# The Rapidly Growing ROI of Search Relevance
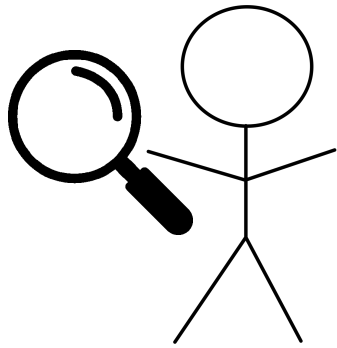
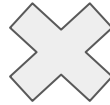## Old Value, New Multipliers

**"How do I pay for search?"**

**"Make it work like Google."**

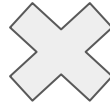**"Make it work like Deep Research."**

# The old ROI model

1 Person : 1 12-hour Task
3 Hours "Discovery"

20 Cycles per year

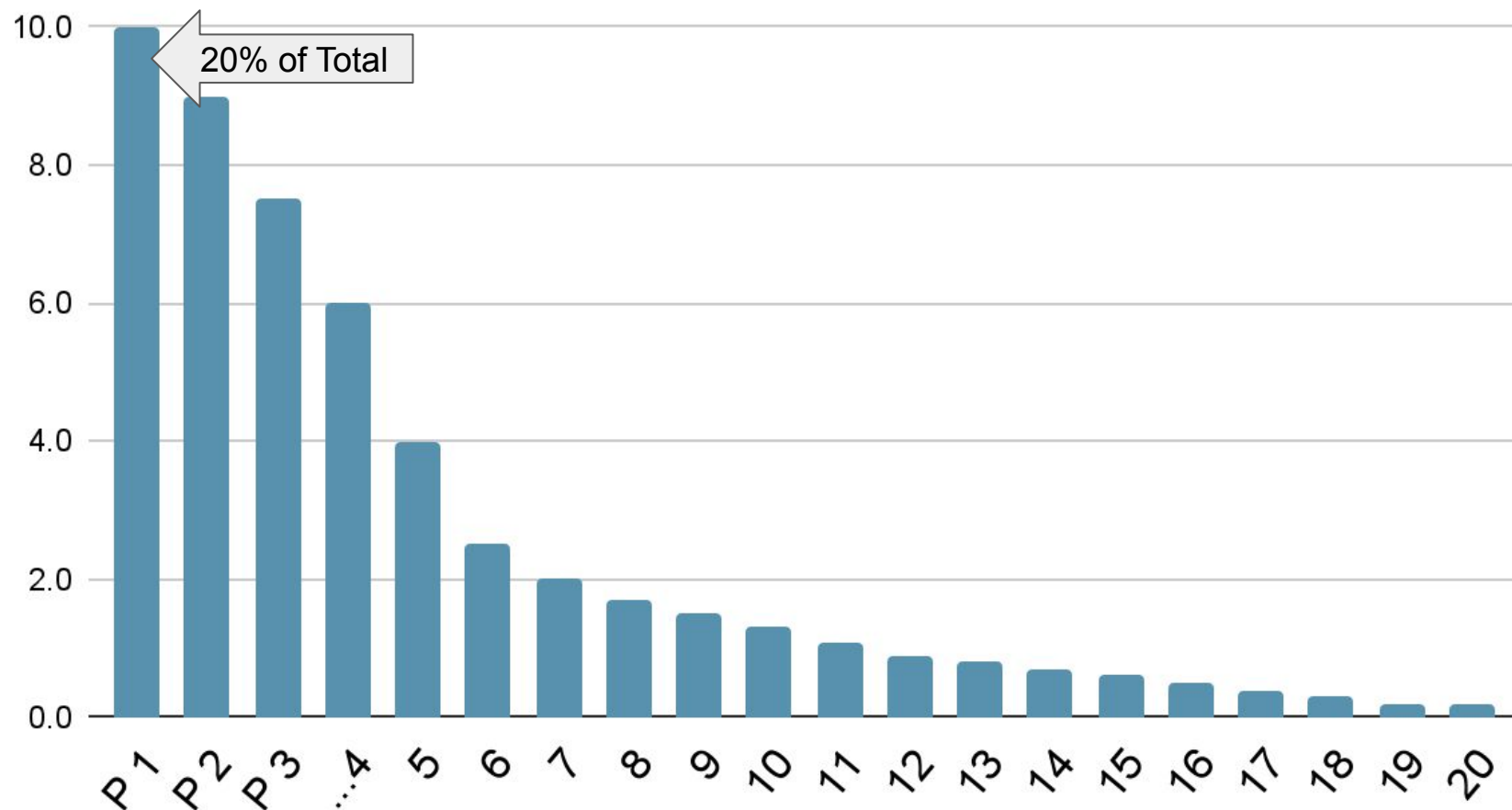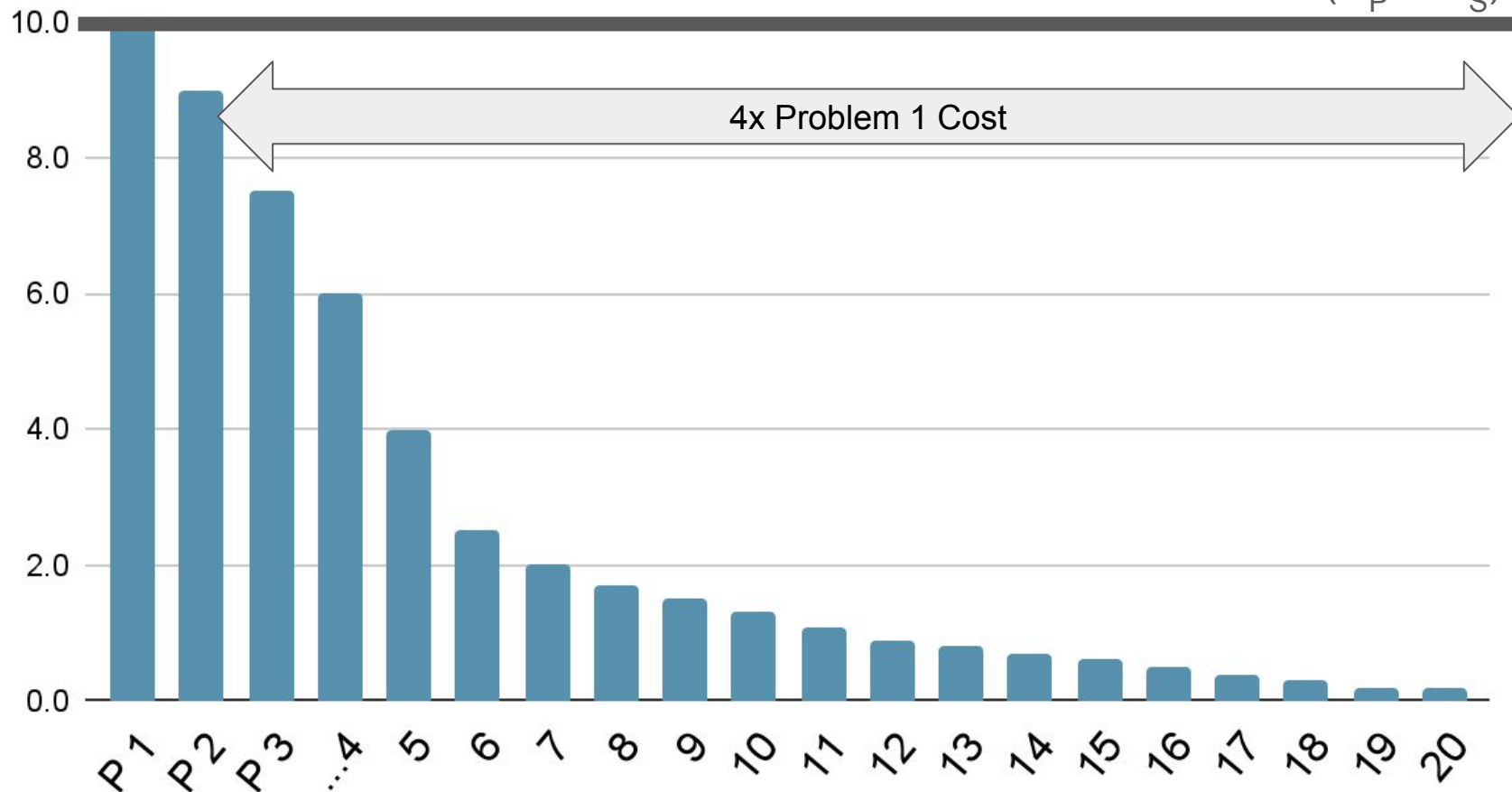25 People

1500 Hours per year

* As long as we find the right information
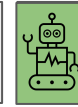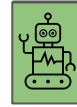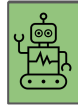
Cost of Problem (P)

20% of Total

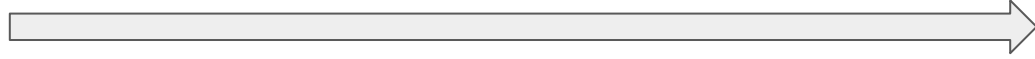Cost of Problem (P)

ROI Line ($C_P > C_S$)

4x Problem 1 Cost

P 1, P 2, P 3, ...4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20

P1 Activities

**OpenAI**
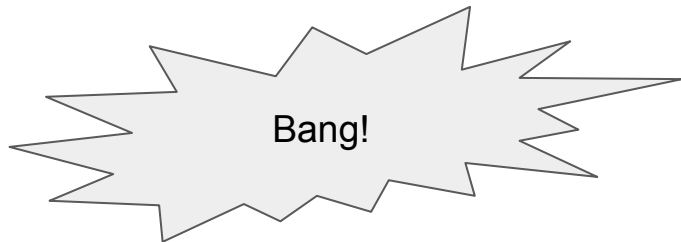
Milestone

# Better language models and their implications

Read paper ↗   View code ↗

"We've trained a large-scale unsupervised language model which generates coherent paragraphs of text, achieves state-of-the-art performance on many language modeling benchmarks, and performs rudimentary reading comprehension, machine translation, question answering, and summarization—all without task-specific training."

2018: **BERT**

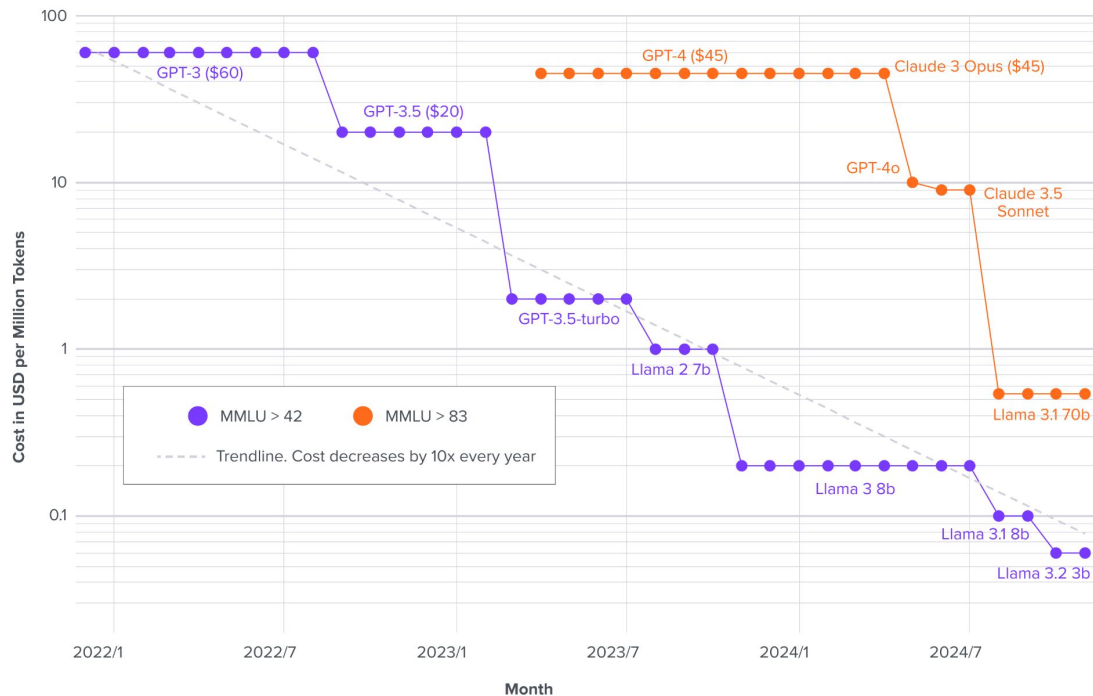T5, GPT2, Megatron-Turning, OPT, BLOOM, …

2022: **ChatGPT**

Bang!

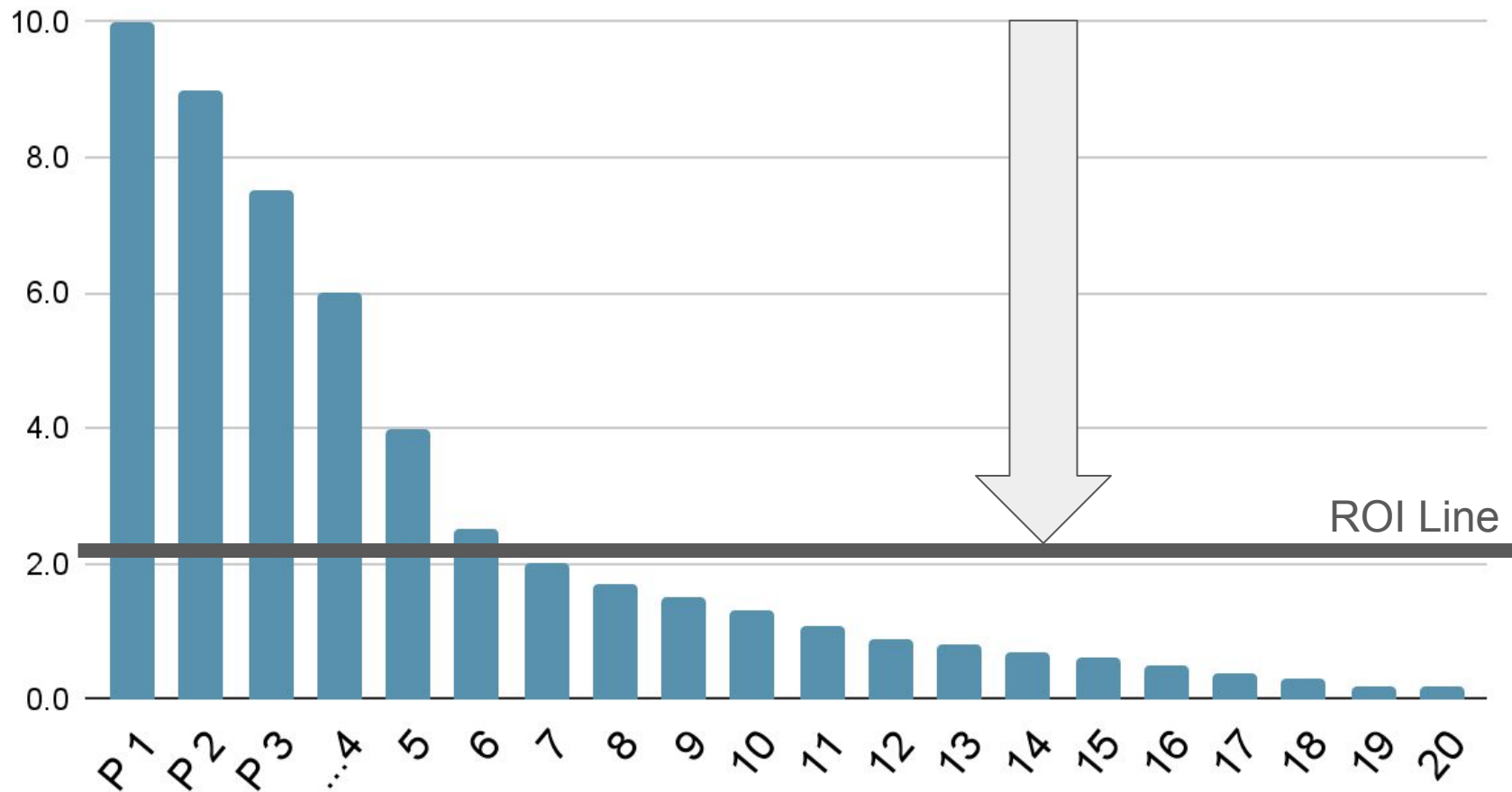# Welcome to LLMflation – LLM inference cost is going down fast ⬇️

Guido Appenzeller

https://a16z.com/llmflation-llm-inference-cost/

**Cost of the Cheapest LLM with a Minimum MMLU Score (Log Scale)**

GPT-3 ($60)

GPT-4 ($45)

Claude 3 Opus ($45)

GPT-3.5 ($20)

GPT-4o

Claude 3.5 Sonnet

GPT-3.5-turbo

Llama 2 7b

Cost in USD per Million Tokens

● MMLU > 42    ● MMLU > 83

--- Trendline. Cost decreases by 10x every year

Llama 3.1 70b

Llama 3 8b

Llama 3.1 8b

Llama 3.2 3b

Month

2022/1    2022/7    2023/1    2023/7    2024/1    2024/7

a16z Infrastructure
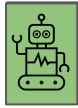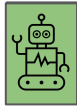
$$C_{LLM}\downarrow \Longrightarrow (C_{Problem} - C_{Solution})\uparrow$$
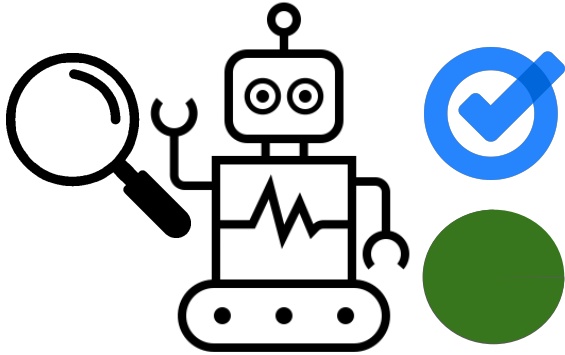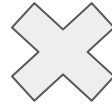
# Cost of Problem (P)



ROI Line

Process Activities

# The new ROI model



**1 Agent : n Tasks**

**n Cycles per year** ✕ = **n Hours per year**
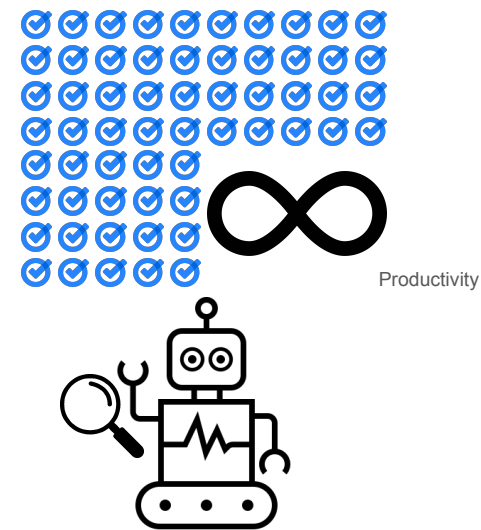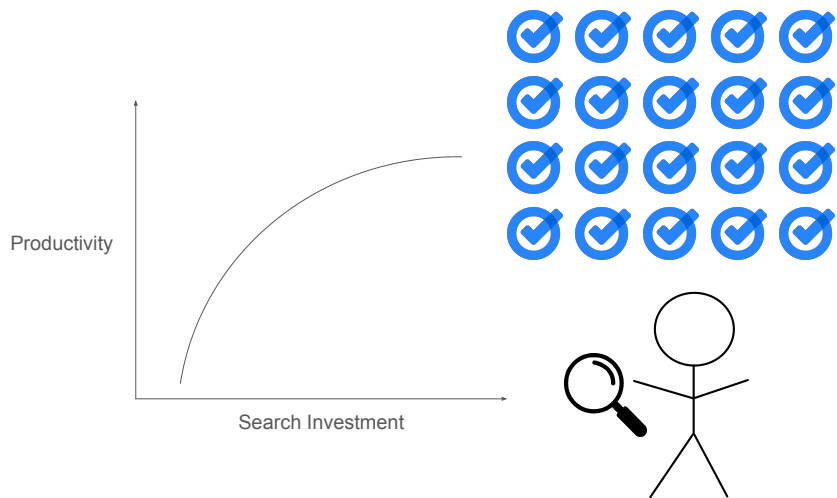
* As long as I find the right information

* As long as we find the right information