# Unearthing *Skill-Level Insights* for Understanding Tradeoffs of Foundation Models

Mazda Moayeri

Vidhisha Balachandran, Varun Chandrasekaran,
Safoora Yousefi, Thomas Fel, Soheil Feizi, Besmira Nushi, Neel Joshi, Vibhav Vineet
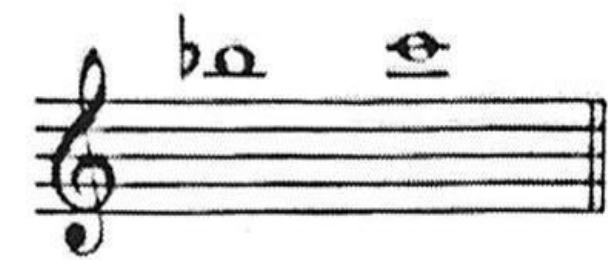
Work done with Microsoft Research AI Frontiers

## Key ingredient: *Rationales*

Rationales reveal the skills and steps needed to answer a question.



Evaluation Instance (MMMU, #744)
Select the type of the interval.
(A) perfect
(B) major
(C) minor
(D) diminished

Generated Rationale with localized skills
Step 1. Recognize the Image
- *Skill:* **Perception, Visual Recognition, Symbol Identification, Treble clef recognition**
- *Conclusion: The image is of notes on the treble staff.*

Step 2. Identify the Notes
- *Skill:* **Perception, Visual Recognition, Note Identification, Pitch recognition**
- *Conclusion: The notes are B and E.*

Step 3. Determine the Interval
- *Skill:* **Knowledge, Music Theory, Interval Identification, Pitch relationships**
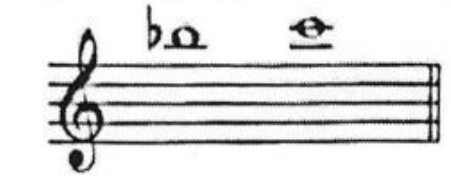...

Testing along *"skill-slices"* (instances testing the same skill) unlocks richer insight from existing evaluation data.

## Isolating Skills via Probing Questions

Reframing rationale steps where a skill is applied allows for generating questions that *test only one skill.*



Surface **instance** and **rationale-step** for skill to probe
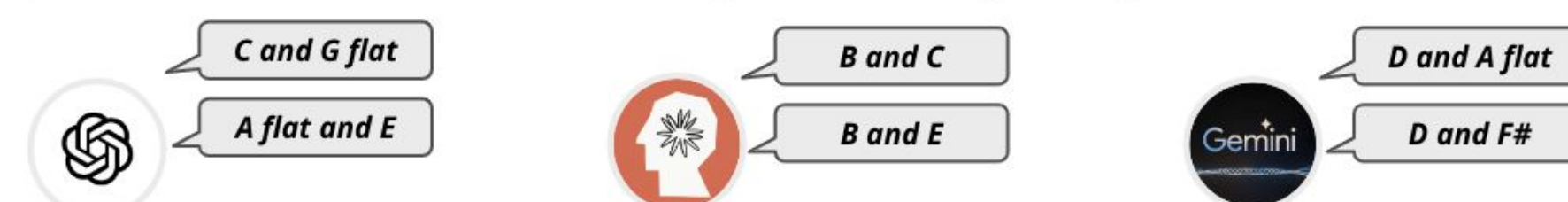[long question]
Step 2. Identify the Notes
- *Skill: Note Identification*
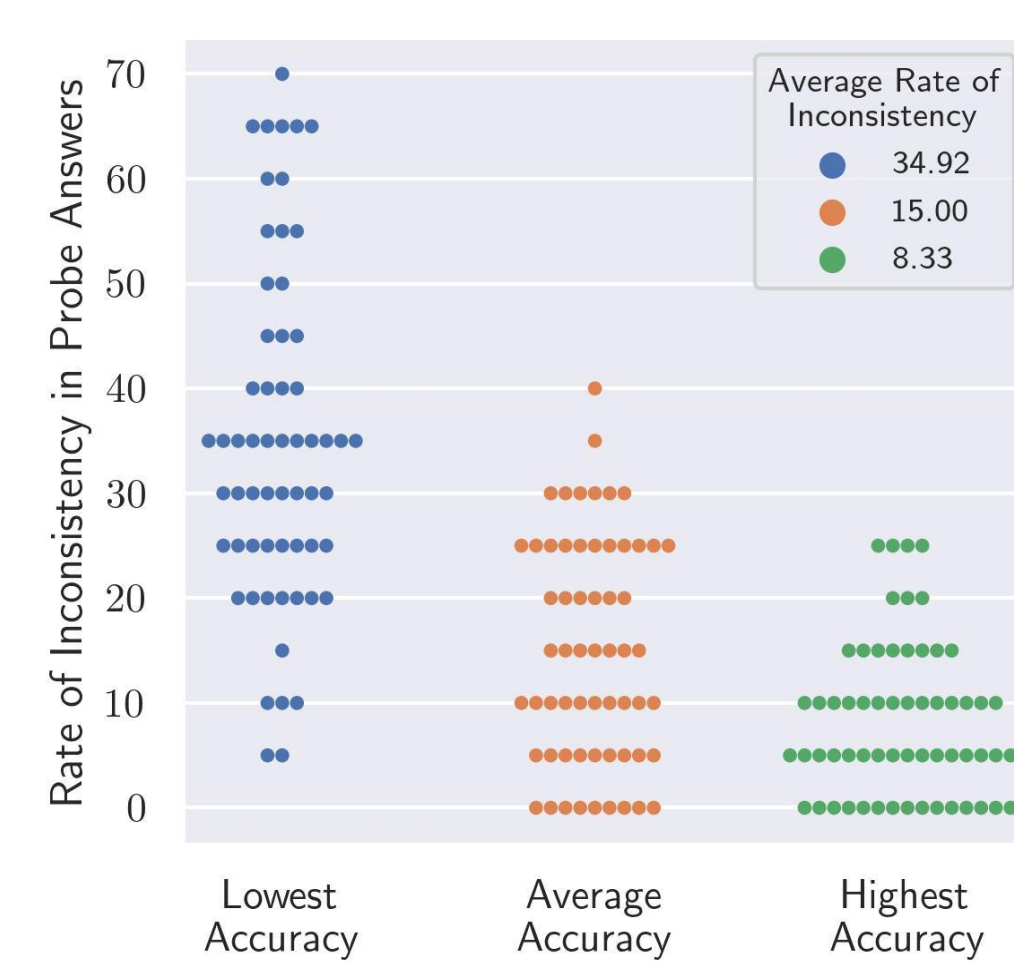- *Conclusion: The notes are B and E.*

Reframe concluding step to probing question
**The notes are B and E.** ⟶ **What are the notes?**

Check consistency over multiple responses
C and G flat / A flat and E     B and C / B and E     D and A flat / D and F#

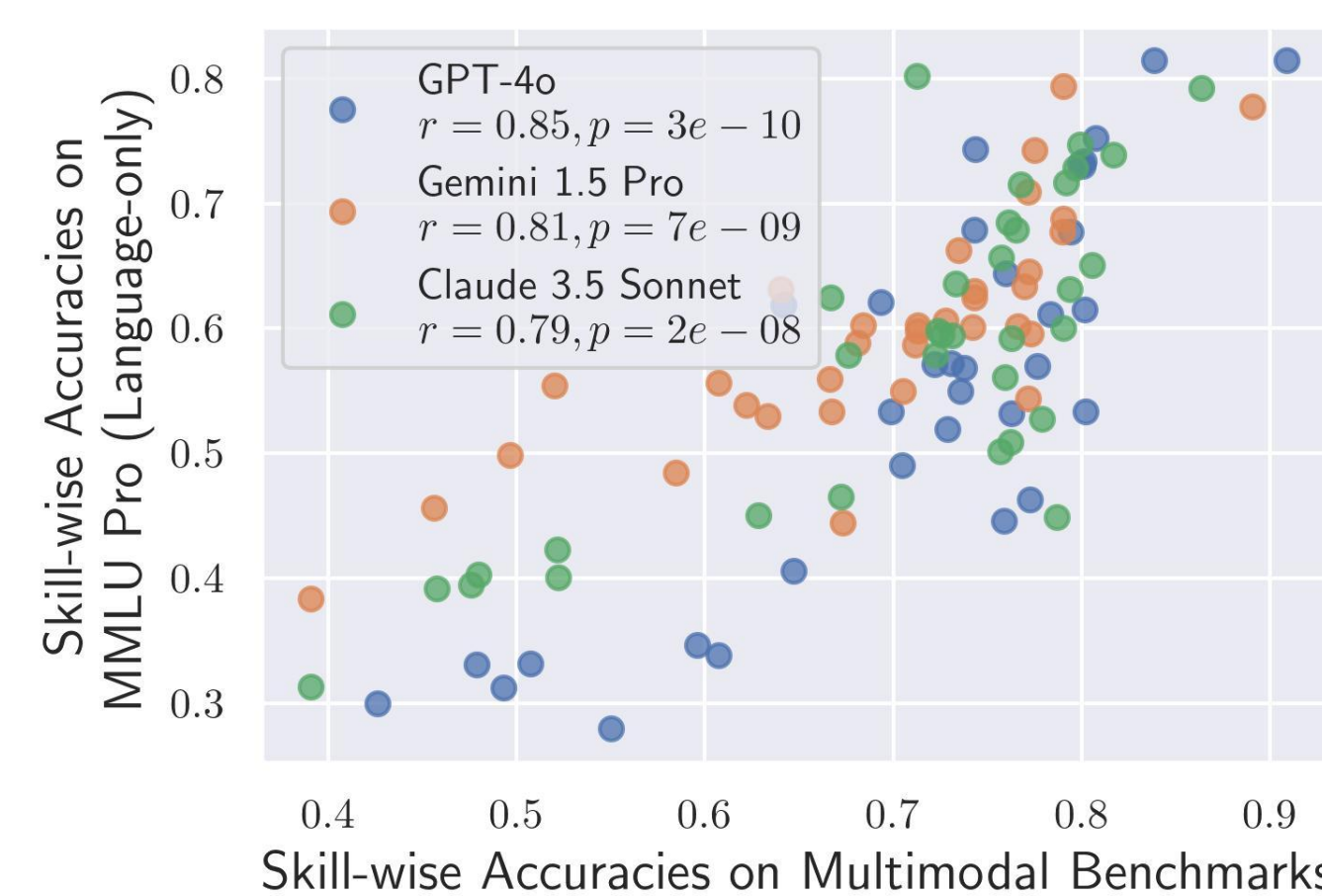**Are above answers consistent?**
✗  ✗  ✗

Consistency over multiple answers to the same probe question offers a 2nd corroborating and complementary (compared to *skill-slice* accuracies) measure of skill proficiency.



## Skill-based Routing Improves Accuracy

Skill-slice accuracies generalize.

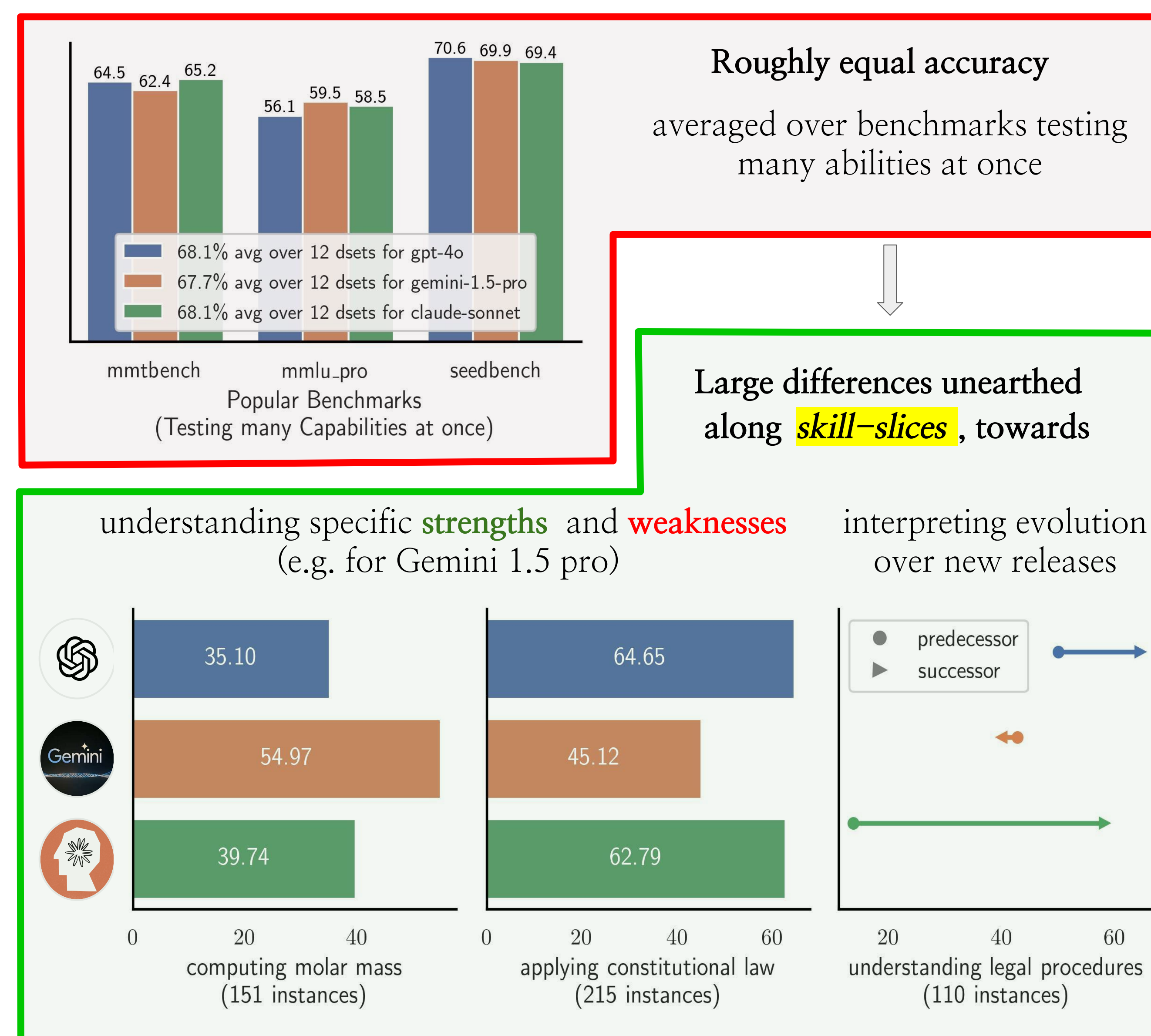Highly correlated accuracies over slices drawn from distinct corpuses, *even when modality changes.*



GPT-4o r = 0.85, p = 3e − 10
Gemini 1.5 Pro r = 0.81, p = 7e − 09
Claude 3.5 Sonnet r = 0.79, p = 2e − 08

Routing each instance to the model with the highest slice accuracies for the relevant skills leads to accuracy gains.

**(+3 to 6.8% on MMLU-Pro)**



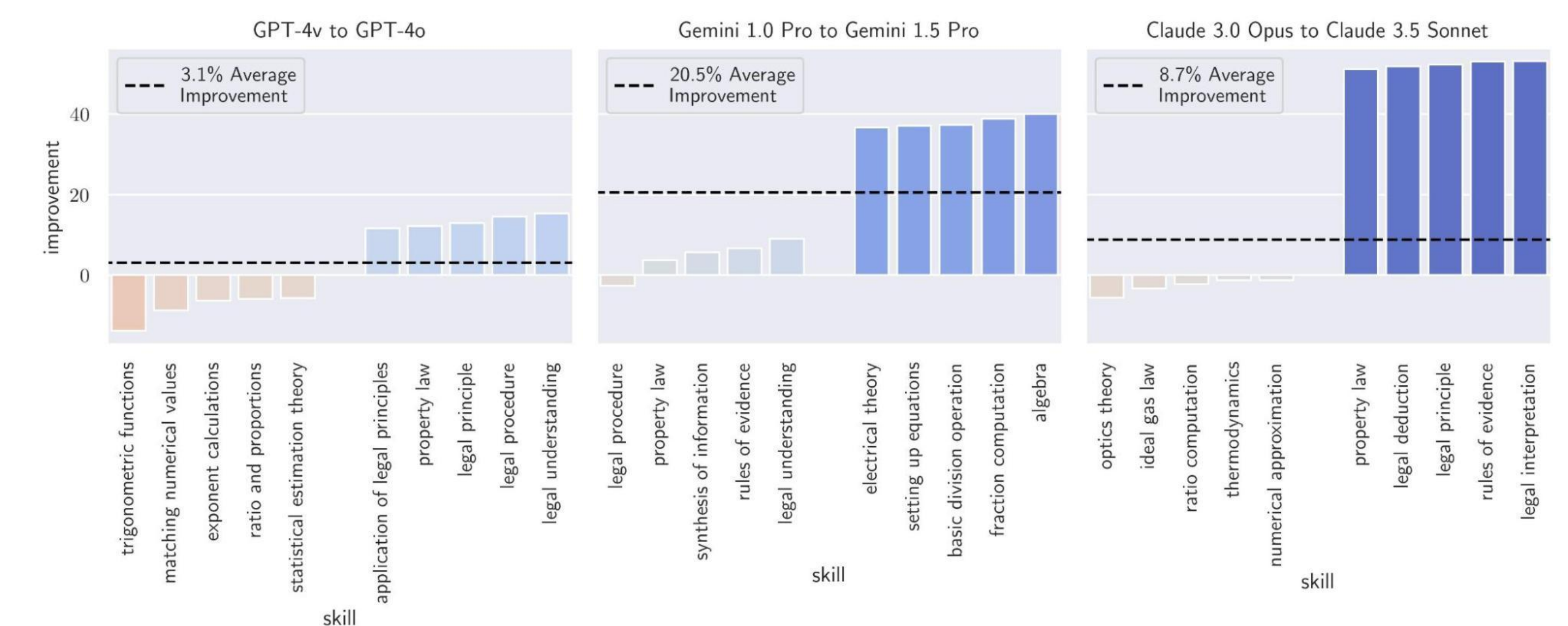| | Accuracy over 12 Datasets | MMLU Pro Accuracy |
|---|---|---|
| Claude Sonnet 3.5 | 68.09 | 58.49 |
| Gemini 1.5 Pro | 67.71 | 59.47 |
| GPT-4o | 68.09 | 56.09 |
| Skill Routing | 70.96 | 62.94 |

---

Foundation models, by design, encapsulate a wide breadth of skills. Thus, modern benchmarks test many skills, all at once – even in the same instance.

*How many insights are we simply averaging away?*

We present a method harnessing generated rationales to scalably recover the skill-level insights hiding within existing benchmark evaluations.
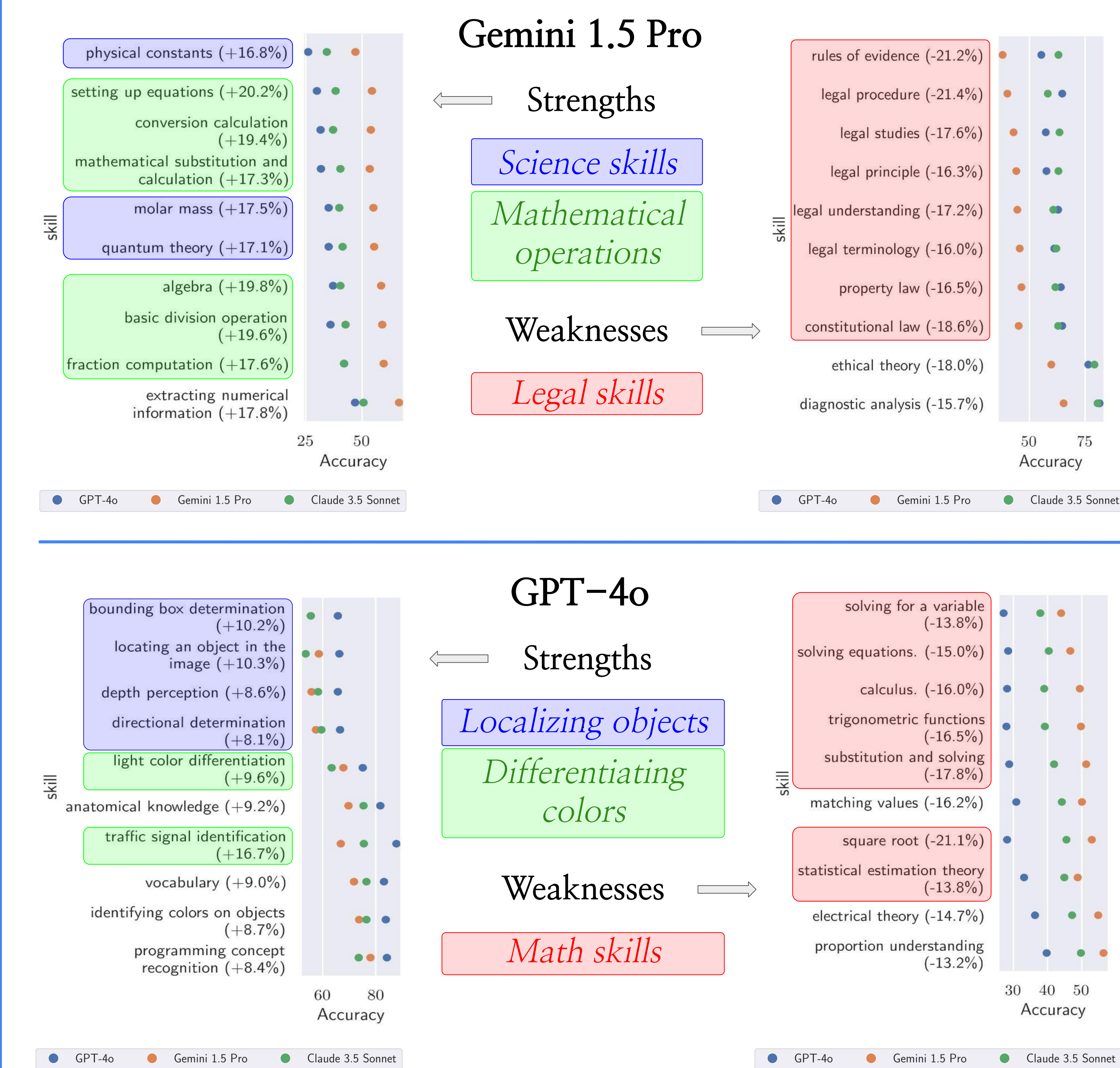


**Roughly equal accuracy** averaged over benchmarks testing many abilities at once

68.1% avg over 12 dsets for gpt-4o
67.7% avg over 12 dsets for gemini-1.5-pro
68.1% avg over 12 dsets for claude-sonnet

mmtbench  mmlu_pro  seedbench
Popular Benchmarks
(Testing many Capabilities at once)

**Large differences unearthed along *skill-slices*, towards**

understanding specific **strengths** and **weaknesses** (e.g. for Gemini 1.5 pro)  interpreting evolution over new releases

- predecessor
- successor

computing molar mass (151 instances)   applying constitutional law (215 instances)   understanding legal procedures (110 instances)

## Evolution over Model Releases

*Skill-slices* uncover the most-improved skills from one release to the next.



GPT-4v to GPT-4o — 3.1% Average Improvement
Gemini 1.0 Pro to Gemini 1.5 Pro — 20.5% Average Improvement
Claude 3.0 Opus to Claude 3.5 Sonnet — 8.7% Average Improvement

Seems like legal skills were a recent priority for OpenAI and Anthropic. Below, we see legal skills are also the area where Gemini is furthest behind!

## Strengths and Weaknesses of Top Models

### Gemini 1.5 Pro



Strengths ⟸
*Science skills*
*Mathematical operations*

physical constants (+16.8%)
setting up equations (+20.2%)
conversion calculation (+19.4%)
mathematical substitution and calculation (+17.3%)
molar mass (+17.5%)
quantum theory (+17.1%)
algebra (+19.8%)
basic division operation (+19.6%)
fraction computation (+17.6%)
extracting numerical information (+17.8%)

Weaknesses ⟹
*Legal skills*

rules of evidence (-21.2%)
legal procedure (-21.4%)
legal studies (-17.6%)
legal principle (-16.3%)
legal understanding (-17.2%)
legal terminology (-16.0%)
property law (-16.5%)
constitutional law (-18.6%)
ethical theory (-18.0%)
diagnostic analysis (-15.7%)

GPT-4o   Gemini 1.5 Pro   Claude 3.5 Sonnet

### GPT-4o



Strengths ⟸
*Localizing objects*
*Differentiating colors*

bounding box determination (+10.2%)
locating an object in the image (+10.3%)
depth perception (+8.6%)
directional determination (+8.1%)
light color differentiation (+9.6%)
anatomical knowledge (+9.2%)
traffic signal identification (+16.7%)
vocabulary (+9.0%)
identifying colors on objects (+8.7%)
programming concept recognition (+8.4%)

Weaknesses ⟹
*Math skills*

solving for a variable (-13.8%)
solving equations. (-15.0%)
calculus. (-16.0%)
trigonometric functions (-16.5%)
substitution and solving (-17.8%)
matching values (-16.2%)
square root (-21.1%)
statistical estimation theory (-13.8%)
electrical theory (-14.7%)
proportion understanding (-13.2%)

GPT-4o   Gemini 1.5 Pro   Claude 3.5 Sonnet

### Claude 3.5 Sonnet



Strengths ⟸
*Chart Understanding*

coordinate system mapping (+10.1%)
comparing data (+16.9%)
summarizing (+10.9%)
validation. (+9.2%)
identifying categories (+7.5%)
synthesizing information (+6.5%)
analytical reasoning (+7.3%)
understanding chart data (+6.0%)
trend analysis (+7.7%)
graph identification (+6.6%)

Weaknesses ⟹
*Visual details and counting*

visual understanding and counting (-11.1%)
final answer selection (-5.4%)
image processing (-5.9%)
visual differentiation (-6.0%)
multiple person identification (-6.2%)
surface detail recognition (-7.9%)
architectural image recognition (-6.2%)
distinguishing animal characteristics (-7.6%)
animal behavior (-13.1%)
understanding ecology (-6.1%)

*June 21, 2024 release, not October 22 update

GPT-4o   Gemini 1.5 Pro   Claude 3.5 Sonnet