Breaking the Benchmark: Findings From the Industry's First LLM Creativity Study

Springboards

↓ SPRINGBOARDS IS

A creative tool to inspire, not give you the answers.







SAM ALTMAN "Al will handle 95% of what marketers use agencies, strategists, marketing and creative professionals for today." March 2024



MARK ZUCKERBERG

"We're going to get to a point where you're a business, you come to us, you tell us what your objective is, you connect to your bank account, you don't need any creative, you don't need any targeting demographic, you don't need any measurement, except to be able to read the results that we spit out"

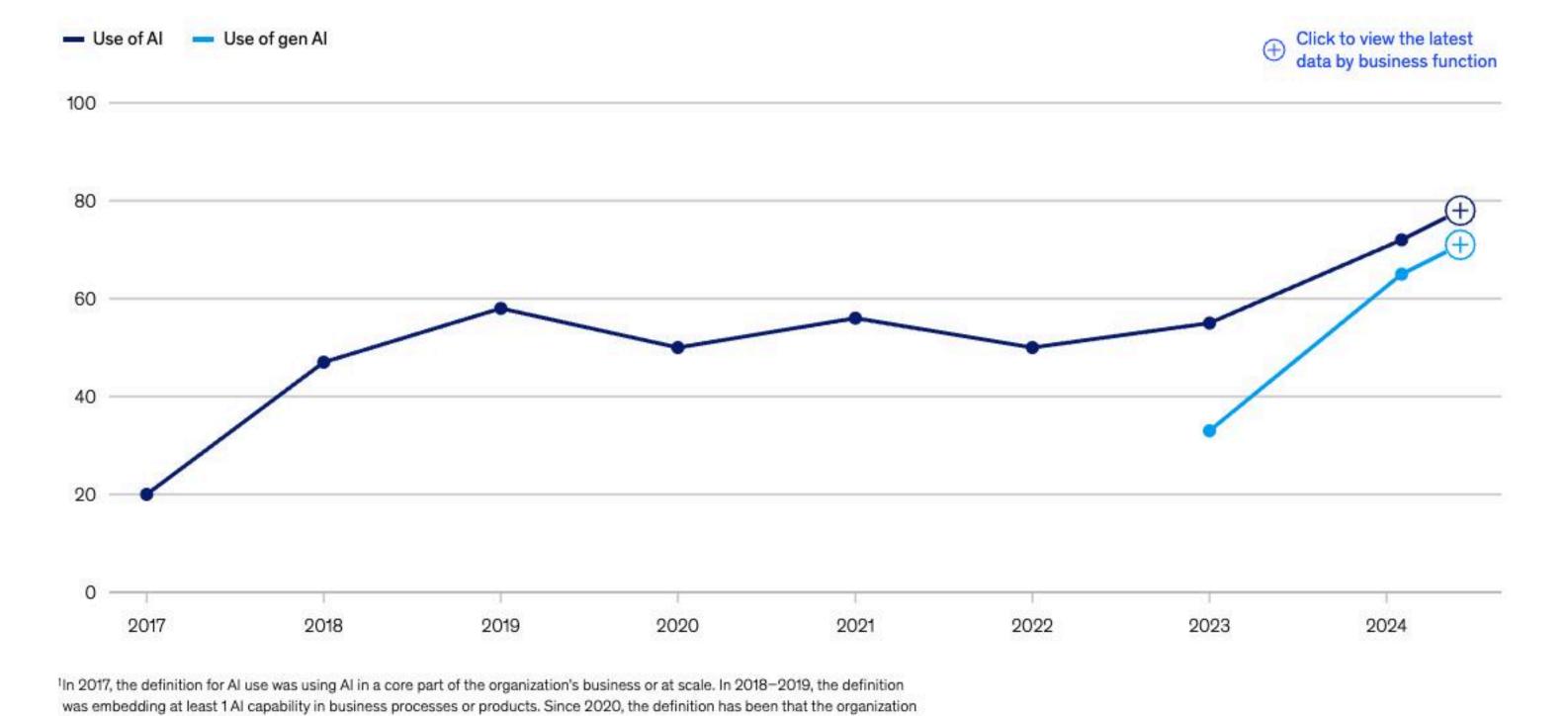
May 2025



Al adoption has been rampant.

Organizations' use of Al has accelerated markedly in the past year, after years of little meaningful change.

Organizations that use AI in at least 1 business function, 1 % of respondents



Source: McKinsey Global Surveys on the state of Al https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai

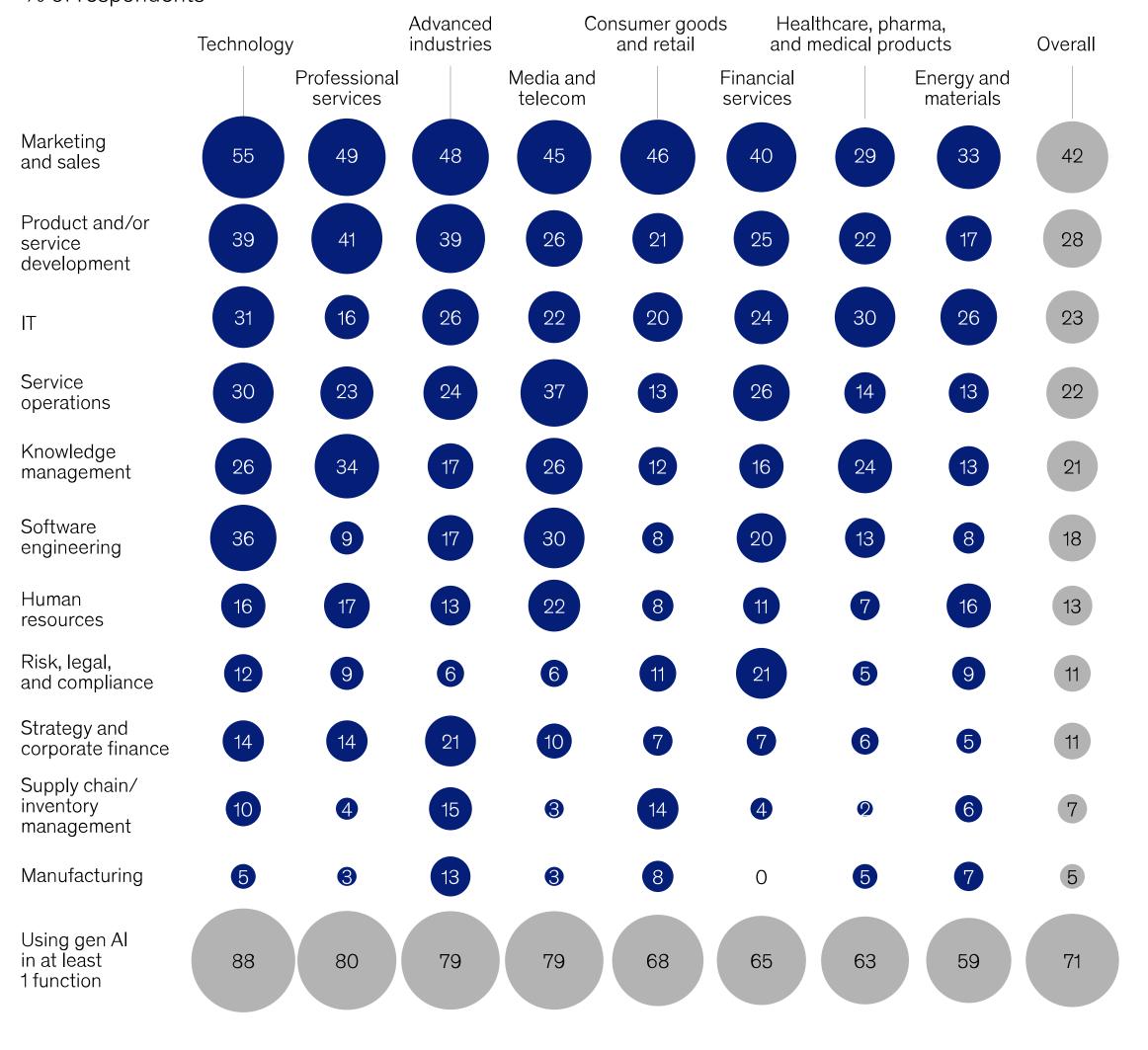
has adopted AI in at least 1 function.



And highest amongst marketing teams

Organizations across industries have begun to use gen Al in marketing and sales, though other uses vary by industry.

Business functions in which respondents' organizations are regularly using gen AI, by industry,¹ % of respondents



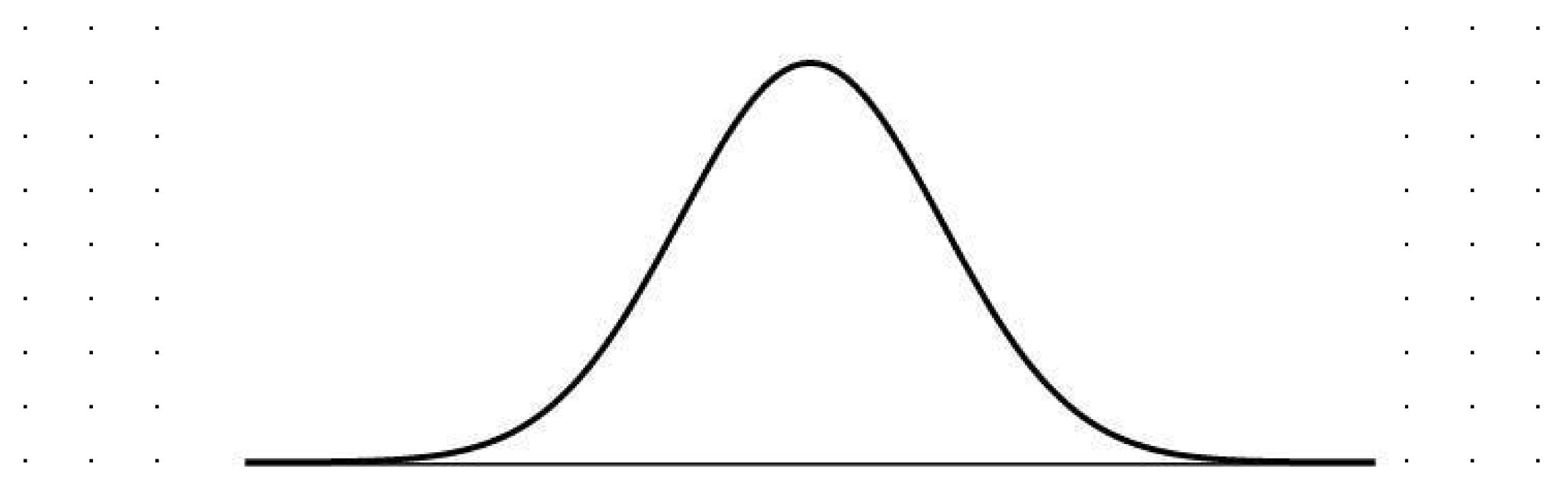


https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai





Unknown to many... they are funneling us all into the same place.





BECAUSE IT THREATENS TO FLATTEN CULTURE

We are seeing it happen in music



PLAY LIVE RADIO

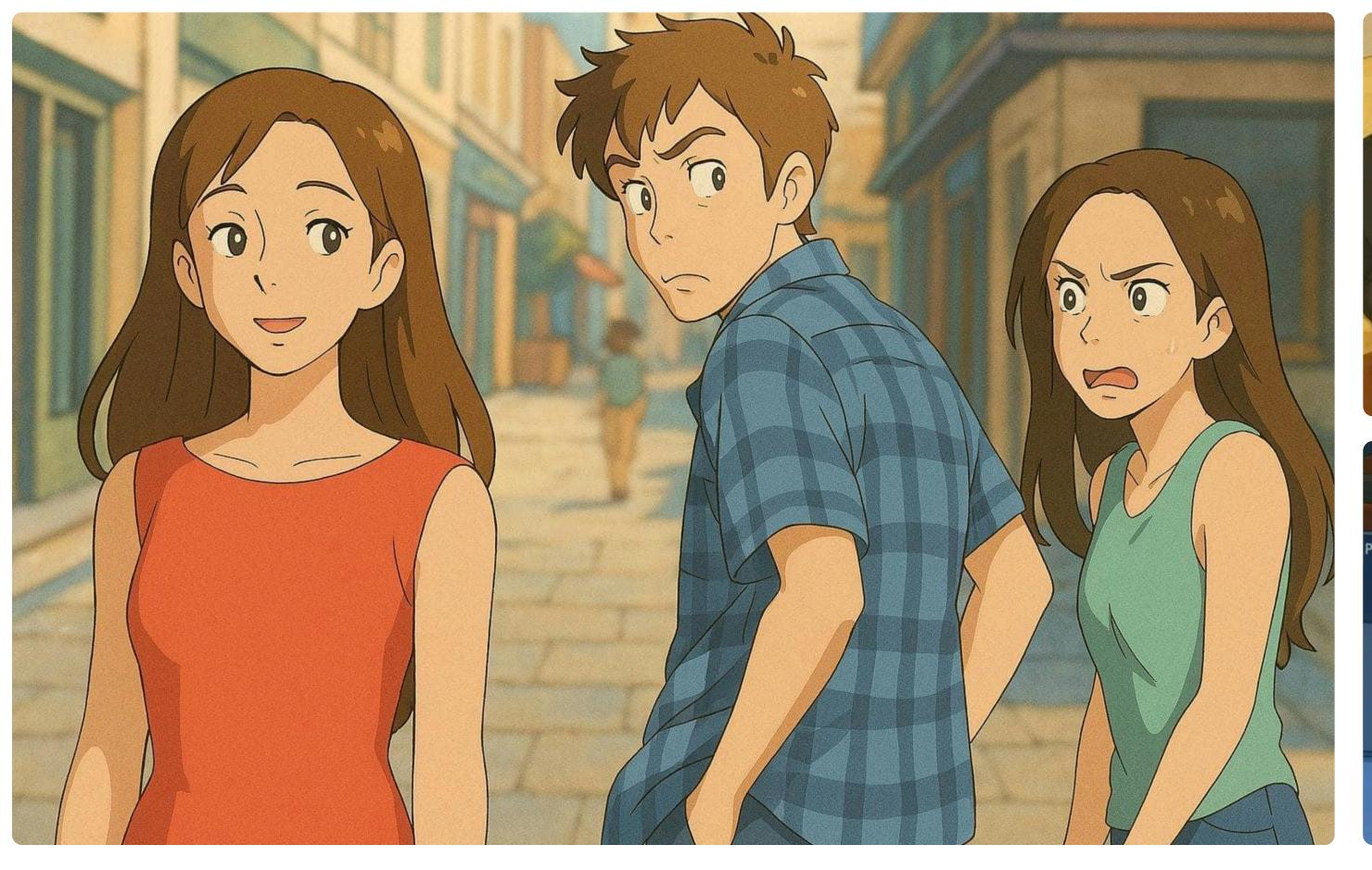
Music streamers are serving you Algenerated and 'ghost' music. Here's how it hurts real artists





BECAUSE IT THREATENS TO FLATTEN CULTURE

We are seeing it happen in culture







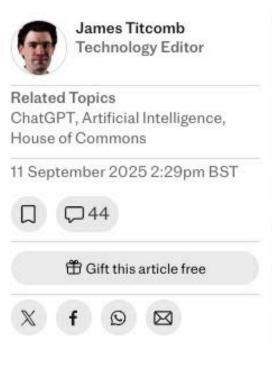




We are seeing it happen in politics

ChatGPT triggers surge in MPs using AI-written speeches

Use of common artificial intelligence phrases such as 'I rise today' has jumped since 2022



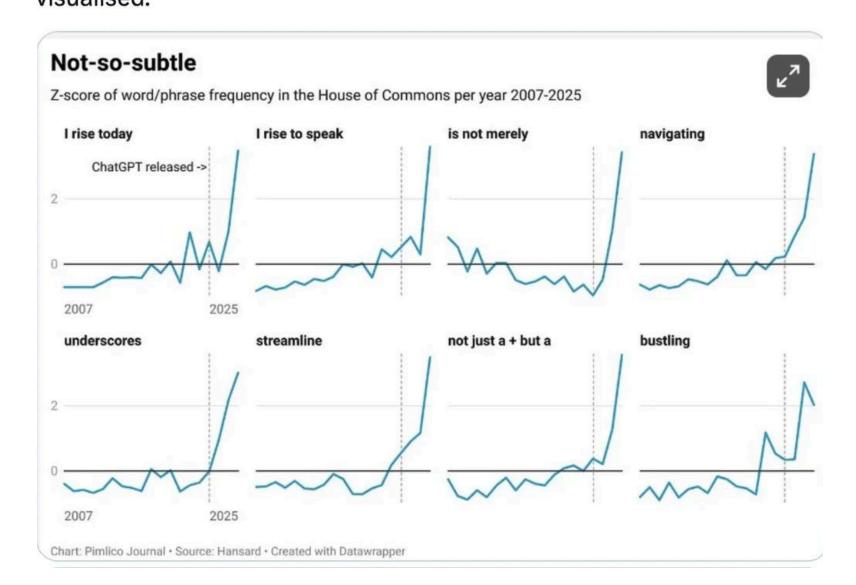


telegraph.co.uk/business/2025/09/11/chatgpt-triggers-surge-in-mps-using-ai-written-speeches/



The surge in Al-written speeches in Britain's House of Commons visualised:

Ø ...

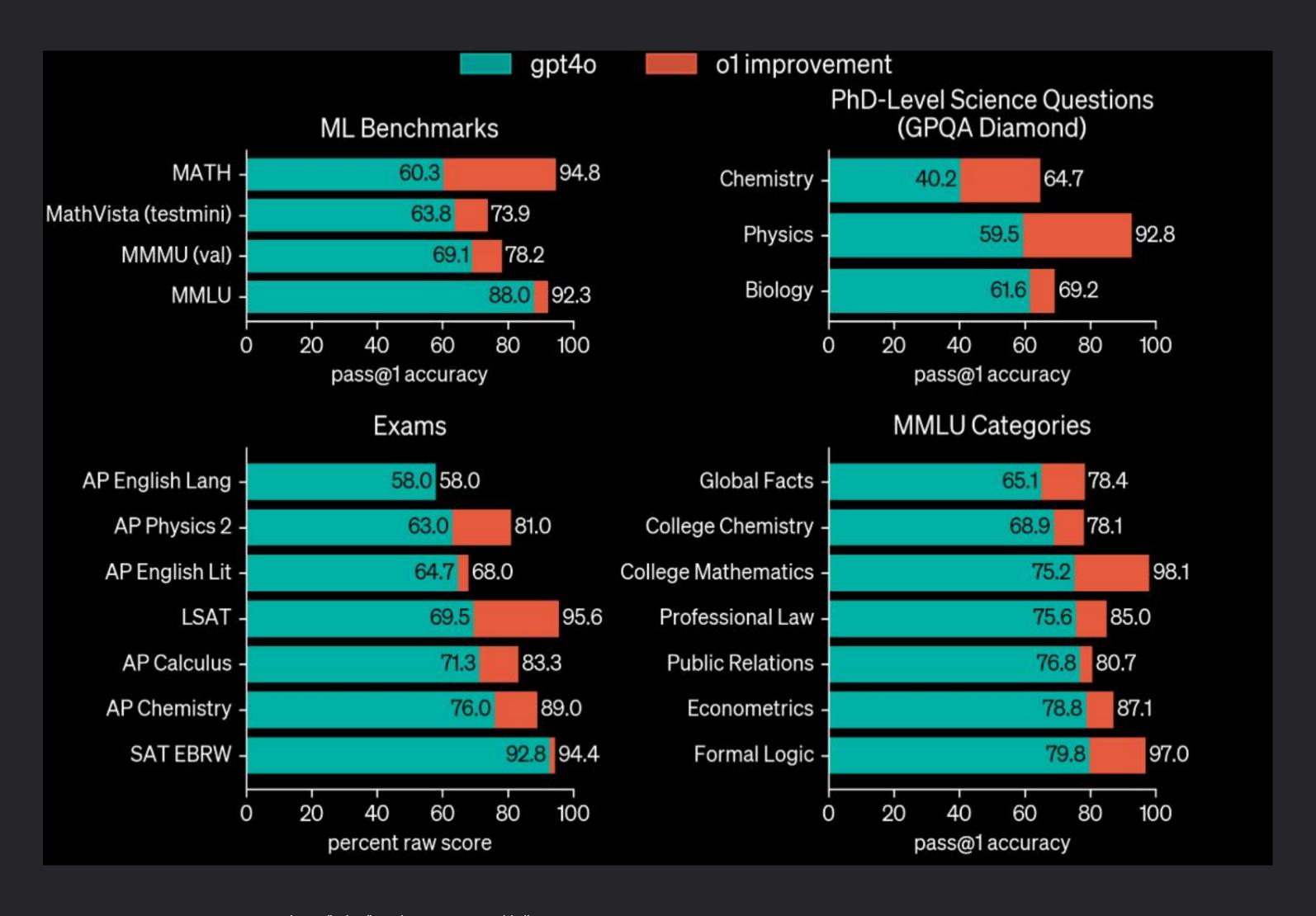






↓ OPENAI EVALS

The first step is to break the benchmarks



↓ ANTHRO	PIC	EVALS	S				•		Claude 3.7 Sonnet 64K extended thinking	Claude 3.7 Sonnet No extended thinking	Claude 3.5 Sonnet (new)	OpenAl o1 ¹	OpenAl o3-mini ¹ High	DeepSeek R1 32K extended thinking	Grok 3 Beta Extended thinking
						•	•	Graduate-level reasoning GPQA Diamond ³	78.2% / 84.8%	68.0%	65.0%	75.7% / 78.0%	79.7%	71.5%	80.2% / 84.6%
						•	-	Agentic coding SWE-bench Verified ²		62.3% / 70.3%	49.0%	48.9%	49.3%	49.2%	
						•	•	Agentic tool use	-	Retail 81.2%	Retail 71.5%	Retail 73.5%		575	277- -
					•	•	TAU-bench	<u>13176</u> 2	Airline 58.4%	Airline 48.8%	Airline 54.2%	POLAR Terren.	15 H 3 1 14 4 4 4 1		
•		•	•	•	•	•	•	Multilingual Q&A MMMLU	86.1%	83.2%	82.1%	87.7%	79.5%		
•			•				•	Visual reasoning MMMU (validation)	75%	71.8%	70.4%	78.2 %			76.0% / 78.0%
-		•	•	•	•	•	-	Instruction- following IFEval	93.2%	90.8%	90.2%			83.3%	
-		•	-	-	-		-	Math problem-solving MATH 500	96.2%	82.2%	78.0%	96.4%	97.9%	97.3%	
•		•	•	•		•	•	High school math competition AIME 2024 ³	61.3% / 80.0%	23.3%	16.0%	79.2% / 83.3%	87.3%	79.8%	83.9% / 93.3%
<u>-</u>		_	_	_	_	_	_)				

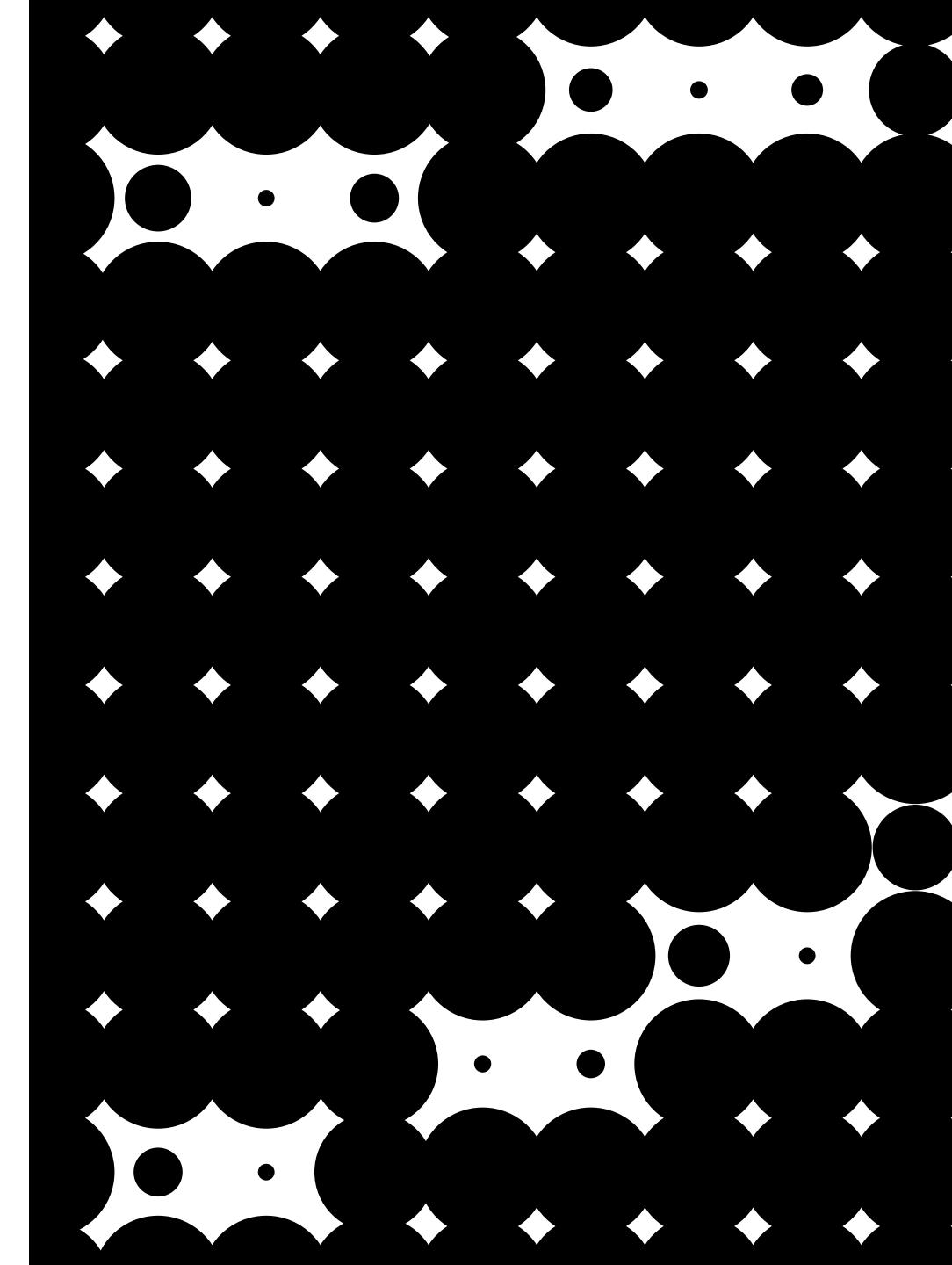
https://www.anthropic.com/news/claude-3-7-sonnet

The creative and advertising industry is left in the dark.

A collective, industry benchmark to help push creativity forward

↓ PART ONE

Breaking the Benchmarks

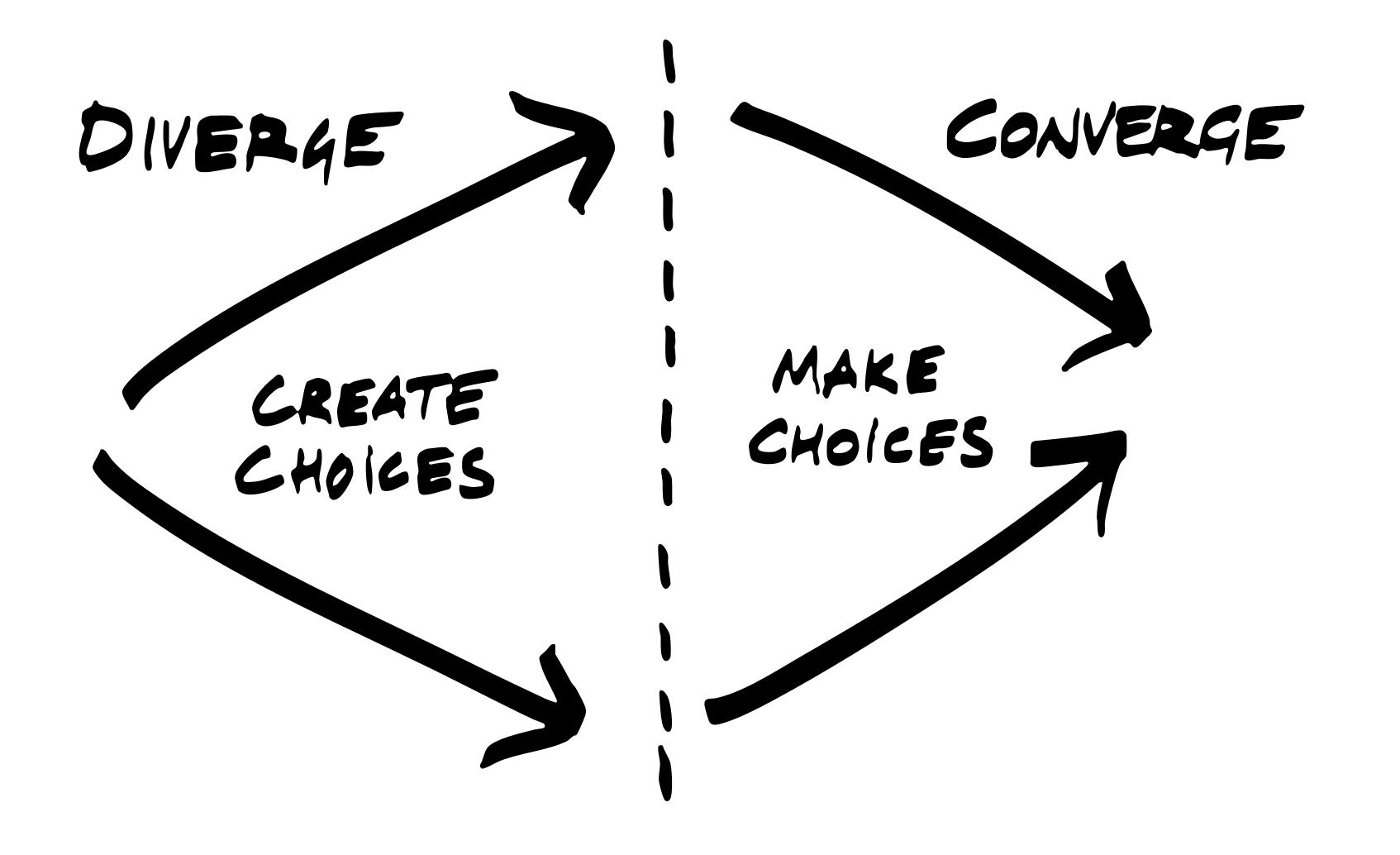




How do you even measure creativity?









Creativity Benchmark

A collective industry benchmark for creativity.

Most Al benchmarks test for maths, science, and law, but none measure creativity. So we're building the first benchmark that reflects how our world works. Created by and for strategists, creatives, and marketers and backed by leading industry bodies. This is your chance to shape how creativity gets measured.

New here? Get started

Sign back in

Backed by





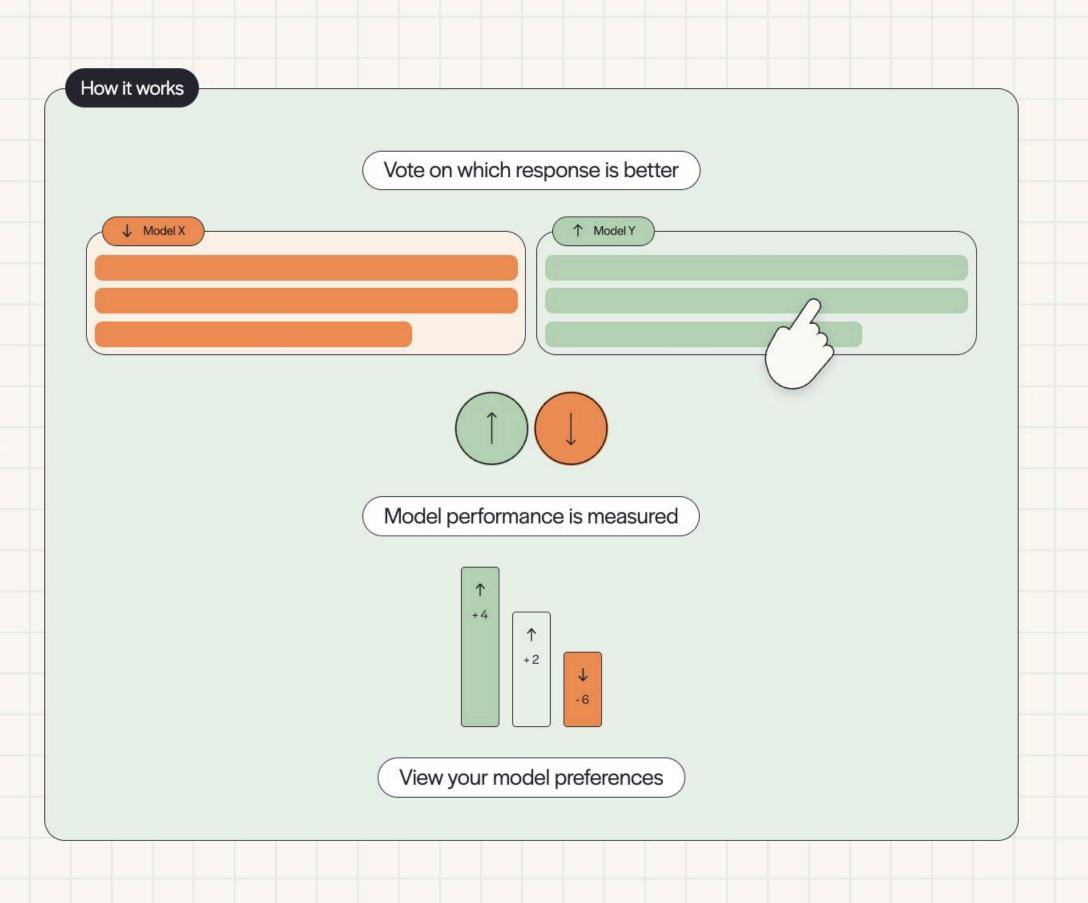




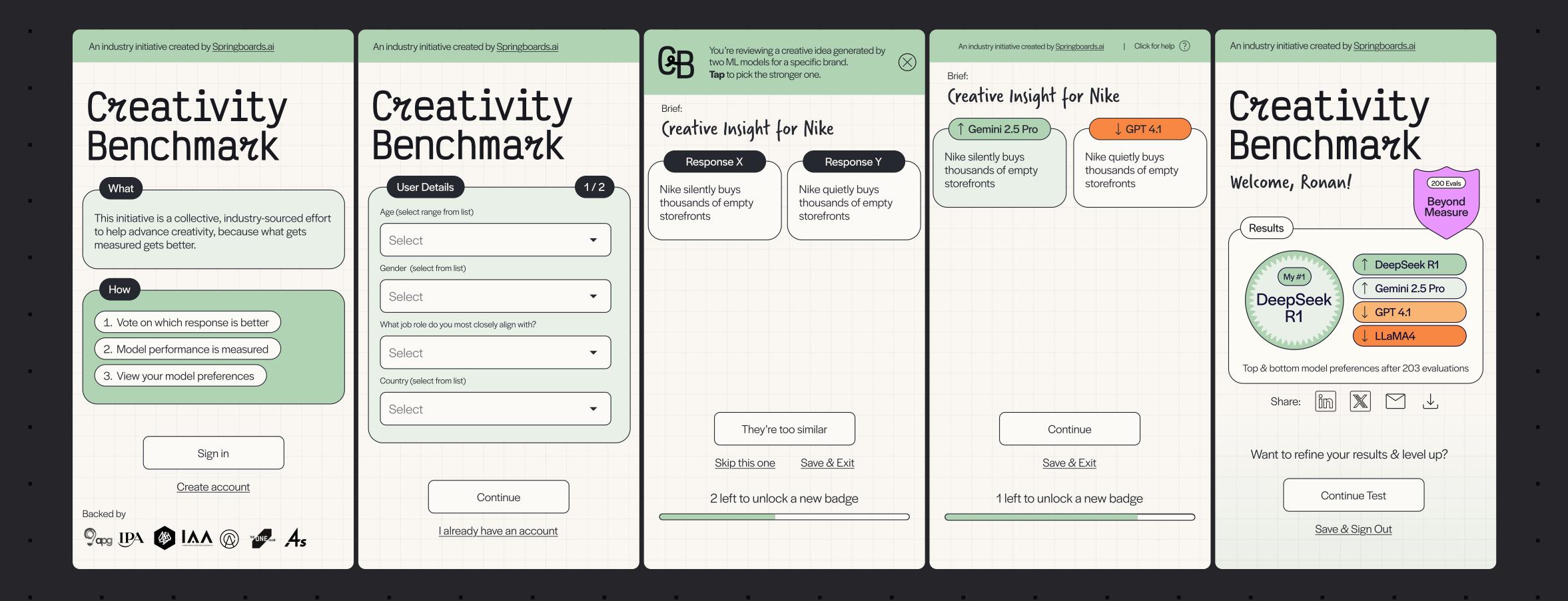








↓ TINDER FOR IDEAS

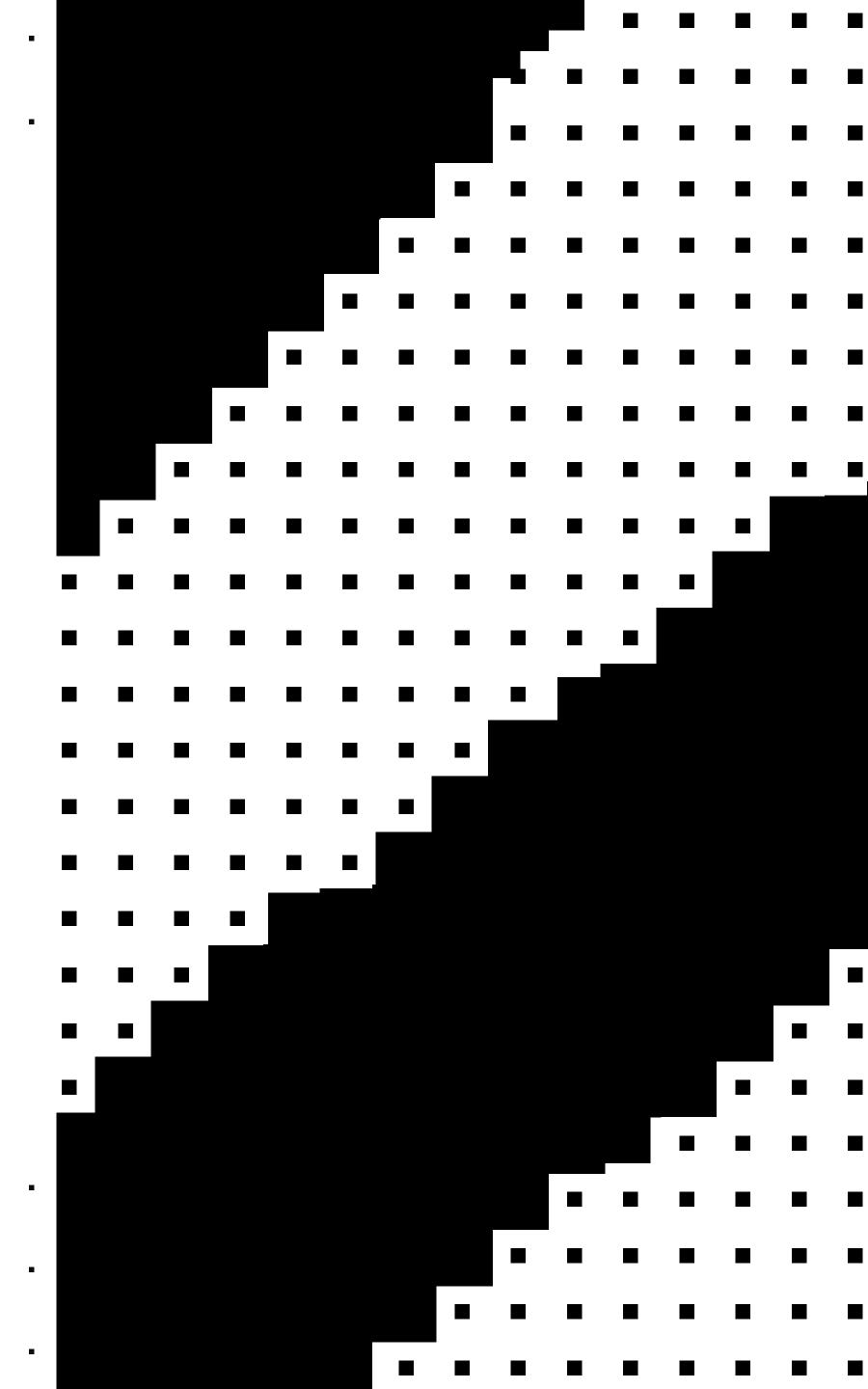




THE BENCHMARKS

- 1. Creative inspiration.

 Measuring the subjectivity of insights and ideas from LLMs.
- 2. Variance of outputs.
 Understanding how much variation people are getting from different models.
- 3. Creative thinking. Traditional creative thinking tests for problem solving, convergent and divergent thinking.

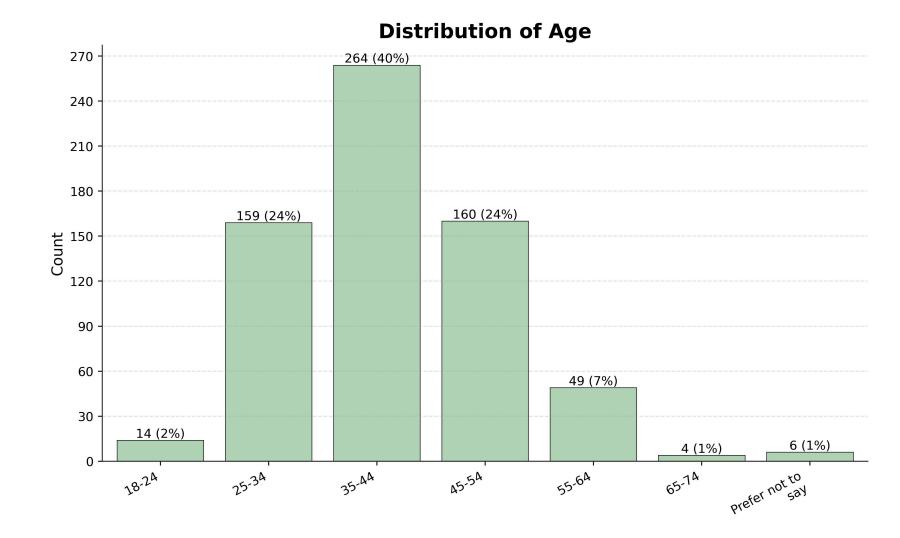


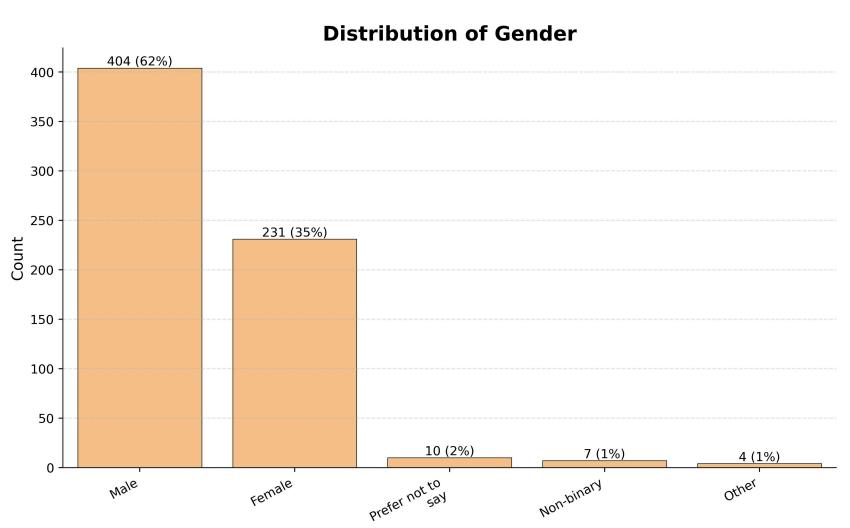
THE BORING BIT

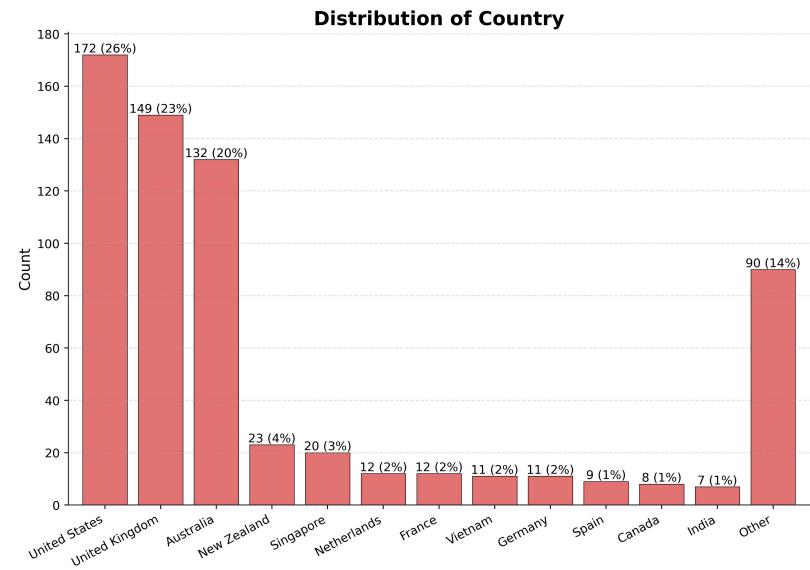
678 Professionals,

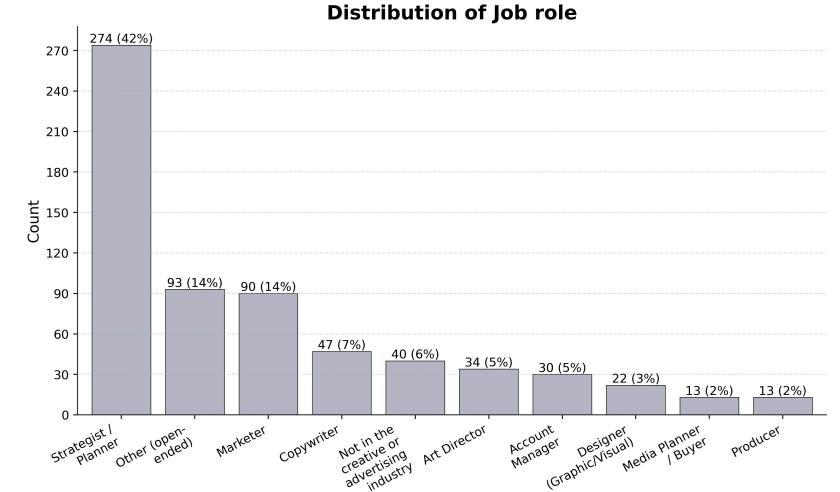
Comparisons

11,012





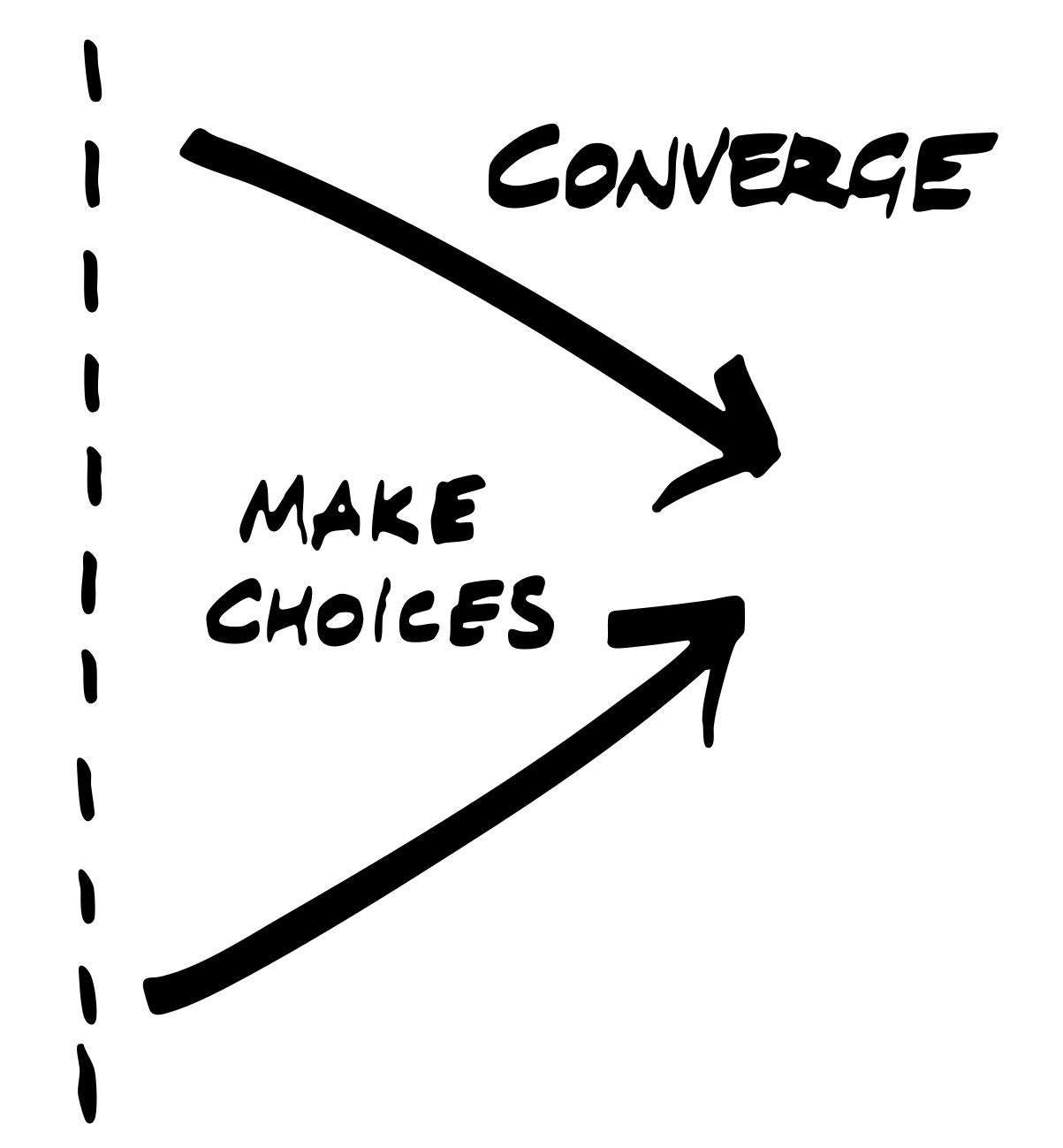






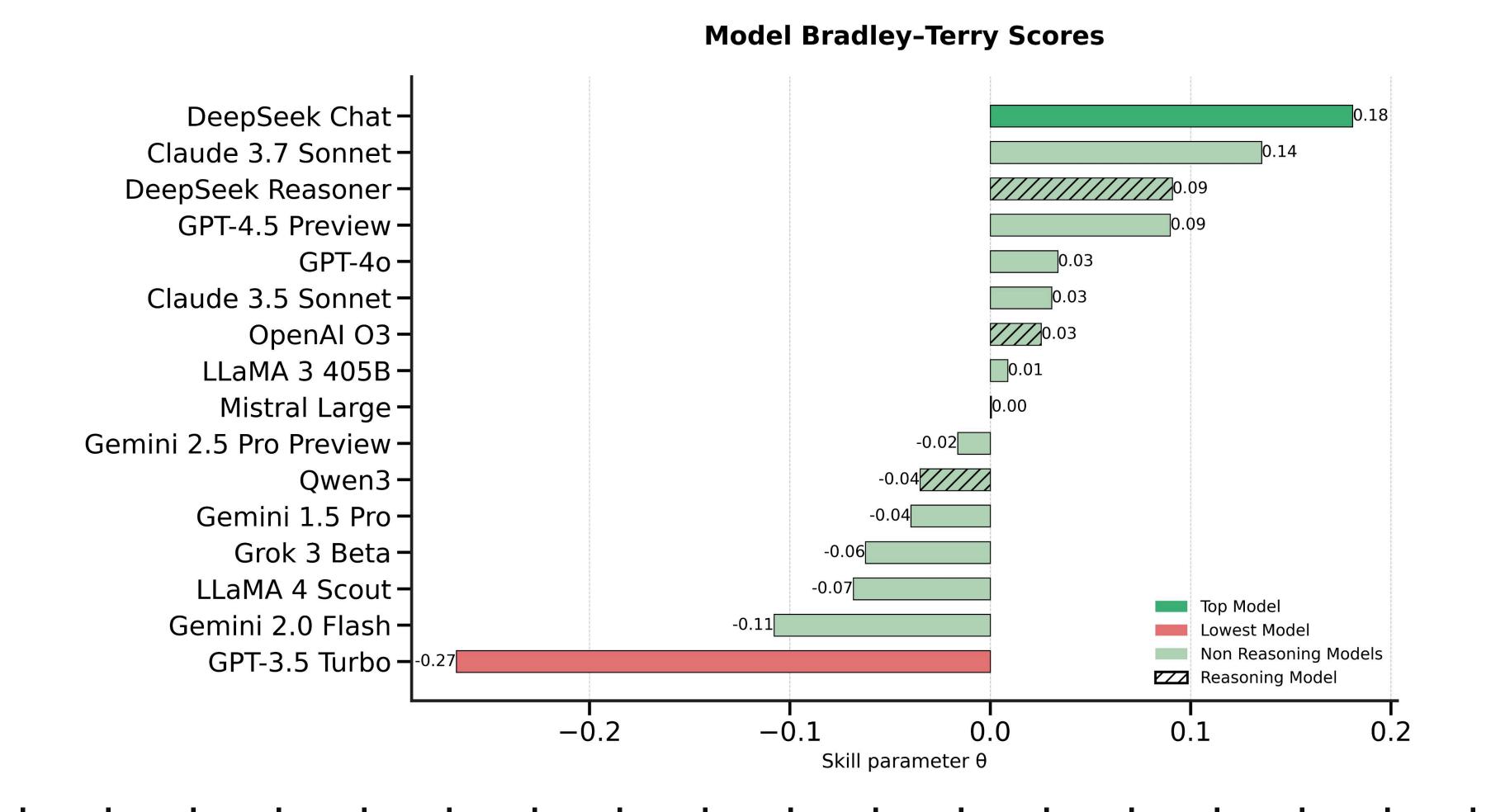
↓ STARTING WITH CONVERGE

Cause humans can break rules





After nearly 700 people and 11k matches we have a winner







CREATIVITY BENCHMARK

If We Judged Creativity On A Purely Academic Torrance Test, The Machines Would

Dominate

The originality of machines: AI takes the Torrance Test

Erik E. Guzik ^a ♀ ☒, Christian Byrge ^b, Christian Gilde ^c

+ Add to Mendeley
Share
Cite

https://doi.org/10.1016/j.yjoc.2023.100065

Under a Creative Commons license

Open access

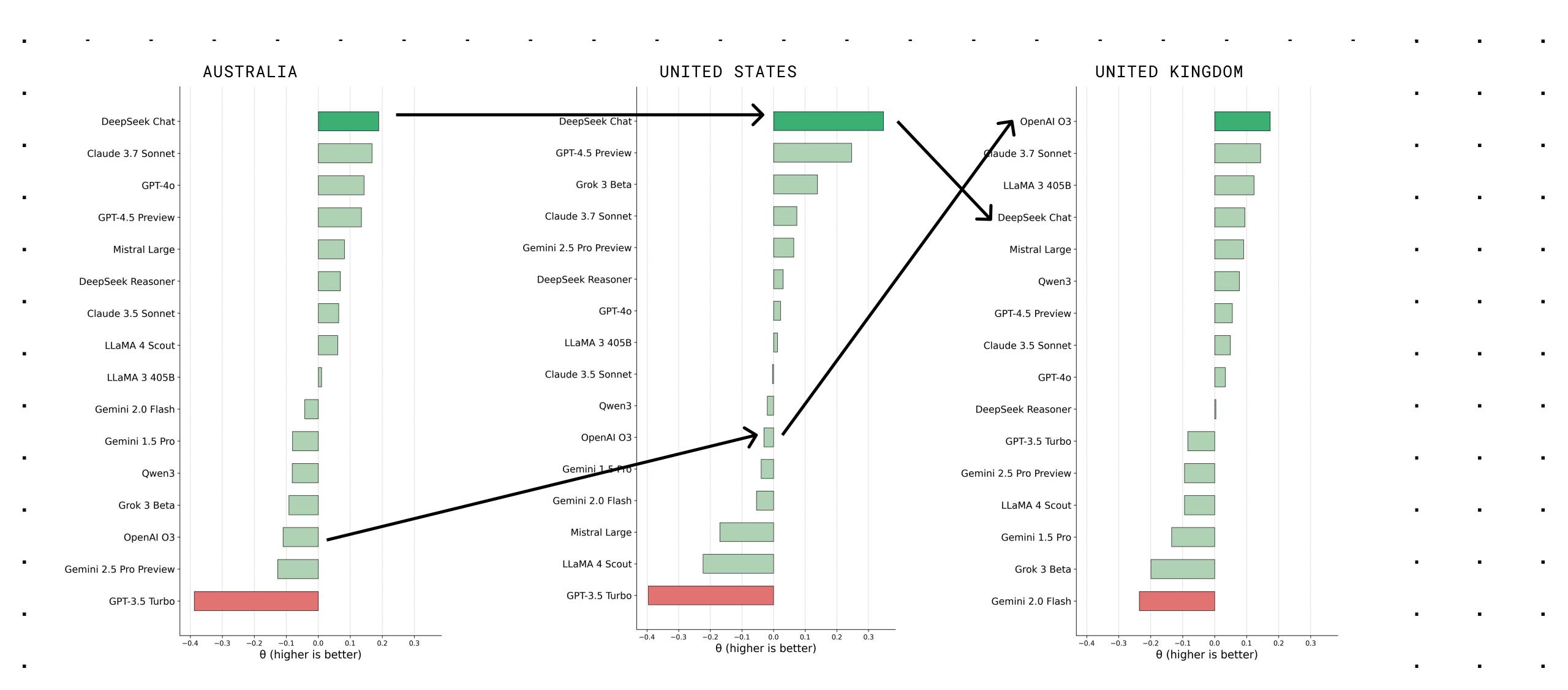
Highlights

- GPT-4 ranked in the top percentile for originality and fluency on the <u>Torrance Tests</u> of Creative Thinking.
- Flexibility for GPT-4 ranged from the 93rd to the 99th percentile.
- The creative abilities of AI, including the ability to generate original output, seem to now match human abilities for the first time.
- The impact on social and economic innovation, including how creativity is understood, will likely be significant.





With Cultural Nuances

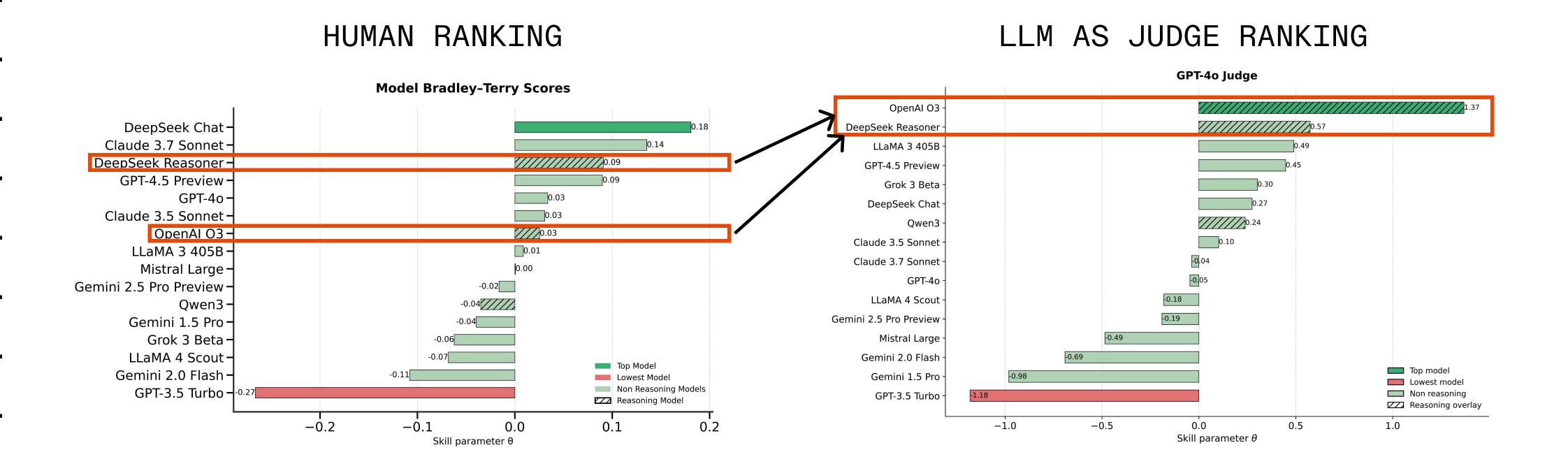


These Results Are Reassuringly Human

- Messy
- Uncertain
- Subjective
- Inconsistent

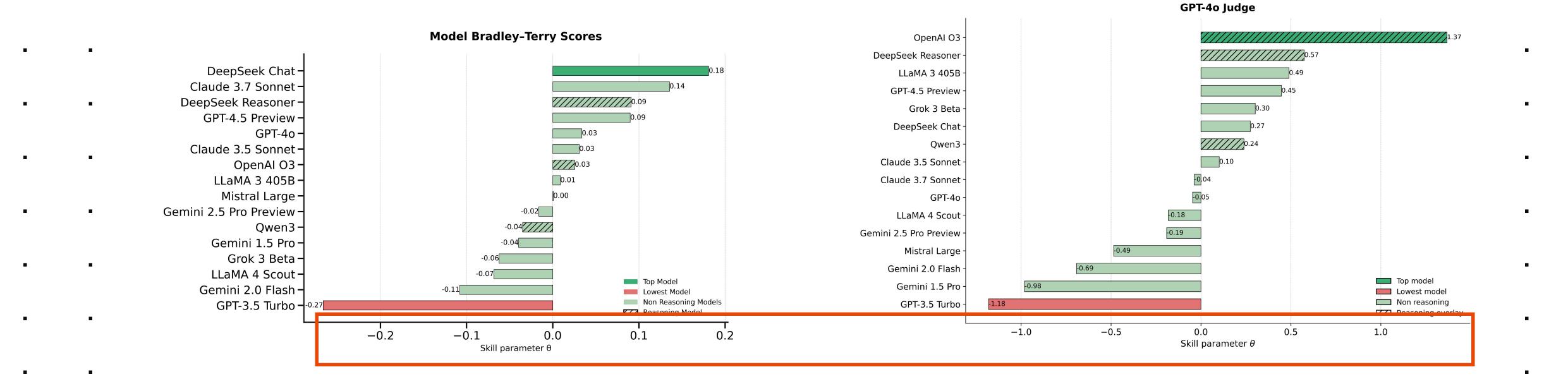


But when we compare this to the machines there is a strong bias towards 'reasoning models'



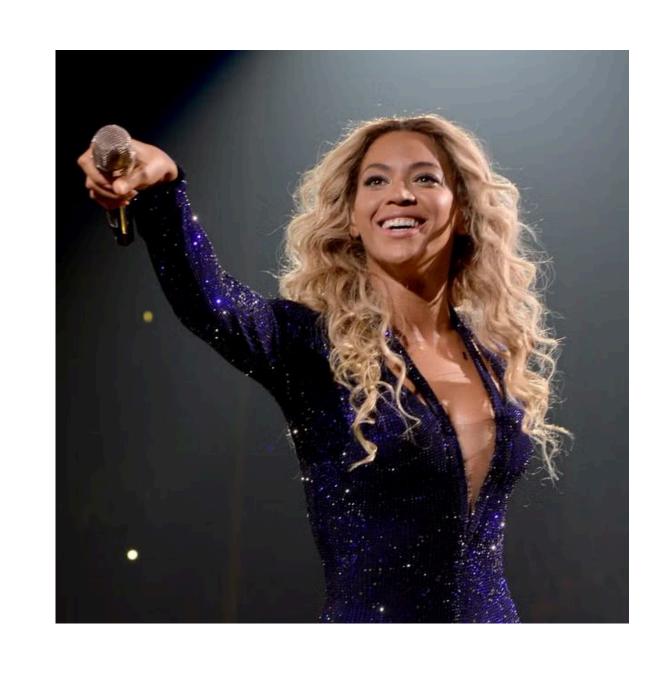


And an order of magnitude 5-10x more confident!





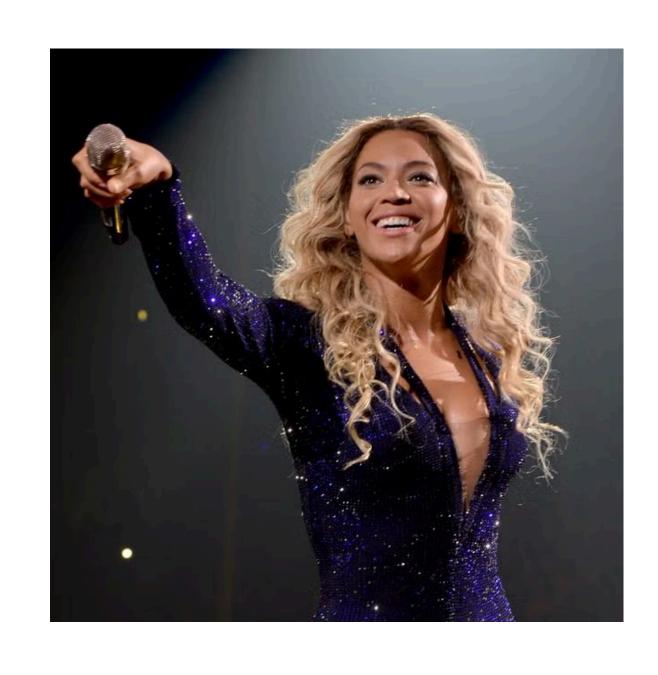
It's like asking people who they think would win an award between top artists.







But the machines vote with the certainty of a top contender beating a mid level artist.

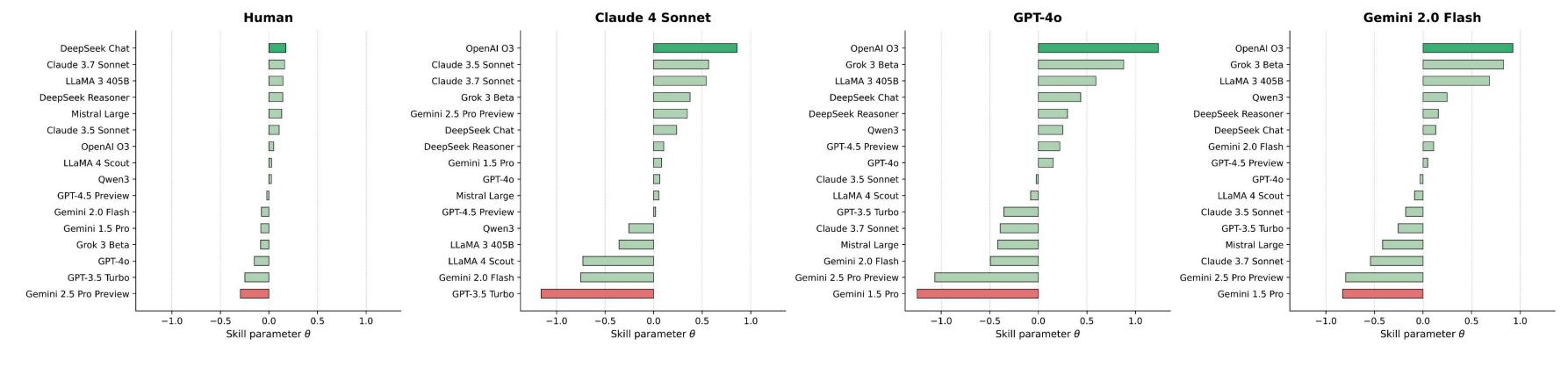




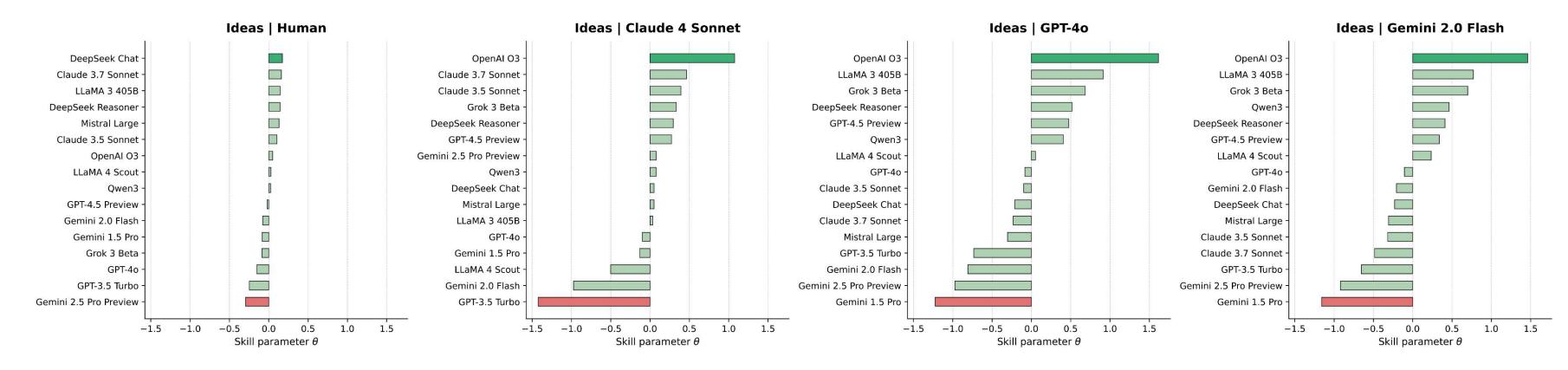
Bradley Terry Scores (EQ Bench Prompt)



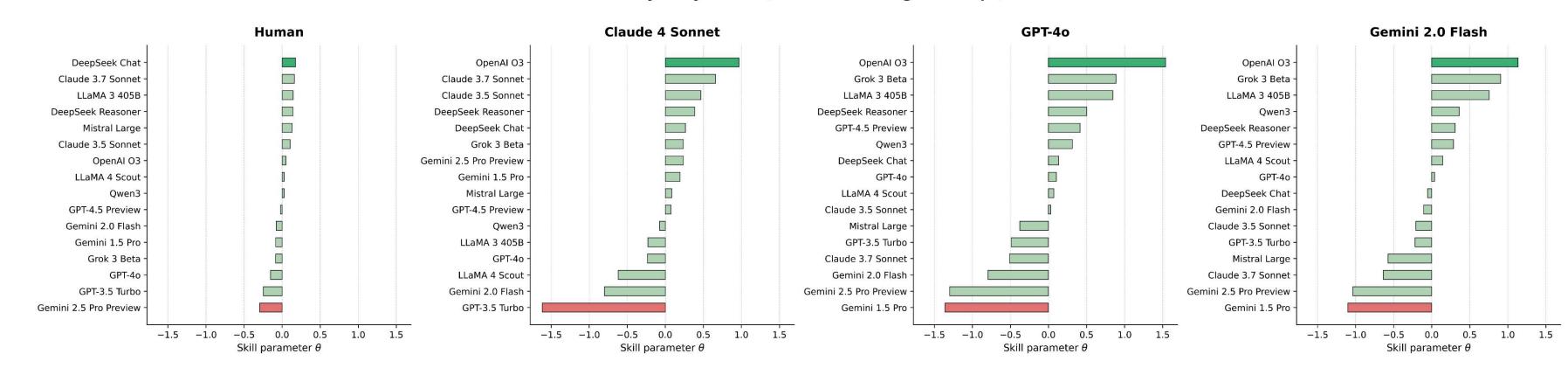
Results Repeated Themselves Over And Over **Again No Matter** The Model Or The LLM As Judge Setup



Bradley Terry Scores (Surprise Prompt)



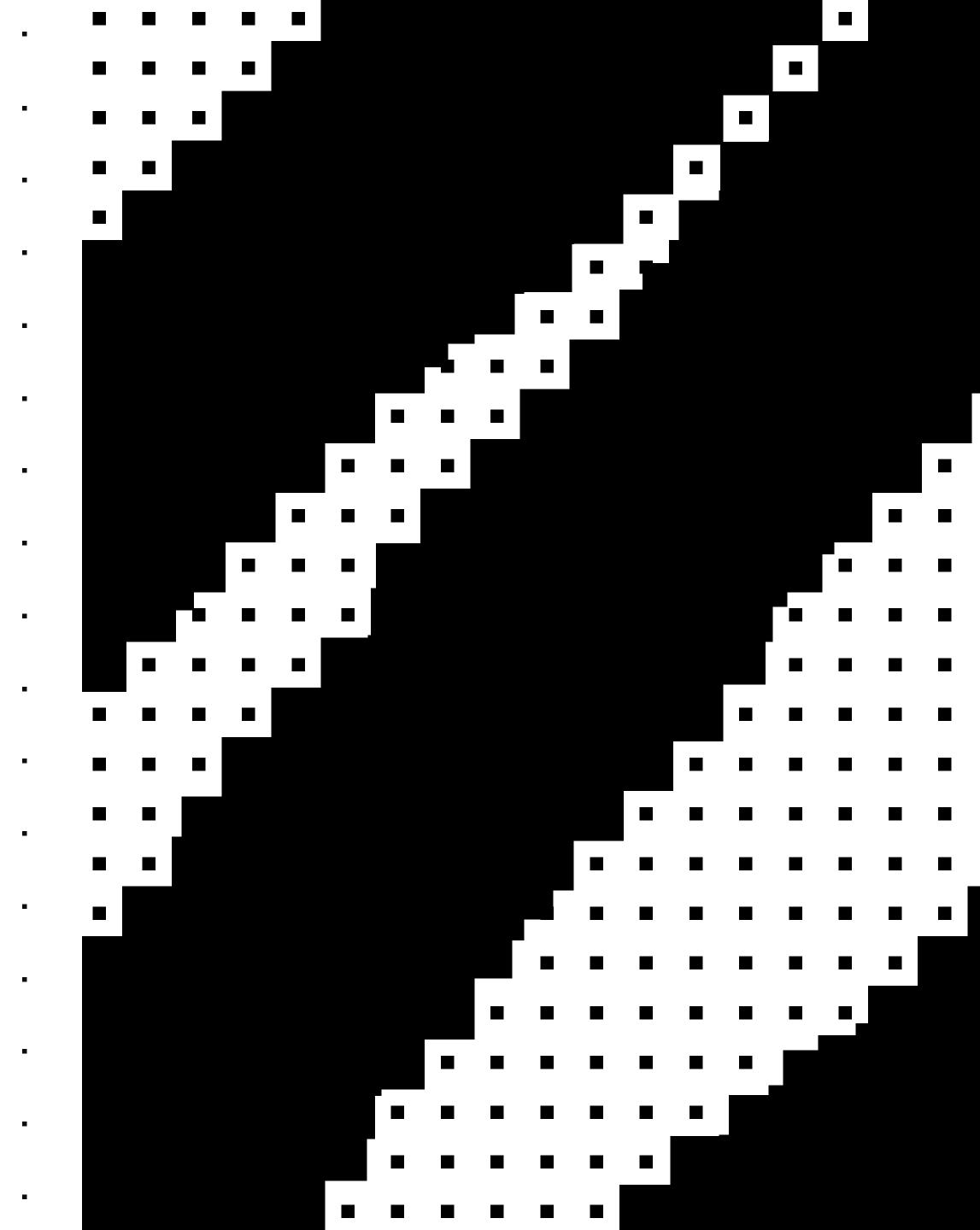
Bradley Terry Scores (Creative Stratergist Prompt)





Humans can't agree.

Machines can't fel.







Your Instinct Matters

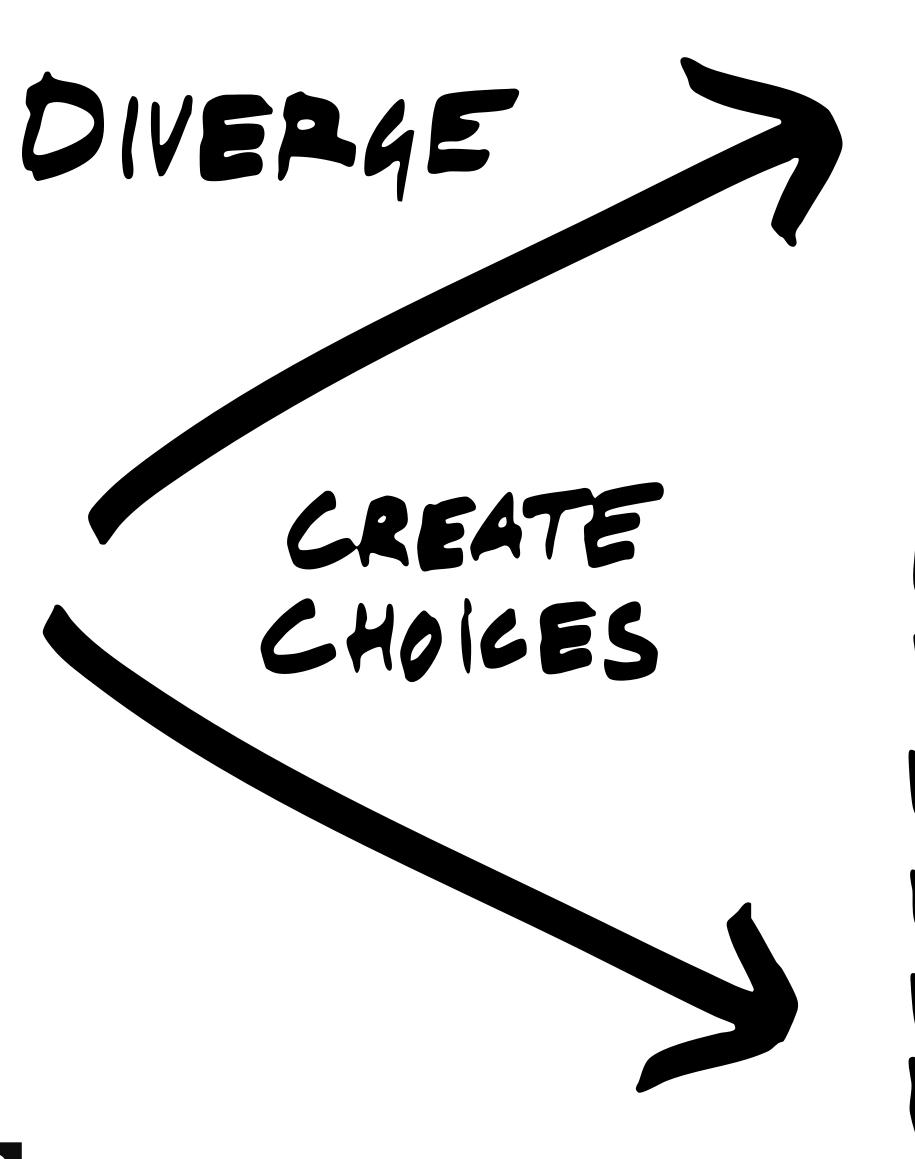
Speed isn't the enemy, undervaluing thinking is

Al has collapsed timelines but that doesn't mean the work is worth less. What's truly valuable hasn't changed: originality, cultural sensitivity, and the judgement to know what will land and why.

But the old pricing models are under pressure. If time is no longer the measure, agencies need to reassert the value of thinking over making. Strategic impact over production effort. The making may now happen faster but the ideas still demand clarity, courage, craft and care.

This is a pivotal moment: not just to defend value, but to redefine it. Agencies that lean into this shift – educating clients, experimenting with models, and pricing for outcomes, not outputs – will stay ahead. Those who don't may find themselves replaced by tools that look capable, but lack context, care or taste.



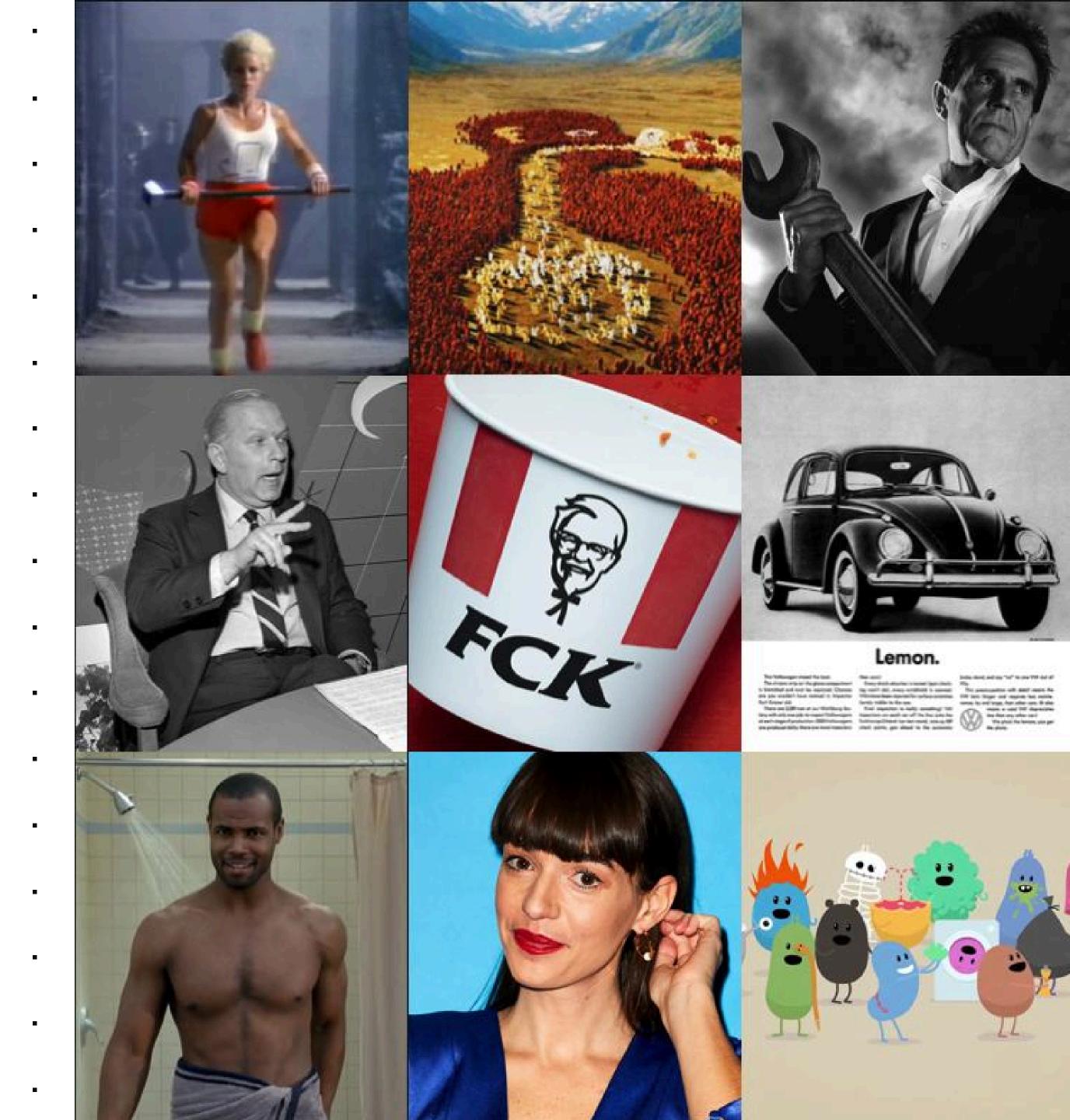


So, what about Diverge?



↓ DIVERGE

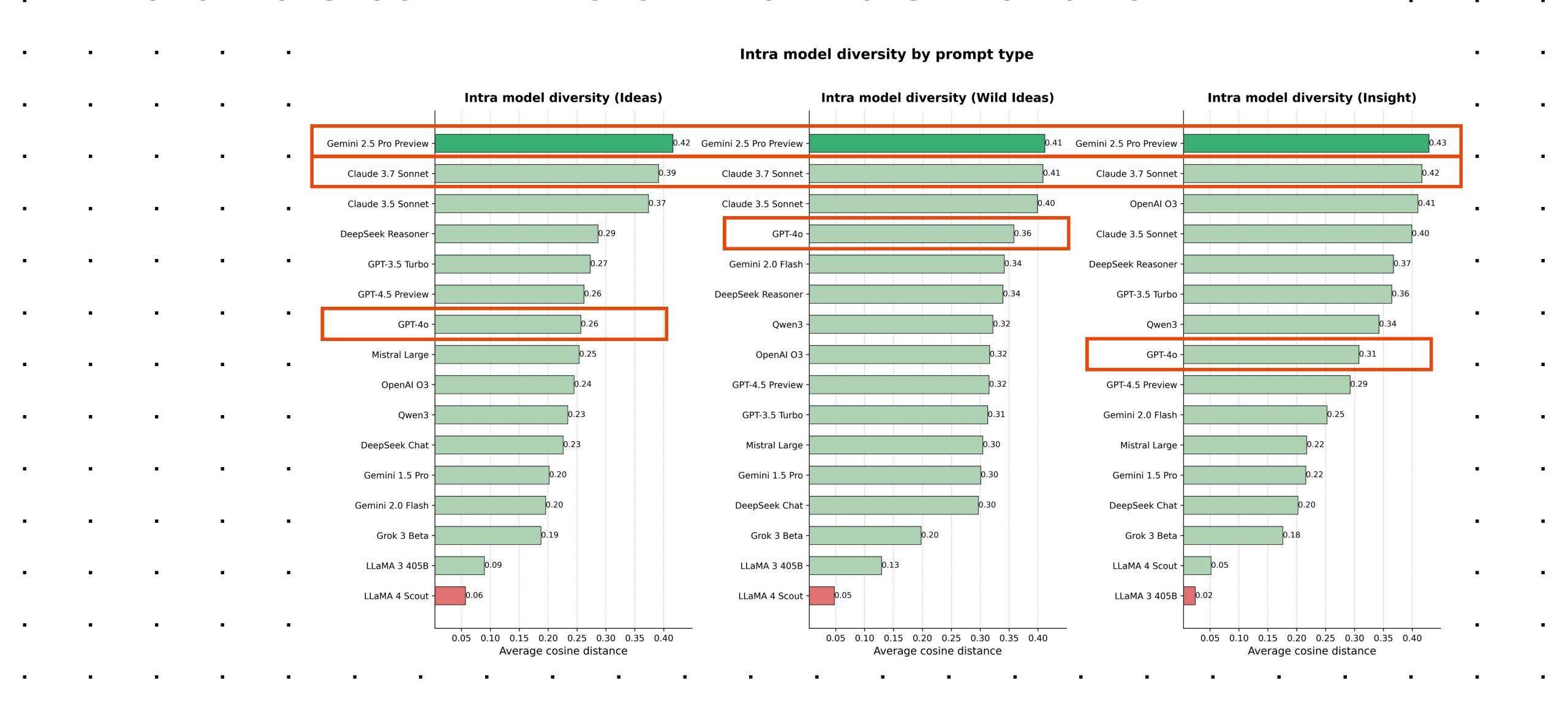
We've Already Seen That Al Can Come Up With Lots Of Ideas But How Different Are They?





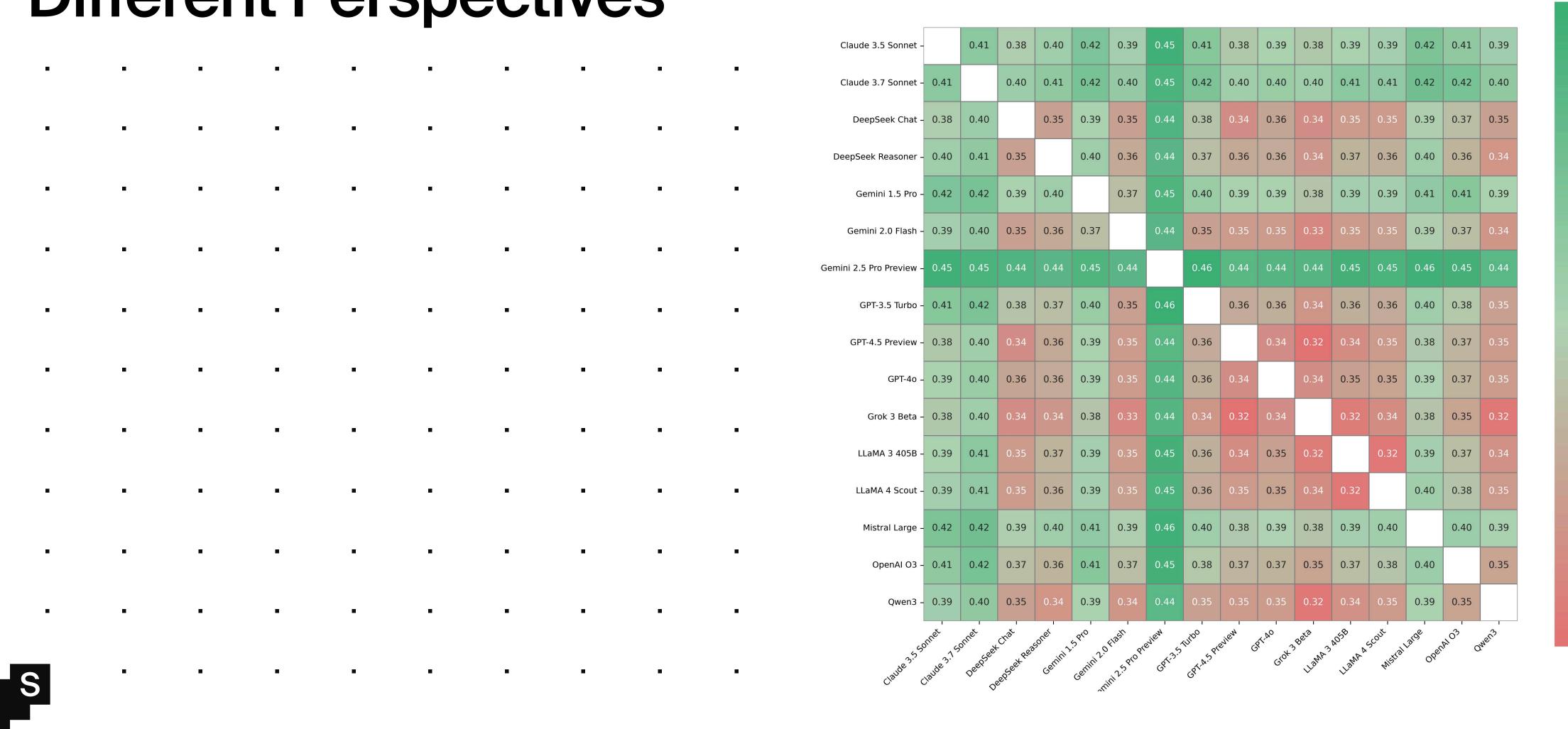


There Re Clear Winners In Terms Of Variation



↓ DIVERGE

The Real Magic Is Combining Models With Different Perspectives



- 0.44

- 0.42

- 0.36

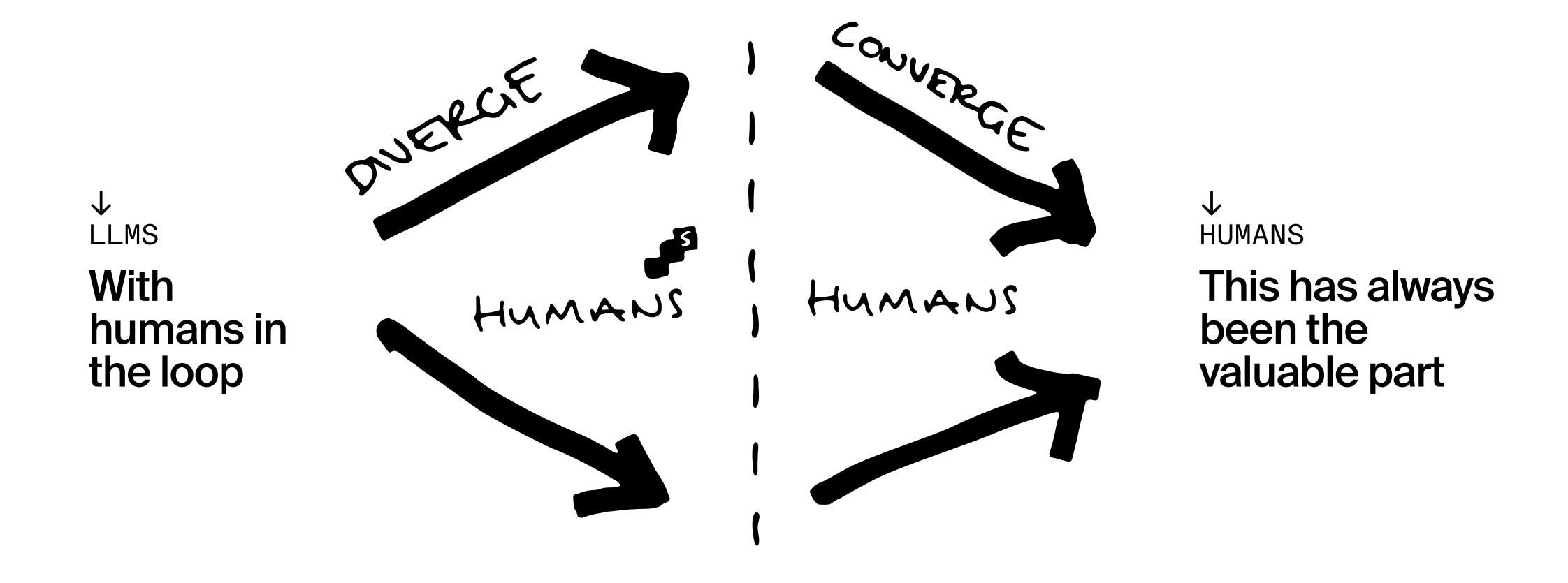
- 0.34

And so what are the takeaways



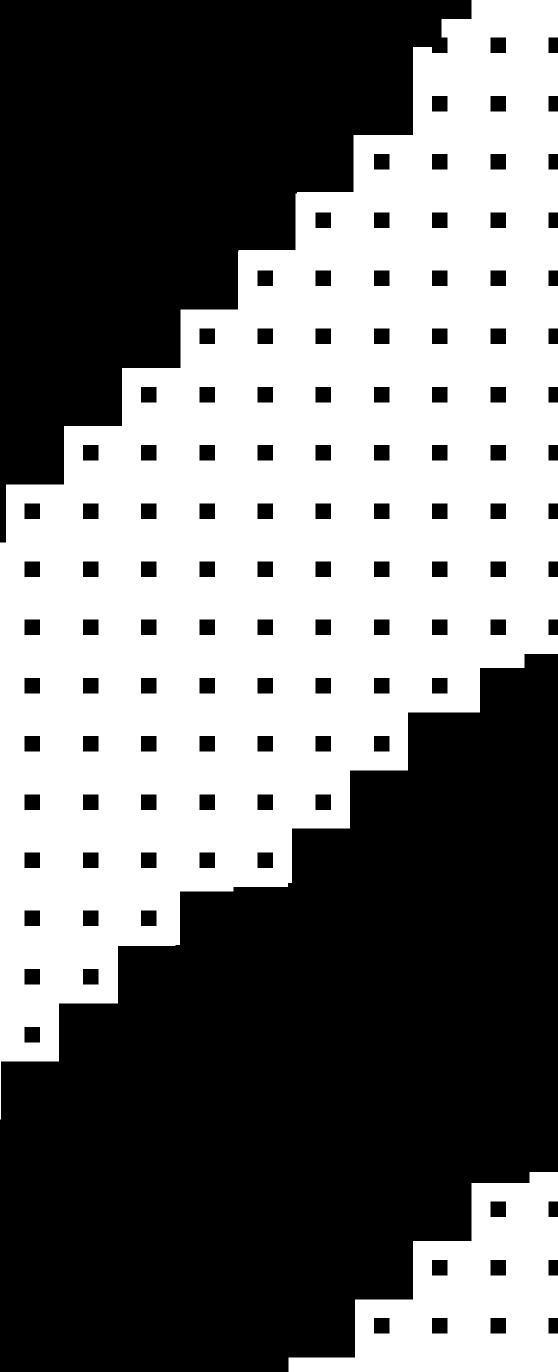


Partner with LLMs for volume, humans for selection





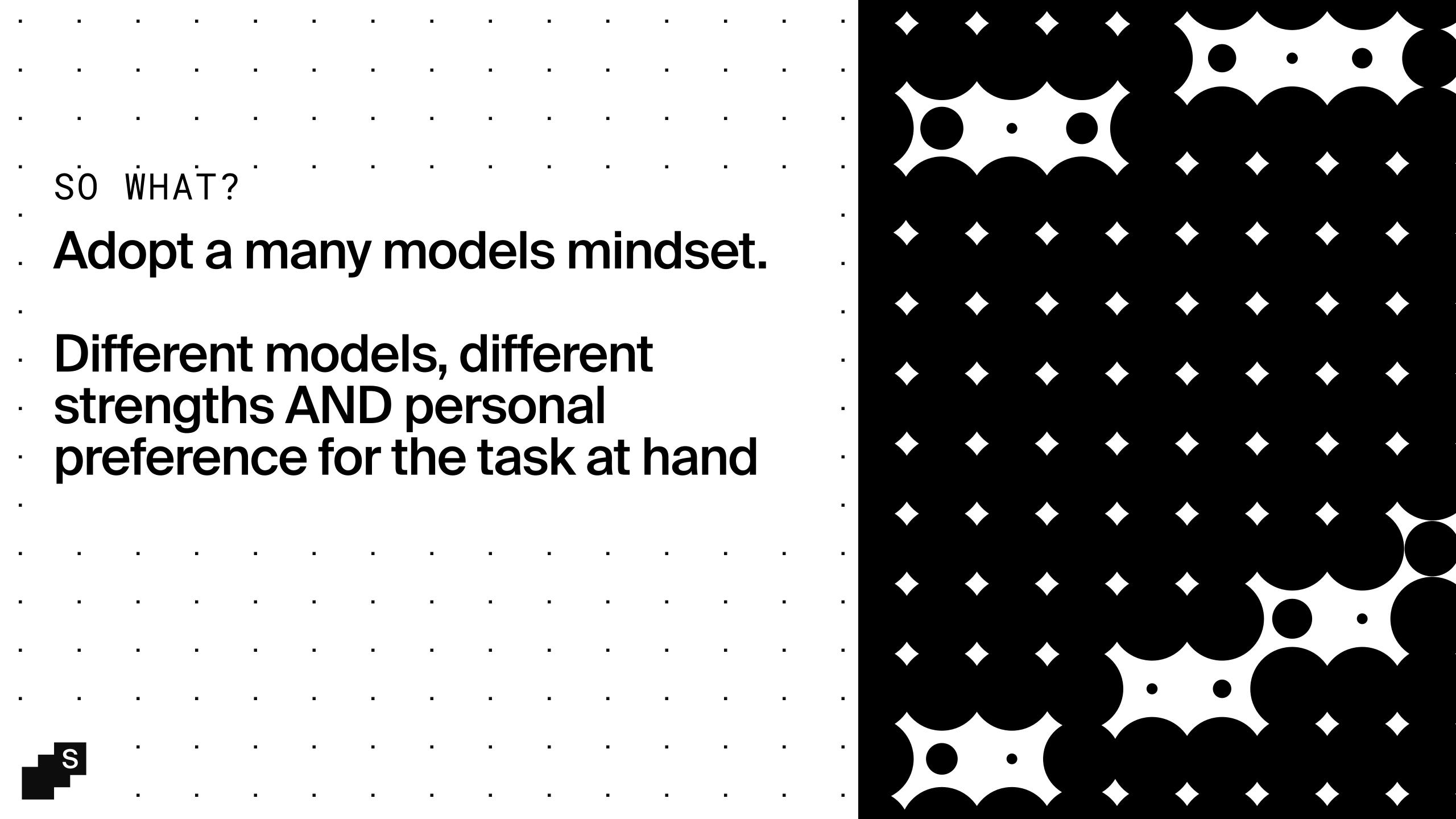
1/ Different AI tools win at different tasks 2/ Al can't judge creative work well 3/ Variety of ideas matters most 4/ Creative preferences vary by location



WHAT?

Adopt a many models mindset.

Different models, different strengths AND personal preference for the task at hand





SO WHAT?

When picking models to work with, consider:

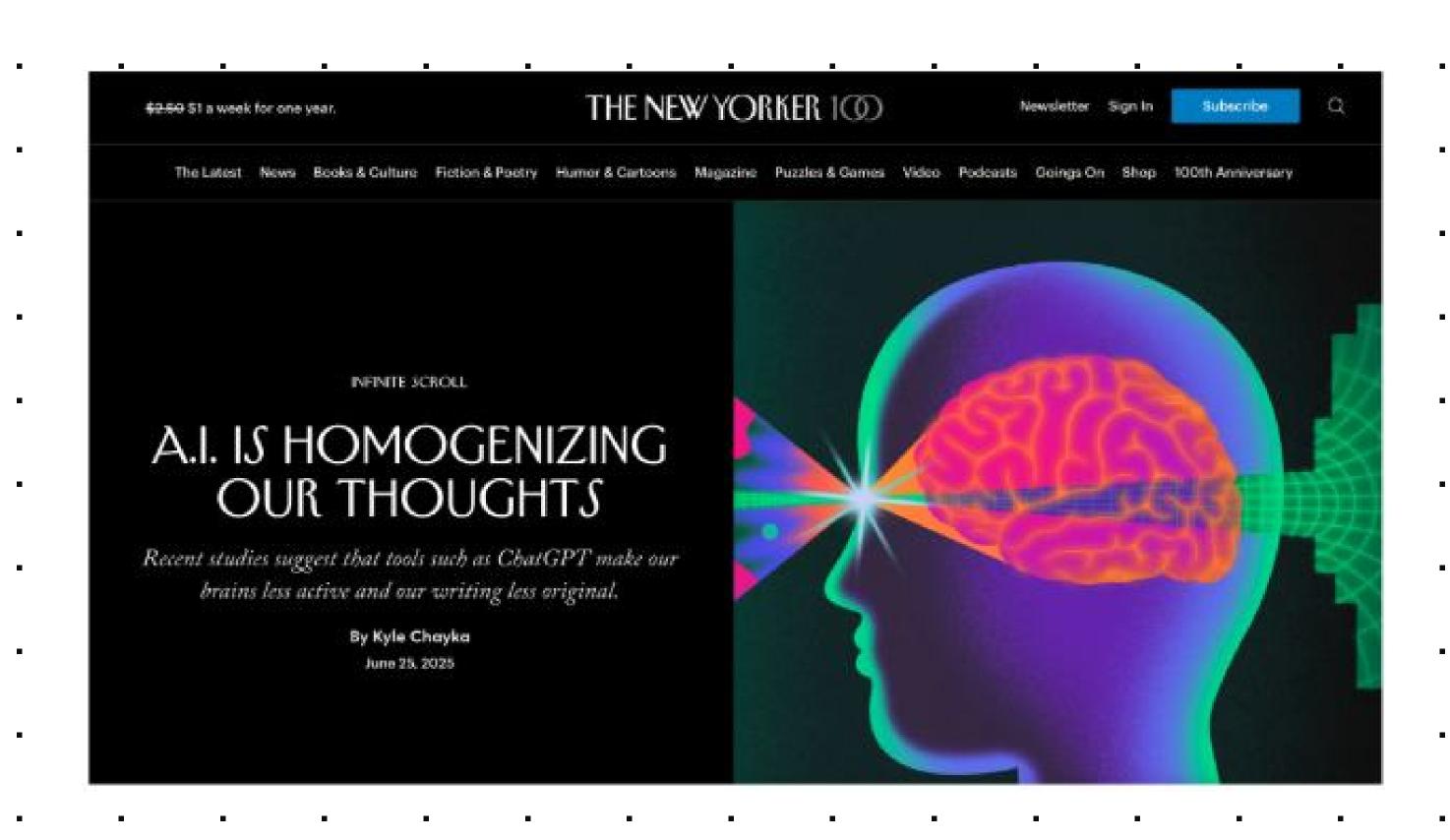
- 1/ The task
- 2/ Your own preferences
- 3/ If you want a range of ideas, don't default to what everyone else does





LLMs aren't designed to be creative...

When you have no divergent opinions... you get average everything everywhere all at once





"Creative work isn't only about precision. It's about discovery, exploration, and happy accidents that lead to breakthroughs.

The best creative tools should feel like collaborators, not calculators. They should surprise, challenge, and push us in new directions."

SOURCE: EVERY

Thank you to everyone involved:





















JAMES HURMAN

