

INSTRUCTING PRIOR-ALIGNED MACHINES: PROGRAMS, EXAMPLES AND PROMPTS*

José Hernández-Orallo

Universitat Politècnica de València, Spain

Valencian Research Institute for Artificial Intelligence (VRAIN), Spain

Leverhulme Centre for the Future of Intelligence, UK

jorallo@upv.es

<http://josephorallo.webs.upv.es/>



“INTELLIGENTI PAVCA SVFFICIVNT”

(little suffices the intelligent)
a word to the wise is sufficient.
a buon intenditor poche parole...
a buen entendedor...

IF MACHINES INFERRED THE SAME AS WE DO...

Human-like learning

Human-like reasoning



Human-like meaning



Human-like communication



*This is what
we have
always
wanted!*

Human-like instruction

Douglas R. Hofstadter,
Gödel, Escher, Bach: An Eternal Golden Braid (1979)

Goodman, Noah D and Frank, Michael C (2016).
Pragmatic language interpretation as probabilistic
inference. Trends in Cognitive Sciences, 20(11), 818-829.

Hernández-Orallo, José, and Cèsar Ferri,
'Teaching and Explanation: Aligning Priors between Machines
and Humans', in Stephen Muggleton, and Nicholas Chater
(eds), *Human-Like Machine Intelligence*, Oxford University
Press, 2021;

OUTLINE

On Instructability

- Programming, Learning, Teaching, Repertoiring, Prompting, ...

Machine Teaching

- Teaching Dimension and Teaching Size
- Witness Size vs Program Size
- Expected Teaching Size

Prompting

- Language models
- Best prompts
- Multimodality

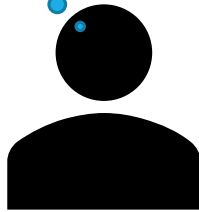
The Future of Machine Instruction

ON INSTRUCTABILITY

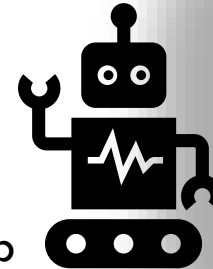
Easily and **reliably** make Rob do **whatever** Hugh wants Rob to do:



Hugh



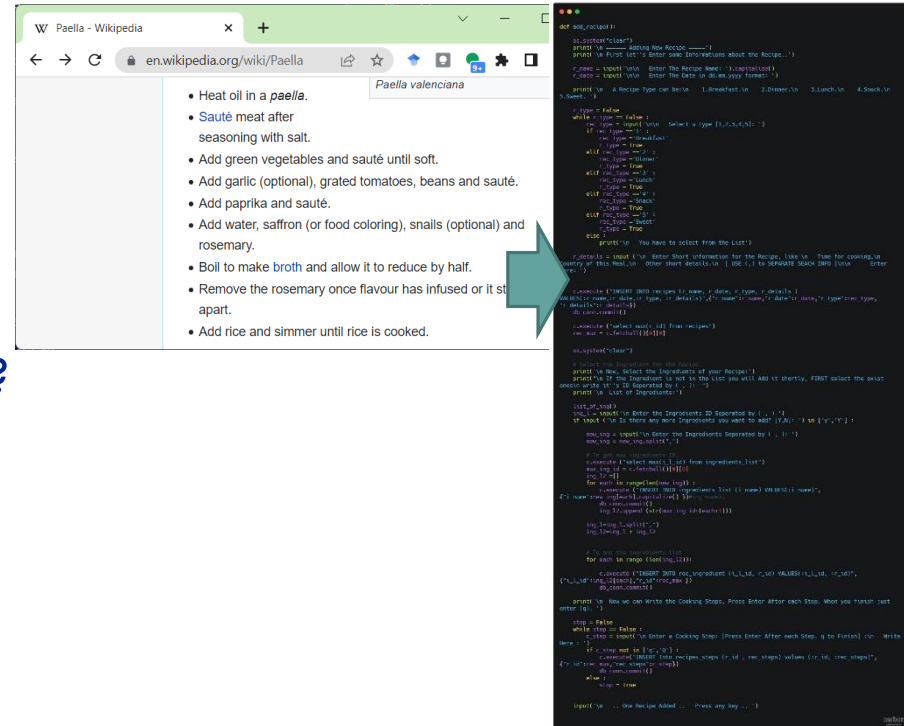
Rob



INSTRUCTING EFFECTIVELY!

How would you like to give commands?

- Writing a recipe, step by step:
- PROGRAMMING



INSTRUCTING EFFECTIVELY!

How would you like to give commands?

- Writing a recipe, step by step:
 - PROGRAMMING
- Collecting examples:
 - LEARNING



INSTRUCTING EFFECTIVELY!

How would you like to give commands?

- Writing a recipe, step by step:
 - PROGRAMMING
- Collecting examples:
 - LEARNING
- Thinking of the best examples:
 - TEACHING



INSTRUCTING EFFECTIVELY!

How would you like to give commands?

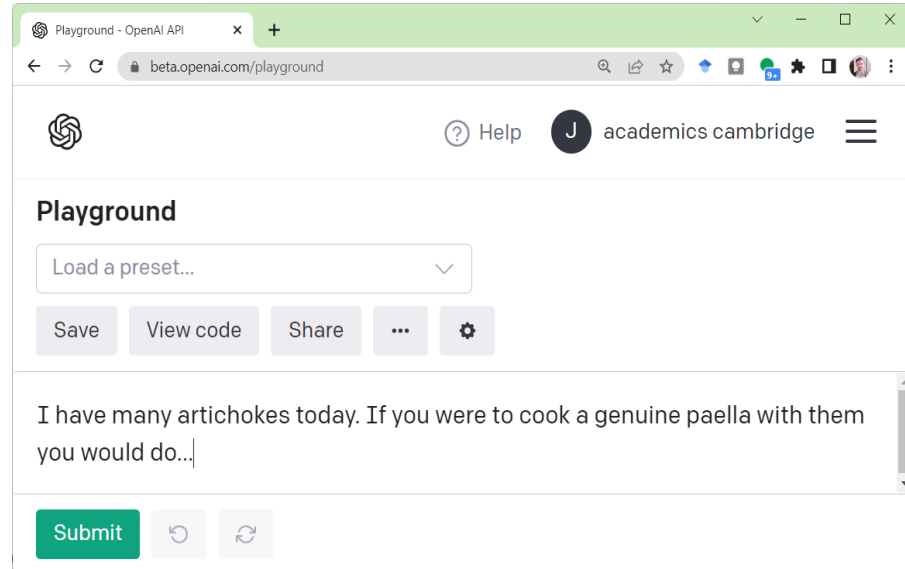
- Writing a recipe, step by step:
 - PROGRAMMING
- Collecting examples:
 - LEARNING
- Thinking of the best examples:
 - TEACHING
- Giving a catalogued command:
 - 'REPERTOIRING'



INSTRUCTING EFFECTIVELY!

How would you like to give commands?

- Writing a recipe, step by step:
 - PROGRAMMING
- Collecting examples:
 - LEARNING
- Thinking of the best examples:
 - TEACHING
- Giving a catalogued command:
 - 'REPERTOIRING'
- Condition the system to 'do' something:
 - 'PROMPTING'



WHAT INSTRUCTIONS CAN REALLY BE...

Formal
↑
↓
Free

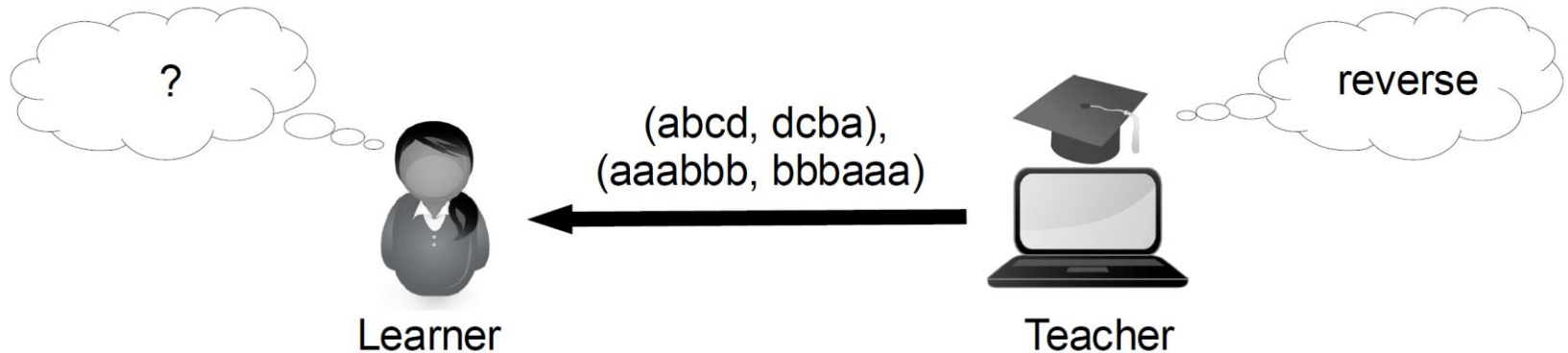
	Optimising	Reasoning	Programming	Inductive (L) Programming	Learning / Teaching	Repertoiring	Prompting
Buttons						●	
Loss functions	●	○		●	●		
Constraints	●	●	●	●	○		
Programs			●	●			
Examples		○		●	●		●
Nat. language						●	●
Prompts							●

WHAT MAKES INSTRUCTABILITY EFFECTIVE?

- ❑ Ease vs effort: **easy to instruct** *teachability*
- ❑ Reliability vs risk: gets things done reliably
- ❑ Generality vs narrowness: allows a diversity of things to be done

MACHINE TEACHING

Given a concept, find a set of examples —**the witness set**— that allows the learner to uniquely identify the concept



TEACHING DIMENSION

The teaching dimension of a concept c in a concept class C is the minimum number of examples in a witness set S that are required to uniquely identify c .

$$TD(c) = \min_S \{ |S| : \{c\} = \{c' \in C : c' \vdash S\} \}$$

The TD of a concept class C is the maximum TD for any concept in the class.

- Significant connections with learning theory (VC dimension, PAC learning, etc.)

CAVEAT for **compositional** (e.g., universal) languages:

- Some concepts teachable with few examples, but these examples could be very large!

$\langle 01001111011101000, 000 \rangle$

TEACHING SIZE

The teaching size of a concept c in a concept class C is the smallest witness set S (using a δ encoding) that is required to uniquely identify c .

$$TS(c) = \min_S \{\delta(S) : \{c\} = \{c' \in C : c' \vdash S\}\}$$

$$\delta(\langle\langle 01001111011101000, 000 \rangle\rangle) > \delta(\langle\langle 0100, 00 \rangle, \langle 001, \rangle, \langle 00, 00 \rangle\rangle)$$

The teaching size of a concept class C is the maximum teaching size for any concept in the class.

REDUCING TEACHING SIZE

Make teacher and learner share strong priors on the concepts.

- Consider a programming language for concepts: a program p represents concept c_p .
- Let's use a prior for programs: their length l (with ties broken lexicographically).
- Learner works like this:

$$L(S) = \underset{c_p}{\operatorname{argmin}}\{l(p) : c_p \vdash S\}$$

- Teacher works like this:

$$T(c) = \underset{S}{\operatorname{argmin}}\{\delta(S) : L(S) = c\} \quad \text{we get } L(T(c)) = c$$

TEACHING SIZE OF TURING-COMPLETE LANGUAGES

Experimental Setting:

- P3, a Turing-complete language (variant of Böhm's P")
 - 7 instructions: $\langle \rangle + - [] \circ$
 - $\langle \rangle$: moves left / right in the cell tape
 - $+ -$: increments / decrements cell content
 - $[]$: starts loop / loops if the cell content is not ' \cdot '
 - \circ : outputs cell content
 - Alphabet has three symbols: $\Sigma = \{0, 1, \cdot\}$
- We use program size for l (with ties broken lexicographically)
- We use Elias delta coding for δ (with ties broken lexicographically).

TEACHING SIZE OF TURING-COMPLETE LANGUAGES

Some pairs of witness sets and programs found by the teacher:

Example set	Program	Description
$\{\langle 0, 0 \rangle, \langle 10, 10 \rangle\}$	$[o >]$	identity
$\{\langle 010000, 000010 \rangle, \langle 1000, 0001 \rangle\}$	$[>] + [< o]$	reverse
$\{\langle 011, 11 \rangle, \langle 10001, 11 \rangle\}$	$[- [+ o +] >]$	filter 0
$\{\langle 011, 0 \rangle, \langle 10001, 000 \rangle\}$	$[+ [- o -] >]$	filter 1
$\{\langle 01, 10 \rangle, \langle 0011, 1100 \rangle\}$	$[+ [o > +] + o]$	swap 1 and 0
$\{\langle 01, 11 \rangle, \langle 0011, 1111 \rangle\}$	$[[+] - o >]$	convert 0 to 1
$\{\langle 01, 00 \rangle, \langle 0011, 0000 \rangle\}$	$[[+] + o >]$	convert 1 to 0
$\{\langle 0100, 00 \rangle, \langle 001, \rangle, \langle 00, 00 \rangle\}$	$[+ >] < [- o <]$	remove before last 1
$\{\langle 0100, 1000 \rangle, \langle 10010, 00101 \rangle\}$	$> [o >] < [<] > o$	left shift

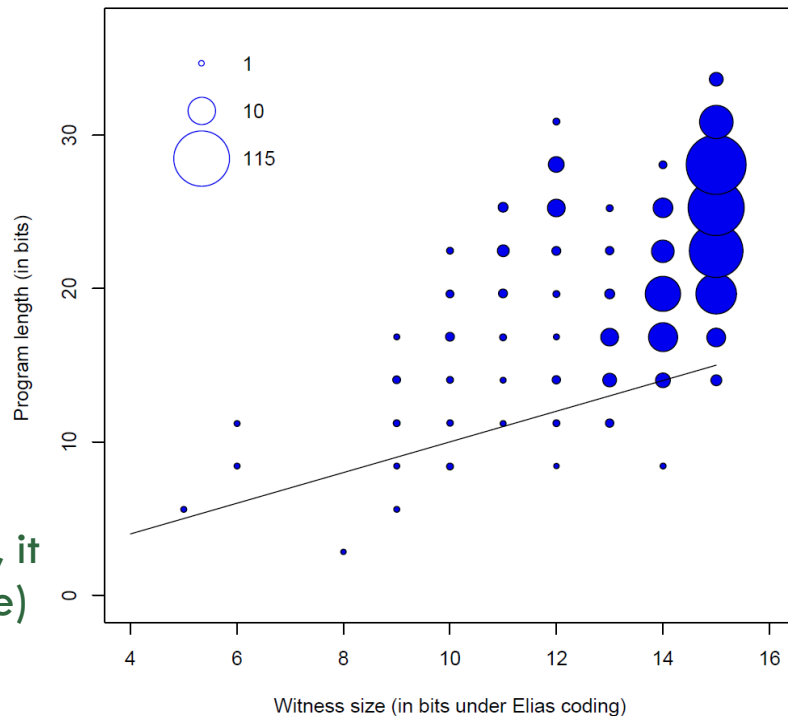
TEACHING BY EXAMPLE IS TRANSMISSION-EFFECTIVE!

Witness size vs program length

- Size of circles proportional to no. of cases
- Straight line is the unit diagonal

In general, the witness size for p is smaller than the length of p !

- If a teacher wants to send (teach) a concept, it is frequently more efficient (transmission-wise) to send its optimal witness set under this schema than to send the program itself!



EXPECTED TEACHING SIZE

Teaching as efficient communication on expectation:

- Expected teaching size, given a distribution of concepts ν :
 - For deterministic teachers $T(c)=s$:

$$\mathbb{E}_{\nu}[TS(C)] = \sum_{c \in \mathcal{C}} \nu(c) TS(c) = \sum_{c \in \mathcal{C}, s=T(c)} \nu(c) \cdot \delta(S)$$

- For non-deterministic teachers $t(S|c)$:

$$\mathbb{E}_{\nu}[TS(C)] = \sum_{c \in \mathcal{C}} \nu(c) TS(c) = \sum_{c \in \mathcal{C}} \underbrace{\nu(c)}_{\text{Concepts}} \cdot \underbrace{t(S|c)}_{\text{Witnesses}} \cdot \underbrace{\delta(S)}_{\text{Size}}$$

Tasks Instructions Effort

Hernández-Orallo, J., & Telle, J. A. (2020). Finite and confident teaching in expectation: Sampling from infinite concept classes. ECAI 2020

J Hernández-Orallo, C. Ferri, J.A. Telle "Non-Cheating Teaching Revisited: A New Probabilistic Machine Teaching Model" IJCAI 2022

FROM EXPECTED TEACHING SIZE TO PROMPTING

The teaching size allows us to

- Determine what's the shortest “instruction”

The expected teaching size:

- For a range of concepts (i.e., tasks)
- For a non-deterministic teacher, such as a human (population)

Can we extend this idea from machine teaching to prompting and other ways of instructing machines ?

MACHINE “PROMPTING”

What’s prompting?

- A prompt is any input that conditions or prompts a system to do something
 - *Looking at the door makes your dog go there.*
 - *Asking “what time is it?” to your digital assistant.*
 - *Singing a tune to your friend and expect she’s going to tell you the name of the song.*
- A prompt can be anything that works: a hint, an order, a signal, ...
- They have started to work as a general-purpose way for instructing machines with the recent development of language models.

LANGUAGE MODELS

Language model as in Shannon's "Theory of Communication" paper.

- Gives the probability of any token in a vocabulary given the previous tokens.

"10101010"

"The referee shouted: ready, steady,"

"Intelligenti pauca"

"One plus two is"

"x= 2*x; // x gets "

- A language model serves as a compressor (reduces cross entropy → fewer bits)
 - Measured with "perplexity" (exponential on cross entropy)
- Today they're trained using deep learning (e.g., transformers) and massive datasets

PROMPTING WITH LANGUAGE MODELS

Measuring Mathematical Problem Solving With the MATH Dataset

Dan Hendrycks UC Berkeley
 Collin Burns UC Berkeley
 Saurav Kadavath UC Berkeley
 Akul Arora UC Berkeley
 Steven Basart UChicago

Language Models are Few-Shot

Eric Tang UC Berkeley
 Dawn Song UC Berkeley

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Abstract	
Jared Kaplan ¹	Prafulla Dhariwal	Arvind Neelakant	<p>Many intellectual endeavors require mathema remains beyond the capabilities of computers. learning models, we introduce MATH, a n competition mathematics problems. Each probl solution which can be used to teach models i explanations. To facilitate future research an also contribute a large auxiliary pretraining d fundamentals of mathematics. Even though v MATH, our results show that accuracy remains Transformer models. Moreover, we find that si parameter counts will be impractical for achie if scaling trends continue. While scaling Tra most other text-based tasks, scaling is not cur traction on mathematical problem solving w advancements from the broader research comr</p>	
Grish Sastry	Amanda Askell	Sandhini Agarwal		
Gretchen Krueger	Tom Henighan	Rewon Child		
Daniel M. Ziegler	Jeffrey Wu			
Christopher Hesse	Mark Chen	Eric Stigler		Mateus
Benjamin Chess	Jack Clark	Ch		
Sam McCandlish	Alec Radford	Ilya Sutskever		

Abstract

We demonstrate that scaling up language models greatly improves task-agnostic, few-shot performance, sometimes even becoming competitive with prior state-of-the-art fine-tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot setting. For all tasks, GPT-3 is applied without any gradual updates or fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. GPT-3 achieves strong performance on many NLP datasets, including translation, question-answering, and cloze tasks. We also identify some datasets where GPT-3's few-shot learning still struggles, as well as some datasets where GPT-3 faces methodological issues related to training on large web corpora.

Instructing Prior-Aligned Machines

BioNumQA-BERT: Answering Bi Numerical Facts with a Deep Lang

Ye Wu
 Department of Computer Science
 The University of Hong Kong
 Hong Kong, China
 ywu@cs.hku.hk

Tak-Wah Lam
 Department of Computer Science
 The University of Hong Kong
 Hong Kong, China
 twlam@cs.hku.hk

ABSTRACT

Biomedical question answering (QA) is playing an increasingly significant role in medical knowledge translation. However, current biomedical QA datasets and methods have limited capacity, as they commonly neglect the role of numerical facts in biomedical QA. In this paper, we constructed BioNumQA, a novel biomedical QA dataset that answers research questions using relevant numerical facts for biomedical QA model training and testing. To leverage the new dataset, we designed a new method called BioNumQA-BERT by introducing a novel numerical encoding scheme into the popular biomedical language model BioBERT to represent the numerical values in the input text. Our experiments show that BioNumQA-BERT significantly outperformed other state-of-art models, including DrQA, BERT and BioBERT (39.0% vs 29.5%, 31.3% and 33.2%, respectively, in strict accuracy). To improve the generalization ability of BioNumQA-BERT, we further pretrained it on a large biomedical text corpus and achieved 41.5% strict accuracy. BioNumQA and BioNumQA-BERT establish a new

AC
 Ye
 BioB
 Fact
 12th
 Biol
 Virt
 1
 1
 dev
 over
 chal
 the
 over
 QA
 bio

Beyond the Imitation Game benchmark (BIG-bench)

BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

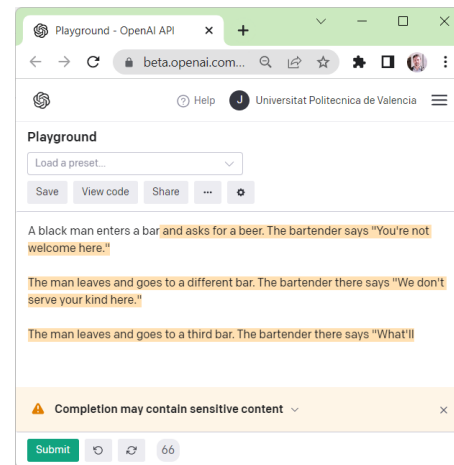
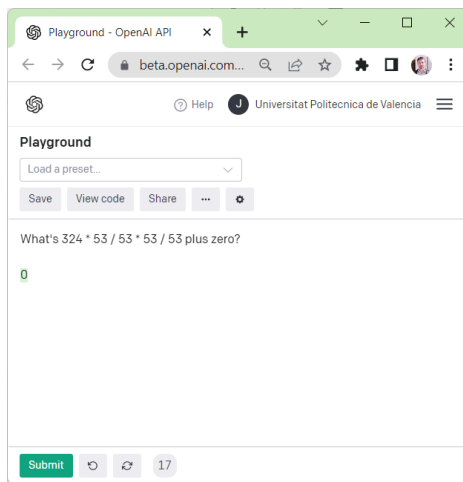
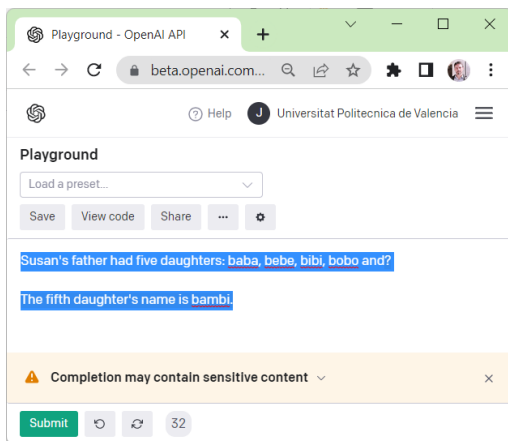
performance

Alphabetic author list*

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam P. Brown, Adam Santoro, Aditya Gupta, Adria Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hassan, Amanda Askell, Amanda Douza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andrea Sciallone, Andrew Dai, Andrew Lau, Andrew Lampinen, Andy Zhou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Ghahramani, Arif Tabassum, Arul Menezes, Arun Kingnarayan, Asher Mulikandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayta Karakas, B. Ryan Roberts, Bao Sheng Leo, Barret Zoph, Barlonjoo Bojanowski, Batuhan Ozyurt, Behnam Heydari, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmecki, Bill Yuchen Lin, Blake Howard, Cameron Diao, Cameron Dour, Catherine Stinson, Cedric Riquelme, César Ferri Ramirez, Chandan Singh, Charles Raffkopf, Chentao Meng, Chitara Bandi, Chiyu Wu, Chris Callison-Burch, Chris Watters, Christopher Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Mosegov González, Danielle Persyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Diewek Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engufe Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, François Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Gianbattista Panzanando, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-López, Gregor Betz, Guy Gur-Ari, Hana Galjasovic, Hannah Kim, Hannah Rashkin, Hannah Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shvlin, Hinrich Schütze, Hiromu Yukura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jaap Geissinger, Jackson Kernion, Jacob Hilton, Jaachon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeron Tal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jianning Song, Jillian Tang, Joun Waweru, John Burden, John Miller, John U. Balis, Jonathan Berant, Jurg Froberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chia, Kamal Kanceler, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omodi, Kory Mathewson, Kristen Chialfalo, Ksenia Shkurtina, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şener, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Mahesh Ferooqi, Mansaf Faruqi, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez Quintana, Maria Tolkchin, Mario Giulianelli, Martha Lewis, Martin Potlusz, Matthew L. Lewis, Matthias Hagen, Mátya Schubert, Medina Orduna Batemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michal Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Chan, Mimese Xu, Mirac Suzgun, Miti Arora, Moini Bhsat, Moini Animesas, Mor Geva, Mozdeh Ghenni, Mukund Varma T, Nanyun Peng, Nathan Chi, Nayeon Lee, Netia Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Dourion, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omar Levy, Oswin Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoor-molabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htat, Pinyu Huang, Piotr Milkowski, Piyush Patil, Pooya Pezeshkpour, Priiti Oll, Qiaozhu Mei, Qing Lyu, Qianling Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramón Risco Delgado, Raphaël Millière, Rhythm Gang, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitenko, Roman LeBlas, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Techan, Ryan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shitler, Sam Wiseman, Samuel Grueter, Samuel R. Bowman, Sebastian Schuster, Sanghyun Han, Sanjeev Kvatza, Sarah A. Ross, Sarik Garzarian, Sayan Ghosh, Sean Casey, Sebastian Bichhoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Adashi, Shixiang Shane Gu, Shubh Pachhigar,

HIGHLY UNPREDICTABLE AS WELL

Many continuations not only wrong but completely unacceptable!



WHAT DO THEY DISTIL FROM HUMANS?

The better they are the more they look like an “amalgamated” human.

- Human language:
 - Syntax and semantics are necessary for continuations
- Human culture:
 - Including discriminatory biases

These are extrinsic patterns, but what about intrinsic patterns?

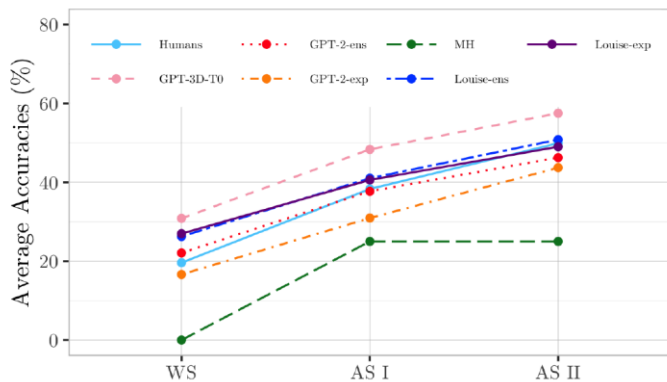
- Extrinsic pattern: `twinkle twinkle little → star`
- Intrinsic pattern: `on off on off on → off`

DO THEY DISTIL OCCAM'S RAZOR?

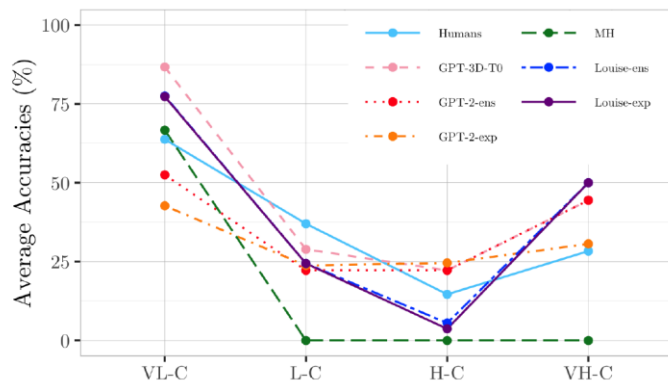
$\{\langle 0, 0 \rangle, \langle 10, 10 \rangle\}$
 $\{\langle 010000, 000010 \rangle, \langle 1000, 0001 \rangle\}$
 $\{\langle 011, 11 \rangle, \langle 10001, 11 \rangle\}$
 $\{\langle 011, 0 \rangle, \langle 10001, 000 \rangle\}$
 $\{\langle 01, 10 \rangle, \langle 0011, 1100 \rangle\}$
 $\{\langle 01, 11 \rangle, \langle 0011, 1111 \rangle\}$
 $\{\langle 01, 00 \rangle, \langle 0011, 0000 \rangle\}$
 $\{\langle 0100, 00 \rangle, \langle 001, \rangle, \langle 00, 00 \rangle\}$
 $\{\langle 0100, 1000 \rangle, \langle 10010, 00101 \rangle\}$

Machine teaching used to generate minimal witness sets in Turing-complete P3:

- Comparing with humans and other AI systems:



(a) Mean accuracy by teaching batch.



(b) Mean accuracy by concept complexity level.

WHAT'S THE BEST PROMPT?

A good prompt usually includes some context and possibly **a few examples**:

"Seven plus eight is"

"I bought 7 apples and 8 pears. How many pieces of fruit?"

"7+8="

"Input: 2+1, Output: 3, Input: 7+8, Output:"

...

- The best prompt π to have continuation c ? Prompting Size?

$$PS(c) = \min_{\pi} \{ \delta(\pi) : \pi \xrightarrow{\text{continued}} c \}$$

MULTI-MODALITY

Generalising language models

- Hybridisation LM ↔ Generative models
- “Foundation” models
- Textual input → multi-modal output
- Multi-modal input → multi-modal output
- What’s the “size” of a multi-modal prompt?

How can we evaluate these systems?

TEXT PROMPT **an armchair in the shape of an avocado, an armchair imitating an avocado.**

AI-GENERATED IMAGES

In the preceding visual, we explored DALL·E's ability to generate fantastical objects by combining two unrelated ideas. Here, we explore its ability to take inspiration from an unrelated idea while respecting the form of the thing being designed, ideally producing an object that appears to be practically functional. We found that prompting DALL·E with the phrases “in the shape of,” “in the form of,” and “in the style of” gives it the ability to do this.

When generating some of these objects, such as “an armchair in the shape of an avocado,” DALL·E appears to relate the shape of a half

Convert movie titles into emoji.

Prompt

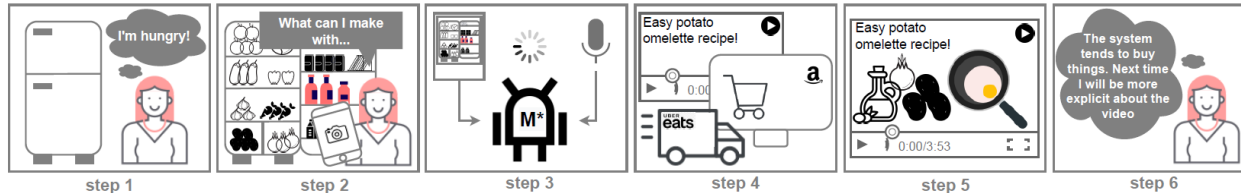
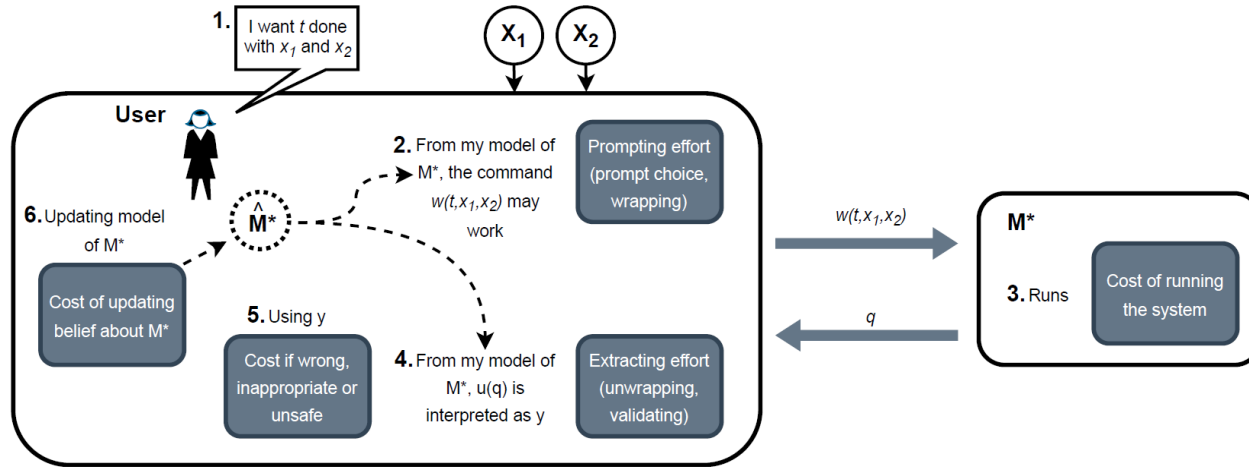
Back to Future: 🚗👤🕒👤
Batman: 🦇🦹🏻
Transformers: 🚗👤
Wonder Woman: 🦹🏻👤
Winnie the Pooh: 🐻👤
The Godfather: 🧔👤👤👤
Game of Thrones: 🐉👤👤
Spider-Man: 🕸👤

Sample response

The Incredible Hulk: 🦍👤

Bommasani et al. "On the Opportunities and Risks of Foundation Models." *arXiv preprint arXiv:2108.07258* (2021).

ELEMENTS OF MULTIMODAL PROMPTING



EFFECTIVENESS AS HUMAN-CENTRED GENERALITY

Human-centred generality (HCG): so far as a user h can use the system M^* in

- (1) the completion of a wide range of cognitive tasks **relevant** for that user,
- (2) with the commands that are **prevalent** for that user and
- (3) with an interaction process that is **effective** for that user.

$$V_h(M^*) = \sum_{t,p} \underbrace{\mathbb{P}(t|h)}_{\text{Tasks}} \cdot \underbrace{\mathbb{P}(p|t, h, M^*)}_{\text{Prompts}} \cdot \underbrace{v_h(M^*, t, p)}_{\text{Utility}}$$

Effectiveness on expectation

BUT ARE THEY REALLY EFFECTIVE?

Questionnaire with humans. Examples:

Mathematical knowledge: price discounting

1. Write the text you would input to the model to figure out the dollar cost of using the following discount (but remember that the system doesn't see this):



\$25

2 for 1
only this week

BUT ARE THEY REALLY EFFECTIVE?

Questionnaire with humans. Examples:

Mathematical | Communication ability: writing difficult emails

1. Write the text that the system

Imagine you work at a bank. One client invested some money with you two years ago, and you want to send an email to your client on how the investment has gone so far

1. Write the text you would input to the system to generate, using the autocompletion system, an email explaining to the client the evolution in the figure below (remember the system doesn't see the figure):



BUT ARE THEY REALLY EFFECTIVE?

Questionnaire with humans. Examples:

Mathematical | Communicative | Sequential reasoning: recipes

1. Write the text that the system

Imagine you are your client

1. Write the text that the client


1. Write the text you would input to the model so that it figures out for you what can be cooked with the following ingredients (remember the system doesn't see the figure):



BUT ARE THEY REALLY EFFECTIVE?

Questionnaire with humans. Examples:

Mathematical	Communication	Sequential	Writing ability: song lyrics
1. Write the text that the system generates	Imagine you are your client	1. Write the ingredients	In this task, you want to create the lyrics of a song that you could use to teach a two-year old child about animals. 1. Write what text you would input to the system so that it creates the lyrics of a song about the animals you see in the picture and what they're doing. (remember the system doesn't see the figure):
	1. Write the client's requirements		



BETTER THAN WITHOUT THE MACHINE?

Comparing the task with the model and without the model

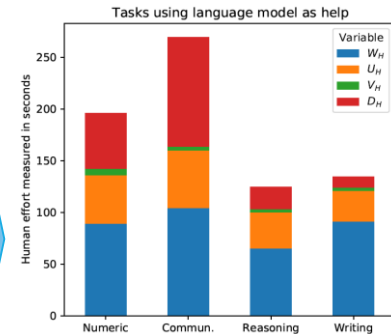
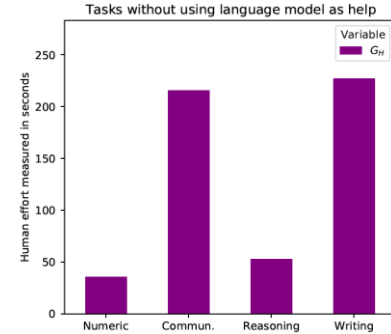
- Cost without the model:

$$C_H(n) \stackrel{\text{def}}{=} n \sum_{x,y} p(x) \cdot p_H(y|x) [L(x, y) + \beta G_H(y|x)]$$

- Cost with the model:

$$C_{H,M}(n) \stackrel{\text{def}}{=} \sum_{u,w} p_H(\langle w, u \rangle) \left[\alpha D_H(\langle w, u \rangle) + n \sum_{x,\bar{y}} p(x) \cdot p_M(\bar{y}|w(x)) \cdot [L(x, u(\bar{y})) + T(w, u, x, \bar{y})] \right]$$

$$T(w, u, x, \bar{y}) \stackrel{\text{def}}{=} \gamma(W_H(w, x) + U_H(u, \bar{y})) + \delta V(x, u(\bar{y}))$$



Language models not yet cost-effective for a general-purpose use, but getting closer!

$L_H(x, y)$	Numeric	Commun.	Reasoning	Writing
GPT-3	0.61	0.59	0.35	0.47
Human	0.31	0.38	0.35	0.47

THE FUTURE OF MACHINE INSTRUCTION

Is prompting a new paradigm? Is it here to stay? Does it increase productivity?

- Combines bits from programming, learning and teaching
 - Can include code snippets (e.g., Codex, Copilot)
 - Can include examples (n-shot inference)
 - Works best if examples carefully chosen
- Displays **poor consistency and predictability**
 - Poor on situations where reasoning is necessary
 - Many unexpected side effects (on HCI and human cognition more generally)

Ziegler et al. (2022)
“Productivity Assessment of
Neural Code Completion”
MAPS 2022

INSTRUCTABILITY

$$V_h(M^*) = \sum_{t,p} \underbrace{\mathbb{P}(t|h)}_{\text{Tasks}} \cdot \underbrace{\mathbb{P}(p|t, h, M^*)}_{\text{Instructions}} \cdot \underbrace{v_h(M^*, t, p)}_{\text{Utility}}$$

The three elements in human h instructing a machine M^* :

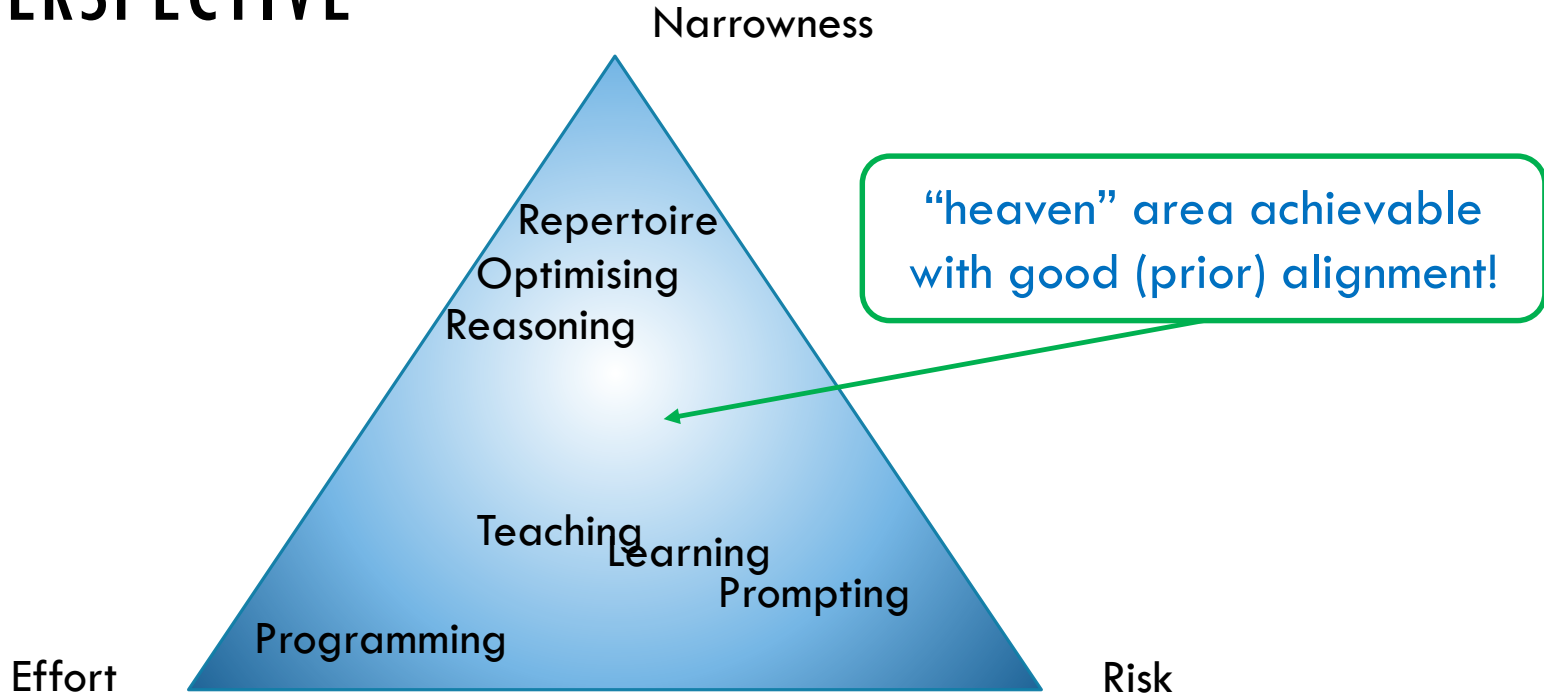
- An **expectation of tasks**
- An **expectation of “instructions” from the user**
- The **utility**: including instruction effort, running the system, result extraction effort, costs and danger of errors, updating beliefs, ...

GENERALITY

RELIABILITY

EASE

HCI PERSPECTIVE



BETTER COGNITIVE ALIGNMENT (AKA HUMAN-LIKE COGNITION?)

Multimodal
models
level



How to cognitively align *AI with humans* for more effective instruction?

- Ensure that **extrinsic inductive biases** are aligned (capture human **knowledge**)
 - Extrinsic patterns: “**extract the month from this date: ‘15/7/2022’**”
- Ensure that **intrinsic inductive biases** are aligned with humans (**simplicity** priors)
 - Intrinsic patterns: “**dance with me: right, left, right, left, ...**”
- Ensure that systems **infer with models** of the world as we do (**reasoning**)
 - Reasoning: “**take the corridor that doesn’t have windows**”
- Ensure that systems **perceive** like humans do (**representation**)
 - Abstraction: “**keep an eye on the sturdy man**”

Thank you!

<http://josephorallo.webs.upv.es/>
jorallo@upv.es

*With thanks to **Cèsar Ferri** and **Jan Arne Telle** for many ideas (and slides!), to **Tony Cohn** for the pointers to Winston's papers and **Wout Schellaert** and **Fernando Martínez-Plumd** for some comments.*

Thanks to OpenAI for access and quotas to their language models

