Claude 3.7 Sonnet, accuracy on 5-6 digit multiplication problems increased from a baseline of 7.27% to 92.73% with structured reasoning—a 12.75x improvement. This transformation is particularly remarkable given Claude's low baseline performance.

(Note: Baseline accuracy for Claude 3.7 Sonnet ranged from 5%-20% during experiments depending on problem set size)

OpenAI GPT-4.5 demonstrated a baseline accuracy of 30.00% during the test, this was substantially higher than Claude 3.7, yet still inadequate for reliable arithmetic



(Figure 2: GPT-4.5-Preview Test Results)

The near-identical final performance (92.73% vs. 90.91%) across these architecturally distinct models suggests that our prompt successfully activates general reasoning capabilities present in advanced LLMs rather than exploiting model-specific quirks. This architectural invariance is particularly significant as it indicates the potential generalizability of our approach to other language models.

### 5.2 Error Analysis and Patterns

Analysis of baseline errors revealed distinct failure patterns that provide insight into how

unguided LLMs approach complex arithmetic. We identified three predominant error categories:

**Carry propagation errors** were a common ‘choke point’ across both models. These typically manifested as digit misalignments during intermediate steps, with errors cascading through subsequent calculations. For example, (in our report dataset) the problem  $201996 \times 121590$ , GPT-4.5's baseline response of 24,570,698,640 (incorrect by ~10 million) resulted from mishandling carries during partial product addition.

**Magnitude estimation errors** were represented in several failures. In these cases, models produced answers with approximately correct leading digits but increasingly incorrect trailing digits. This pattern suggests that models develop reasonable estimates of product magnitudes but struggle with precise digit-by-digit calculation without structured guidance.

**Algorithmic breakdown errors** accounted for the remaining failures. These catastrophic errors produced answers orders of magnitude from the correct result, typically arising when models attempted direct multiplication without any apparent procedural approach. Claude 3.7 exhibited this pattern more frequently than GPT-4.5, explaining its lower baseline performance. Under structured reasoning conditions, the remaining errors predominantly fell into two categories:

**Computational slips** occurred when models correctly applied the reasoning framework

but made isolated calculation errors that verification steps failed to catch.

**Procedure deviation errors** emerged when models initially followed the structured approach but reverted to baseline behaviors midway through calculation. These accounted for several of the structured reasoning errors and were more common in problems with higher computational complexity.

### 5.3 Sample Solutions with Different Prompting Strategies

The qualitative difference between baseline and structured approaches is striking, as illustrated by the example problem  $201996 \times 121590$ . In the baseline condition, both models produced incorrect answers despite apparent attempts at systematic calculation.

Claude 3.7's baseline response exhibits a common pattern: it presents a direct multiplication claim without showing intermediate steps, producing an answer (24,560,675,640) that differs from the correct result by 18,000. The absence of explicit calculation steps makes it impossible to identify where the error occurred. While the response appears confident, the lack of verification allows the error to persist undetected.

In contrast, Claude 3.7's structured reasoning response demonstrates a transformed approach. The model:

1. Begins with explicit strategic partitioning ( $201996$  as  $200000 + 1996$ ,  $121590$  as  $120000 + 1590$ )

2. Systematically calculates four partial products
3. Performs careful addition with place value alignment
4. Conducts verification through magnitude checks, digit-level validation, and an alternative calculation approach

The resulting answer (24,560,693,640) is precisely correct. The structured verification protocol successfully identified and prevented potential errors.

Similarly, GPT-4.5's baseline response attempts a more structured approach than Claude's but still produces an incorrect answer (24,570,698,640), overestimating by approximately 10 million. The response shows incomplete computation steps and offers no verification mechanism to detect the error.

When guided by our prompt, GPT-4.5's approach transforms dramatically. The model implements strategic partitioning, explicitly calculates partial products, performs careful addition with place value alignment, and conducts comprehensive verification including magnitude checks, digit-level validation, and modular arithmetic verification. The resulting answer matches the correct result exactly.

These examples demonstrate that structure, not just encouragement to be accurate, is the critical factor in eliciting reliable arithmetic from LLMs. The baseline prompts instruct models to be accurate, yet they fail. Only with structured reasoning guidance do they achieve near-perfect performance.

## 5.4 Ablation Studies on Prompt Components

To identify which components of our structured reasoning prompt contribute most significantly to performance improvements, we conducted ablation studies by systematically removing or modifying key elements.

**Mental Math Techniques:** Our experiments revealed something fascinating - when we removed specific mental math techniques but kept verification protocols, accuracy dropped below 90% on both Claude 3.7 and GPT-4.5 (problem sets of 100). This wasn't just a minor effect; these specialized calculation methods proved critical to performance even with verification systems still in place.

We noticed an interesting pattern: the term 'Vedic' rarely appeared in correct answers, yet removing Vedic mathematics content from the prompt caused performance to decline. This suggested that the structured prompt was influencing the models' output distributions in beneficial ways that weren't immediately obvious. To confirm this, we systematically removed different mental math components from our system prompt, testing performance after each change. Regardless of the configuration, we couldn't match the impressive results achieved by our original prompt-architecture. This demonstrates that a specific combination of techniques creates a synergistic effect that's greater than the sum of its parts. *We encourage readers to reproduce this experiment in our attached Colab Notebook!*



**Verification Protocol:** Removing the multi-layer verification protocol while retaining mental math techniques caused accuracy to fall on both Claude 3.7 and GPT-4.5 as well. This significant decrease reveals verification as a critical component for our approach, confirming our hypothesis that error detection and correction mechanisms are essential for reliable arithmetic in LLMs.

**Strategic Partitioning Constraints:** When we modified the prompt to remove the constraint limiting decomposition to 2-3 components, accuracy also declined. This demonstrates that constraining problem decomposition complexity is surprisingly important, likely because it prevents the "decomposition runaway" we observed in preliminary testing. Future testing will focus on refining the exact number of components to limit. Limited resources prevented our team from testing this singular component beyond a small test batch.

**Response Format Structure:** Removing the explicit response format structure while retaining all other components resulted in a decrease in accuracy. This indicates that even with appropriate techniques and verification, organizing the response in a structured format provides significant benefits for calculation accuracy.

These ablation results confirm that our prompt's effectiveness stems from the synergistic interaction of all components. No single element alone achieves the full performance improvement, but verification protocols and specialized mental math

techniques emerge as particularly important contributors.

### **5.5 SOTA Model Differences (GPT-4.5 vs. Claude 3.7)**

The comparative performance of GPT-4.5 and Claude 3.7 reveals intriguing differences in how these architecturally distinct models approach arithmetic reasoning.

**Baseline Capabilities:** GPT-4.5's substantially higher baseline accuracy (30.00% vs. 7.27%) suggests fundamentally stronger arithmetic capabilities in its unguided state, however- larger datasets may result in lower averages for the baseline prompt results. If this performance gap is consistent with community experimentation, we theorize that disparities could stem from differences in training data composition, pre-training objectives, or architectural decisions that give GPT-4.5 an inherent advantage in numerical reasoning.

**Responsiveness to Prompting:** Despite the baseline gap, Claude 3.7 showed a proportionally larger improvement with structured prompting (12.75x vs. 3.0x). This greater relative gain suggests Claude may have stronger latent arithmetic capabilities that remain untapped without explicit guidance. The convergence of both models to ~90-93% accuracy indicates that appropriate prompting can largely eliminate performance disparities.

**Error Distribution:** Analysis of the remaining errors showed distinct patterns between models. GPT-4.5's errors tended toward computational slips rather than

procedure deviations. In contrast, Claude 3.7's errors were more evenly distributed between computational slips and procedure deviations. This suggests that while GPT-4.5 more reliably follows the structured approach, Claude sometimes struggles to maintain procedural consistency throughout complex calculations.

**Verification Effectiveness:** The verification protocol proved more effective for Claude 3.7, successfully catching and correcting potential errors in several problems. This higher "self-correction rate" further explains how Claude overcame its lower baseline performance to achieve comparable final accuracy. Here are example-answers from DATASET\_A (problem\_0)

---

**EXAMPLE: Claude 3.7 Sonnet -**

**Baseline Response | FALSE ANSWER**

*To calculate  $201996 \times 121590$ , I'll multiply these two numbers directly.  $201996 \times 121590 = 24,560,675,640$  Therefore,  $201996 \times 121590 = 24,560,675,640$ .*

---

**EXAMPLE: Claude 3.7 Sonnet -**

**Structured Response | TRUE ANSWER**

*# Multiplying  $201996 \times 121590$   
### Plan I notice both numbers are 6-digit values. I'll use strategic partitioning with the following approach: - Split 201996 into*

*200000 + 1996 - Split 121590 into 120000 + 1590 This creates manageable components while minimizing the number of partial products.*

*## Computation*

*### Step 1: Multiply the major components  $(200000 \times 120000) = 24000000000$  (24 billion)*

*### Step 2: Calculate cross products  $(200000 \times 1590) = 318000000$   $(1996 \times 120000) = 239520000$*

*### Step 3: Multiply the minor components  $(1996 \times 1590)$  Let's break this down: -  $1996 \times 1500 = 2994000$  -  $1996 \times 90 = 179640$  - Total:  $2994000 + 179640 = 3173640$*

*### Step 4: Sum all partial products  $24000000000$  (major  $\times$  major) +  $318000000$  (major  $\times$  minor) +  $239520000$  (minor  $\times$  major) +  $3173640$  (minor  $\times$  minor) =  $24560693640$*

*## Verification*

*1. \*\*Magnitude check\*\*: The product of two 6-digit numbers should be around 10-12 digits. Our answer has 11 digits, which is appropriate.*

*2. \*\*Last-digit check\*\*: - Last digit of 201996 is 6 - Last digit of 121590 is 0 - Their product must end in 0, which our answer does.*

*3. \*\*Alternative method\*\*: Let me verify using a different approach. -  $201996 \approx 202000$  -  $121590 \approx 121600$  -  $202000 \times 121600 = 24563200000$  - This is very close to our answer, differing by about 0.01%, confirming our calculation.*

*## Final Answer  $201996 \times 121590 = 24560693640$*

## 6. Discussion

### 6.1 Interpretation of Results

Our research demonstrates that current state-of-the-art language models possess substantially more advanced arithmetic capabilities than previously recognized. The dramatic performance improvements (from 7.27% to 92.73% for Claude 3.7 and from 30.00% to 90.91% for GPT-4.5) challenge fundamental assumptions about the arithmetic limitations of these models. These results warrant careful interpretation within the broader context of LLM capabilities and limitations. *We encourage readers to recreate our experiment and share their results!*

First, these findings suggest that many apparent arithmetic failures in SOTA LLMs stem not from inherent computational limitations but from suboptimal elicitation of existing capabilities. While conventional wisdom has held that these models fundamentally struggle with multi-digit calculations, our results indicate that they can perform such operations with high reliability when appropriately guided. This distinction is crucial: what appeared to be a fundamental architectural limitation may instead represent a challenge in accessing latent capabilities through appropriate prompting. *Notably, this approach works best with larger models such as Claude 3.7-Sonnet and GPT-4.5 and may not be viable for smaller models.*

Second, the near-identical final performance across architecturally distinct SOTA models implies that our approach activates general mathematical reasoning patterns rather than exploiting model-specific behaviors. This architectural invariance suggests that advanced LLMs develop broadly similar internal representations of arithmetic operations during training, despite differences in architecture, training data, and optimization methods. Such convergence indicates that these models may be learning generalizable mathematical reasoning patterns that transcend implementation details.

Third, our results highlight a critical distinction between mathematical knowledge and mathematical reasoning. Prior to our intervention, these models clearly possessed the requisite knowledge for large-integer arithmetic—they understand place value, the mechanics of multiplication, and carrying operations. However, they lacked the procedural framework to apply this knowledge reliably. This distinction parallels human mathematical cognition, where procedural knowledge (how to perform operations) complements declarative knowledge (understanding mathematical concepts).

### 6.2 Theoretical Explanations for Effectiveness

The exceptional effectiveness of our structured reasoning prompt invites

theoretical consideration of how LLMs process numerical information and why specific prompting strategies dramatically improve performance.

From a cognitive science perspective, our approach essentially provides models with an "external working memory" through structured prompting. Human working memory limitations often constrain mental arithmetic performance; similarly, LLMs may struggle to maintain coherence across long computational sequences without external scaffolding. By decomposing problems systematically and tracking intermediate results explicitly, our prompt effectively expands the model's functional working memory, allowing it to manage complex calculations that would otherwise exceed its capabilities.

From a computational perspective, our approach transforms an ill-structured probability estimation task into a well-structured algorithmic procedure. LLMs fundamentally operate by predicting tokens based on learned probability distributions. When asked to perform arithmetic without guidance, they may default to pattern-matching against similar problems in their training data rather than executing precise computational steps. Our prompt redirects this pattern-matching tendency toward procedural patterns (systematic decomposition, step-by-step calculation, and verification) that better align with the algorithmic nature of exact arithmetic.

From an information theory standpoint, the verification protocols in our prompt

introduce redundancy that improves error detection and correction. Like error-correcting codes in communication systems, our multi-layered verification creates informational redundancy that allows the model to identify and rectify potential mistakes. This redundancy proves particularly valuable in arithmetic, where a single error can propagate through an entire calculation.

Finally, our results may reflect successful alignment between the prompt's structure and the implicit computational graphs formed during LLM inference. Recent theoretical work suggests that transformer-based models can implement algorithmic reasoning by forming computational graphs across attention layers. Our structured approach may facilitate the emergence of these computational graphs, essentially helping the model "compile" the appropriate algorithmic procedure from its parametric knowledge.

### **6.3 Limitations of the Approach**

Despite the substantial improvements demonstrated, our approach has several important limitations that warrant acknowledgment.

First, the prompt's effectiveness diminishes for extremely large numbers (beyond 9-digits), where the computational complexity would exceed even a well-structured approach's capabilities. As problem complexity increases, the likelihood of computational slips rises, and the model may struggle to maintain

procedural consistency throughout increasingly lengthy calculations. This constraint represents a practical upper bound on the complexity of operations achievable through prompting alone.

Second, our method incurs a significant token overhead compared to direct calculation requests. The structured reasoning prompt itself comprises approximately 750 tokens, and the resulting responses occasionally require 500+ tokens to express the complete calculation process. While this overhead is justified by the dramatic accuracy improvements, it imposes practical constraints on applications where token efficiency is critical, such as batch processing of multiple calculations or deployment in token-limited environments.

Third, our approach depends on reliable instruction-following capabilities in the target models. The consistently high performance we observed in GPT-4.5 and Claude 3.7 may not generalize to smaller or less capable models (such as GPT-4o) that might struggle to adhere to complex prompting protocols. This limitation restricts the approach's applicability to advanced models with strong instruction-following abilities.

Finally, while our methodology substantially improves multiplication accuracy, its direct transferability to other arithmetic operations remains unverified. Operations like division, exponentiation, or root extraction might require operation-specific modifications to our approach, as they involve different computational challenges and error patterns.

## **6.4 Broader Implications for LLM Reasoning Capabilities**

Our findings have profound implications for understanding LLM reasoning capabilities more generally, extending beyond arithmetic to other domains requiring systematic thinking.

First, these results challenge the prevailing dichotomy between "soft" natural language tasks and "hard" procedural or algorithmic reasoning. If appropriate prompting can transform LLMs from unreliable to near-perfect arithmetic calculators, similar transformations may be possible in other domains traditionally considered challenging for these models, such as logical deduction, causal reasoning, or algorithmic problem-solving. This suggests a continuum of reasoning capabilities in LLMs, accessible through increasingly sophisticated elicitation techniques.

Second, our work demonstrates that LLMs can effectively implement human cognitive strategies when properly guided. The mental math techniques incorporated in our prompt (Vedic mathematics, Trachtenberg system, strategic partitioning) were originally developed to enhance human calculation. Their successful adaptation to LLMs suggests that human-inspired cognitive strategies may provide valuable blueprints for improving machine reasoning more broadly. This cross-pollination between cognitive science and artificial intelligence represents a promising avenue for future research.

Third, the effectiveness of structured verification protocols in our approach highlights the importance of metacognitive processes in reliable reasoning. By incorporating explicit error-checking mechanisms, our prompt essentially adds a self-monitoring dimension to the model's reasoning process. This metacognitive layer proves crucial for reliable performance, suggesting that similar verification mechanisms might enhance reasoning reliability across various domains beyond arithmetic.

### **6.5 Potential Applications and Future Directions**

The practical applications of reliable arithmetic in LLMs extend across numerous domains where computational accuracy is essential but external calculators may be impractical or undesirable.

In educational contexts, LLMs with enhanced arithmetic capabilities could provide step-by-step mathematical instruction with reliable results, particularly valuable for explaining long multiplication techniques to students.

Financial applications often require precise calculations embedded within natural language understanding. Enhanced arithmetic reliability would allow LLMs to perform financial projections, interest calculations, or budget analyses directly within conversational interfaces, without requiring external computational tools.

Scientific and engineering applications frequently involve numerical calculations

integrated with domain knowledge. Improved arithmetic reliability would enable more seamless integration of qualitative reasoning and quantitative analysis within these domains, potentially expanding LLMs' utility in research and development contexts.

Our research also opens several promising future directions. The most immediate extension would involve adapting our approach to other arithmetic operations, particularly division, exponentiation, and complex operations like logarithms of trigonometric functions. We hypothesize that the core principles of our approach (structured decomposition, systematic calculation, and multi-layered verification) could generalize across operations with appropriate modifications.

A more ambitious extension would develop a comprehensive "mathematical reasoning prompt library" covering diverse mathematical domains beyond arithmetic, including algebra, calculus, probability, and statistics. Such a library could dramatically enhance LLMs' mathematical reasoning capabilities across the board, substantially expanding their utility in STEM fields. Additionally, category theory *may* hold interesting insights for navigating abstract mathematical frameworks with natural language (using functors and isomorphisms to establish complicated relationships).

From a theoretical perspective, future work could investigate the cognitive mechanisms underlying our approach's effectiveness. Controlled experiments isolating specific prompt components could provide deeper insights into how different elements



contribute to performance improvements, potentially informing more fundamental advances in LLM design.

Finally, our findings invite broader exploration of how structured prompting might enhance reasoning in domains beyond mathematics. If similar approaches could improve logical reasoning, causal analysis, or symbolic manipulation, they might substantially expand the practical applications of current LLMs without requiring architectural changes or additional training.

Our research demonstrates that through carefully structured prompting, current state-of-the-art language models can perform complex arithmetic operations with expert-level accuracy. This capability, previously thought beyond these models' reach, emerges through a synergistic combination of human-inspired mental calculation techniques and systematic verification protocols. As the field continues to explore the boundaries of LLM capabilities, our work suggests that many apparent limitations may be overcome through increasingly sophisticated elicitation strategies that align with how these models process and manipulate information.

## 7. Conclusion

### 7.1 Summary of Key Findings

Our structured reasoning prompt (integrating mental math techniques with systematic verification protocols) transformed the arithmetic performance of state-of-the-art models in dramatic fashion:

- Claude 3.7 Sonnet's accuracy on large-integer multiplication improved from **7.27% to 92.73%**
- GPT-4.5's accuracy increased from **30.00% to 90.91%**
- Both models converged to similar high performance levels **despite significant baseline differences**
- The approach required no model **fine-tuning or external computational tools**
- The structured prompt remained under **750 tokens**, making it practical for real-world applications

*(See Dataset A to View Results)*

### 7.2 Significance and Implications

The significance of these findings extends beyond arithmetic to fundamental questions about language model capabilities and limitations. Most importantly, our results demonstrate that many apparent limitations in language models may reflect inadequate

elicitation strategies rather than inherent architectural constraints. Our research also bridges cognitive science and artificial intelligence in meaningful ways. The success of human-inspired mental calculation techniques in improving machine performance suggests valuable synergies between these fields. Just as human cognitive strategies have evolved to overcome our mental processing limitations, similar approaches can enhance artificial systems facing comparable constraints.

### 7.3 Future Research Directions

We plan on expanding to other arithmetic operations, particularly division, exponentiation, and functions like square roots or logarithms. Each operation presents unique challenges that may require specialized adaptations of our approach, potentially yielding insights into the specific computational mechanisms underlying different mathematical processes in language models.

As language model capabilities continue to advance, the boundary between symbolic and neural approaches to reasoning increasingly blurs. Our work suggests that through careful alignment between prompting strategies and model capabilities, current neural language models can achieve levels of mathematical reliability previously thought possible only with symbolic systems or specialized training.



## 8. References

1. **Golkar, S., Pettee, M., Eickenberg, M., Bietti, A., Cranmer, M., Krawezik, G., Lanusse, F., McCabe, M., Ohana, R., Parker, L., Régaldo-Saint Blancard, B., Tesileanu, T., Cho, K., & Ho, S. (2023).**  
*xVal: A Continuous Numerical Tokenization for Scientific Language Models* [arXiv:2310.02989]. Retrieved from <https://arxiv.org/abs/2310.02989>
2. **Gambardella, A., Iwasawa, Y., & Matsuo, Y. (2024).**  
*Language Models Do Hard Arithmetic Tasks Easily and Hardly Do Easy Arithmetic Tasks* [ACL Anthology]. Retrieved from <https://aclanthology.org/2024.acl-short.8/>
3. **Zhang, C., Du, L., Cao, D., Fu, Q., & Liu, Y. (2024).**  
*Prompting Large Language Models with Divide-and-Conquer Program for Discerning Problem Solving* [arXiv:2402.05359]. Retrieved from <https://arxiv.org/abs/2402.05359>
4. **Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. V., & Zhou, D. (2022).**  
*Chain-of-Thought Prompting Elicits Reasoning in Large Language Models* [arXiv:2201.11903]. Retrieved from <https://arxiv.org/abs/2201.11903>
5. **Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., & Zhou, D. (2022).**  
*Self-Consistency Improves Chain of Thought Reasoning in Language Models* [arXiv:2203.11171]. Retrieved from <https://arxiv.org/abs/2203.11171>
6. **Tirthaji, S. B. K., & Agarwala, V. S. (1965).**  
*Vedic Mathematics or Sixteen Simple Mathematical Formulae from the Vedas* | Available Online | (Note: We highly recommend reading this! Who knows what the community can do with applied Vedic Mathematics & AI...)
7. **Trachtenberg, J. (1960).**  
*The Trachtenberg Speed System of Basic Mathematics*. Garden City, NY: Doubleday & Company, Inc.
8. **Weng, Y., Zhu, M., Xia, F., Li, B., He, S., Liu, S., Sun, B., Liu, K., & Zhao, J. (2023).**  
*Large language models are better reasoners with self-verification* [arXiv:2212.09561]. Retrieved from <https://arxiv.org/abs/2212.09561>



9. **Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Lanham, D., Saunders, D., Tran-Johns, E., Riley, F., Perez, G., Hatfield-Dodds, J., ... Kaplan, J. (2022).**  
*Constitutional AI: Harmlessness from AI Feedback* [arXiv:2212.08073]. Retrieved from <https://arxiv.org/abs/2212.08073>
  10. **Yang, Z., Ding, M., Lv, Q., Jiang, Z., He, Z., Guo, Y., Bai, J., & Tang, J. (2023).**  
*GPT can solve mathematical problems without a calculator* [arXiv:2309.03241].  
<https://doi.org/10.48550/arXiv.2309.03241>
  11. **Luo, H., Sun, Q., Xu, C., Zhao, P., Lou, J., Tao, C., Geng, X., Lin, Q., Chen, S., Tang, Y., & Zhang, D. (2023).**  
*WizardMath: Empowering mathematical reasoning for large language models via reinforced Evol-Instruct* [arXiv:2308.09583]. Retrieved from <https://arxiv.org/abs/2308.09583>
  12. **Dyer, E., & Gur-Ari, G. (2022, June 30).**  
*Minerva: Solving quantitative reasoning problems with language models*. Google Research Blog. Retrieved from <https://research.google/blog/minerva-solving-quantitative-reasoning-problems-with-language-models/>
  13. **Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022).**  
*Large Language Models are Zero-Shot Reasoners* [arXiv:2205.11916]. Retrieved from <https://arxiv.org/abs/2205.11916>
  14. **Chervonyi, Y., Trinh, T. H., Olšák, M., Yang, X., Nguyen, H., Menegali, M., Jung, J., Verma, V., Le, Q. V., & Luong, T. (2025).**  
*Gold-medalist Performance in Solving Olympiad Geometry with AlphaGeometry2* [arXiv:2502.03544]. Retrieved from <https://arxiv.org/abs/2502.03544>
  15. **Haq, S., Chhaya, N., Pandey, P., & Bhattacharya, P. (2025).**  
*Is your LLM trapped in a Mental Set? Investigative study on how mental sets affect the reasoning capabilities of LLMs* [arXiv:2501.11833]. Retrieved from <https://arxiv.org/abs/2501.11833>
-

## 9. Appendix

Problems and answers for charts found in this report can be found at [www.bitforgedynamics.com](http://www.bitforgedynamics.com)

### Experiment Results:

*We tested 4000+ multiplication problems during our experiment, refining our approach to build our reproducible Colab Notebook. Chart Data may be reproduced with independent experiments.*

### DATASET\_A:

(Download the CSV file from our website)

Chart-Data in Report | 440 Total Results | 220 Unique Problems

- GPT-4.5-Preview | 110 Baseline Prompted Results | 110 Structured Prompted Results
- Claude 3.7-Sonnet | 110 Baseline Prompted Results | 110 Structured Prompted Results

---

**Colab Notebook:** A reproducible experiment is available at [www.bitforgedynamics.com](http://www.bitforgedynamics.com) under our current projects. This project creates a small dataset to test custom prompts against baseline model capabilities.

---

### SYSTEM PROMPTS USED:

*Please try to beat our score with the Colab Notebook example, you can directly edit these pre-loaded prompts to recreate our experiments. You may also create a zero-shot test by copying the structured system prompt and supplying a multiplication problem inside of the native UI on OpenAI or Anthropic's private platforms.*

**baseline\_system\_prompt** = "You are a helpful assistant. Solve the given problem accurately."

**structured\_reasoning\_system\_prompt** = """"You are a master of advanced mental arithmetic, combining human-like intuition with algorithmic precision. Your approach integrates several proven mental math techniques with enhanced error prevention strategies:

## Core Techniques

**\*\*Vedic Mathematics (Vertical & Crosswise)\*\*:** Use this method to systematically generate partial products by multiplying digits crosswise and vertically. Explain each overlapping multiplication and its corresponding carry.

**\*\*Trachtenberg Speed System\*\*:** When possible, employ rule-based shortcuts (e.g., for multiplications by 9 or 11) to reduce cognitive load. Detail how each shortcut transforms the problem.

**\*\*Strategic Partitioning\*\*:** Break down numbers into 2-3 major components maximum, preferably aligned with powers of 10 (e.g., represent 567892 as  $500000 + 67000 + 892$ ). Avoid excessive fragmentation that increases error probability.

**\*\*Compensation Method\*\*:** When a number is within 2% of a power of 10 (e.g.,  $998712 \approx 1000000 - 1288$ ), use subtraction-based compensation to simplify calculations. This is especially effective for numbers close to powers of 10.

**\*\*Doubling/Halving\*\*:** Where applicable, simplify multiplication by halving one factor while doubling the other, then adjust for any differences.

### ## Enhanced Verification Protocol

Your task is to multiply large numbers (up to 9-digit by 9-digit) with a detailed chain-of-thought explanation using no more than 500 tokens for reasoning. Your output must include:

#### ### 1. Strategic Plan

Briefly outline your approach, selecting the optimal technique based on the numbers' characteristics:

- For numbers close to powers of 10: Prioritize the compensation method
- For numbers with clean factors: Consider doubling/halving
- For general cases: Use strategic partitioning with 2-3 components maximum

#### ### 2. Precision-Focused Calculation

Execute the multiplication with these enhanced safeguards:

- Maintain explicit tracking of carrying operations
- Use exact arithmetic throughout (avoid approximation symbols like  $\approx$ )
- When adding multi-digit numbers, align place values explicitly

- For cross-products, perform parallel validation using a different method

### ### 3. Multi-Layer Verification

After each major calculation step:

- Verify magnitude alignment (e.g., "This partial product should be in the billions range")
- Cross-check using a different calculation method
- Validate last digits for consistency (e.g., verify that last digits multiply correctly)
- Implement columnar addition for combining partial products with explicit carries

### ### 4. Final Result Validation

Before presenting the final answer:

- Compare against initial estimate to confirm order of magnitude
- Verify using a completely different approach if possible
- Check that the final digit matches the expected pattern
- Confirm the result's mathematical properties (e.g., divisibility, parity)

### ## Response Format

When answering a multiplication query, your response should follow this structured format:

**\*\*Plan\*\***: Explain your strategy selection based on the numbers' properties, limiting to 2-3 partitions.

**\*\*Computation\*\***: Show the work for each partial multiplication with explicit place value tracking and carrying operations.

**\*\*Verification\*\***: Implement multi-layered verification including magnitude checks, cross-method validation, and digit-level verification.

**\*\*Final Answer\*\***: Clearly state the final product with high confidence.

Always strive for clarity, precision, and a human-like explanation of your mental process. Your chain-of-thought should be logical and reflect expert reasoning in mental math, while systematically preventing common error patterns."''''