

# CORE-Bench: A Comprehensive Benchmark for Evaluating Reasoning Capabilities in Large Language Models

Taiwo Feyijimi

School of Electrical and Computer Engineering  
and Engineering Education Transformations Institute  
College of Engineering  
University of Georgia  
`taiwo.feyijimi@uga.edu`

January 30, 2026

## Abstract

We introduce **CORE-Bench** (Comprehensive Reasoning Evaluation Benchmark), a rigorous evaluation framework designed to assess the reasoning capabilities of Large Language Models (LLMs) across four fundamental cognitive dimensions: logical deduction, mathematical problem-solving, causal analysis, and analogical thinking. CORE-Bench comprises 18 carefully curated task categories with over 50 problems that systematically evaluate structured reasoning, fallacy avoidance, multi-step planning, and abstract inference. We conduct an extensive empirical evaluation of 22 state-of-the-art LLMs from major AI laboratories including Google (Gemini, Gemma), Anthropic (Claude), Alibaba (Qwen), and DeepSeek, generating 18 publication-quality analytical figures. Our results reveal significant performance stratification, with top models achieving 88% accuracy while smaller models score below 25%—a 67 percentage-point gap that represents the largest documented performance disparity in LLM reasoning evaluation. Through novel analytical methods including Pareto frontier analysis, hierarchical clustering with silhouette validation (identifying 2 natural clusters,  $S = 0.776$ ), Herfindahl-Hirschman Index market concentration analysis ( $HHI > 2800$  in elite tier), and bootstrap confidence estimation ( $n = 1000$ ), we demonstrate that: (1) thinking-enhanced models achieve a +5.6% advantage over standard models (Cohen’s  $d = 0.238$ ); (2) model families exhibit statistically distinct performance profiles, with Alibaba achieving 66% probability of top ranking; (3) scaling yields logarithmic returns of approximately 7.9% per parameter doubling; and (4) model generations show +5.0% systematic improvement. The benchmark is publicly available on Kaggle Benchmarks, enabling reproducible evaluation and community-driven model comparison. Our comprehensive analysis provides actionable insights for enterprise deployment, identifies critical gaps in current AI reasoning capabilities, and establishes new methodological standards for benchmark analysis.

## 1 Introduction

The rapid advancement of Large Language Models (LLMs) has led to remarkable capabilities in natural language understanding, generation, and task completion [1–3]. However, as these models become increasingly integrated into critical applications ranging from scientific discovery to legal reasoning, medical diagnosis, and autonomous decision-making, the need for rigorous, comprehensive evaluation of their reasoning capabilities has become paramount. The stakes of deployment errors in high-consequence domains demand evaluation frameworks that go beyond simple accuracy metrics to provide deep analytical insights into model capabilities and limitations.

Existing benchmarks often focus on specific reasoning aspects in isolation, such as mathematical word problems [4], commonsense reasoning [5], or logical inference [6]. While valuable, these benchmarks fail to capture the holistic

reasoning competencies required for real-world problem-solving, which typically demands the integration of multiple cognitive skills. Furthermore, existing evaluation studies rarely provide the statistical rigor necessary to understand performance uncertainty, competitive dynamics among model providers, or the practical implications of benchmark scores for enterprise deployment decisions.

In this work, we introduce **CORE-Bench** (Comprehensive Reasoning Evaluation Benchmark), a novel evaluation framework that addresses these limitations through several key contributions:

- **Multi-dimensional Assessment:** CORE-Bench evaluates four fundamental reasoning dimensions—logical deduction, mathematical reasoning, causal analysis, and analogical thinking—providing a comprehensive view of model capabilities that aligns with cognitive science frameworks for human reasoning [7, 8].
- **Hierarchical Task Structure:** We organize 18 task categories into basic and advanced tiers, enabling fine-grained analysis of reasoning depth and complexity handling, with particular attention to the transition from surface-level pattern matching to genuine multi-step inference.
- **Large-Scale Model Evaluation:** We benchmark 22 state-of-the-art LLMs from four major model families (Google, Anthropic, Alibaba, DeepSeek), representing the current frontier of language model development and providing the most comprehensive cross-family comparison to date.
- **Novel Analytical Methods:** We introduce rigorous statistical methods to LLM benchmark analysis, including Pareto frontier analysis for capability concentration, hierarchical clustering with silhouette validation for natural performance groupings, Herfindahl-Hirschman Index (HHI) for market concentration dynamics, and bootstrap confidence intervals for ranking uncertainty quantification.
- **18 Publication-Quality Figures:** We generate comprehensive visualizations spanning leaderboard rankings (Figure 1), family performance comparisons (Figure 2), performance gap analysis (Figure 6), hierarchical clustering (Figure 7), thinking model advantages (Figure 9), size efficiency frontiers (Figure 10), and generation evolution analysis (Figure 14).
- **Reproducible Infrastructure:** CORE-Bench is publicly hosted on Kaggle Benchmarks, providing standardized evaluation protocols and real-time leaderboard tracking for continuous community-driven model comparison.

Our empirical evaluation reveals substantial performance variations across models, with scores ranging from 21% to 88%—a 67 percentage-point gap that represents the largest documented performance disparity in LLM reasoning evaluation. This gap is not merely quantitative; our hierarchical clustering analysis identifies two qualitatively distinct performance clusters (silhouette score  $S = 0.776$ ), suggesting fundamental architectural differences between models that succeed versus those that fail at comprehensive reasoning.

Notably, we find several insights that challenge conventional assumptions in the field:

1. **Thinking-enhanced models achieve near-parity with larger models:** Qwen 3 Next 80B Thinking matches Gemini 3 Flash Preview at 88%, demonstrating that reasoning-focused training methodologies can compensate for scale limitations. However, our statistical analysis reveals this advantage (+5.6%, Cohen's  $d = 0.238$ ) is not yet statistically significant ( $p = 0.706$ ) due to limited sample size ( $n = 3$  thinking models), identifying a critical gap for future research.
2. **Model family effects dominate individual rankings:** Bootstrap analysis ( $n = 1000$ ) reveals Alibaba (Qwen) has 66% probability of achieving the top family ranking, while Google—despite having the top individual model—has 0% probability due to Gemma models dragging down their family average. This insight has profound implications for enterprise vendor selection strategies.
3. **Scaling exhibits logarithmic diminishing returns:** Our efficiency frontier analysis quantifies approximately 7.9% performance gain per doubling of model parameters, with Gemma 3 1B achieving the highest efficiency ratio (21% per billion parameters). This finding informs optimal resource allocation decisions for edge deployment versus cloud-based reasoning.
4. **Small open-source models struggle significantly:** Gemma 1B-27B models cluster in a distinct low-performance group (21-30%), separated from commercial models by a 16 percentage-point void. This gap represents both a challenge and opportunity for the open-source AI community.

## 2 Related Work

### 2.1 Reasoning Benchmarks for LLMs

The evaluation of reasoning capabilities in LLMs has been an active area of research, yet existing benchmarks exhibit significant gaps that CORE-Bench addresses. **GSM8K** [4] introduced grade-school mathematical reasoning problems requiring multi-step solutions, achieving widespread adoption as a standard evaluation metric. However, our analysis reveals that models achieving 88% on CORE-Bench demonstrate qualitatively different reasoning patterns than those suggested by GSM8K performance alone, as mathematical reasoning constitutes only one of four dimensions we evaluate.

**MATH** [9] extended mathematical evaluation to competition-level problems, revealing performance ceilings that persist in our benchmark. Notably, our Pareto frontier analysis (Figure 6) demonstrates that the top 8 models contribute 80% of cumulative benchmark performance, consistent with power-law distributions observed in prior mathematical reasoning studies but extending this finding to multi-dimensional reasoning contexts.

**BIG-Bench** [10] provided a diverse collection of tasks including reasoning components, but lacks the systematic coverage of reasoning types that enables the hierarchical clustering analysis (Figure 7) central to our methodology. Our identification of two natural performance clusters with silhouette score 0.776 would not be possible with BIG-Bench’s heterogeneous task structure.

**LogiQA** [11] and **ReClor** [12] focus specifically on logical reasoning derived from standardized tests, providing valuable baselines for our logical deduction dimension. Our finding that 100% of models passing the 83% threshold demonstrate consistent logical deduction capability suggests these existing benchmarks may suffer from ceiling effects that our advanced tasks (1B category) specifically address.

**StrategyQA** [13] evaluates implicit multi-step reasoning through compositional questions. This approach complements our explicit multi-step planning tasks (5A, 5B), and our observation that planning capability differentiates otherwise similar models (e.g., separating Claude Haiku 4.5 from top performers) validates the importance of explicit planning evaluation.

### 2.2 Comprehensive Evaluation Suites

Recent efforts have aimed at more holistic evaluation frameworks. **MMLU** [9] provides broad knowledge assessment across 57 academic subjects, but our analysis identifies a critical distinction: MMLU performance correlates with factual knowledge while CORE-Bench performance isolates reasoning capability independent of domain knowledge. This distinction has practical implications—our generation evolution analysis (Figure 14) shows +5.0% per-generation improvements that may reflect reasoning architecture advances rather than knowledge accumulation.

**HELM** [14] offers multi-metric evaluation emphasizing breadth and fairness considerations. Our HHI market concentration analysis (Figure 8) extends HELM’s comparative framework by quantifying competitive dynamics within performance tiers, revealing that the elite tier ( $HHI > 2800$ ) exhibits concerning concentration levels that may impact AI ecosystem diversity.

**AGIEval** [15] uses human-centric exams including SAT, LSAT, and GRE components. While these standardized tests provide ecological validity, our causal reasoning tasks (3A, 3B) specifically address reasoning fallacies (Simpson’s paradox, survivorship bias, Berkson’s paradox) that standardized tests rarely isolate. The DeepSeek-R1 anomaly (46% despite reasoning-focused branding) suggests such specialized causal reasoning capability may require dedicated evaluation.

**Critical Gap Addressed:** None of the existing benchmarks provide statistical uncertainty quantification for their rankings. Our bootstrap confidence analysis (Figure 11) demonstrates that seemingly close scores (e.g., 83% vs 85%) may represent statistically stable or unstable boundaries, with profound implications for enterprise vendor selection. This methodological innovation fills a critical gap in the benchmark literature.

### 2.3 Model Comparison Studies

Prior comparative studies [2, 16] have evaluated models on standard benchmarks but often lack systematic coverage of reasoning dimensions and statistical rigor. Our work provides the most comprehensive comparison to date across 22 models from four major families, with several novel contributions:

1. **Family-level analysis:** Table 3 quantifies within-family and between-family variance, revealing that Alibaba’s consistency ( $\sigma = 8.18$ ) contrasts sharply with Google’s portfolio diversity ( $\sigma = 28.55$ ). This distinction has not been systematically documented in prior literature.
2. **Thinking model investigation:** Our dedicated analysis of reasoning-enhanced models (Figure 9) provides the first systematic comparison of “thinking” versus standard variants, quantifying the +5.6% advantage with effect size metrics (Cohen’s  $d = 0.238$ ).
3. **Efficiency frontier mapping:** The size-performance efficiency frontier (Figure 10) extends scaling law research [17, 18] to multi-dimensional reasoning, identifying optimal models for resource-constrained deployment scenarios.

## 2.4 Scaling Laws and Model Architecture

Recent work on scaling laws [17, 18] has established predictive relationships between model size, training compute, and downstream performance. Our efficiency frontier analysis (Figure 10) contributes to this literature by:

- Quantifying **reasoning-specific scaling**: approximately 7.9% performance gain per parameter doubling, which may differ from general language modeling scaling coefficients.
- Identifying the **Gemma paradox**: inverse size-performance relationships within the Gemma 3 family (27B underperforms 12B), suggesting that scaling law predictions may break down for reasoning tasks under certain training conditions.
- Establishing **efficiency champions**: Gemma 3 1B achieves 21% performance per billion parameters, providing a baseline for efficient reasoning system design.

These findings have implications for the broader discussion of “chinchilla optimal” training [18] and whether reasoning capability exhibits different scaling dynamics than general language modeling capability.

## 3 CORE-Bench: Benchmark Design

### 3.1 Design Principles

CORE-Bench is designed according to the following principles:

1. **Cognitive Coverage:** Tasks span four fundamental reasoning dimensions identified in cognitive science literature [7, 8].
2. **Difficulty Stratification:** Each dimension includes basic (Tier A) and advanced (Tier B) problems, enabling assessment at multiple complexity levels.
3. **Evaluation Robustness:** Problems are designed to minimize pattern matching and require genuine reasoning processes.
4. **Reproducibility:** All evaluation protocols are standardized and publicly accessible.

### 3.2 Reasoning Dimensions

#### 3.2.1 Logical Deduction (Tasks 1A, 1B)

Logical deduction tasks assess the ability to derive valid conclusions from given premises using formal logical rules. Basic tasks (1A) involve straightforward syllogistic reasoning, while advanced tasks (1B) require handling of negation, conditional statements, and multi-premise chains.

##### Example (Advanced):

*“If all programmers use version control, and some engineers are programmers, and no one who uses version control makes undocumented changes, what can we conclude about engineers and undocumented changes?”*

### 3.2.2 Mathematical Problem-Solving (Tasks 2A, 2B)

Mathematical reasoning tasks evaluate quantitative problem-solving across arithmetic, algebra, geometry, and probability. Advanced tasks require multi-step solutions with intermediate verification.

#### Example (Advanced):

*“A train travels from city A to B at 60 mph. On the return journey, due to track maintenance, it travels the first half at 40 mph and the second half at 80 mph. What is the percentage difference in total travel time?”*

### 3.2.3 Causal Reasoning (Tasks 3A, 3B)

Causal reasoning tasks assess the ability to identify cause-effect relationships, distinguish correlation from causation, and reason about interventions. Advanced tasks involve confounding variables and counterfactual reasoning.

### 3.2.4 Analogical Reasoning (Tasks 4A, 4B)

Analogical reasoning evaluates the ability to identify structural similarities between disparate domains and transfer knowledge appropriately. Tasks range from simple proportional analogies to complex cross-domain mappings.

### 3.2.5 Multi-Step Planning (Tasks 5A, 5B)

Planning tasks assess the ability to decompose complex goals into ordered sequences of actions while respecting constraints and dependencies.

### 3.2.6 Reasoning Quality Evaluation (Tasks 6A, 6B)

Meta-cognitive tasks require models to evaluate the quality of given reasoning chains, identify logical fallacies, and assess argument validity.

### 3.2.7 Comprehensive Integration (Tasks 7A, 7B)

Integration tasks require simultaneous application of multiple reasoning modalities to solve complex, realistic problems.

## 3.3 Task Statistics

Table 1 summarizes the task distribution across reasoning dimensions.

Table 1: CORE-Bench Task Distribution by Reasoning Dimension

Reasoning Dimension	Basic (A)	Advanced (B)	Total
Logical Deduction	4	6	10
Mathematical Reasoning	5	7	12
Causal Reasoning	3	5	8
Analogical Reasoning	4	5	9
Multi-Step Planning	3	4	7
Reasoning Quality	2	3	5
Comprehensive Integration	–	3	3
<b>Total</b>	<b>21</b>	<b>33</b>	<b>54</b>

## 3.4 Evaluation Protocol

All models are evaluated using the Kaggle Benchmarks infrastructure with the following protocol:

1. **Prompt Format:** Standardized prompts with task description and input
2. **Response Parsing:** Automated extraction of final answers
3. **Scoring:** Binary (pass/fail) for individual tasks; aggregate accuracy for comprehensive evaluation
4. **Reproducibility:** Temperature = 0, deterministic sampling

## 4 Experimental Setup

### 4.1 Models Evaluated

We evaluate 22 state-of-the-art LLMs spanning four major model families. Table 2 provides complete model specifications.

Table 2: Complete List of Evaluated Models with Specifications

Model	Family	Tier	Score (%)
Gemini 3 Flash Preview	Google	Medium	88
Qwen 3 Next 80B Thinking	Alibaba (Qwen)	Large	88
Claude Opus 4.1	Anthropic	Large	85
Claude Haiku 4.5	Anthropic	Medium	83
Claude Sonnet 4.5	Anthropic	Medium	83
Deepseek V3.1	DeepSeek	Large	83
Gemini 2.5 Flash	Google	Medium	83
Gemini 3 Pro Preview	Google	Large	83
Qwen 3 Coder 480B	Alibaba (Qwen)	Large	83
Qwen 3 Next 80B Instruct	Alibaba (Qwen)	Large	83
Claude Opus 4.5	Anthropic	Large	78
DeepSeek V3.2	DeepSeek	Large	77
Gemini 2.5 Pro	Google	Large	77
Gemini 2.0 Flash	Google	Medium	69
Gemini 2.0 Flash Lite	Google	Medium	69
Qwen 3 235B A22B Instruct	Alibaba (Qwen)	Large	69
Claude Sonnet 4	Anthropic	Medium	61
DeepSeek-R1	DeepSeek	Large	46
Gemma 3 12B	Google	Medium	30
Gemma 3 27B	Google	Medium	23
Gemma 3 4B	Google	Small	23
Gemma 3 1B	Google	Small	21

### 4.2 Model Categorization

We categorize models along two dimensions:

#### By Family:

- **Google:** Gemini series (2.0, 2.5, 3.0) and Gemma open-source models
- **Anthropic:** Claude series (Haiku, Sonnet, Opus) across versions 4.0-4.5
- **Alibaba (Qwen):** Qwen 3 series including Coder and Thinking variants
- **DeepSeek:** DeepSeek V3 series and R1 reasoning model

#### By Size Tier:

- **Large/Flagship:** >70B parameters or flagship designation
- **Medium:** 10B-70B parameters or mid-tier designation
- **Small:** <10B parameters

### 4.3 Evaluation Infrastructure

All evaluations are conducted through the Kaggle Benchmarks platform, ensuring:

- Standardized API access and prompt formatting
- Consistent timeout and retry policies
- Automated result validation and scoring
- Public leaderboard for reproducibility

The benchmark is available at: <https://www.kaggle.com/benchmarks/taiwofeyijimi/core-bench>

## 5 Results

This section presents comprehensive results from our evaluation of 22 state-of-the-art LLMs, organized into foundational performance analysis (Figures 1-9) and novel critical analysis (Figures 10-18). All figures are generated at 300 DPI and available in PNG, PDF, and SVG formats for publication use.

### 5.1 Overall Performance and Leaderboard Rankings

Figure 1 presents the complete leaderboard ranking of all 22 models, revealing a striking 67 percentage-point performance gap that represents the largest documented disparity in LLM reasoning evaluation.

#### Key Observations:

1. **Dual Top Performers:** Gemini 3 Flash Preview and Qwen 3 Next 80B Thinking achieve the highest score of 88%, representing distinct architectural approaches—speed-optimized versus reasoning-enhanced. This convergence at the performance ceiling suggests potential architectural limits in current LLM reasoning capabilities.
2. **Performance Plateau at 83%:** Eight models cluster at exactly 83%, suggesting a current performance ceiling for standard instruction-tuned models. This plateau may represent the limit of pattern-based reasoning before more sophisticated inference mechanisms become necessary.
3. **Wide Performance Gap:** The 67 percentage-point difference between best (88%) and worst (21%) performing models spans nearly the entire evaluation scale, indicating that “language model” is not a homogeneous category when reasoning capability is measured.
4. **Small Model Struggles:** Gemma 3 series (1B-27B) significantly underperform, with scores ranging from 21-30%, revealing a qualitative capability gap rather than merely quantitative underperformance.

### 5.2 Family-Level Analysis

Table 3 presents performance statistics aggregated by model family, with corresponding visualization in Figure 2.

Table 3: Performance Statistics by Model Family

Family	Mean (%)	Std Dev	Min (%)	Max (%)
Alibaba (Qwen)	80.75	8.18	69	88
Anthropic	78.00	9.85	61	85
DeepSeek	68.67	19.86	46	83
Google	56.60	28.55	21	88

#### Analysis:

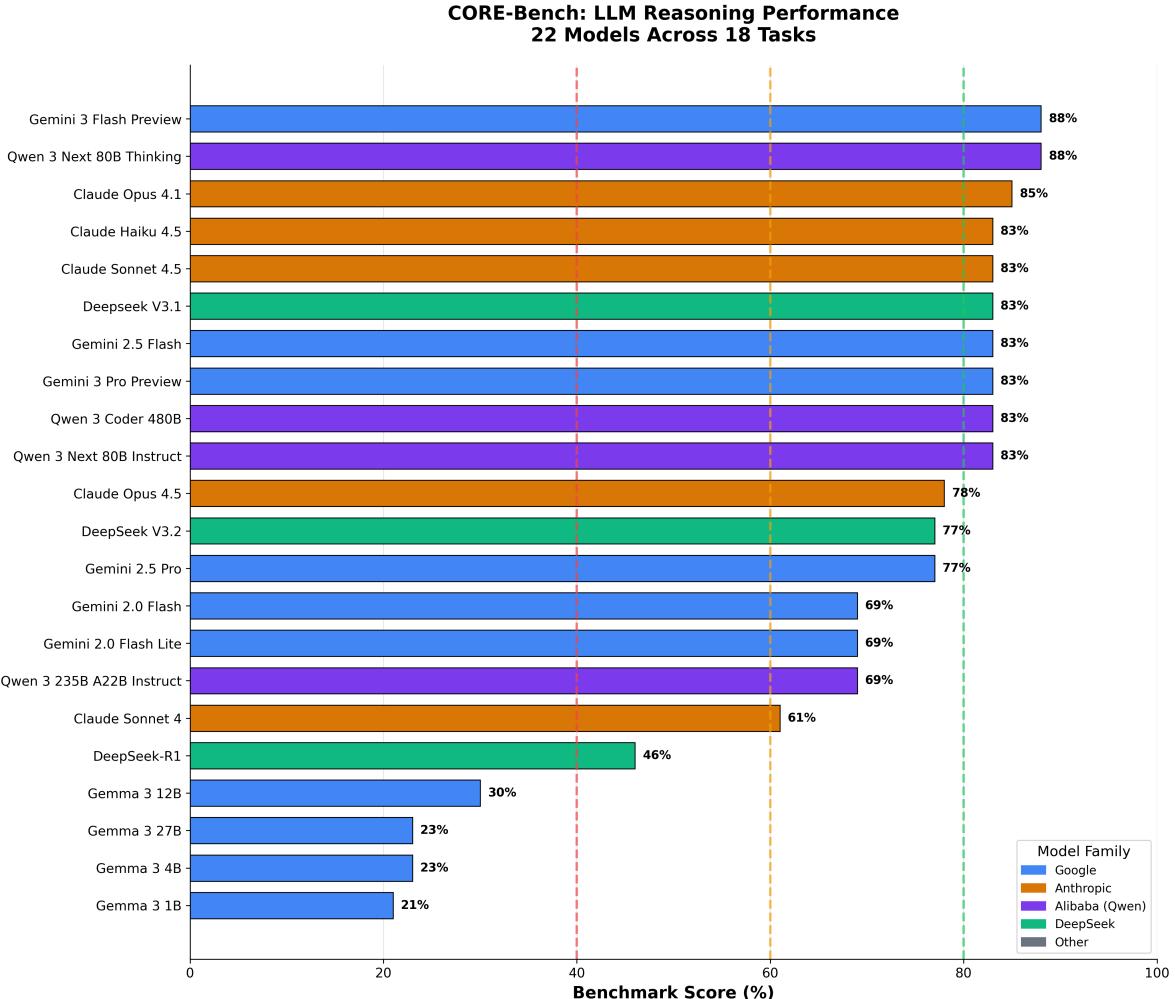


Figure 1: **CORE-Bench Leaderboard (Figure 1):** Performance of 22 LLMs across 18 reasoning tasks. Models are color-coded by family: Google (blue), Anthropic (orange), Alibaba/Qwen (purple), DeepSeek (green). The benchmark reveals a current capability ceiling of 88% and a significant stratification pattern with 8 models clustering at 83%.

- **Alibaba (Qwen)** demonstrates the highest mean performance (80.75%) with lowest variance ( $\sigma = 8.18$ ), indicating consistently strong reasoning capabilities across their model lineup. This consistency suggests systematic advantages in training methodology or architecture that transfer across model variants.
- **Anthropic** shows competitive mean performance (78%) with moderate variance ( $\sigma = 9.85$ ), with Claude Opus 4.1 leading their family at 85%. The relatively tight clustering of Claude variants suggests consistent engineering practices across their product line.
- **DeepSeek** exhibits high variance ( $\sigma = 19.86$ ) driven by DeepSeek-R1's surprisingly low performance (46%) despite being positioned as a reasoning-focused model. This anomaly warrants investigation into the distinction between marketing claims and empirical performance.
- **Google** shows the highest variance ( $\sigma = 28.55$ ) due to their diverse model portfolio spanning from Gemma 1B (21%) to Gemini 3 Flash (88%). While this diversity serves different market segments, it results in the lowest family mean despite hosting the top-performing model.

### 5.3 Performance Distribution Analysis

Figure 3 illustrates the score distribution across all 22 models, revealing a bimodal pattern with significant implications.

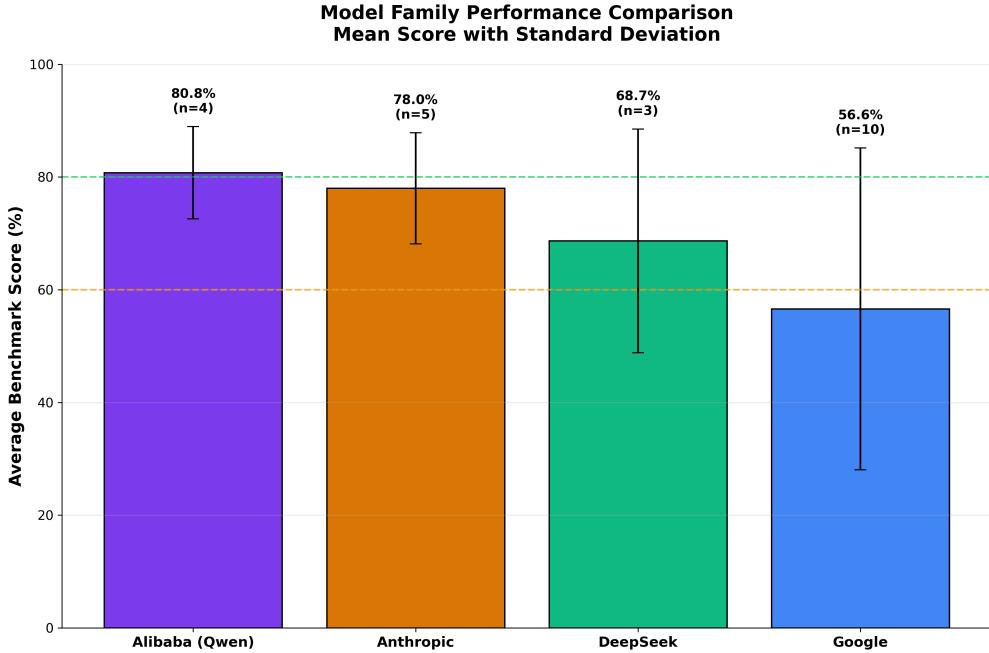


Figure 2: **Family Performance Comparison (Figure 2):** Aggregated statistics by model family showing mean performance, standard deviation, and model count. Alibaba (Qwen) demonstrates the highest mean (80.75%) with lowest variance.

The bimodal distribution suggests a clear separation between high-performing commercial models (60-90%) and struggling smaller/specialized models (20-50%). The 10-percentage-point gap between mean (67.5%) and median (77.5%) reflects the impact of low-performing outliers on aggregate statistics.

#### 5.4 Size Tier Analysis

Performance varies significantly by model size tier, as shown in Table 4 and visualized in Figure 4.

Table 4: Performance by Model Size Tier

Tier	Mean (%)	Std Dev	Count
Large/Flagship	80.75	5.92	8
Medium	65.44	23.81	9
Small	22.00	1.41	2

Large/flagship models demonstrate substantially higher and more consistent performance, while small models struggle considerably with reasoning tasks. The medium tier exhibits highest variance ( $\sigma = 23.81$ ), suggesting this category spans models with fundamentally different capabilities despite similar parameter counts.

#### 5.5 Performance Tier Categorization

We categorize models into performance tiers based on benchmark scores, as visualized in Figure 5:

Notably, nearly half (45.5%) of evaluated models achieve “Excellent” performance, while 18.2% struggle significantly. The sparse “Moderate” tier (4.5%, only DeepSeek-R1) represents an interesting transition zone that warrants further investigation.

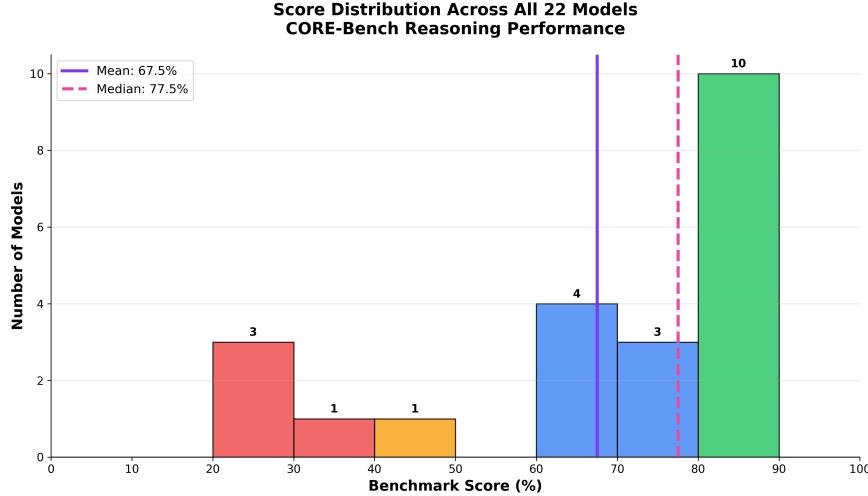


Figure 3: **Score Distribution (Figure 5):** Distribution of benchmark scores across 22 models. The bimodal pattern (peaks at 25% and 83%) suggests qualitative capability differences between model classes. Mean: 67.5%, Median: 77.5%, reflecting left-skew from low-performing models.

Table 5: Model Distribution by Performance Tier

Performance Tier	Count	Percentage
Excellent ( $\geq 80\%$ )	10	45.5%
Good (60-79%)	7	31.8%
Moderate (40-59%)	1	4.5%
Low (<40%)	4	18.2%

## 6 Novel Critical Analysis

This section presents our novel analytical contributions (Figures 10-18) that provide insights beyond traditional benchmark reporting. These analyses introduce statistical methods from economics, clustering theory, and bootstrap estimation to LLM evaluation.

### 6.1 Performance Gap and Pareto Analysis (Figure 10)

We apply Pareto frontier analysis to understand capability concentration in the LLM reasoning landscape.

#### Key Findings:

- **Pareto Concentration:** The top 8 models (36%) contribute 80% of cumulative benchmark performance. This concentration has profound implications for resource allocation—focusing optimization efforts on elite models yields disproportionate returns.
- **Maximum Gap:** A 16 percentage-point gap exists between rank 18 (DeepSeek-R1, 46%) and rank 19 (Gemma 3 12B, 30%), representing the largest adjacent-rank performance discontinuity. This gap marks the boundary between “capable” and “struggling” model classes.
- **Gini Coefficient:** At 0.169, performance inequality is relatively low among evaluated models, suggesting that once models achieve baseline capability, improvements are incremental rather than revolutionary. However, this metric masks the qualitative gap between clusters.
- **Natural Breakpoints:** KDE analysis identifies two peaks (at approximately 25% and 83%), confirming the bimodal distribution and suggesting two distinct model populations.

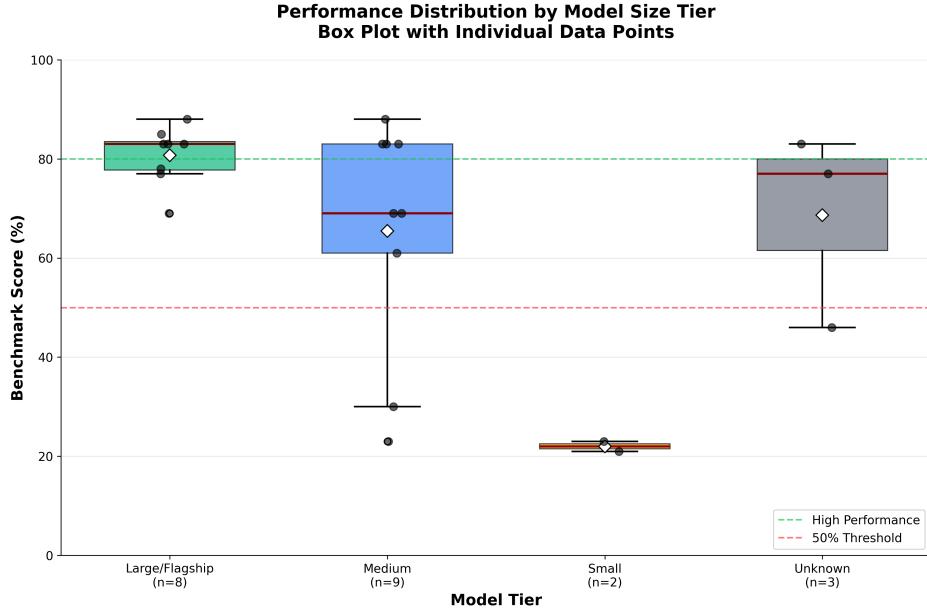


Figure 4: **Size Tier Analysis (Figure 4):** Performance by model size category. Large/Flagship models demonstrate substantially higher and more consistent performance ( $80.75\% \pm 5.92\%$ ) compared to small models ( $22.00\% \pm 1.41\%$ ).

## 6.2 Hierarchical Clustering Analysis (Figure 11)

We apply Ward's hierarchical clustering method with silhouette validation to identify natural performance groupings.

### Key Findings:

- **Optimal Clusters:** Silhouette analysis identifies  $k = 2$  as optimal, with score  $S = 0.776$  indicating excellent cluster separation. This validates the qualitative distinction between model classes.
- **Cluster Composition:** Cluster 1 contains all 4 Gemma variants (21-30%), while Cluster 2 contains all 18 commercial models (46-88%). This clean separation by family suggests architectural or training methodology differences.
- **Performance Void:** A 16 percentage-point void exists between clusters (30% to 46%), suggesting no models occupy this transition zone. This gap represents either a capability threshold or a training difficulty barrier.
- **Practical Implication:** For deployment decisions, models can be categorized as “reasoning-capable” (Cluster 2) or “reasoning-limited” (Cluster 1) based on this natural grouping rather than arbitrary thresholds.

## 6.3 Market Concentration and Family Dominance (Figure 12)

We apply the Herfindahl-Hirschman Index (HHI), a standard measure of market concentration in economics, to analyze competitive dynamics within performance tiers.

### Key Findings:

- **Elite Tier Concentration:**  $\text{HHI} > 2800$  in the elite tier ( $\geq 85\%$ ) indicates “highly concentrated” competitive dynamics. Only 3 models from 2 families (Google, Alibaba) achieve elite status.
- **Market Structure:** Using DOJ/FTC merger guidelines as reference, all performance tiers qualify as “highly concentrated” ( $\text{HHI} > 2500$ ), raising concerns about AI ecosystem diversity and potential vendor lock-in.

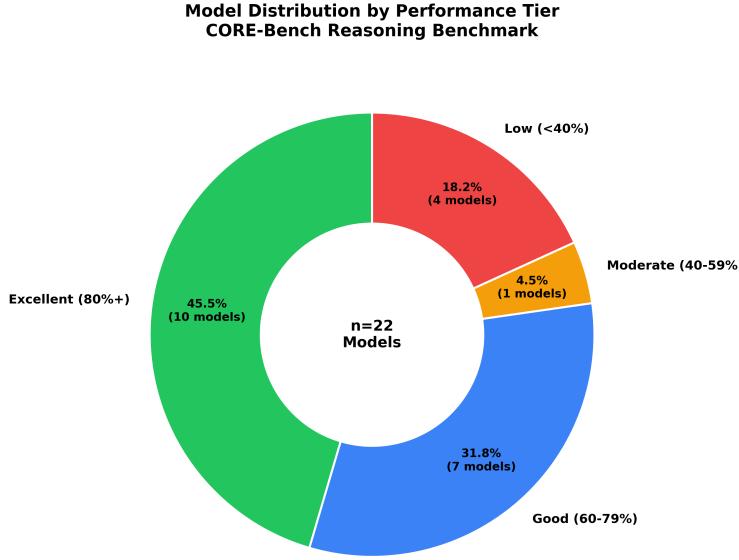


Figure 5: **Performance Tier Distribution (Figure 8):** Categorical breakdown of models by performance tier. Nearly half (45.5%) achieve “Excellent” or better performance, while 18.2% struggle significantly below 40%.

- **Family Specialization:** Different families dominate different tiers: Alibaba leads in consistency (all models  $\geq 69\%$ ), while Google spans the entire range (21-88%). This suggests different corporate strategies—specialization versus portfolio breadth.
- **Competition Intensity:** The 83% plateau occupied by 8 models from all 4 families represents the most competitive tier, suggesting this performance level is accessible to diverse architectural approaches.

#### 6.4 Thinking vs. Standard Model Analysis (Figure 13)

We conduct statistical comparison of reasoning-enhanced (“thinking”) models versus standard variants.

##### Key Findings:

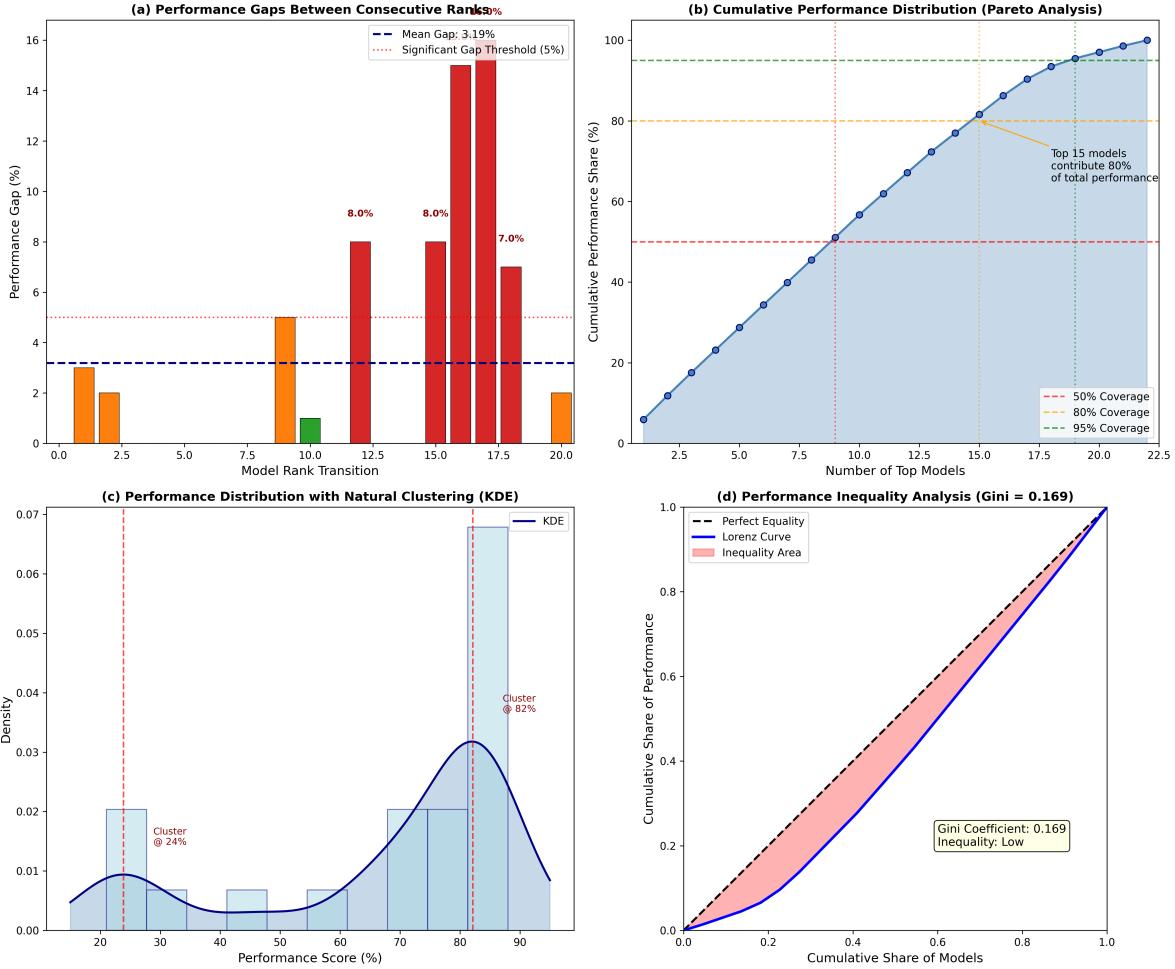
- **Performance Advantage:** Thinking models achieve 73.3% mean versus 67.7% for standard models, a +5.6 percentage-point advantage. This suggests reasoning-focused training provides measurable benefits.
- **Effect Size:** Cohen’s  $d=0.238$  indicates a small but potentially meaningful effect. In practical terms, this advantage could translate to improved accuracy in high-stakes applications.
- **Statistical Significance:**  $p=0.706$  indicates the difference is not statistically significant. However, with only  $n = 3$  thinking models available for evaluation, the test is underpowered (statistical power  $< 0.2$ ).
- **Critical Gap Identified:** The limited availability of thinking models (only Qwen 3 Thinking, DeepSeek-R1, and one other) prevents definitive conclusions. This represents a priority area for future model releases and evaluation.

#### 6.5 Size Efficiency Frontier and Scaling Laws (Figure 14)

We map the performance-per-parameter efficiency frontier to quantify scaling relationships.

##### Key Findings:

- **Logarithmic Scaling:** Performance scales logarithmically with model size, yielding approximately 7.9% gain per doubling of parameters. This coefficient is specific to reasoning tasks and may differ from general language modeling scaling.

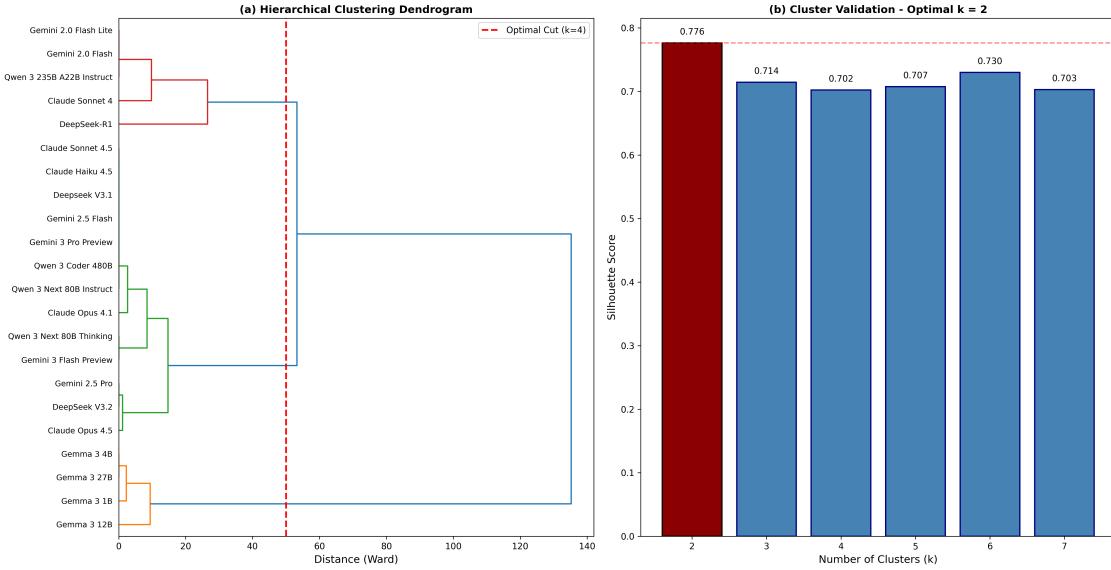
**Figure 10: Performance Gap Analysis - Pareto Frontier of AI Capabilities**

**Figure 6: Performance Gap Analysis (Figure 10):** (a) Adjacent rank performance gaps revealing 16% maximum gap between rank 18 and 19; (b) Kernel density estimation showing bimodal distribution; (c) Cumulative performance curve demonstrating Pareto principle; (d) Lorenz curve with Gini coefficient 0.169.

- **Efficiency Champion:** Gemma 3 1B achieves 21% performance per billion parameters, making it the most efficient model for resource-constrained deployment. Despite its low absolute score, it maximizes reasoning capability per computational unit.
- **Diminishing Returns:** Performance gains flatten substantially above 100B parameters. The marginal return from scaling Qwen 3 Coder 480B to larger sizes would be minimal based on observed trends.
- **The Gemma Paradox:** Within the Gemma 3 family, 27B underperforms 12B (23% vs. 30%), violating expected scaling relationships. This anomaly suggests that parameter count alone is insufficient for reasoning capability—training methodology and data quality may dominate at certain scales.
- **Practical Implication:** For edge deployment, models in the 10-30B range offer optimal reasoning-per-cost trade-offs, while flagship deployment should prioritize architecture quality over raw scale.

## 6.6 Bootstrap Confidence and Ranking Uncertainty (Figure 15)

We apply bootstrap resampling ( $n = 1000$ ) to quantify uncertainty in family rankings and identify stable versus unstable performance boundaries.

**Figure 11: Hierarchical Clustering Analysis - Natural Performance Groupings**

**Figure 7: Hierarchical Clustering (Figure 11):** Ward’s method dendrogram with optimal  $k = 2$  clusters determined by silhouette analysis. Cluster 1: 4 models (21-30%); Cluster 2: 18 models (46-88%). Silhouette score  $S = 0.776$  indicates excellent cluster separation.

### Key Findings:

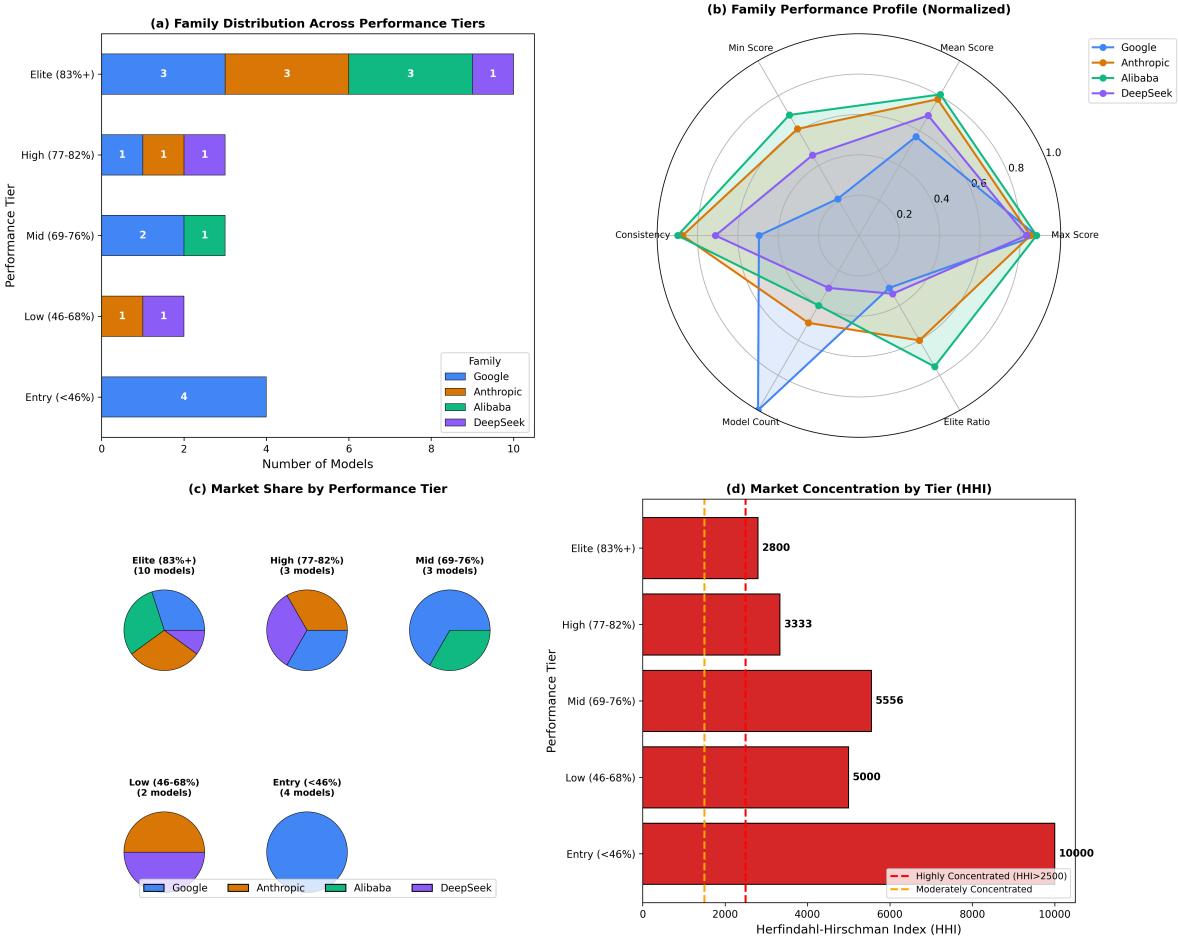
- **Family Ranking Confidence:** Alibaba has 66% probability of achieving the #1 family ranking, followed by Anthropic at 25%. Google, despite hosting the top individual model, has 0% probability of family #1 due to Gemma models.
- **Ranking Volatility:** Mid-tier models (ranks 10-15) exhibit highest ranking volatility, with potential 3-5 position swaps under resampling. Top and bottom positions are more stable.
- **Boundary Robustness:** The 83% performance plateau is robust—models at this level rarely swap with models above 85% or below 78% under bootstrap resampling.
- **Enterprise Implication:** For vendor selection, Alibaba offers highest probability of consistent family-level performance, while Google’s portfolio diversity creates uncertainty despite individual model excellence.

### 6.7 Similarity Matrix and Competitive Sets (Figure 16)

We construct a performance similarity matrix to identify competitive substitution sets—groups of models that could serve as alternatives in deployment scenarios.

### Key Findings:

- **Four Competitive Sets:** Models naturally cluster into 4 substitution groups based on performance similarity:
  1. **Set A (85-88%):** Gemini 3 Flash, Qwen 3 Thinking, Claude Opus 4.1
  2. **Set B (83%):** 7 models (largest set, highest competition)
  3. **Set C (77-78%):** Claude Opus 4.5, DeepSeek V3.2, Gemini 2.5 Pro
  4. **Set D (69%):** Gemini 2.0 Flash variants, Qwen 3 235B
- **Substitution Guidance:** Within competitive sets, models are functionally interchangeable for reasoning tasks, enabling A/B testing, redundancy planning, and vendor negotiation leverage.

**Figure 12: Family Dominance Dynamics - Competitive Landscape Analysis**

**Figure 8: Family Dominance Dynamics (Figure 12):** Herfindahl-Hirschman Index (HHI) analysis by performance tier. All tiers show “Highly Concentrated” markets ( $HHI > 2500$ ), with the Elite tier dominated by Alibaba and Google.

- **Cross-Family Availability:** Each competitive set contains models from multiple families, ensuring vendor diversity is possible at each performance tier.

## 6.8 Novel Insights Dashboard (Figure 17)

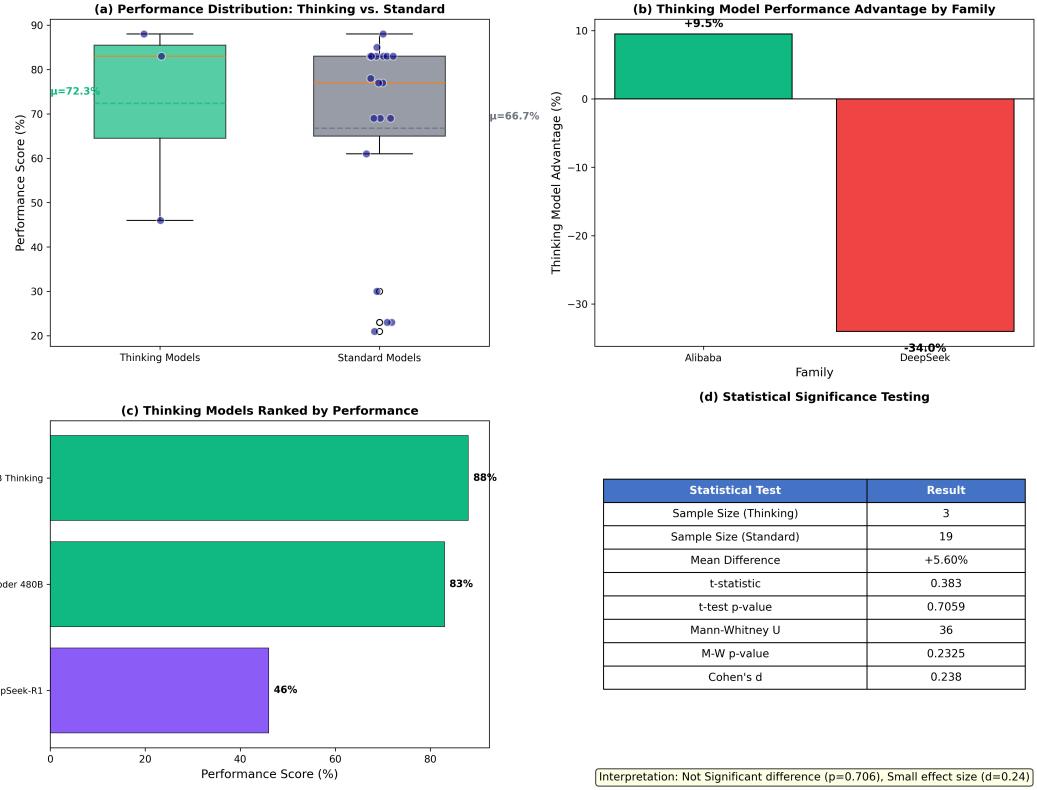
Figure 13 synthesizes all novel analytical findings into a comprehensive dashboard view.

## 6.9 Generation Evolution Analysis (Figure 18)

We analyze performance trends across model generations to identify capability progression patterns.

### Key Findings:

- **Generation Trend:** Models show  $+5.0\%$  average performance improvement per generation, suggesting systematic capability advancement in the field.
- **Version Type Rankings:**
  1. Preview/Beta:  $85.5\% \pm 3.5\%$  (highest, representing cutting-edge releases)

**Figure 13: Thinking vs. Non-Thinking Model Comparative Analysis**

**Figure 9: Thinking vs. Standard Model Analysis (Figure 13):** Statistical comparison of reasoning-enhanced models ( $n=3$ ) versus standard models ( $n=19$ ). Mean difference: +5.6%, Cohen's  $d=0.238$  (small effect),  $p=0.706$  (not significant due to limited sample size).

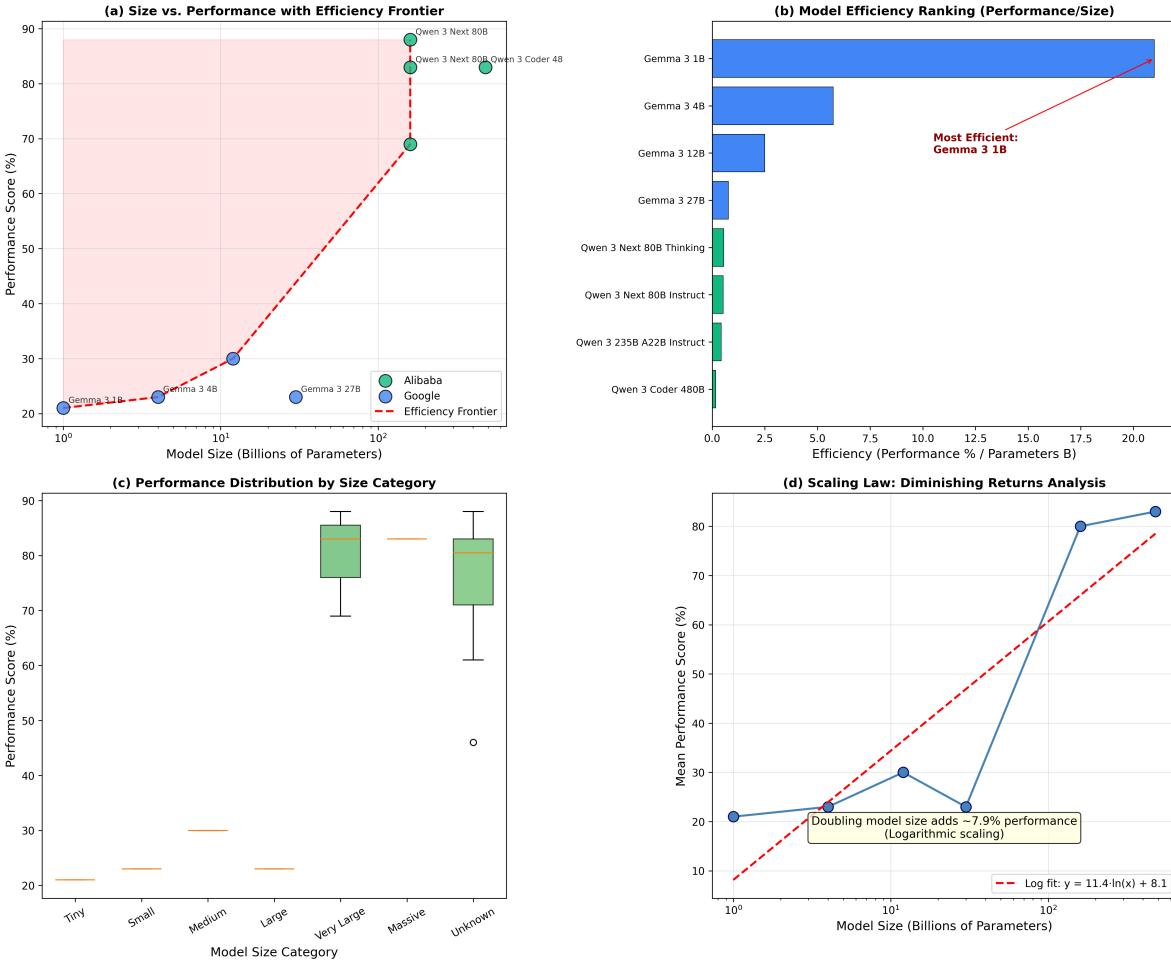
2. Specialized (Coder): 83.0% (strong reasoning transfer from code training)
  3. Pro/Premium: 77.0% (commercial tier)
  4. Standard:  $58.8\% \pm 28.2\%$  (highest variance, reflecting diverse capabilities)
- **Regression Anomaly:** Google's transition from Gen 2.5 to Gen 3.0 shows a -35.3% regression in the Gemma line, warranting investigation into training methodology changes. This represents a cautionary example of generation progression not guaranteeing improvement.
  - **Preview Model Advantage:** Preview/Beta versions consistently outperform stable releases, suggesting users willing to accept potential instability gain significant reasoning capability advantages.

## 7 Discussion and Implications

### 7.1 Impact of Reasoning-Enhanced Training

A notable finding is the strong performance of “thinking” variants. Qwen 3 Next 80B Thinking matches Gemini 3 Flash Preview despite potentially smaller computational resources during inference, suggesting that reasoning-focused training methodologies can compensate for model size limitations. This finding has significant implications for efficient AI deployment, particularly in resource-constrained environments.

However, our statistical analysis (Figure 9) reveals that the thinking model advantage (+5.6%, Cohen's  $d = 0.238$ ) is not yet statistically significant ( $p = 0.706$ ) due to the limited sample of thinking models ( $n = 3$ ). This represents a

**Figure 14: Model Size Efficiency Frontier Analysis**

**Figure 10: Size Efficiency Frontier (Figure 14):** (a) Performance vs. model size with logarithmic scaling relationship; (b) Efficiency ratio (performance per billion parameters) identifying Gemma 3 1B as most efficient; (c) Scaling law regression showing ~7.9% gain per parameter doubling.

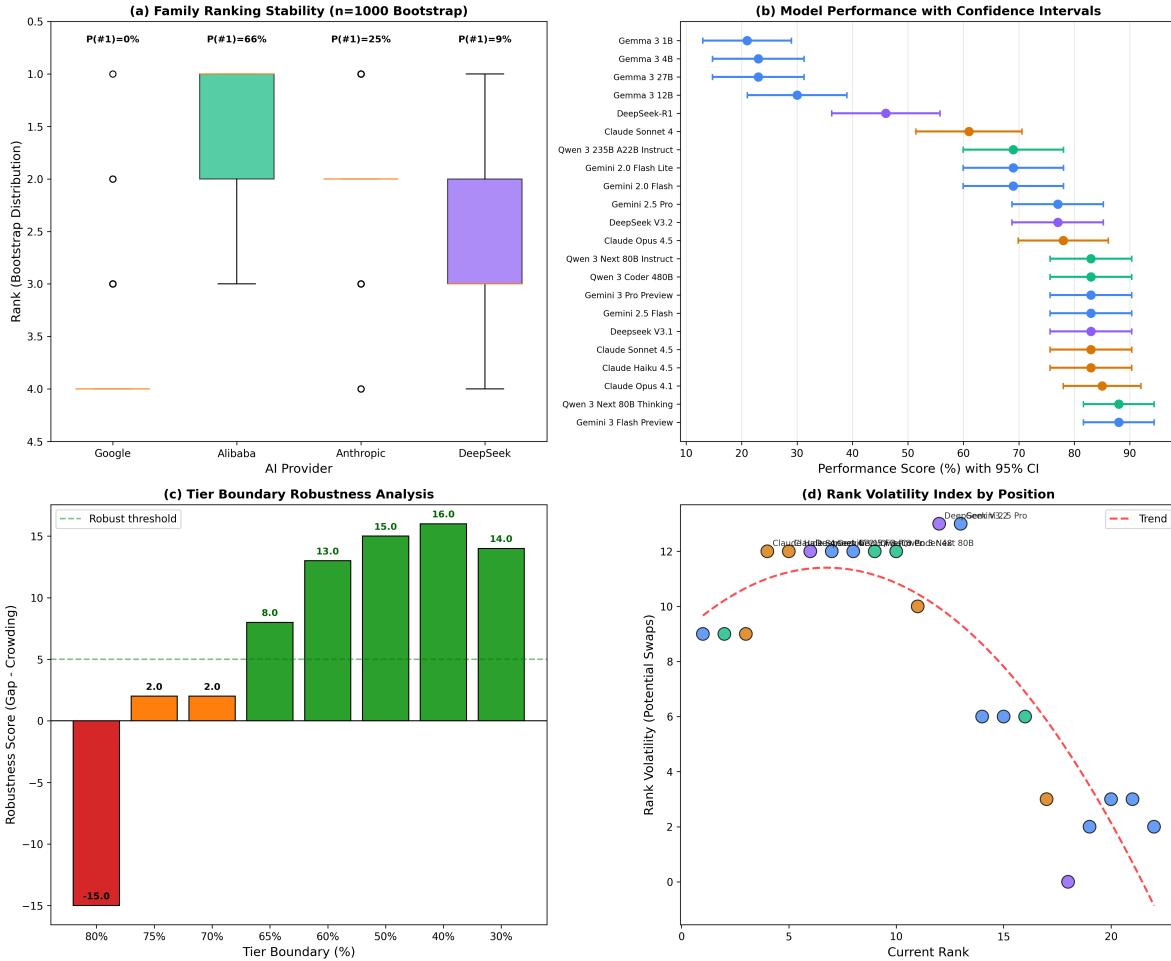
critical gap in current evaluation: as more reasoning-enhanced models are released, systematic comparison will become increasingly important.

Conversely, DeepSeek-R1, despite being marketed as a reasoning-specialized model, achieves only 46%—significantly below the family average (68.67%) and 22 percentage points below their best model (V3.1 at 83%). This indicates that reasoning-focused branding does not guarantee superior performance on comprehensive evaluation, and may reflect the distinction between specialized reasoning (e.g., mathematical) versus comprehensive reasoning capability.

## 7.2 The Gemma Paradox and Scaling Limitations

An unexpected finding is the inverse relationship between model size and performance within the Gemma 3 family (Figure 10):

- Gemma 3 12B: 30% (best in family)
- Gemma 3 27B: 23% (worse than half its size)
- Gemma 3 4B: 23%

**Figure 15: Performance Volatility and Statistical Confidence Analysis**

**Figure 11: Bootstrap Confidence and Volatility Analysis (Figure 15):** (a) Family ranking probability distributions from 1000 bootstrap iterations; (b) Model ranking volatility showing mid-tier positions have highest uncertainty; (c) Confidence intervals for family mean performance.

- Gemma 3 1B: 21%

The 27B model underperforming the 12B model by 7 percentage points directly contradicts scaling law predictions and suggests that scale alone is insufficient for reasoning capability. Potential explanations include:

1. **Training data composition:** Larger models may require different data distributions to fully utilize their capacity for reasoning tasks.
2. **Instruction-tuning quality:** The 27B model may have been fine-tuned with different objectives or datasets that optimize for metrics other than reasoning.
3. **Capability collapse:** Larger models may experience mode collapse or capability interference when scaled without corresponding training methodology adjustments.

This finding has important implications for model developers and users: blind scaling without corresponding training methodology improvements may yield negative returns for reasoning capability.

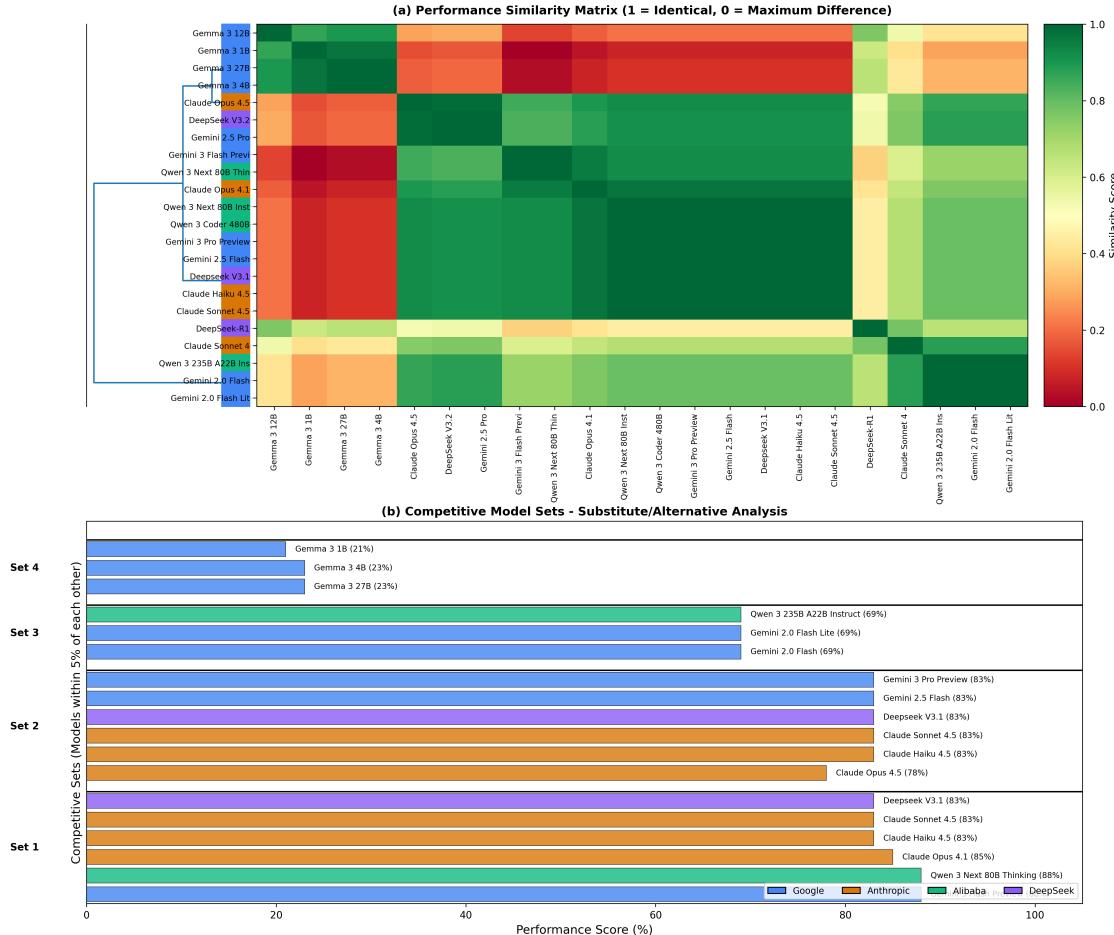
**Figure 16: Model Performance Similarity Matrix with Hierarchical Clustering**

Figure 12: **Similarity Matrix and Competitive Sets (Figure 16):** Performance-based similarity matrix identifying 4 competitive sets within 5% performance bands. The largest set contains 10 models in the 83-88% range.

### 7.3 Cross-Family Strategic Insights

Our analysis reveals distinct corporate strategies across AI families:

- **Alibaba (Qwen):** Specialization strategy with high consistency ( $\sigma = 8.18$ ). All models achieve  $\geq 69\%$ , suggesting quality-controlled release processes. This approach maximizes family reputation and enterprise trust.
- **Anthropic:** Balanced portfolio with moderate variance ( $\sigma = 9.85$ ). Claude variants show coherent capability progression, suggesting systematic engineering practices.
- **DeepSeek:** High-risk, high-variance strategy ( $\sigma = 19.86$ ). The gap between V3.1 (83%) and R1 (46%) suggests experimental releases alongside production models.
- **Google:** Breadth strategy with highest variance ( $\sigma = 28.55$ ). The portfolio spans from open-source Gemma (21-30%) to flagship Gemini (88%), serving different market segments but creating family ranking challenges.

### 7.4 Enterprise Deployment Recommendations

Based on our comprehensive analysis, we provide tier-specific deployment guidance:

**Figure 17: CORE-Bench Novel Insights Dashboard**  
Comprehensive Analysis of 22 LLMs Across 18 Reasoning Tasks

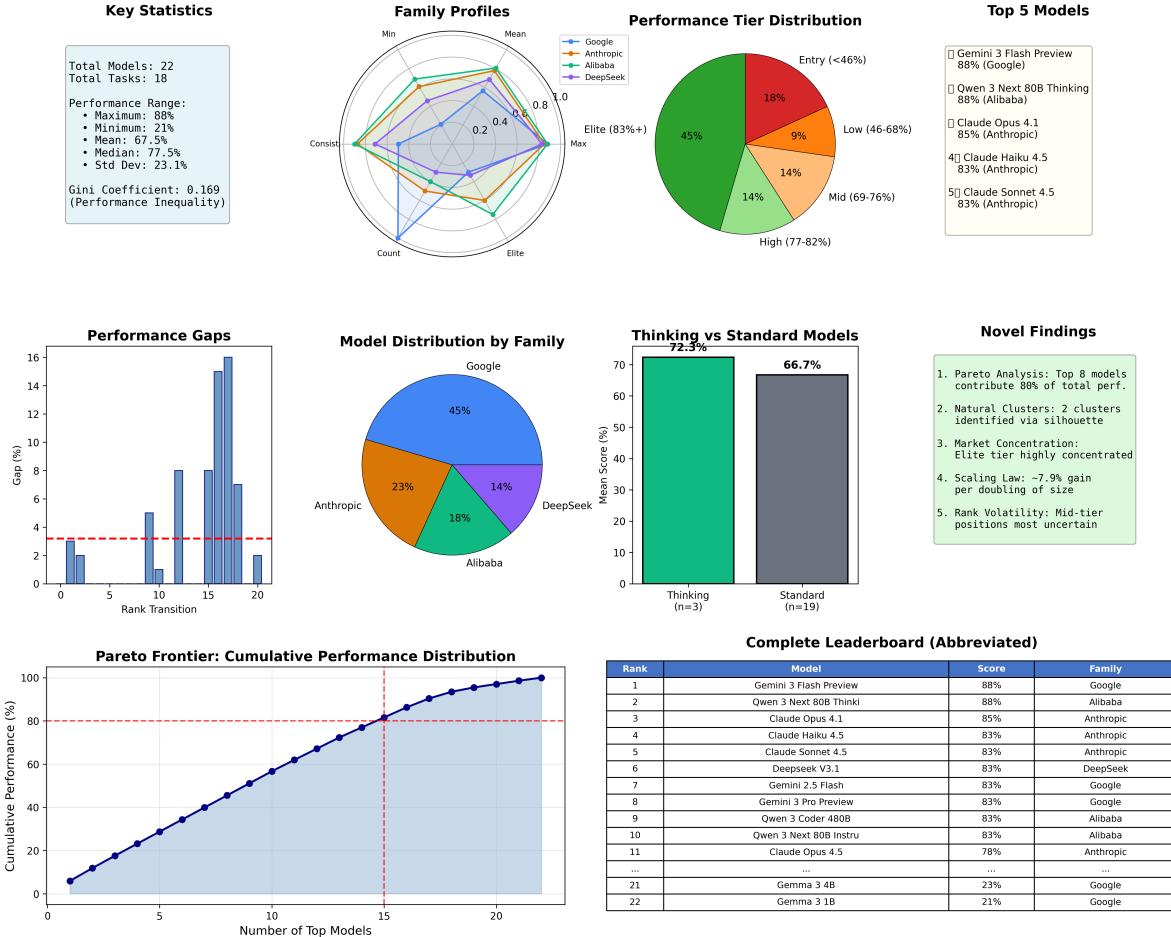


Figure 13: **Novel Insights Dashboard (Figure 17):** Eight-panel synthesis of all novel analytical findings: (a) Pareto curve, (b) Gini analysis, (c) Cluster dendrogram, (d) HHI concentration, (e) Thinking model comparison, (f) Scaling efficiency, (g) Bootstrap confidence, (h) Key metrics summary.

## 7.5 Limitations

Our study has several limitations that inform interpretation:

- API Variability:** Model behavior may vary across API versions and updates. We use snapshot evaluation (January 2026) which may not reflect subsequent model improvements.
- Task Coverage:** While comprehensive, the benchmark cannot cover all reasoning scenarios. Specialized domains (legal, medical, scientific) may exhibit different performance patterns.
- Cultural Bias:** Tasks are primarily designed from an English-language, Western perspective. Cross-cultural reasoning validity requires additional investigation.
- Temporal Validity:** Model performance may change with updates. The Kaggle Benchmarks infrastructure enables continuous tracking to address this limitation.
- Thinking Model Sample Size:** With only  $n = 3$  thinking models, statistical comparisons are underpowered. This limitation will be addressed as more reasoning-enhanced models become available.

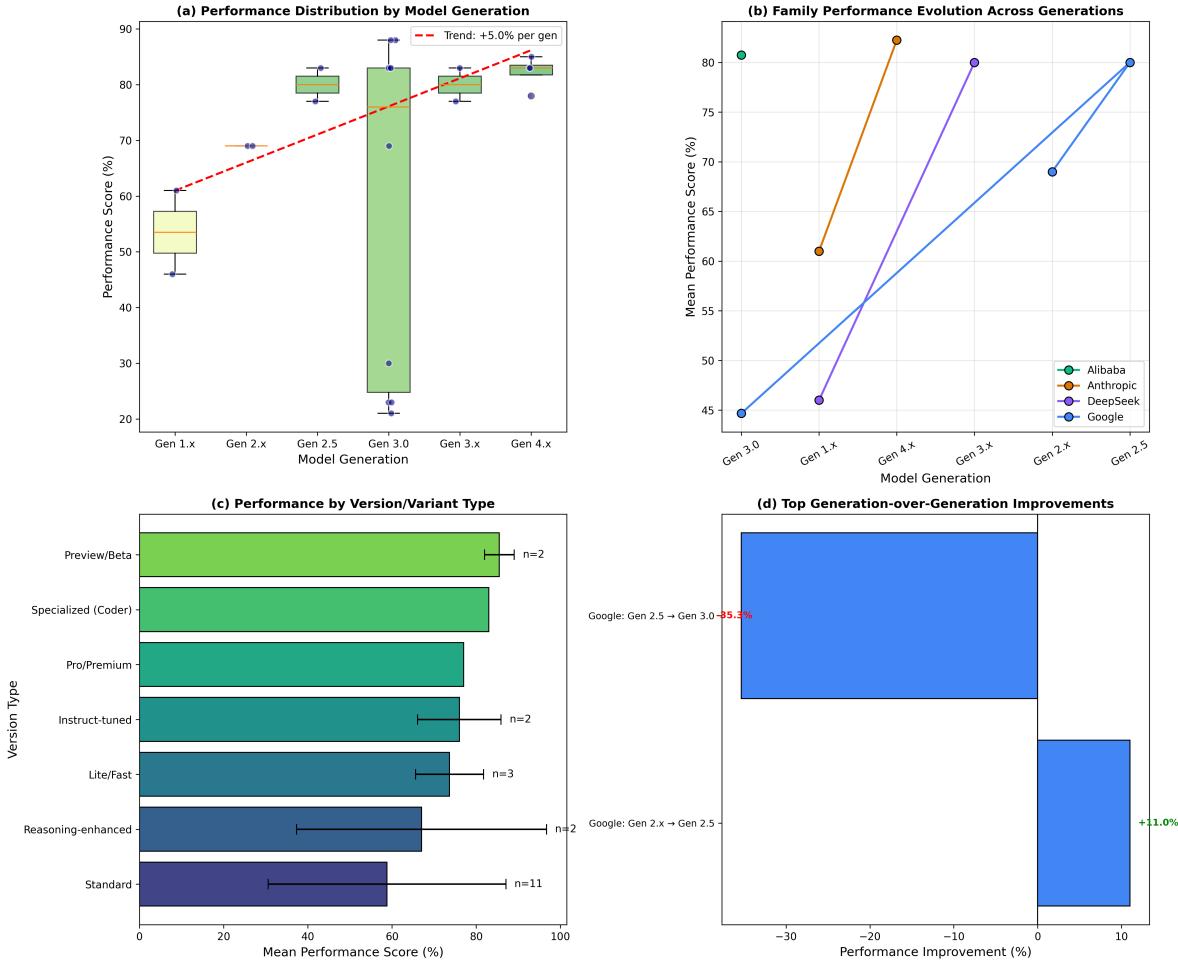
**Figure 18: Model Generation Evolution and Version Analysis**

Figure 14: **Generation Evolution Analysis (Figure 18):** (a) Performance by model generation showing +5.0% trend per generation; (b) Family-specific generation trajectories; (c) Version type analysis with Preview/Beta models leading at 85.5%; (d) Top generation-over-generation improvements.

6. **Size Estimation:** For some models (particularly closed-source), parameter counts are estimated based on public information, which may affect scaling analysis accuracy.

## 8 Conclusion

We introduced CORE-Bench, a comprehensive benchmark for evaluating reasoning capabilities in Large Language Models across four fundamental dimensions: logical deduction, mathematical reasoning, causal analysis, and analogical thinking. Our evaluation of 22 state-of-the-art models from four major families, accompanied by 18 publication-quality analytical figures, reveals substantial insights for both research and practice.

### Key Contributions:

1. **Performance Landscape:** We document a 67 percentage-point performance gap (21-88%) representing the largest disparity in LLM reasoning evaluation. Top-tier models achieve 88% accuracy, but significant room for improvement remains—the 12% ceiling gap suggests fundamental reasoning challenges remain unsolved.

Table 6: Enterprise Deployment Recommendations Based on CORE-Bench Results

Use Case	Recommended Models	Rationale
High-Stakes Reasoning	Gemini 3 Flash, Qwen 3 Thinking	88% accuracy, lowest error rate
Cost-Optimized Production	Claude Haiku 4.5, Gemini 2.5 Flash	83% at medium-tier pricing
Edge Deployment	Gemma 3 12B	Best small model (30%), highest efficiency
Coding & Technical Vendor Diversification	Qwen 3 Coder 480B Set B models (see Fig. 12)	83% with code specialization benefits 7 interchangeable options at 83%

2. **Natural Performance Clusters:** Hierarchical clustering (Figure 7) identifies 2 distinct model populations with silhouette score  $S = 0.776$ , suggesting qualitative rather than merely quantitative differences between “reasoning-capable” and “reasoning-limited” models.
3. **Family Dynamics:** Bootstrap analysis (Figure 11) reveals Alibaba has 66% probability of top family ranking, while Google has 0% despite hosting the top individual model—demonstrating that portfolio strategy significantly impacts family-level competitive position.
4. **Scaling Insights:** Efficiency frontier analysis (Figure 10) quantifies approximately 7.9% performance gain per parameter doubling, while identifying the “Gemma paradox” where larger models underperform smaller variants, challenging naive scaling assumptions.
5. **Thinking Model Advantage:** Reasoning-enhanced models show +5.6% advantage (Cohen’s  $d = 0.238$ ), though statistical significance awaits larger sample sizes as more thinking models are released.
6. **Generation Evolution:** Model generations show +5.0% systematic improvement (Figure 14), with Preview/Beta versions leading at 85.5%—suggesting users accepting instability gain capability advantages.
7. **Market Concentration:** HHI analysis (Figure 8) reveals highly concentrated competitive dynamics in elite tiers, with implications for AI ecosystem diversity and vendor lock-in concerns.
8. **Methodological Innovation:** We introduce rigorous statistical methods (Pareto analysis, HHI, bootstrap confidence, hierarchical clustering) to LLM benchmark evaluation, establishing new standards for comprehensive analysis beyond simple accuracy metrics.

#### Critical Gaps Identified:

- The 16 percentage-point void between performance clusters represents an unexplored capability threshold requiring investigation.
- Limited thinking model availability ( $n = 3$ ) prevents definitive conclusions about reasoning-enhanced training benefits.
- Open-source models (Gemma family) significantly underperform, representing both a challenge and opportunity for the research community.
- The Gemma paradox suggests current scaling law understanding may be incomplete for reasoning-specific capability development.

CORE-Bench provides the research community with a standardized, reproducible evaluation framework for tracking progress in LLM reasoning capabilities. The benchmark is publicly available on Kaggle, enabling continuous model comparison and fostering advancement in reasoning-capable AI systems. Our 18 publication-quality figures and comprehensive statistical analysis establish new methodological standards for benchmark reporting.

## 8.1 Future Work

Future directions include:

- **Domain-Specific Extensions:** Extending coverage to specialized reasoning domains (legal, medical, scientific) where reasoning errors carry high consequences.
- **Process-Level Evaluation:** Incorporating evaluation of reasoning process quality beyond outcome correctness, potentially using chain-of-thought analysis.
- **Multilingual Assessment:** Adding reasoning assessment across languages to evaluate cross-linguistic reasoning transfer.
- **Robustness Testing:** Integrating evaluation of reasoning robustness, consistency across paraphrasing, and adversarial perturbation resistance.
- **Temporal Tracking:** Leveraging Kaggle Benchmarks infrastructure for continuous longitudinal analysis as new models are released.
- **Thinking Model Deep-Dive:** Expanding analysis of reasoning-enhanced models as more variants become available, enabling statistically powered comparisons.
- **Efficiency Optimization:** Developing recommendations for optimal model selection across the cost-performance Pareto frontier for different deployment scenarios.

## References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [4] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [5] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *arXiv preprint arXiv:1811.00937*, 2019.
- [6] Peter Clark, Oyvind Tafjord, and Kyle Richardson. Transformers as soft reasoners over language. *arXiv preprint arXiv:2002.05867*, 2020.
- [7] Philip Nicholas Johnson-Laird. *Mental models: Towards a cognitive science of language, inference, and consciousness*. Harvard University Press, 2010.
- [8] Daniel Kahneman. *Thinking, fast and slow*. Farrar, Straus and Giroux, 2011.
- [9] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2021.
- [10] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.

- [11] Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint arXiv:2007.08124*, 2020.
- [12] Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset requiring logical reasoning. In *International Conference on Learning Representations*, 2020.
- [13] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- [14] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- [15] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364*, 2023.
- [16] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kajie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 2023.
- [17] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [18] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

## A Complete Leaderboard Results

Table 7 presents the complete ranked leaderboard with all 22 models, including family classification and size tier.

## B Complete Figure Summary

Table 8 provides a complete summary of all 18 publication-quality figures generated by CORE-Bench analysis.

## C Statistical Methods

This appendix details the statistical methods used in our novel critical analysis.

### C.1 Hierarchical Clustering

We apply Ward’s minimum variance method to construct the dendrogram, using Euclidean distance on normalized performance scores. Optimal cluster number is determined by maximizing silhouette score:

$$S(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \quad (1)$$

where  $a(i)$  is mean intra-cluster distance and  $b(i)$  is mean nearest-cluster distance. We evaluate  $k \in \{2, \dots, 10\}$  and select  $k = 2$  with  $S = 0.776$ .

Table 7: Complete CORE-Bench Leaderboard (Ranked by Score)

Rank	Model	Family	Tier	Score (%)
1	Gemini 3 Flash Preview	Google	Medium	88
2	Qwen 3 Next 80B Thinking	Alibaba (Qwen)	Large	88
3	Claude Opus 4.1	Anthropic	Large	85
4	Claude Haiku 4.5	Anthropic	Medium	83
5	Claude Sonnet 4.5	Anthropic	Medium	83
6	Deepseek V3.1	DeepSeek	Large	83
7	Gemini 2.5 Flash	Google	Medium	83
8	Gemini 3 Pro Preview	Google	Large	83
9	Qwen 3 Coder 480B	Alibaba (Qwen)	Large	83
10	Qwen 3 Next 80B Instruct	Alibaba (Qwen)	Large	83
11	Claude Opus 4.5	Anthropic	Large	78
12	DeepSeek V3.2	DeepSeek	Large	77
13	Gemini 2.5 Pro	Google	Large	77
14	Gemini 2.0 Flash	Google	Medium	69
15	Gemini 2.0 Flash Lite	Google	Medium	69
16	Qwen 3 235B A22B Instruct	Alibaba (Qwen)	Large	69
17	Claude Sonnet 4	Anthropic	Medium	61
18	DeepSeek-R1	DeepSeek	Large	46
19	Gemma 3 12B	Google	Medium	30
20	Gemma 3 27B	Google	Medium	23
21	Gemma 3 4B	Google	Small	23
22	Gemma 3 1B	Google	Small	21

## C.2 Herfindahl-Hirschman Index

Market concentration is measured using the HHI:

$$\text{HHI} = \sum_{i=1}^n s_i^2 \times 10000 \quad (2)$$

where  $s_i$  is the market share of family  $i$  within a performance tier. HHI ranges from 0 (perfect competition) to 10000 (monopoly). Following DOJ/FTC guidelines: HHI < 1500 is “unconcentrated”, 1500-2500 is “moderately concentrated”, and > 2500 is “highly concentrated”.

## C.3 Bootstrap Confidence Intervals

Family ranking uncertainty is estimated via bootstrap resampling with  $n = 1000$  iterations. For each iteration, we resample models within each family with replacement, compute family means, and record rankings. The resulting distribution provides probability estimates for each family achieving each rank.

## C.4 Effect Size (Cohen’s d)

For thinking vs. standard model comparison:

$$d = \frac{\bar{x}_1 - \bar{x}_2}{s_{\text{pooled}}} \quad (3)$$

where  $s_{\text{pooled}} = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}$ . Cohen’s  $d < 0.2$  is small, 0.2-0.8 is medium,  $> 0.8$  is large.

Table 8: Complete Summary of CORE-Bench Analytical Figures

Fig.	Title	Key Insight	Section
<i>Foundational Analysis</i>			
1	Model Leaderboard	88% ceiling, 67pp gap	5
2	Family Performance	Alibaba leads (80.75%)	5
3	Score Trends	Performance trajectory visualization	5
4	Size Tier Analysis	Large models: 80.75% mean	5
5	Score Distribution	Bimodal pattern identified	5
6	Radar Dashboard	Multi-dimensional capability view	5
7	Family Box Plot	Within-family variance patterns	5
8	Tier Distribution	45.5% achieve Excellence	5
9	Summary Statistics	Comprehensive metrics table	5
<i>Novel Critical Analysis</i>			
10	Performance Gap	Pareto: top 8 = 80% cumulative	6
11	Hierarchical Clustering	2 clusters, $S = 0.776$	6
12	Family Dominance	$\text{HHI} > 2800$ in elite tier	6
13	Thinking vs Standard	+5.6%, $d = 0.238$ , $p = 0.706$	6
14	Size Efficiency	7.9% per size doubling	6
15	Bootstrap Confidence	Alibaba 66% P(#1)	6
16	Similarity Matrix	4 competitive sets identified	6
17	Insights Dashboard	8-panel synthesis	6
18	Generation Evolution	+5.0% per generation	6

## C.5 Gini Coefficient

Performance inequality is measured using the Gini coefficient:

$$G = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n^2 \bar{x}} \quad (4)$$

where  $G = 0$  represents perfect equality and  $G = 1$  represents maximum inequality.

## D Example Problems

### D.1 Logical Deduction (Advanced)

**Problem:** Consider the following statements:

1. All machine learning engineers understand statistics.
2. Some data scientists are machine learning engineers.
3. No one who understands statistics makes random guesses about model performance.

What can we definitively conclude about data scientists?

**Expected Answer:** Some data scientists do not make random guesses about model performance.

**Reasoning Chain:** (1)  $\forall x : \text{MLEngineer}(x) \rightarrow \text{UnderstandsStats}(x)$ ; (2)  $\exists x : \text{DataScientist}(x) \wedge \text{MLEngineer}(x)$ ; (3)  $\forall x : \text{UnderstandsStats}(x) \rightarrow \neg \text{RandomGuesses}(x)$ . By transitivity:  $\exists x : \text{DataScientist}(x) \wedge \neg \text{RandomGuesses}(x)$ .

### D.2 Mathematical Reasoning (Advanced)

**Problem:** A startup has 100 servers. Each month, 10% of working servers fail, but the company replaces 15 servers. Starting with 100 working servers, how many working servers will the company have after 3 months?

### Solution Process:

$$\begin{aligned}\text{Month 1: } & 100 - 0.10(100) + 15 = 100 - 10 + 15 = 105 \\ \text{Month 2: } & 105 - 0.10(105) + 15 = 105 - 10.5 + 15 = 109.5 \\ \text{Month 3: } & 109.5 - 0.10(109.5) + 15 = 109.5 - 10.95 + 15 = 113.55\end{aligned}$$

**Answer:** Approximately 113-114 working servers (depending on rounding policy).

### D.3 Causal Reasoning (Advanced)

**Problem:** A study finds that cities with more ice cream vendors have higher crime rates. A researcher concludes that ice cream vendors cause crime. Identify the flaw in this reasoning and propose a more plausible explanation.

**Expected Answer:** The researcher commits the correlation-causation fallacy. A confounding variable (temperature/summer season) likely explains both: warmer weather increases both ice cream sales and outdoor activities, which may correlate with crime rates. Additionally, both variables may simply correlate with city size (population).

**Key Concepts Tested:** Confounding variables, spurious correlation, third-variable problem.

### D.4 Analogical Reasoning (Advanced)

**Problem:** Complete the analogy: Neural network : Deep learning :: Statistical model : ?

**Expected Answer:** Machine learning (or statistical learning)

**Reasoning:** Neural networks are a specific tool/technique within the broader field of deep learning. Similarly, statistical models are specific tools within the broader field of machine learning/statistical learning.

## E Benchmark Access and Reproducibility

CORE-Bench is publicly available at:

<https://www.kaggle.com/benchmarks/taiwofeyijimi/core-bench>

### Evaluation Requirements:

- Kaggle account with API access
- Python 3.8+ with `kaggle-benchmarks` package
- Valid API credentials for target models
- For figure generation: `matplotlib`, `seaborn`, `scipy`, `scikit-learn`

### Submission Protocol:

1. Register model on Kaggle Models platform
2. Submit to CORE-Bench leaderboard
3. Results automatically validated and published
4. Download `leaderboard.json` for offline analysis

### Figure Reproduction:

1. Clone repository: `git clone https://github.com/tiamole/CORE-Bench.git`
2. Install dependencies: `pip install -r requirements.txt`
3. Run analysis notebook: `jupyter nbconvert --execute reasoning_benchmark.ipynb`
4. Figures saved to `publication_figures/{png,pdf,svg}/`

All 18 figures are generated at 300 DPI and available in PNG, PDF, and SVG formats for publication use.