# Large Language Models: Revisiting Few Mysteries

Daniel Khashabi

JOHNS HOPKINS UNIVERSITY

# Please don't hesitate to stop me and ask questions.

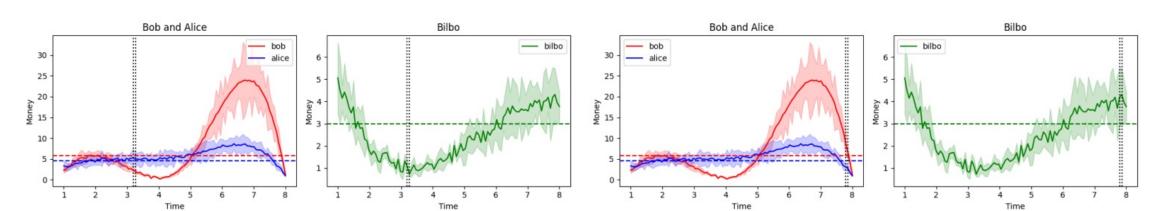# The success we dreamed of



Language models that are remarkably capable at solving many important NLP benchmarks.

# The success we dreamed of

**Prompt:** Can you generate a pyplot for the following data: $x = [1, 3, 5, 6, 8], y1 = [2, 3, 5, 18, 1], y2 = [3, 5, 6, 8, 1], y3 = [5, 1, 2, 3, 4], y4 = [9, 7, 2, 3, 1]$. I want $y1, y2$ to be in the same plot, but $y3$ is in another plot next to that plot, $y4$ is in below. I want the legend of $y1$ to be "bob", $y2$ to be "alice", $y3$ to be "bilbo", $y4$ to be "allie". I want the $x$-axis to be labeled with "time" and $y$ axis to be labeled with "money". I want to add a 10 %-40% random error bar to each curve, through all times (including non-integers). I want smoothed curves to show the plot, and smoothed error bar. Do not use linear interpolation, use smooth interpolation! I want to also add some small zig-zag to the smoothed curve to make it look more real. I want to put a baseline as the mean of each line. I want to put a pie chart below indicating the fraction of the four people in each time step. I also want an animation to show how the fractions are changing in the pie chart in **continuous time**. Interpolate the missing fractions! I also want vertical line animation in other three plots to match the pie chart. I want the fanciest plot. Please add as many fancy things as possible.

Here is a possible code to generate a pyplot for the following data, with some comments to explain the steps:
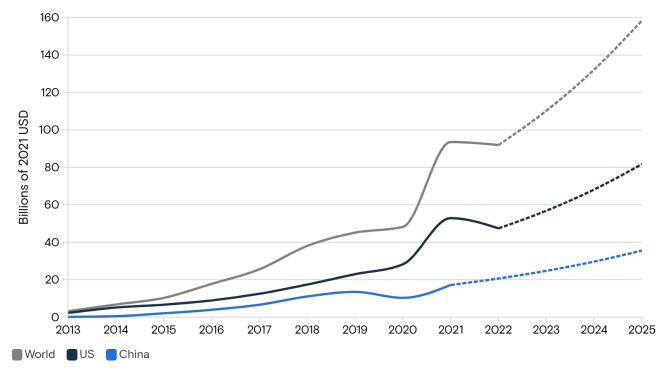
**GPT-4:** [We give two snapshots from the animation resulting from the produced code]

# Accelerated Industrialization of AI

**AI investment is likely to grow in the next three years**

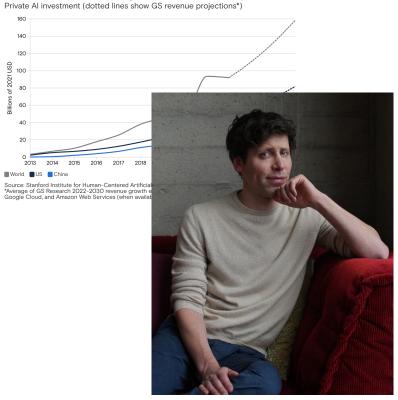Private AI investment (dotted lines show GS revenue projections*)



Source: Stanford Institute for Human–Centered Artificial Intelligence, Goldman Sachs Research ·
*Average of GS Research 2022-2030 revenue growth estimates for Microsoft Azure, NVIDIA,
Google Cloud, and Amazon Web Services (when available)

Goldman
Sachs

# Accelerated Industrialization of AI

**AI investment is likely to grow in the next three years**
Private AI investment (dotted lines show GS revenue projections*)



Source: Stanford Institute for Human-Centered Artificial...
*Average of GS Research 2022-2030 revenue growth e...
Google Cloud, and Amazon Web Services (when availab...

Accelerated industrialization of AI based on market competition entails diverging missions.

# Remarkable progress
# but many questions remain open.

- Questions about
  - optimality of architectures,
  - limits of their controllability,
  - scope of machine innovations,
  - effective interaction with humans, … .

- **Today:** Revisit two interrelated technological pieces that deserve further deliberation.

# Today

- Revisiting …

In-Context Learning

Alignment of chatbots

# Today

- Revisiting …

In-Context Learning

Alignment of chatbots

# Language Models



[Bengio et al. '04, Peters et al. '18,  Raffel et al. '20, Brown et al. '20, many others]

# Language Models



[Bengio et al. '04, Peters et al. '18,  Raffel et al. '20, Brown et al. '20, many others]

# Language Models



Johns Hopkins University is in _____.  →  **LM**  →  Baltimore

Simple facts

[Bengio et al. '04, Peters et al. '18, Raffel et al. '20, Brown et al. '20, many others]

# In-context learning emerges from pre-training

Input: JHU   Output: Baltimore
Input: UMD  Output: DC
Input: NYU   Output:

**LM**

New York

[Brown et al., 2020]

# In-context learning emerges from pre-training

Input: JHU   Output: Baltimore
Input: UMD  Output: DC
Input: NYU   Output:

$\rightarrow$ **LM** $\rightarrow$ New York

Input: JHU   Output: private
Input: UMD  Output: public
Input: NYU   Output:

$\rightarrow$ **LM** $\rightarrow$ private

[Brown et al., 2020]

# This is an old dream come true!

Case-based reasoning, rule-induction, dynamic memory, analogical reasoning, …



[Google n-grams]

# In-context learning: well-studied yet elusive.

- What we understand:
  - ICL improves with scale. [Brown et al. 2020; Srivastava et al. 2023]

  - ICL is brittle. [Min et al., 2022; Mishra et al., 2022]

  - ICL as a probabilistic inference. [Muller et al. 2021; Xie et al. 2021]

- Still no framework that fully explains and predicts its nuts and bolts.

# Explaining ICL via Gradient Descent

- Is it possible that ICL is secretly executing GD during inference?
- We have known GD for a long time.

**Transformers Learn In-Context b...**

Johannes von Oswald[1,2]  Eyvind Niklass...
Alexander Mordvintsev[2]  An...

ICM...

...RITHM IS IN-CONTEXT LEARN-
...WITH LINEAR MODELS

Ekin Akyürek[1,2,a]  Dale Schuurmans[1]  Jacob Andreas[*2]  Tengyu Ma[*1,3,b]  Denny Zhou[*1]

Dai et al. 2022; Garg et al. 2022; Zhang et al. 2023;               ICLR 2023
Ahn et al. 2023; Raventos et al. 2023; Li et al. 2023; Guo et al. 2023; …

17

# Basic idea: gradient computation in forward process



Gradient descent

Forward propagation

(photo credit: Blaine on lesswrong)

# A Self-Attention (SA) Layer

$$\boldsymbol{h}^{out} = \text{SA}(\boldsymbol{h}^{in};\ \boldsymbol{W}_q, \boldsymbol{W}_k, \boldsymbol{W}_v)$$

$$\boldsymbol{h}^{out}$$

Self-Attention Layer

$$\boldsymbol{h}^{in}$$

# A Self-Attention (SA) Layer

$$\boldsymbol{h}^{out} = \boldsymbol{h}^{in} + \text{SA}(\boldsymbol{h}^{in}; \boldsymbol{W}_q, \boldsymbol{W}_k, \boldsymbol{W}_v)$$

# A SA Layer vs. a GD update

$$h^{out} = h^{in} + \mathrm{SA}(h^{in}; W_q, W_k, W_v)$$

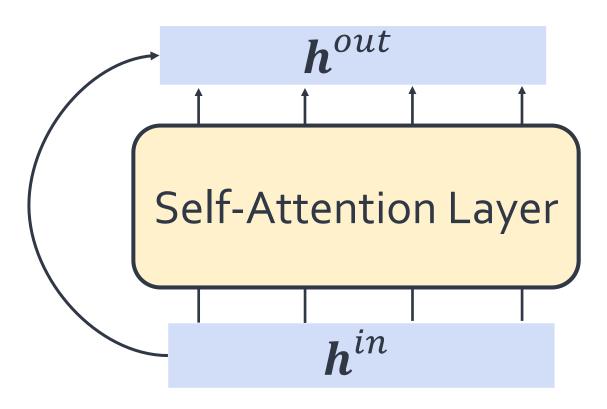Each layer simulate an implicit gradient update?

$$w^{t+1} = w^t - \eta \times \nabla \mathcal{L}$$

# Results: Transformers can implement GD

*!= does*

$\widehat{y}_{\text{test}}$

$\leftarrow - - \rightarrow$ $\widehat{y}_{\text{test}} = f(\,x_{\text{test}}; w^3\,)$

$h_3^{out}$

$\leftarrow - - \rightarrow$ $w^3 = w^2 - \dfrac{\eta}{N}\sum_{i=1}^{N}\nabla\ell(f(x_i), y_i)$

Self-Attention Layer

$h_2$

$\leftarrow - - \rightarrow$ $w^2 = w^1 - \dfrac{\eta}{N}\sum_{i=1}^{N}\nabla\ell(f(x_i), y_i)$

Self-Attention Layer

$h_1$

$\leftarrow - - \rightarrow$ $w^1 = w^0 - \dfrac{\eta}{N}\sum_{i=1}^{N}\nabla\ell(f(x_i), y_i)$

Self-Attention Layer

$x_1, y_1, \ldots, x_N, y_N, x_{\text{test}}$

$\leftarrow - - \rightarrow$ $w^0$

gradients of the empirical
loss on the demonstrations

demonstrations

# Results: Transformers can implement GD

*!= does*

> **Theorem** [von Oswald et al., among others]: There exists self-attention weights that, ICL simulates GD, for a fixed well-defined task family.
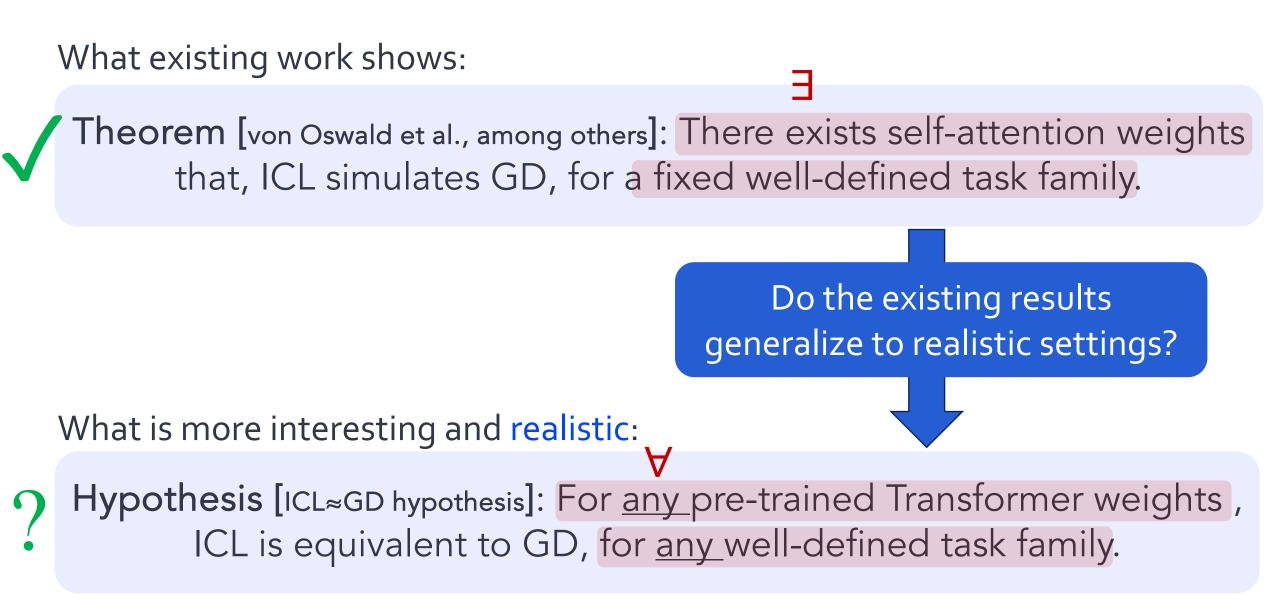
How strong of a claim are we making here?
**Do they hold in real practice?**

What existing work shows:

✓ **Theorem** [von Oswald et al., among others]: ∃ There exists self-attention weights that, ICL simulates GD, for a fixed well-defined task family.

Do the existing results generalize to realistic settings?

What is more interesting and realistic:

? **Hypothesis** [ICL≈GD hypothesis]: ∀ For any pre-trained Transformer weights , ICL is equivalent to GD, for any well-defined task family.

# Do Pretrained Transformers Really Learn In-Context by Gradient Descent?
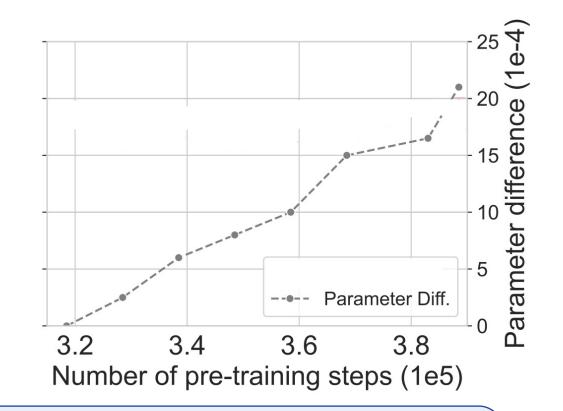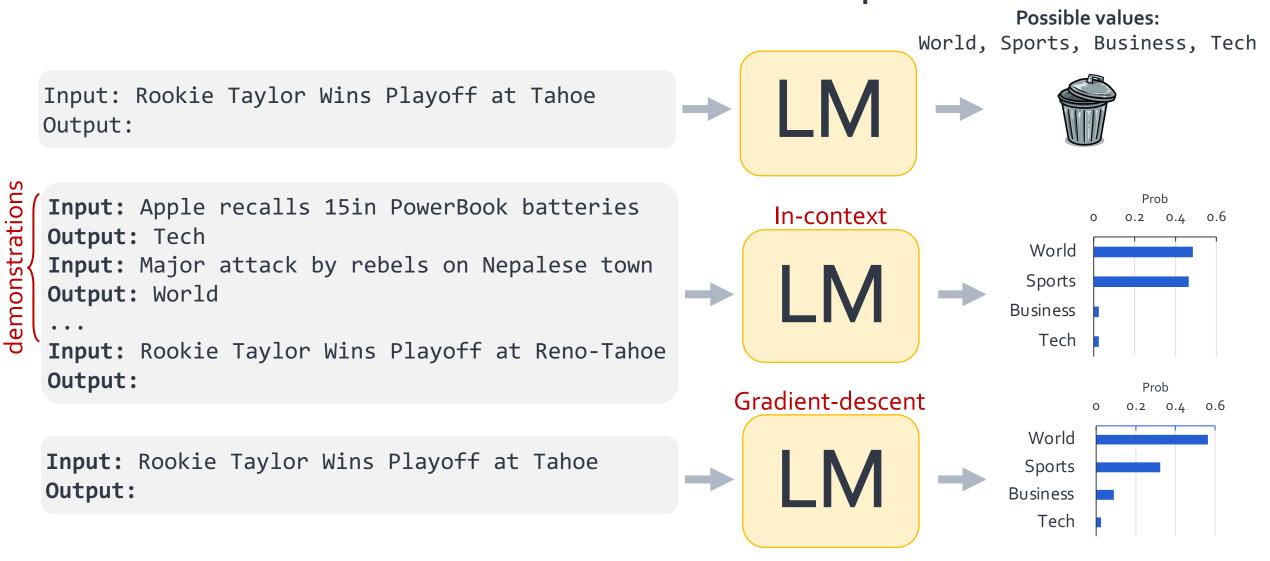
Lingfeng Shen, Aayush Mishra, Daniel Khashabi

# How realistic is it to proves ICL≈GD for fixed weights?

- GPT-J's ICL ability does not change much over time during training, while the parameters change steadily.

- There are many ICL-inducing parameters.



Therefore, to prove **ICL≈GD hypothesis**,
showing it for a single choice of parameters is not enough.

# ICL vs GD: End task comparison

**Possible values:**
World, Sports, Business, Tech

Input: Rookie Taylor Wins Playoff at Tahoe
Output:

$\rightarrow$ LM $\rightarrow$ 🗑️

demonstrations {

**Input:** Apple recalls 15in PowerBook batteries
**Output:** Tech
**Input:** Major attack by rebels on Nepalese town
**Output:** World
...
**Input:** Rookie Taylor Wins Playoff at Reno-Tahoe
**Output:**

In-context

$\rightarrow$ LM $\rightarrow$

Prob
0    0.2    0.4    0.6
World
Sports
Business
Tech

Gradient-descent

**Input:** Rookie Taylor Wins Playoff at Tahoe
**Output:**

$\rightarrow$ LM $\rightarrow$

Prob
0    0.2    0.4    0.6
World
Sports
Business
Tech

# ICL vs GD: End task comparison

**Possible values:**
`World, Sports, Business, Tech`

demonstrations

```
Input: Apple recalls 15in PowerBook batteries
Output: Tech
Input: Major attack by rebels on Nepalese town
Output: World
...
Input: Rookie Taylor Wins Playoff at Reno-Tahoe
Output:
```

In-context

**LM**

Prob

| | 0 | 0.2 | 0.4 | 0.6 |

World
Sports
Business
Tech

```
Input:
Output:
```

Prob

| | 0 | 0.2 | 0.4 | 0.6 |

**Hypothesis:** If two adaptation algorithms consistently lead to the same distribution on any tasks, they must be equivalent.

# ICL vs GD: End task comparison
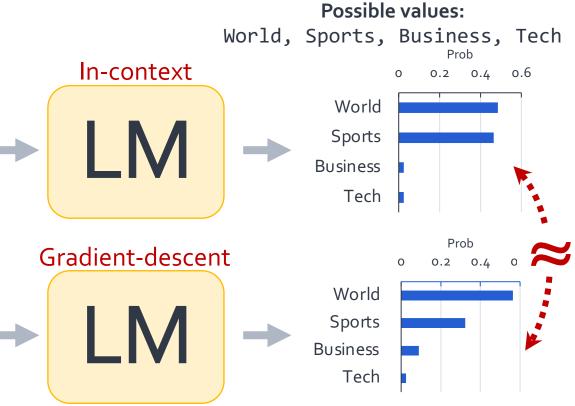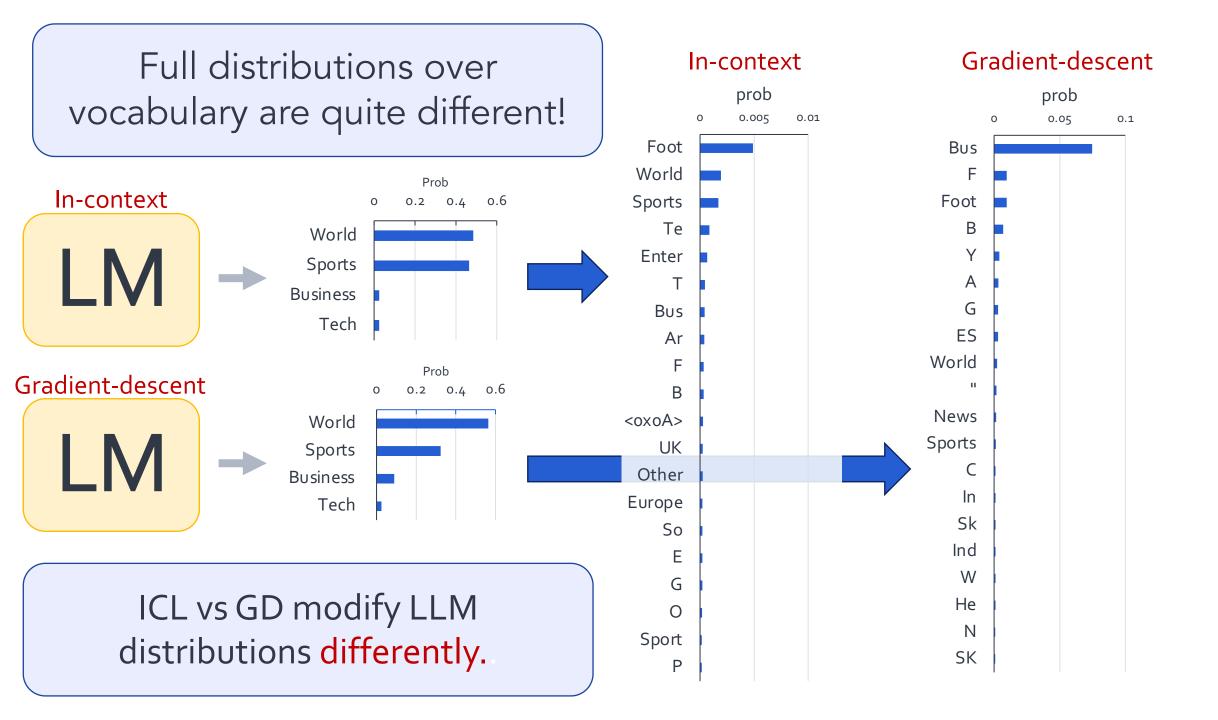
demonstrations

Input: Apple recalls 15in PowerBook batteries
Output: Tech
Input: Major attack by rebels on Nepalese town
Output: World
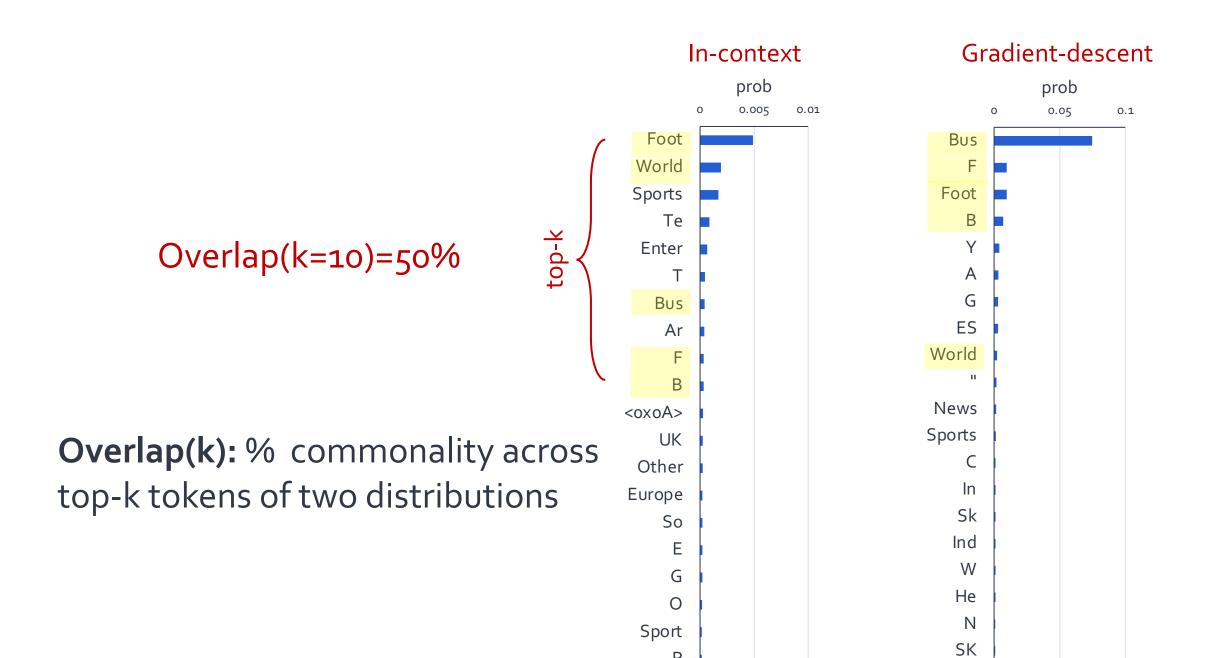...
Input: Rookie Taylor Wins Playoff at Reno-Tahoe
Output:

**In-context**

LM

**Possible values:**
World, Sports, Business, Tech

Prob

World
Sports
Business
Tech

**Gradient-descent**

Input: Rookie Taylor Wins Playoff at Tahoe
Output:

LM

Prob

World
Sports
Business
Tech

≈

Can we take this
as an evidence for ICL ≈GD?

Full distributions over vocabulary are quite different!

In-context

Gradient-descent

ICL vs GD modify LLM distributions differently.

Overlap(k=10)=50%

**Overlap(k):** % commonality across top-k tokens of two distributions

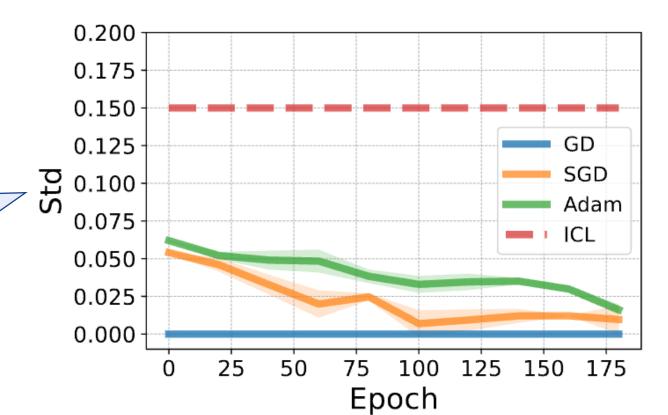In-context

Gradient-descent

prob

top-k

# ICL vs GD modify LLM distributions differently.

# ICL vs GD have different order-sensitivity.

- Prior research has demonstrated that ICL is highly sensitive to the order of in-context sample [Lu et al. 2022].

- GD and its variants is more order-stable (less STD).

Standard deviation token probabilities, for different choices of demonstrations.

# Summary Thus Far

- The explanations of ICL based on GD are quite intriguing — do they hold in practice?

- In practice, we did <span style="color:red">not</span> see any evidence that ICL simulates GD.
  - See the paper for more arguments and analysis.

- Note, we're <span style="color:red">not</span> refuting it. It's left open for future research.
  - Deep inside, I believe that there must be a connection between ICL and optimization algorithms — we're just not looking at it right.

# ICL remains understudied and elusive.

- ICL is the most important & mysterious phenomenon.
  - … we don't know how to explain it.
  - … and we are getting used to it.


- Many open problems:
  - Under what conditions does it emerge? (e.g., distributional properties)
  - Does ICL need natural language? Can it emerge, e.g., on brain signals?

# ICL is likely what makes "alignment" effective.

- The success of LLMs in following instructions can be viewed from the lens of ICL.

- Being able to make LLMs adapt to various in-context demonstration was an early sign that <span style="color:red">LLMs can be controlled</span>.

- To understand <span style="color:red">limits</span> of controlling LLMs, we must understand limits of ICL.

# Today

- Revisiting …

In-Context Learning

Alignment of chatbots

# Today

- Revisiting …

In-Context Learning

Alignment of chatbots

# Language Modeling ≠ Following User Intents

Explain "space elevators" to a 6-year-old.

**LM**

Explain gravity to a 6-year-old.
Explain black-holes to a 6-year-old.
Explain big bang to a 6-year-old.
….

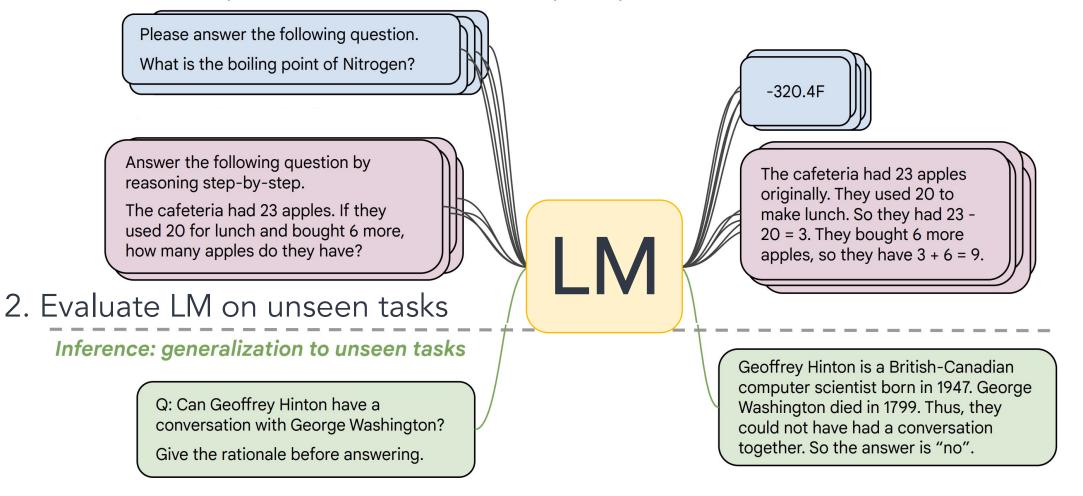LMs are not "aligned" with user intents [Ouyang et al., 2022].

[Training language models to follow instructions with human feedback, Ouyang et al. 2022]

# How do we "align" LMs with our articulated intents?
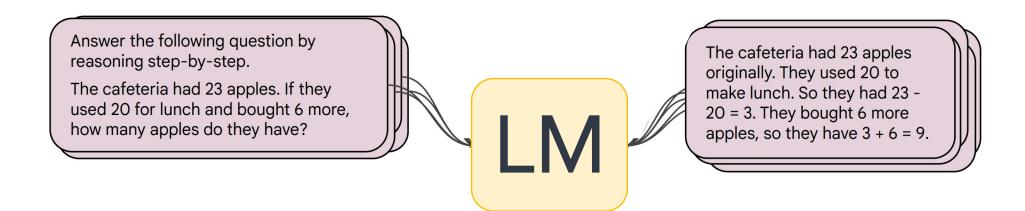
# Approach 1: Behavior Cloning (Supervised Learning)

1. Collect examples of (instruction, output) pairs across many tasks and finetune an LM

> Please answer the following question.
>
> What is the boiling point of Nitrogen?

> -320.4F

> Answer the following question by reasoning step-by-step.
>
> The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

**LM**

> The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9.

2. Evaluate LM on unseen tasks

*Inference: generalization to unseen tasks*

> Q: Can Geoffrey Hinton have a conversation with George Washington?
>
> Give the rationale before answering.

> Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".
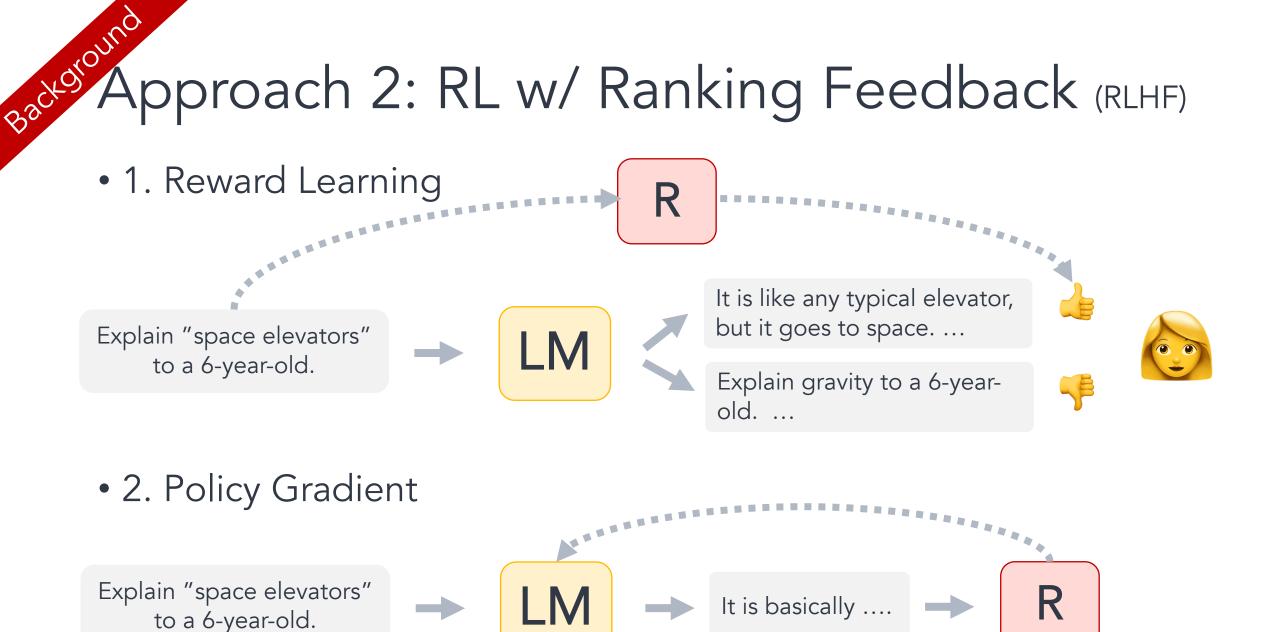
[McCann et al., 2019, Weller et al. 2020. Mishra et al. 2021; Wang et al. 2022, Sanh et al. 2022; Wei et al., 2022, [Chung et al. 2022, many others ] 42

# Approach 1: Behavior Cloning (Supervised Learning)

- Incentivizes word-by-word rote learning => limits creativity

- => The resulting models' generality/creativity is bounded by that of their supervision data.
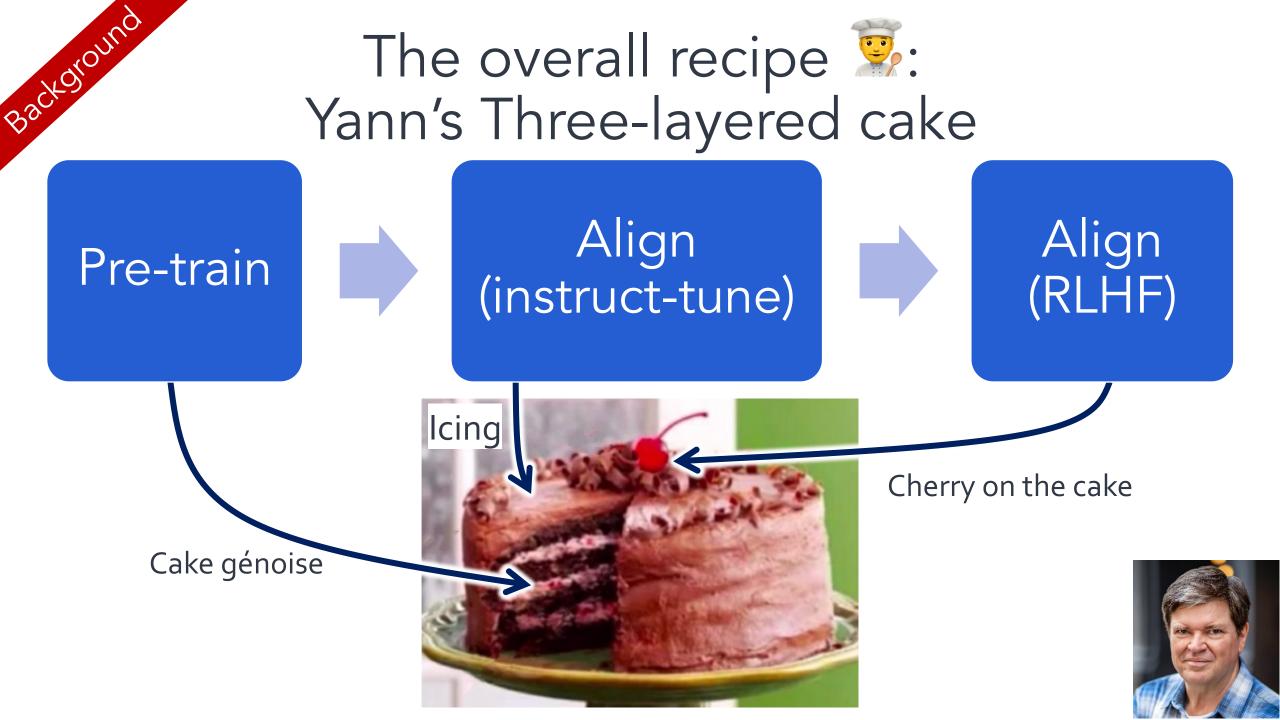


Answer the following question by reasoning step-by-step.

The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

**LM**

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9.

[McCann et al., 2019, Weller et al. 2020. Mishra et al. 2021; Wang et al. 2022,  Sanh et al. 2022; Wei et al., 2022, Chung et al. 2022, many others ] 43

# Approach 2: RL w/ Ranking Feedback (RLHF)

- **1. Reward Learning**



| Explain "space elevators" to a 6-year-old. | → | LM | ⟶ | It is like any typical elevator, but it goes to space. … | 👍 👩 |
|  |  |  | ⟶ | Explain gravity to a 6-year-old. … | 👎 |

R

- **2. Policy Gradient**

| Explain "space elevators" to a 6-year-old. | → | LM | → | It is basically …. | → | R |

[Christiano et al. 2017; Stiennon et al. 2020; Ouyang et al., 2022]

44

# The overall recipe 👨‍🍳:

Pre-train → Align (instruct-tune) → Align (RLHF)

# The overall recipe 👨‍🍳:

**Pre-train** → **Align (instruct-tune)** → **Align (RLHF)**

# The overall recipe 👨‍🍳:
# Yann's Three-layered cake

**Pre-train** → **Align (instruct-tune)** → **Align (RLHF)**

Cake génoise

Icing

Cherry on the cake

# Are these steps equally important?

Pre-train → Align (instruct-tune) → Align (RLHF)

# Are these steps equally important?

Pre-train → Align (instruct-tune) → Align (RLHF)

Who should care?

1. **Product designers:** If you have $X million to build your best chatbot, how would you allocate it?

2. **Scientists:** Fundamentally, is this the ultimate pipeline?

[Brown et al., 2020. GPT3, Ouyang et al., 2022. InstructGPT]

# Are these steps equally important?

Pre-train → Align (instruct-tune) → Align (RLHF)

How far can we reduce the human annotations?

# How far can we reduce the human annotations?

- **Idea:** we can bootstrap "instruction" from off-the-shelf LMs.
  - LMs have seen humans talk about their needs and goals.

Pretraining
(GPT3*: 499 Billion
tokens)

LLMs should know
a lot of tasks!

Self-Instruct:

# Aligning Language Models w/ Self-Generated Instructions

Warning: the paper is a year old!!

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, Hannaneh Hajishirzi

https://arxiv.org/abs/2212.10560

# Get humans to write "seed" tasks ✍️

- I am planning a 7-day trip to Seattle. Can you make a detailed plan for me?
- Is there anything I can eat for breakfast that doesn't include eggs, yet includes protein and has roughly 700-1000 calories?
- Given a set of numbers find all possible subsets that sum to a given number.
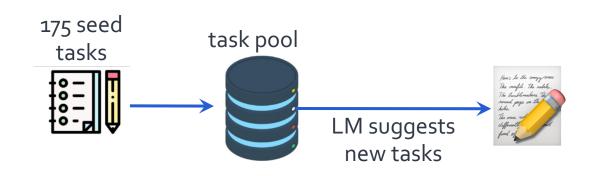- Give me a phrase that I can use to express I am very happy.

175 seed
tasks

# Put them your task bank 📦

- I am planning a 7-day trip to Seattle. Can you make a detailed plan for me?
- Is there anything I can eat for breakfast that doesn't include eggs, yet includes protein and has roughly 700-1000 calories?
- Given a set of numbers find all possible subsets that sum to a given number.
- Give me a phrase that I can use to express I am very happy.

175 seed
tasks

task
pool

# Sample and get LLM to expand it

- I am planning a 7-day trip to Seattle. Can you make a detailed plan for me?
- Is there anything I can eat for breakfast that doesn't include eggs, yet includes protein and has roughly 700-1000 calories?
- Given a set of numbers find all possible subsets that sum to a given number.
- Give me a phrase that I can use to express I am very happy.

**LM** Pre-trained, but **not aligned yet**

- Create a list of 10 African countries and their capital city?
- Looking for a job, but it's difficult for me to find one. Can you help me?
- Write a Python program that tells if a given string contains anagrams.

175 seed tasks

task pool

LM suggests new tasks

# Get LLM to answers the new tasks

- Task: Convert the following temperature from Celsius to Fahrenheit.
- Input: 4 °C
- Output: 39.2 °F

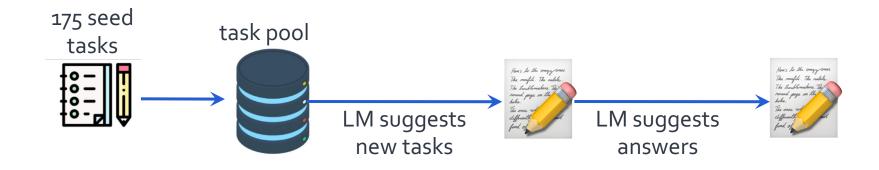- Task: Write a Python program that tells if a given string contains anagrams.

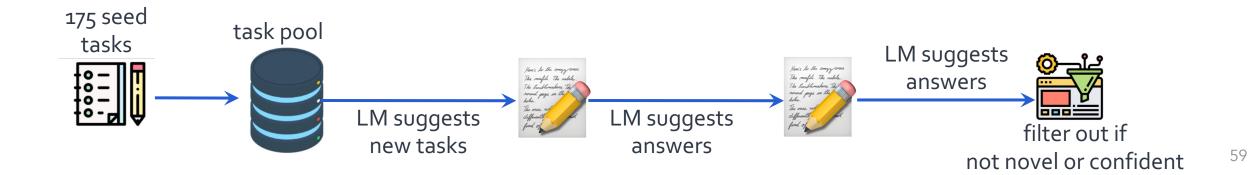**LM** Pre-trained, but **not aligned yet**

- Input: -
- Output:
    def isAnagram(str1, str2): ...

175 seed tasks

task pool

LM suggests new tasks

LM suggests answers

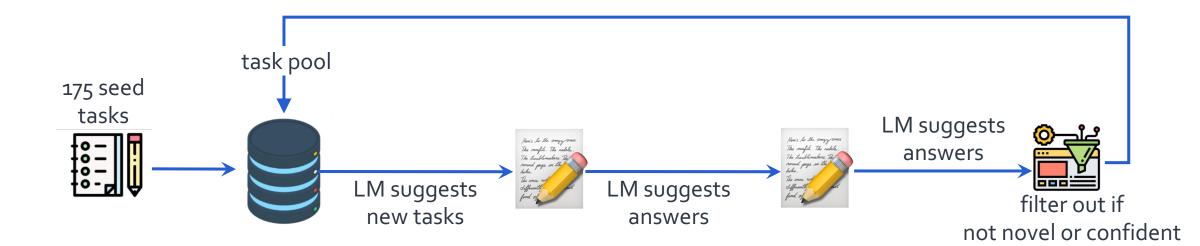# Filter tasks

- Drop tasks if LM assigns <span style="color:red">low probability</span> to them.

- Drop tasks if they have a <span style="color:red">high overlap</span> with one of the existing tasks in the task pool.
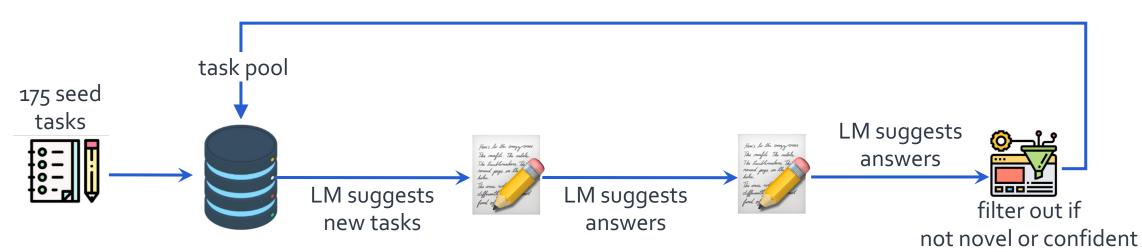  - Otherwise, common tasks become more common — <span style="color:red">tyranny of majority.</span>



175 seed tasks

task pool

LM suggests new tasks

LM suggests answers

LM suggests answers

filter out if not novel or confident

# Close the loop

- Add the filtered tasks to the task pool.
- Iterate this process (generate, filter, add) until yield is near zero.



task pool

175 seed
tasks

LM suggests
new tasks

LM suggests
answers

LM suggests
answers

filter out if
not novel or confident

# Self-Instructing GPT3 (base version)
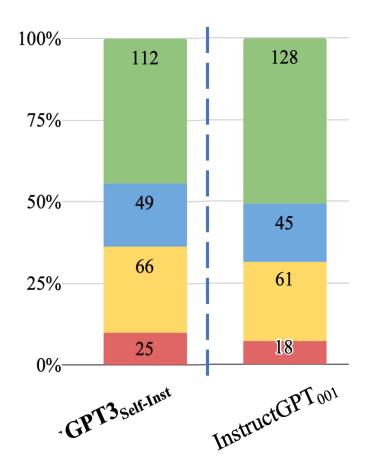
- ## Generate:
  - GPT3 ("davinci" engine).
  - We generated 52K instructions and 82K instances.
  - API cost ~$600

- ## Align:
  - We finetuned GPT3 with this data via OpenAI API (2 epochs). **
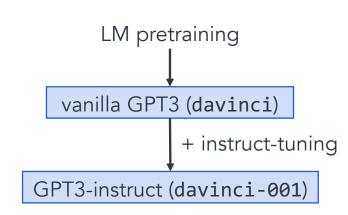  - API cost: ~$338 for finetuning



task pool

175 seed
tasks

LM suggests
new tasks

LM suggests
answers

LM suggests
answers

filter out if
not novel or confident

# Evaluation on User-Oriented Instructions

■ **A**: correct and satisfying response　　■ **B**: acceptable response with minor imperfections

■ **C**: responds to the instruction but has significant errors　■ **D**: irrelevant or invalid response



LM pretraining

vanilla GPT3 (`davinci`)

+ instruct-tuning

GPT3-instruct (`davinci-001`)

Diverse, "self-instruct" data ~ thousands of human-written data

[Self-Instruct: Aligning Language Model with Self-Generated Instructions, Wang et al. 2023]

# Summary Thus Far

- There is a lot of room to reduce the reliance on <span style="color:red">human</span> annotations in the "alignment" stage.

  - Well-read LLMs know a lot of our needs and demands.

  - Magic of "in-context learning" can surface these.

- Self-Instruct: Rely on creativity induced by LLMs themselves.

  - Lots of open-source adoption, but that's not the point …

(* See also concurrent work: Unnatural-Instructions [Honovich et al. 2022] and Self-Chat [Xu et al. 2023] )

# The weight of "alignment" step

Fundamentally, what is the role of post hoc alignment (step #2/3)?

| Step #1: Pre-train | → | Step #2/3: Align (RLHF or instruction-tune) |

It's playing a small role — Lightly modify LM so it can articulate its existing knowledge of tasks.

(+ put guardrails for what it can articulate)

It's playing a big role — Teaching LM knowledge of new tasks.

# Implications for what comes out

Fundamentally, what is the role of post hoc alignment (step #2/3)?

Step #1:
Pre-train

→

Step #2/3: Align
(RLHF or instruction-tune)

Unexpected behaviors may "emerge".

It will be as good as the alignment supervision.

It's playing a small role —
Lightly modify LM so it can articulate its existing knowledge of tasks.

It's playing a big role —
Teaching LM knowledge of new tasks.

(+ put guardrails for what it can articulate)

# The weight of "alignment" step: My 2 cents

- Most of the heavy lifting is done via pre-training (unlabeled).

- Alignment to "instructions" (tuning/RLHF) is a light touch on LLMs.
    - Can (and should) be done more efficiently.

It's playing a small role —
Lightly modify LM so it
can articulate its existing
knowledge of tasks.

It's playing a big role —
Teaching LM knowledge
of new tasks.

(+ put guardrails for what it can articulate)

# RLHF is patchwork for lack of grounding

- RLHF teach LMs (ground) the communicative intent of users.

  - For example, what is intended by "summarize"? The act of producing a summary grounded in the human concept of "summary".

- Not a panacea, but a short-term "band-aid" solution.



Intents and norms

RLHF or instruct-tuning

LM

[Some remarks on Large Language Models, Goldberg 2023]

# Alignment as a social process

• Can alignment emerge as a social experience?

• Internet also captures a subset of the world's interactive experiences.

# The future is a cheesecake



- Future: A unifying process that encompasses various steps that are done separately today.

- The margins between alignment stages are getting murkier.
  - Using model itself for feedback and verification
  - Alignment during pre-training (Korbak et al. 2023)
  - Building bridges between supervised learning and RL (see DPO vs. RLHF)
  - …

Pre-train → Align (instruct-tune) → Align (RLHF)

# The future is a cheesecake



- Future: A unifying process that encompasses various steps that are done separately today.

- Yann's framework was good for getting a system off the ground.

- Now that we are moving to interactive setups, alignment and pre-training will be a continual process. Systems that :
  - Adaptively change to our needs and habits;
  - Seamlessly pick up on implicit reward;
  - ....

# Thanks!