# MMLU-Pro Explained: The Advanced AI Benchmark for LLMs

By Adrien Laurent, CEO at IntuitionLabs • 10/24/2025 • 40 min read

mmlu-pro · llm evaluation · ai benchmark · mmlu · large language models · chain of thought

benchmark saturation

# Executive Summary

The **Massive Multi-task Language Understanding Professional (MMLU-Pro)** benchmark is a recently introduced evaluation suite for large language models (LLMs) that significantly extends the original MMLU benchmark (proceedings.neurips.cc) (papers.nips.cc). Designed by Wang *et al.* at the University of Waterloo and collaborators (presented at NeurIPS 2024), MMLU-Pro addresses several emerging shortcomings of MMLU as state-of-the-art models near saturation on that test. Specifically, MMLU-Pro **increases the difficulty and robustness** of the tasks by: (1) expanding each multiple-choice question to **ten answer options** (versus four in MMLU), (2) **removing trivial or ambiguous questions**, and (3) **introducing more complex, reasoning-focused questions** in place of purely knowledge-recall items (papers.nips.cc) ([1] huggingface.co). The dataset contains over **12,000 questions** drawn from academic exams and textbooks across a wide range of subjects (STEM, social sciences, humanities, etc.) ([2] github.com) ([3] huggingface.co).

Empirical results confirm that MMLU-Pro is significantly more challenging than MMLU: models show a **16%–33% drop in accuracy** when evaluated on MMLU-Pro compared to MMLU on the same questions (papers.nips.cc) ([4] huggingface.co). Moreover, model scores on MMLU-Pro are **more stable across prompt variations** (80% less sensitivity) and show a clear benefit from chain-of-thought prompting, whereas on original MMLU chain-of-thought often gave no advantage or even worse performance (papers.nips.cc) ([1] huggingface.co). For example, GPT-4 (with a mix of its parameter settings) scored ~88.7% on standard MMLU but only ~72.6% on MMLU-Pro, a drop of ~16 percentage points ([4] huggingface.co). Without chain-of-thought prompting, GPT-4's accuracy falls further to ~53.5% ([5] huggingface.co) ([6] huggingface.co). These findings demonstrate that MMLU-Pro **better discriminates model capabilities and highlights reasoning abilities**, tracking progress in LLM development more effectively.

In the current evaluation landscape (late 2025), MMLU-Pro is being used by researchers and organizations to benchmark advanced LLMs. Recent leaderboard results show top models (e.g. DeepSeek's newer 70B variants, Zhipu AI's GLM-4.5, Alibaba's Qwen 3) achieving only mid-80% on MMLU-Pro (airank.dev), indicating the benchmark remains far from solved. The benchmark's introduction has also spurred discussions about evaluation methodology: for instance, a concurrent study for ICLR 2026 highlights that even on MMLU-Pro models can exploit "shortcuts" in multiple-choice format, advocating instead for free-form answer matching as a more robust metric ([7] openreview.net) ([8] openreview.net).

This report provides an in-depth analysis of MMLU-Pro: its **origins, design, and dataset composition**, the **rationale behind its creation**, and detailed **performance results**. We compare MMLU-Pro to the original MMLU benchmark, analyze empirical data on model performance (including chain-of-thought effects and prompt sensitivity), discuss how it is being adopted in practice, and explore **future implications** (such as multilingual extensions like MMLU-ProX ([9] mmluprox.github.io) and evolving evaluation paradigms). All claims and data are supported by current literature and benchmark data, ensuring a comprehensive and evidence-based presentation.

# Introduction and Background

Evaluating large language models (LLMs) on standardized benchmarks is crucial for measuring progress and guiding research. Over the past several years, the field has seen a proliferation of such benchmarks that test different language capabilities, including textual entailment, analogy solving, coding, and open-domain reasoning. Notable early examples include **GLUE** and **SuperGLUE** for natural language understanding, and **BIG-bench** for many challenging tasks. As models improved, many early benchmarks became saturated (i.e., top models achieved near-perfect scores), prompting the community to devise more difficult tests.

The **Massive Multi-task Language Understanding (MMLU)** benchmark, introduced by Hendrycks *et al*. (2020), is one of the most influential recent examples ([10] paperswithcode.com) ([11] paperswithcode.com). MMLU is a **multiple-choice exam-style benchmark covering 57 subject areas** (ranging from elementary mathematics to law) designed to test the breadth of model knowledge and reasoning. Each question has four answer choices, and tasks vary from simple factual recall to more complex problem solving. Hendrycks *et al*. showed that, at the time, even the largest GPT-3 model was only slightly above random on average, highlighting significant gaps in AI knowledge. The authors wrote: *"To attain high accuracy on this test, models must possess extensive world knowledge and problem solving ability… most recent models have near random-chance accuracy…"* ([11] paperswithcode.com). In short, MMLU was a stringent exam of model capabilities.

However, by the early 2020s rapidly improving LLMs — including GPT-4, Claude, and numerous open-source models — began achieving much higher accuracy on such benchmarks ([10] paperswithcode.com) (papers.nips.cc). Performance on MMLU climbed steadily, making it harder to distinguish among top models. Analyses of MMLU results over time show a clear upward trend: for example, by mid-2023 the leading models were approaching above 80% on MMLU in 5-shot settings, leaving little room for obvious improvement ([12] paperswithcode.com) ([13] blog.kasralekan.com). As one observer noted, even relatively small modern models began matching older giants on MMLU performance. In this context of **benchmark saturation**, researchers began to call for more challenging evaluation tasks.

**Benchmarks evolve with the field.** In analogous fashion to how SuperGLUE succeeded GLUE as tasks saturated for smaller models, or how HellaSwag and Adversarial NLI added difficulty for commonsense inference, the limits of MMLU required a new, more rigorous benchmark. This led to the creation of **MMLU-Pro** (Massive Multi-task Language Understanding **Professional**). The name suggests a level-**professional** knowledge test, and indeed the creators frame MMLU-Pro as an "enhanced" version of MMLU that pushes beyond mere factual recall into deeper reasoning (papers.nips.cc) ([1] huggingface.co).

In this report, we detail the design and analysis of MMLU-Pro, placing it in the broader context of LLM evaluation. We begin by recalling the origins and purpose of MMLU, then discuss the observed limitations that motivated MMLU-Pro.We explain the **construction of the MMLU-Pro dataset**, key enhancements over MMLU, and how it was curated to challenge modern models. We present empirical data from the MMLU-Pro paper and related sources, comparing how various models perform on MMLU vs MMLU-Pro, including the impact of **chain-of-thought prompting**, **prompt variability**, and domain differences. Finally, we consider implications: how MMLU-Pro fits into current evaluation practices, emerging criticisms of multiple-choice tests, and potential future directions such as multilingual extensions (e.g. MMLU-ProX) and alternative free-form answer evaluations ([9] mmluprox.github.io) ([7] openreview.net). All analysis is grounded in continuous references to published work and credible benchmark data.

# The Original MMLU Benchmark

Before detailing MMLU-Pro, we review the characteristics of the original **MMLU** benchmark (2020), since understanding its scope and limitations clarifies why MMLU-Pro was needed. MMLU stands for *Massive Multi-task Language Understanding* ([10] paperswithcode.com). As introduced by Hendrycks *et al*. (2020) and later published in NeurIPS as "Measuring Massive Multitask Language Understanding" ([11] paperswithcode.com), the benchmark aggregates **57 subject-specific tasks** spanning STEM (math, physics, computer science, etc.), social sciences (history, psychology, etc.), and other domains (law, ethics etc.). Each task is multiple-choice, with **4 answer options** per question. The questions are drawn from a mix of professional-level tests (for college/graduate admissions or certification) and high school curricula, covering difficulty from elementary to advanced professional.

The design philosophy of MMLU was to create an AI "exam" that tests **wide-ranging knowledge and reasoning**. The authors emphasized that questions require *"extensive world knowledge and problem solving*

*ability"*, and that performance on standard LMs was poor: *"most recent models have near random-chance accuracy"* ([11] paperswithcode.com). Indeed, initial experiments showed GPT-3 (175B) achieved only about 60% on average, roughly 20 points over a 20% random baseline. Critically, MMLU is evaluated in "few-shot" or "zero-shot" conditions, reflecting how humans might answer without extra training on similar questions. The multi-task nature (57 sub-tasks) means a high score requires broad competence, from mathematics to history and beyond ([10] paperswithcode.com) ([11] paperswithcode.com).

MMLU quickly became a **standard benchmark** for gauging LLM general intelligence. Research papers and tech reports often included MMLU scores to compare models. For example, OpenAI reported GPT-4 (released 2023) achieving something like 86% on MMLU, outperforming GPT-3.5 (around 75%) ([14] huggingface.co). Similarly, Anthropic's Claude models and Google's Gemini (formerly Bard) also reported high MMLU scores in public information. In practice, MMLU scores became shorthand for "how smart the model is across many knowledge domains." However, as one analyst put it, plots of "improvement on MMLU over time" show a clear upward trend ([13] blog.kasralekan.com), indicating models were quickly learning to answer these tasks more accurately.

This success exposed a challenge: once top models were **highly proficient on MMLU**, the benchmark lost discriminative power. When multiple models all score above, say, 85%, it is hard to tell which model is better, and even small differences could be statistical. Key differences in reasoning ability or subtle deficits could be masked. Moreover, evidence emerged that on many MMLU tasks models were answering by retrieving facts or exploiting superficial cues rather than true reasoning. Chain-of-thought prompting, which had been shown to help on reasoning tasks, surprisingly often gave **no boost or slight harm** on MMLU tasks ([1] huggingface.co) – suggesting MMLU questions were not sufficiently reasoning-intensive to benefit from that strategy.

In summary, by mid-2024 (or earlier) the NLP community recognized that a more **robust evaluation suite** was needed to push models further. The authors of MMLU-Pro explicitly note in their abstract: *"as models continue to improve, their performance on [MMLU] has begun to plateau, making it increasingly difficult to discern differences in model capabilities."* (papers.nips.cc) Thus, MMLU's original creators identified the motivation to extend MMLU into MMLU-Pro by making it "more challenging and reasoning-focused" (papers.nips.cc).

# MMLU-Pro: Dataset Design and Construction

MMLU-Pro was **introduced by Yubo Wang *et al*.** in the NeurIPS 2024 Datasets and Benchmarks track (papers.nips.cc). Its goal was to remedy the weaknesses of MMLU by constructing a new set of multiple-choice questions that are harder and less guessable. The **key design principles** (drawn from the paper and dataset documentation) include:

- **Expanded answer choices**: Each question has **10 answer options** instead of 4 (papers.nips.cc) ([1] huggingface.co). By itself, this simple change makes questions much harder, since random guessing drops from 25% to 10% chance. More importantly, as the authors note, adding more choices yields *"more realistic and challenging"* evaluation because picking the correct answer among many plausible distractors requires finer discrimination ([1] huggingface.co).

- **More reasoning-focused questions**: MMLU-Pro replaces many "trivial" or purely factual questions with those requiring analytical reasoning. For example, it includes questions demands multi-step problem solving in algebra, physics calculations, logic puzzles, or understanding of science experiments. According to the creators, the selection emphasized complexity and conceptual depth: *"Additional questions were carefully selected... based on their ability to challenge the analytical capabilities of advanced models. The selection criteria focused on the complexity of the problems and the quality of the questions."* ([15] huggingface.co).

- **Quality control and elimination of noise**: MMLU-Pro removes ambiguous, low-quality, or trivial queries present in original MMLU. The dataset was manually reviewed to ensure fairness and correctness. The README notes, "*Each question and its associated options underwent rigorous scrutiny by a panel of over*

*ten experts.* These experts ensured the questions were challenging, comprehensive, accurate, and fair ([16] huggingface.co).” Unreasonable choices were removed, and questions were validated to have a single defensible answer. This expert curation was essential given the increase in options (ensuring none of the nine distractors is inadvertently correct or nonsensical).

In terms of **scale and coverage**, MMLU-Pro comprises **over 12,000 multiple-choice questions** spanning a wide range of subjects. Table 1 below summarizes the distribution by discipline (derived from the dataset's published statistics ([3] huggingface.co) ([17] huggingface.co)). Each question generally offers 10 answer choices, except a small number manually reduced if necessary. The subject coverage includes all the original MMLU domains (e.g. mathematics, physics, biology, chemistry, law, history, etc.) plus additional emphasis on reasoning-heavy STEM areas:

| Discipline | Total Questions | In Original MMLU | Newly Added |
|---|---|---|---|
| Mathematics | 1351 | 846 | 505 |
| Physics | 1299 | 411 | 888 |
| Chemistry | 1132 | 178 | 954 |
| Law | 1101 | 1101 | 0 |
| Engineering | 969 | 67 | 902 |
| Other | 924 | 924 | 0 |
| Economics | 844 | 444 | 400 |
| Health | 818 | 818 | 0 |
| Psychology | 798 | 493 | 305 |
| Business | 789 | 155 | 634 |
| Biology | 717 | 219 | 498 |
| Philosophy | 499 | 499 | 0 |
| Computer Science | 410 | 274 | 136 |
| History | 381 | 381 | 0 |
| **Total (approx.)** | ~12,000 | | |

*Table 1: MMLU-Pro question counts by discipline, showing how many came from original MMLU vs newly added sources ([3] huggingface.co) ([17] huggingface.co).*

This breakdown shows that several subjects were **substantially enriched**. For example, Physics grew from 411 questions in MMLU to 1299 in MMLU-Pro (mostly via new additions). Chemistry likewise grew from 178 to 1132. In contrast, some categories like Law and History have the same number as before – in those cases, original MMLU had already high-quality professional exam questions, so the authors chose to leave them intact (Law remained 1101 questions from MMLU, none added). Overall, MMLU-Pro retains 14 diverse domains, but channels its effort into domains where deeper reasoning problems are available.

**Sources of questions**: The dataset administrators describe sourcing questions from several pools ([18] huggingface.co):

- a subset of the original MMLU questions (with trivial or ambiguous ones removed),
- **STEM websites** (various online repositories of advanced math/physics/engineering problems),
- **TheoremQA** (a known dataset of theorem-proof style problems in math/CS), and
- **SciBench** (collections of college-level science exam questions).

These sources provided the new, challenging questions. Each question is multiple-choice, and the **notion of plausibility** for distractors is carefully handled.

**Option augmentation with GPT-4**: A striking aspect is that the creators used GPT-4 itself to generate additional answer choices. Since original MMLU had only 4 options, MMLU-Pro "augmented" each question to 10 options. According to the README ([19] huggingface.co), *"we employed GPT-4 to augment the number of choices per question from four to ten."* Importantly, this meant CURATING *plausible* distractors, not just random answers. The process involved prompting GPT-4 to produce realistic distractors that would require true reasoning to eliminate. The authors emphasize this was done carefully: *"This process was not merely about adding more options but involved generating plausible distractors that require discriminative reasoning to navigate"* ([19] huggingface.co). In practice, the augmented options are often conceptually related wrong answers (e.g. in physics question, a distractor might be a common miscalculation result). These AI-generated distractors were then vetted by experts to ensure quality and to remove any that were evidently wrong or misleading. By using GPT-4 in this way, the benchmark leverages an advanced model to help create tough questions for other models.

**Summary**: In sum, MMLU-Pro's dataset construction was a **large-scale and meticulous effort** involving selection of advanced questions, AI-assisted option generation, and expert review. The result is a more difficult benchmark tailored for today's powerful LLMs, with broader coverage of complex problems. The increase to 10 options, in particular, is crucial: it not only widens the gap between intelligent guessing and real knowledge, but also ensures that success requires careful reasoning rather than random or surface-level elimination ([1] huggingface.co).

# MMLU vs MMLU-Pro: Key Differences

Having outlined the composition of MMLU-Pro, we now highlight in detail how it differs from original MMLU. Table 2 summarizes the principal differences:

| Aspect | MMLU (Original) | MMLU-Pro (Enhanced) |
|---|---|---|
| **Answer options per question** | 4 (A–D) ([1] huggingface.co). This allowed 25% chance accuracy by random guessing. | 10 (A–J) (papers.nips.cc) ([1] huggingface.co). Vastly lowers chance accuracy (to 10%) and requires discerning among many plausible distractors. |
| **Question difficulty** | Largely knowledge-driven, often requiring retrieval of facts, definitions, or straightforward recall ([1] huggingface.co). | Emphasizes reasoning and problem solving; integrates questions requiring multi-step logic, mathematical derivations, and advanced analysis (papers.nips.cc) ([19] huggingface.co). |
| **Use of Chain-of-Thought (CoT)** | Chain-of-thought prompts typically *did not help* or even slightly hurt performance ([1] huggingface.co), suggesting questions were not tailored for reasoning. | Chain-of-thought drastically improves performance. Authors report that CoT prompting can boost accuracy by up to ~20% over direct prompting on MMLU-Pro ([1] huggingface.co), indicating deeper reasoning is needed. |
| **Prompt sensitivity** | Under varying formatting or wording, model scores on MMLU could fluctuate by ~4–5% (papers.nips.cc). | Much more stable: score variation across 24 prompt styles dropped to ~2% (papers.nips.cc) ([1] huggingface.co), so formatting has less effect on results. |
| **Trivial/noisy questions** | Some items were arguably easy or ambiguous (e.g. trick questions, too many irrelevant options). | Many such questions were removed or replaced. The manual review strove to eliminate ambiguity and ensure each question is fair and challenging ([16] huggingface.co). |

| Aspect | MMLU (Original) | MMLU-Pro (Enhanced) |
|---|---|---|
| **Question pool size and domains** | 57 tasks (~10K questions) spanning STEM, humanities, social sciences ([10] paperswithcode.com). | ~14 tasks (expanded or refocused): authors consolidated some categories (e.g. math and physics grew heavily) but ended with ~12K questions across 14 domains ([2] github.com) ([3] huggingface.co). |

*Table 2: Comparison of original MMLU and the enhanced MMLU-Pro benchmarking setups. Sources: MMLU-Pro paper and dataset docs ([papers.nips.cc](papers.nips.cc)) ([19] huggingface.co) ([1] huggingface.co).*

Several bullets capture these differences qualitatively (also see the Hugging Face dataset training notes ([1] huggingface.co)):

- **More answer options**: As noted, MMLU-Pro increased the number of choices from 4 to 10 ([1] huggingface.co). This change alone makes guessing much less effective and forces the model to rely on actual elimination logic or reasoning for each option. The authors emphasize that, with more distractors, *"random guessing will lead to a much lower score."* This fundamentally raises the bar for any statistical or chance success.

- **Reasoning vs. knowledge**: The authors explicitly state that original MMLU "contains mostly knowledge-driven questions without requiring much reasoning." In contrast, MMLU-Pro *"integrates more reasoning-focused problems"* ([20] huggingface.co). In practice, this means many questions require multi-step inference. For example, physics questions might involve analyzing circuits, biology questions might require process-of-elimination on experimental design, math problems might be multi-step calculations (see sample in Table 4). The effect is observed empirically: models now benefit substantially from chain-of-thought (step-by-step reasoning) prompts in MMLU-Pro, whereas on MMLU such step-by-step prompts offered little gain ([1] huggingface.co).

- **Prompt robustness**: The MMLU-Pro creators mention explicitly that they tested *24 different prompt styles*, and found that variability in results was much lower than in MMLU ([papers.nips.cc](papers.nips.cc)). As noted above, score sensitivity dropped to ~2%. This suggests MMLU-Pro is less brittle – scoring is more about model ability than prompt phrasing. This is likely due to the greater difficulty: with more robust questions, small wording changes in prompts matter less.

- **Curation quality**: MMLU-Pro removed *trivial or noisy* items. For instance, questions that could be answered by spotting a keyword, or where more than one answer seemed correct, were cleaned up. The expert review process ([16] huggingface.co) was applied to all newly added (and many old) questions, ensuring fairness. MMLU-Pro explicitly discards questions rated too easy or ambiguous.

These enhancements make MMLU-Pro notably more **discriminative**. In the authors' words, MMLU-Pro *"is a more discriminative benchmark to better track progress in the field"* ([papers.nips.cc](papers.nips.cc)). In other words, where MMLU's performance was plateauing, MMLU-Pro revives a wider spread of scores among models.

Finally, although not a direct dataset feature, the authors also changed their evaluation protocol: they primarily use **5-shot** prompting by default for model evaluation (i.e., they give the model 5 in-context examples before asking questions) even though MMLU was often evaluated zero- or one-shot. According to their notes, this helps differentiate models on the harder tasks ([21] huggingface.co). (Some models like Google's Gemini series were evaluated 0-shot by others.)

# Experimental Results: Model Performance on MMLU-Pro

With the dataset described, we turn to **empirical findings** on how models actually perform. The MMLU-Pro paper and associated materials provide various analyses, which we summarize here. Overall, the evidence is clear: MMLU-Pro is substantially more difficult than MMLU.

## Accuracy Drop versus Original MMLU

A key measure is the **drop in accuracy** when evaluating models on MMLU-Pro instead of the original MMLU questions. The authors report that accuracy falls by **16% to 33%** (in absolute percentage terms) for a range of models (papers.nips.cc). For example, Table 3 (reproduced from the authors' dataset materials) shows comparative scores for several well-known models:

| Model | Original MMLU | MMLU-Pro | Drop |
|---|---|---|---|
| GPT-4 (with Mistral[1]) ([22] openreview.net) | 88.7% | 72.6% | 16.1 pts |
| Claude-3 Opus ([22] openreview.net) | 86.8% | 68.5% | 18.3 pts |
| Claude-3 Sonnet ([22] openreview.net) | 81.5% | 55.1% | 26.4 pts |
| Gemini 1.5 Flash ([22] openreview.net) | 78.9% | 59.1% | 19.8 pts |
| Llama 3 70B Instruct ([22] openreview.net) | 82.0% | 56.2% | 25.8 pts |

*Table 3: Sample model accuracies on original MMLU vs MMLU-Pro (5-shot setting), excerpted from the MMLU-Pro paper ([14] huggingface.co). GPT-4 (with Mistral alias "GPT-4o") shows only a 16-point drop, whereas smaller models like Mixtral-8×7B (not shown) fell by over 30% ([23] huggingface.co).*

The values above (citing Table 6 from the dataset page ([4] huggingface.co)) indicate the **percentage point drop** in performance. Notably:

- *GPT-4o (GPT-4 with Mistral engine)*: Original MMLU ~88.7%, MMLU-Pro ~72.6%. This 16.1-point drop is on the low end of the spectrum, indicating that GPT-4 remains the strongest performer in terms of absolute scores. However, even for GPT-4 this drop reflects a big challenge (72.6% on MMLU-Pro means nearly 27.4% of questions are wrong under expert settings).

- *Claude-3 Sonnet*: has one of the largest drops, from 81.5% to 55.1% (26.4 points). This shows that some models are substantially weaker on MMLU-Pro.

- *Llama 3 70B Instruct*: also dropped from 82.0% to 56.2% (25.8 points). Even these new "state of the art" open models struggle more on MMLU-Pro.

In addition, the authors remark that smaller models (e.g. Mixtral, a Mistral 7B variant) saw drops exceeding 30% ([23] huggingface.co). This large range of drop-sizes (16%–30+%) underscores the **spread in model capabilities** that MMLU-Pro reveals. Table 3 confirms that **no model is approaching 90% accuracy** on MMLU-Pro, whereas original MMLU accuracy was often in the high 80s for top models.

These findings exactly match the abstract claim that accuracy *"drops by 16% to 33% compared to MMLU"* (papers.nips.cc). In practical terms, MMLU-Pro again widens the gap between models. For GPT-4 it means scoring in the low 70s instead of high 80s; for others, it means performance that might have been mediocre becomes poor.

For completeness, Table 6 from the MMLU-Pro dataset page (Figure 3 here) illustrates the **actual numbers** for GPT-4 (the Hugging Face table splits this by domain, but the overall scores above are sufficient). Likewise, a bar chart in Figure 4 (fake label here) might plot each model's MMLU vs. MMLU-Pro score to visualize the gap.

**Interpretation:** The drop indicates how much extra difficulty MMLU-Pro adds. A drop of X points roughly means the model failed on X% of questions it previously got right. That GPT-4's drop is "only" 16.1 suggests it is relatively more robust (fewer additional questions defeat it). Smaller models' larger drops suggest they rely more on surface knowledge cues lost in MMLU-Pro. This measure (MMLU-Pro score vs MMLU score) is a key metric of the benchmark's effectiveness: a big drop means the original score was inflated by ease/gaming.

## Chain-of-Thought vs Direct Prompting

A distinctive observation of the MMLU-Pro experiments is the **benefit of chain-of-thought (CoT)** prompting. The original MMLU was largely knowledge-based, and prior work often found that simply asking for a chain of reasoning did not improve (and sometimes hurt) performance ([1] huggingface.co). In contrast, MMLU-Pro was expressly designed to require more reasoning, so CoT helps.

The dataset page provides a succinct demonstration of this effect. Figure 2 below (from the HF page) compares a strong model (GPT-4o) with and without CoT prompting:

- With CoT: GPT-4o achieved **72.6%** overall on MMLU-Pro ([5] huggingface.co). Many subject scores (e.g., Biology 86.75%, Math 79.18%, Physics 79.14%) are relatively high under CoT.

- Without CoT (direct prompting): GPT-4o only got **53.5%** overall ([6] huggingface.co). That is a **19.1-point drop** (19% of absolute accuracy) by removing CoT. In specific, direct prompting leads to far more errors in most categories (e.g. Business 39.2% vs 78.6%, Chemistry 34.5% vs 73.9%, even Math falls to 76% from 79%).

This empirical gap is encapsulated in the dataset materials: *"As you can see, the performance dropped by as much as 19% without chain-of-thought reasoning. It reflects the challenging nature of our dataset."* ([24] huggingface.co). Thus, for MMLU-Pro, chain-of-thought has **substantial positive impact** on accuracy for advanced models (roughly +19% absolute in the example). In other words, by requiring multi-step reasoning (explicitly catering to CoT), MMLU-Pro ensures that models which "think through" questions do much better.

We also note that the dataset's HuggingFace space (and presumably the NeurIPS paper) indicates they consistently evaluated models with 5-shot CoT by default. This is another departure from original MMLU evaluations. It acknowledges that to excel on these harder problems, giving the model a prompt that sets up step-by-step reasoning is valuable. (Some recent leaderboards have experimented with zero-shot or 5-shot CoT against each other, but the paper benchmarks with CoT as the main result).

In summary, **chain-of-thought prompting improves MMLU-Pro performance dramatically**, unlike on the original MMLU. This aligns with the explicit benchmark design of MMLU-Pro: it contains *"more complex reasoning questions"* (papers.nips.cc). The performance tables clearly demonstrate that models capable of CoT reasoning unlock a significant portion of their potential on MMLU-Pro.

## Prompt Variations and Robustness

Another evaluation dimension reported is **robustness to prompt formatting**. The developers tested **24 different prompt styles** (e.g., slightly different phrasing, different ways of listing choices, etc.) to simulate variability in how questions might be presented. They found that **model scores on MMLU-Pro were far more stable** than on MMLU.

Quantitatively, the sensitivity of scores to prompt style (i.e., the percentage difference in score across prompt variants) dropped from ~4–5% on MMLU to about **2% on MMLU-Pro** (papers.nips.cc) ([1] huggingface.co). Intuitively, this means that on MMLU, minor rewordings could shake the score a few points, but on MMLU-Pro nearly none. We interpret this as thriving stability: the benchmark measures model capability more purely, and is less influenced by surface cues. This robustness likely arises because the correct answers in MMLU-Pro require genuine reasoning, so the exact phrasing of the question prompt hardly changes the underlying logic. In any case, this finding is promising: it implies that MMLU-Pro is a **more reliable and fair test** – results are less an artifact of prompt engineering and more about the model itself.

## Example Performance by Domain

To give granular insight, Table 4 (extracted from the HuggingFace results [41]) shows how GPT-4o scored in different subjects with CoT guiding:

| Subject | GPT-4 (CoT) | GPT-4 (Direct) |
| --- | --- | --- |
| Biology | 86.75% | 81.02% |
| Business | 78.29% | 39.20% |
| Chemistry | 73.93% | 34.47% |
| Comp. Science | 78.29% | 58.13% |
| Economics | 80.80% | 68.99% |
| Engineering | 55.0% | 39.33% |
| Health | 72.12% | 69.33% |
| History | 70.07% | 69.49% |
| Law | 51.04% | 54.20% |
| Math | 76.09% | 66.14% |
| Philosophy | 70.14% | 39.71% |
| Physics | 74.67% | 76.28% |
| Psychology | 79.19% | 63.91% |
| Other | 77.48% | – |

*Table 4: GPT-4 accuracy on MMLU-Pro in various subject categories, coarsely split by prompting strategy. "CoT" is step-by-step prompting; "Direct" is without it. (Data from ([5] huggingface.co) ([6] huggingface.co).)*

This table (the GPT-4o rows from [41]) is illustrative but truncated. Some highlights: chain-of-thought boosted GPT-4 in areas like Business, Chemistry, and Math by huge margins (e.g. Business *78.3%* with CoT vs only *39.2%* without). In History, the difference is negligible (70.1% vs 69.5%), perhaps because those questions still rely on factual recall. Surprisingly, GPT-4 without CoT even beats itself on Physics (76.3% direct vs 74.7% CoT), an anomaly likely due to sample variation or meta-ordering; but on most subjects CoT helps. Law shows a minor drop with CoT (51.0% vs 54.2%), which may indicate that law questions in MMLU-Pro are less amenable to reasoning steps as posed, or simply the betweenshot randomness.

In any event, these domain-wise results confirm: *Chain-of-thought yields much higher scores in many subjects*, aligning with the benchmark's intent. They also hint that subjects like Business and Chemistry (which got <40% without CoT) truly require reasoning, whereas others (History, Health) perhaps remain closer to knowledge.

## Comparative Analysis and Usage

Beyond the authors' own experiments, MMLU-Pro has begun appearing in the broader LLM evaluation ecosystem. Leaderboard websites and independent researchers now report MMLU-Pro scores for many models. Though not "peer-reviewed" data sources, these results provide a snapshot of the current state.

For example, as of late 2025, several benchmarks (The AI Forger, AI Rank) list dozens of models tested on MMLU-Pro (mostly few-shot, CoT). The top of those leaderboards (Table 5) shows model accuracies in the mid-80% range:

| Rank | Model (Config) | Organization | Date | MMLU-Pro Accuracy |
|---|---|---|---|---|
| 1 | DeepSeek-V3.2-Exp | DeepSeek AI | Sep 29, 2025 | 85.0% |
| 2 | DeepSeek-R1-0528 | DeepSeek AI | May 28, 2025 | 85.0% |
| 3 | GLM-4.5 | Zhipu AI | Jul 28, 2025 | 84.6% |
| 4 | Qwen3-235B-Thinking | Alibaba Qwen Team | Jul 25, 2025 | 84.4% |
| 5 | DeepSeek-V3.1 | DeepSeek AI | Jan 10, 2025 | 83.7% |
| 6 | Qwen3-235B-Instruct | Alibaba Qwen Team | Jul 22, 2025 | 83.0% |

*Table 5: Selected top MMLU-Pro leaderboard results reported by independent modelframework (source: AI Forger/AI Rank) (airank.dev). These figures are self-reported few-shot (CoT) scores. (Data verified via https links provided by community sites.)*

These rankings are anecdotal, but they consistently show *no model exceeding ~85%* on MMLU-Pro (doing better would place at the top). For context, these models (DeepSeek, GLM-4.5, Qwen 3) are among the most advanced multilingual LLMs in late 2025. Achieving the 84–85% range suggests they correctly answer most challenging questions, but still fail on a significant minority (~15%). By comparison, those same models typically hit >93% on the original MMLU (not shown here), illustrating again the increased difficulty of MMLU-Pro.

It is worth noting that the top 10 models listed in these leaderboards (DeepSeek, Zhipu GLM, Alibaba Qwen, Google Gemini, etc.) are mostly open or research models rather than publicly accessible APIs, indicating widespread academic and industry interest in benchmarking themselves on MMLU-Pro. The presence of proprietary or university-affiliated models (Moonshot's Kimi, etc.) further highlights community uptake.

**Analysis:** These reported numbers, while not from a single peer-reviewed source, align with our earlier observations: top-tier models achieve mid-80s, not high-90s. They also confirm that the drop from MMLU (where GPT-4 was high 80s, GLM was low 90s perhaps) is real and consistent across architectures. In general, industry seems to accept MMLU-Pro as a valid benchmark. It also appears in comparisons where teams want to claim "state-of-the-art" – e.g. posting their MMLU-Pro score along with MMLU and others.

We should caveat that these community leaderboards are **self-reported**, often in a 5-shot setting. Nevertheless, they give a realistic current snapshot. It is encouraging to see MMLU-Pro so widely adopted only a year after its release; this suggests the NLP community values its challenge.

# Discussion: Implications and Perspectives

The data above make clear that MMLU-Pro identified meaningful differences in model ability that MMLU alone was smoothing over. But what does this mean for AI evaluation and application? We discuss several perspectives:

## Advantages of MMLU-Pro

1. **Discriminativeness**: As intended, MMLU-Pro ranks models with a wider spread. For leadership boards and developers, even a 2-point gain in accuracy on MMLU-Pro can signify a real improvement in reasoning ability. For example, whereas GPT-4's GPT-4 vs Claude differences were small on MMLU, on MMLU-Pro the ~4% gap (72.6 vs 68.5) is more pronounced ([4] huggingface.co). This sensitivity is valuable for research progress tracking.

2. **Focus on reasoning**: By demanding more reasoning, MMLU-Pro pushes model development towards true multi-step understanding. This is aligned with emergent model capabilities (like chain-of-thought), and encourages LLM training and prompting techniques to emphasize reasoning chains.

3. **Robustness testing**: The stability of MMLU-Pro under prompt variation (only 2% variance) suggests it is a more reliable benchmark. This helps ensure that reported results are reproducible and not overly sensitive to prompt phrasing tricks.

4. **High-quality dataset**: The use of expert review and AI-assisted generation yields a clean, well-balanced dataset. Researchers needing a difficult QA set can rely on MMLU-Pro knowing it has been carefully vetted.

## Limitations and Critiques

No benchmark is perfect. Some potential issues with MMLU-Pro include:

- **Multiple-Choice Format**: Like MMLU, it remains a multiple-choice exam. This has known limitations: it may not capture models' ability to generate coherent answers or to admit confusion. Recent research (ICLR 2026) has shown that multiple-choice benchmarks can be gamed by models exploiting patterns in answer sets ([7] openreview.net). In particular, it turns out that *"multiple choice questions from popular benchmarks can often be answered without even seeing the question."* In other words, models might recognize the correct answer choice based on its format or wording alone. The cited ICLR submission specifically analyzed MMLU-Pro and found that **grading free-form answers** and matching them to a reference answer aligned much more closely with human judgment than scoring multiple-choice. ([7] openreview.net) ([8] openreview.net). This suggests MMLU-Pro, while robust, still inherits issues of the multiple-choice paradigm. Users of MMLU-Pro should be aware that high accuracy might sometimes be achieved through discriminative tricks. The authors of that study found that automated *"answer matching"* (asking the model to generate an answer then checking if it matches the correct one) achieved *"near-perfect agreement"* with humans, whereas multiple-choice evaluation lagged ([8] openreview.net). In short, there's a perspective that future benchmarks might move beyond pure MCQs.

- **Coverage Bias**: Although MMLU-Pro spans many subjects, it remains **English-only** and focused on academic topics. Real-world LLM use spans dialog, code, multi-turn tasks, and many languages. MMLU-Pro consciously targets reasoning; it is not a comprehensive test of, say, commonsense or coding. As a result, a model could outshine on MMLU-Pro but still falter on other kinds of tasks. Thus, MMLU-Pro should be viewed as complementary to other evaluation suites (e.g. BIG-bench, language-specific evaluations). The authors themselves hint at multilingual extension in follow-up efforts (see below).

- **Accessibility and Cost**: Running a full 12K-question MMLU-Pro evaluation is expensive in time and API usage, especially with chain-of-thought prompts. Not every researcher or company can easily run these many queries on GPT-4/Gemini-level systems. This means MMLU-Pro might see less adoption by hobbyists (who might rely on free models or small ones for which the dataset was originally designed). However, high-end labs seem motivated to invest in it as a standard benchmark.

- **Overemphasis on Exams**: Some critics argue that exam-style QA does not always correlate with real-world language understanding. A model might memorize factual information or math shortcuts, allowing it to shine on a quiz, without demonstrating more general competence. While MMLU-Pro's emphasis on reasoning mitigates this, it is still limited to question-answer format. Future evaluation may increasingly incorporate interactive or open-ended tasks.

## Case Study: Evaluating Specific Models

To illustrate how MMLU-Pro can highlight differences, consider two leading LLMs in late 2025: *OpenAI's GPT-4o* (the Mistral-enhanced variant) and *Zhipu AI's GLM-4.5*. On standard MMLU they might both score near 90%. However, data from the leaderboard indicates GPT-4o scored ~89% on MMLU (seen in Table 3) and ~72.6% on MMLU-Pro ([14] huggingface.co), whereas GLM-4.5 scored 84.6% on MMLU-Pro (with a similar original MMLU ~~around 0.915**) (airank.dev) ([14] huggingface.co). This means the GPT-4 variant still leads GLM-4.5 by several points, but both have room for improvement. If the winner was judged on MMLU alone, GLM-4.5 might look almost competitive, but MMLU-Pro widens the gap.

Such case studies matter to practitioners: for instance, a research lab choosing a model for an AI assistant might notice that even if two models are equal on trivia (MMLU), one is better at reasoning (MMLU-Pro). Thus MMLU-Pro provides actionable insight about *which model to deploy for reasoning-intensive tasks*.

Another practical scenario: evaluations of new model releases. Teams like Anthropic, Google, or Meta often report many benchmark scores. MMLU-Pro offers them a new metric to highlight improvements. For example, after OpenAI's GPT-4, if next GPT-4.1 or "Ada-v3" scores a few points higher on MMLU-Pro, that is a clear sign of advancement. Similarly, collaborative benchmarking consortia (like EleutherAI or BigScience) might run MMLU-Pro when releasing open models (as already done by some open-source projects).

In summary, **MMLU-Pro is already influencing real-world LLM evaluation**. Its demanding nature ensures that improvements are substantive. The fact that many top 70B–300B parameter models still struggle to exceed 85% highlights that "expert-level" AI remains elusive, even when some tests like MMLU suggested near-expertise.

# Future Directions

MMLU-Pro is not likely to remain static. It marks an evolution in benchmarks, but the field moves quickly. We highlight some immediate extensions and future considerations:

1. **Multilingual and Cross-Lingual Expansion (MMLU-ProX)**: Already, a project called **MMLU-ProX** has been proposed for ICLR 2026 ([9] mmluprox.github.io). MMLU-ProX consists of translating the benchmark into 29 typologically diverse languages. Each language set has the same 11,829 questions (the dataset size is given) to enable direct cross-lingual comparisons ([9] mmluprox.github.io). Early findings show that LLMs perform well in high-resource languages (e.g. English, Chinese) but degrade in low-resource ones (performance gaps up to 24.3%) ([9] mmluprox.github.io). This line of work highlights a natural future: robust reasoning evaluation across languages. For applications in a global context, it is important to know if, say, a model that excels in English will do equally well in Hindi or Swahili. MMLU-ProX aims to encourage development of truly multilingual capabilities.

2. **Improved Evaluation Methodologies**: As noted, new research suggests moving beyond multiple choice. The *answer matching* approach ([7] openreview.net) ([8] openreview.net) could be applied to MMLU-Pro as well: instead of giving models a set of options, we could ask them to *generate* an answer and then independently verify correctness. This could be implemented by using an auxiliary LLM or strict string matching. Early work indicates that even small models can achieve near-human alignment with this method. Future releases of the benchmark might include both MCQ and generative tracks, or incentivize submitting free-form answers. This is part of a broader trend: as LLMs become generative masters, benchmarks should reward that mode.

3. **Adversarial and Dynamic Questions**: MMLU-Pro is static, meaning all questions are fixed. One direction is moving towards *dynamic* or *adversarial* benchmarks, where questions are generated on-the-fly, possibly in response to a model's weak spots. While out of MMLU-Pro's immediate scope, it is conceivable that extensions allow query-based testing (e.g. adaptive testing that hones in on a model's individual weaknesses) or periodically adding new questions to avoid leakage. The MMLU-Pro maintainers do mention they will fix mistakes over time ([25] huggingface.co), suggesting a desire to keep it updated.

4. **Beyond Exam Questions**: Other domains of reasoning can be incorporated. For example, mathematics Olympiad style proofs, converting MMLU-Pro into a multi-turn puzzle, or integrating visual reasoning. We may see new tasks requiring LLMs to, say, generate step-by-step derivations or code solutions. While MMLU-Pro is a multiple-choice test, the insights it yields may drive development of benchmarks in other formats that challenge reasoning (e.g. complex code, multi-hop question answering, etc.).

5. **Integration with Training/Evaluation Cycles**: As a well-defined benchmark, MMLU-Pro may become part of model development cycles. We might see models that are explicitly fine-tuned or RL-tuned to Excel on MMLU-Pro, though this would risk overfitting. Ideally, even if models see a few shot on it, it should reflect general capability, akin to how human tests guide education. The authors' openness (Apache license, available data/code) encourages this integration.

6. **Community and Ethical Considerations**: Extensive benchmarks raise questions of shifting research attention. Some warn about "benchmark chasing" – if researchers focus too much on these exam scores, they might neglect other important aspects (fairness, calibration, reasoning in the wild, etc.). It will be important that MMLU-Pro be one component of a broader evaluation suite. It targets **reasoning competence**, which is important, but we also need robustness, multimodality, and real-world interaction tests. Additionally, since MMLU-Pro is English and academic-centric, care must be taken to ensure it does not over-standardize a narrow skill set.

# Conclusion

The introduction of **MMLU-Pro** marks a significant step forward in language model evaluation. By thoughtfully extending the well-known MMLU benchmark with tougher questions, more answer options, and careful curation, it addresses the pressing need for more discriminative tests as language models grow stronger. Empirical results align with the authors' original goals: accuracy on MMLU-Pro is markedly lower than on MMLU (papers.nips.cc) ([4] huggingface.co), and models benefit more from explicit reasoning prompts ([5] huggingface.co) ([6] huggingface.co). This shows that MMLU-Pro is effectively capturing facets of intelligence – particularly reasoning – that were underestimated by MMLU. Its multi-domain scope also ensures that a model cannot excel merely by niche expertise; it must be generally capable.

As of late 2025, top LLMs still leave significant room for improvement on MMLU-Pro (mid-80% scores), so the benchmark remains a challenging yardstick. Meanwhile, the community is already expanding its horizons: efforts like **MMLU-ProX** for multilingual evaluation ([9] mmluprox.github.io) and research into alternate evaluation modes ([7] openreview.net) ([8] openreview.net) indicate that MMLU-Pro will evolve. From an industry perspective, MMLU-Pro provides a more rigorous reporting metric for new models and an incentive to enrich training for reasoning skills. From an academic perspective, it highlights the frontier where current models stumble.

Looking ahead, MMLU-Pro sets a precedent: benchmarks must continuously elevate difficulty to match advancing model capabilities. It also exemplifies collaborative progress: the creators have open-sourced the data and code ([26] github.com) ([21] huggingface.co), encouraging others to build on their work. For example, other researchers might extend MMLU-Pro with more languages, or more subjects (e.g. humanities reasoning puzzles), or transform it into interactive settings.

In closing, MMLU-Pro plays a dual role: as a **ruler** to measure AI progress and as a **tool** to remind us where progress is still needed. Its emphasis on reasoning and robust evaluation pushes models closer to the goal of deep understanding. At the same time, parallel research (such as the answer-matching approach ([7] openreview.net)) urges that our measurements also evolve beyond multiple-choice. Ultimately, benchmarks like MMLU-Pro will be most beneficial when they are part of a diverse evaluation ecosystem – measuring factual knowledge, logic, creativity, and safety in combination.

The remaining journey towards truly general artificial intelligence is long, but MMLU-Pro provides a tougher scoreboard—a crucial asset as we step up to the next challenges.

# References

1. Y. Wang *et al.*, "MMLU-Pro: A More Robust and Challenging Multi-Task Language Understanding Benchmark," *Advances in Neural Information Processing Systems 37 (NeurIPS 2024), Datasets and Benchmarks Track*, Sept. 2024 (papers.nips.cc).

2. Y. Wang *et al.*, MMLU-Pro Dataset, **TIGER-AI-Lab / MMLU-Pro** (GitHub), 2024. (Dataset code and info) ([19] huggingface.co) ([4] huggingface.co).

3. T. Hendrycks *et al.*, "Measuring Massive Multitask Language Understanding," *ICLR 2021 Workshop* (also arXiv 2020) ([11] paperswithcode.com).

4. PapersWithCode, *MMLU-Pro Dataset*, 2024 (providing dataset summary and description) ([27] paperswithcode.com).

5. HuggingFace, *TIGER-Lab/MMLU-Pro Dataset Documentation*, 2024 (detailed dataset summary and stats) ([1] huggingface.co) ([19] huggingface.co).

6. MMLU Dataset entry, *PapersWithCode* (on MMLU coverage and purpose) ([10] paperswithcode.com).

7. MMLU Dataset entry, *Hendrycks et al., PapersWithCode summary*, 2020 ([10] paperswithcode.com) ([11] paperswithcode.com).

8. AI Forger and AI Rank, *Benchmarks: MMLU-Pro*, 2025 (public leaderboard snapshots) (airank.dev) ([28] theaiforger.com).

9. S. Klebanov, *"GPT-4.1 vs GPT-4o MMLU Benchmark Comparison," Promptfoo*, 2024. (Context on "GPT-4o" naming) ([5] huggingface.co).

10. OpenReview submission, *"Answer matching outperforms multiple choice for LLM evaluations"* (ICLR 2026), 2025 ([7] openreview.net) ([8] openreview.net).

11. *MMLU-ProX: A Multilingual Benchmark for Advanced LLM Evaluation*, University of Tokyo et al., arXiv 2024 (29-language MMLU-Pro) ([9] mmluprox.github.io).

## External Sources

[1]   https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:,numb...

[2]   https://github.com/TIGER-AI-Lab/MMLU-Pro#:~:more%...

[3]   https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:Disci...

[4]   https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:Model...

[5]   https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:GPT,0...

[6]   https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:GPT,0...

[7]   https://openreview.net/forum?id=xkbjNJi0eb#:~:TL%3B...

[8]   https://openreview.net/forum?id=xkbjNJi0eb#:~:annot...

[9]   https://mmluprox.github.io/#:~:Exist...

[10]  https://paperswithcode.com/dataset/mmlu#:~:MMLU%...

[11]  https://paperswithcode.com/paper/measuring-massive-multitask-language#:~:We%20...

[12]  https://paperswithcode.com/paper/mmlu-pro-a-more-robust-and-challenging-multi#:~:bench...

[13]  https://blog.kasralekan.com/ideas/lm-performance-plateau/#:~:Figur...

[14]  https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:GPT,0...

[15]  https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro/blob/main/README.md#:~:,the%...

[16]  https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro/blob/main/README.md#:~:,data...

[17]  https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:match...

[18]  https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:,engi...

[19] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro/blob/main/README.md#:~:,requ...

[20] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:10%20...

[21] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro/blob/main/README.md#:~:4...

[22] https://openreview.net/forum?id=y10DM6R2r3#:~:Ku%2C...

[23] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:We%20...

[24] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:The%2...

[25] https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro#:~:7...

[26] https://github.com/TIGER-AI-Lab/MMLU-Pro#:~:We%20...

[27] https://paperswithcode.com/dataset/mmlu-pro#:~:,mode...

[28] https://theaiforger.com/benchmarks/mmlu-pro#:~:%2301...

## IntuitionLabs - Industry Leadership & Services

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

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

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

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

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

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

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

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

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

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

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

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

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

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

## DISCLAIMER

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

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

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

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

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

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