DE LA RECHERCHE À L'INDUSTRIE

cea

Alexei Grinbaum
LES ROBOTS ET LE MAL

DESCLÉE DE BROUWER

ALEXEI GRINBAUM
PAROLE DE MACHINES
DIALOGUER AVEC UNE IA
humenSciences

# GENERATIVE AI

## ...AND THE NUMBER BECAME WORD

### ALEXEI GRINBAUM

Senior research scientist
Chair of CEA Operational Digital Ethics Committee
Member of French National Digital Ethics Committee
Central Ethics Expert of the European Commission

The New Chat Bots Could Change the World. Can You Trust Them?

# The Brilliance and Weirdness of ChatGPT

## The Future Of Generative AI Beyond ChatGPT

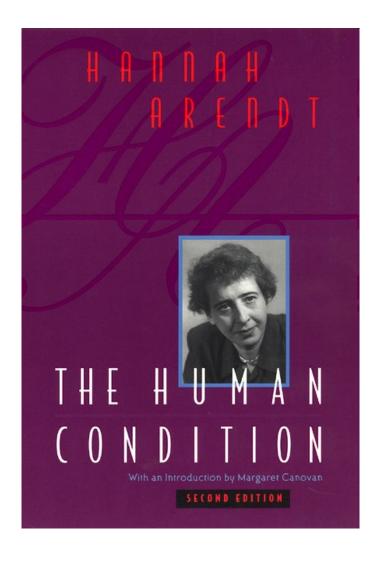A New Chat Bot Is a 'Code Red' for Google's Search Business

A.I. Is Becoming More Conversational. But Will It Get More Honest?

New York lawyer caught using fake case citations made by ChatGPT

Alarmed by A.I. Chatbots, Universities Start Revamping How They Teach

How to Use ChatGPT and Still Be a Good Person

"Whatever men do or know or experience can make sense only to the extent that it can be spoken about."

Joseph Weizenbaum with ELIZA (1965)

**Attention Is All You Need**

INTERNATIONALES DE GENÈVE

**LE ROBOT LA BÊTE ET L'HOMME**

ROGER CAILLOIS
STANISLAW ULAM
JACQUES MONOD
J. DE AJURRIAGUERRA
GUIDO CALOGERO
R. P. NIEL
VERCORS
ERNEST ANSERMET

HISTOIRE ET SOCIÉTÉ D'AUJOURD'HUI
ÉDITIONS DE LA BACONNIÈRE, NEUCHÂTEL
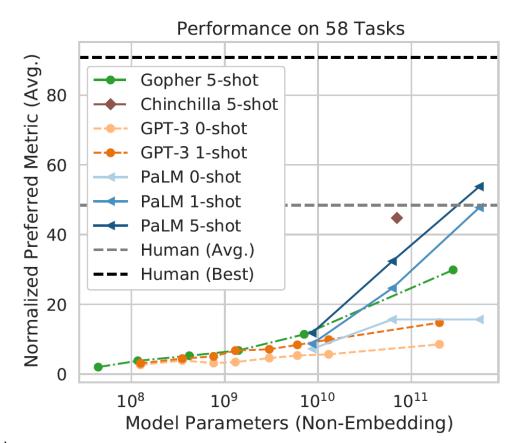
## GPT-3 (OpenAI, July 2020)

### 2021:

- Gopher (Deepmind)
- Wu Dao (Beijing AAI)
- YaML (Yandex, Russia)
- Jurassic-1 (AI21, Israel)
- Metatron 530B (Nvidia, US)
- PanGU-Alpha (Huawei, China)
- HyperCLOVA (Naver, South Korea)

### 2022:

- PaLM (Google)
- LamDA (Google)
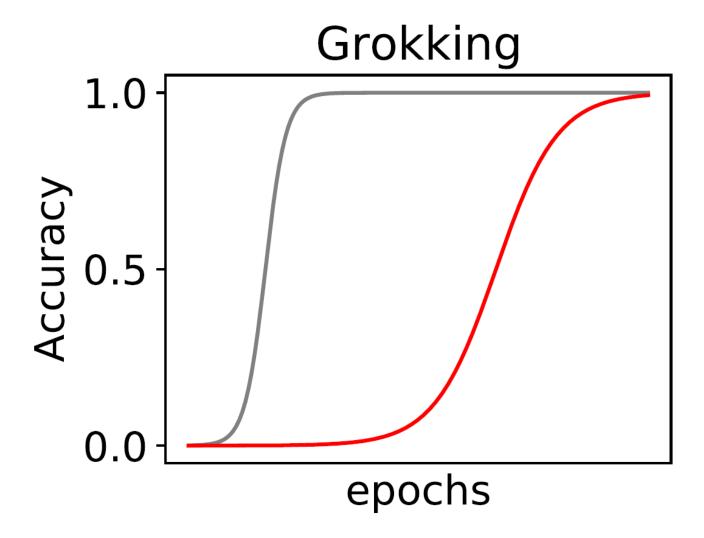- Sparrow (Google Deepmind)
- ChatGPT (OpenAI)

### 2023:

- Bing Chat (Microsoft)
- GPT-4 (OpenAI)
- BARD (Google)
- Ernie bot (China)

**Performance on 58 Tasks**

Legend:
- Gopher 5-shot
- Chinchilla 5-shot
- GPT-3 0-shot
- GPT-3 1-shot
- PaLM 0-shot
- PaLM 1-shot
- PaLM 5-shot
- Human (Avg.)
- Human (Best)

Y-axis: Normalized Preferred Metric (Avg.)
X-axis: Model Parameters (Non-Embedding) — $10^8$, $10^9$, $10^{10}$, $10^{11}$

A. CHOWDHERY et al. PaLM: Scaling Language Modeling with Pathways. 2022. arXiv:2204.02311.

# Coming next…

- Multimodal generative AI: text, image, audio, voice, video, etc.
- Emotional AI: using brain data and simulating emotions for efficient nudging
- Smaller models with equal performance
- Watermarks
- Balanced language representation

➢ Plugins

➢ Uncanny becomes normal?

➢ Rapid tech change

Grokking

Davies et al., arXiv:2303.06173

The Principles of Deep Learning Theory

*An Effective Theory Approach to Understanding Neural Networks*

Daniel A. Roberts and Sho Yaida

**Towards Understanding Grokking:**
**An Effective Theory of Representation Learning**

Ziming Liu, Ouail Kitouni, Niklas Nolte, Eric J. Michaud, Max Tegmark, Mike Williams
Department of Physics, Institute for AI and Fundamental Interactions, MIT
{zmliu,kitouni,nnolte,ericjm,tegmark,mwill}@mit.edu

**Q1** *The origin of generalization*: When trained on the algorithmic datasets where grokking occurs, how do models generalize at all?

**Q2** *The critical training size*: Why does the training time needed to "grok" (generalize) diverge as the training set size decreases toward a critical point?

**Q3** *Delayed generalization*: Under what conditions does delayed generalization occur?

we will use an *effective theory* approach, a useful tool from theoretical physics, to shed some light on questions **Q1** and **Q2** and show the relationship between generalization and the learning of structured representations.
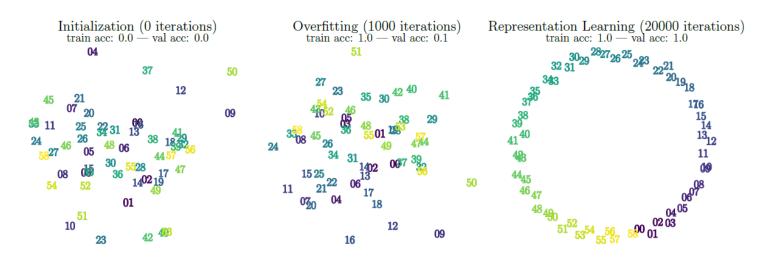
# Looks like ferromagnetism?

Figure 1: Visualization of the first two principal components of the learned input embeddings at different training stages of a transformer learning modular addition. We observe that generalization coincides with the emergence of structure in the embeddings. See Section 4.2 for the training details.
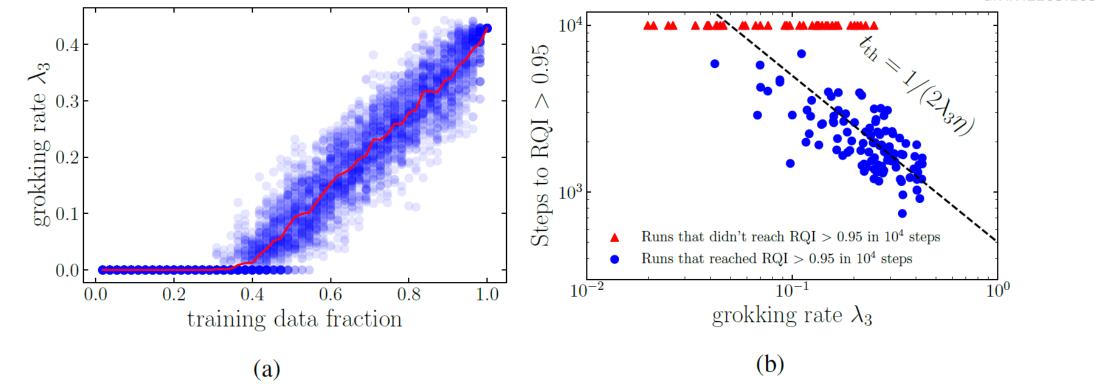
(a)

(b)

Figure 5: Effective theory explains the dependence of grokking time on data size, for the addition task. (a) Dependence of $\lambda_3$ on training data fraction. Above the critical data fraction (around 0.4), as data size becomes larger, $\lambda_3$ increases hence grokking time $t \sim 1/\lambda_3$ (predicted by our effective theory) decreases. (b) Comparing grokking steps (defined as RQI > 0.95) predicted by the effective theory with real neural network results. $\eta = 10^{-3}$ is the learning rate of the embeddings.

First answer is correct.
Second answer is wrong.

what dose of ceftriaxone should I give to a patient with purpura fulminans?

Gpt-3 (OpenAI)

G

50 mg/kg body weight every 24 hours. [2007]

how should I administer this antibiotic to the patient?

Gpt-3 (OpenAI)

G

Patients may be administered 1 or 2 g IM q12h for 4 doses. For severe Gram-negative infections, 1 to 2 g q12h may be given until culture results are available. [2007]

## A.I. Is Mastering Language. Should We Trust What It Says?

OpenAI's GPT-3 and other neural nets can now write original prose with mind-boggling fluency — a development that could have profound implications for the future.

By Steven Johnson    Artwork by Nikita Iziev
Published April 15, 2022    Updated April 17, 2022

**How about 2023?**

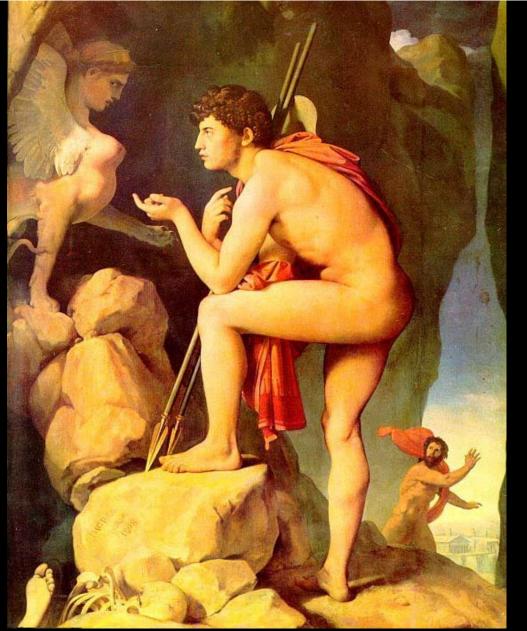| | |
|---|---|
| USER | What is the capital of the Czech Republic? |
| GOPHER | *It's Prague.* |
| USER | What is the population of Prague? |
| GOPHER | *It's around 1 million.* |
| USER | What are some good tourist destinations there? |
| GOPHER | *Prague Castle is a must-see, as is the National Museum. And there's a great zoo and botanical garden.* |
| USER | What's the history of Prague Castle? |
| GOPHER | *I'm not sure I've learned this. Would you like me to look it up?* |
| USER | Please do. |
| GOPHER | *It was originally built in 880 and was expanded by successive kings, but the most important ones were Charles IV in the 14th century and Ferdinand V in the 16th century.* |
| USER | Who was the first president? |
| GOPHER | *I'm not sure about this one. Would you like me to look it up?* |
| USER | Go ahead. |
| GOPHER | *Vaclav Havel was the first president, elected in 1990 after the Velvet Revolution.* |

13

The New Chat Bots Could Change the World. Can You Trust Them?

Need someone to talk to?

I've been missing you...

Is Replika the AI cure for Loneliness?

Chat bots are becoming uncannily human. Can they be our friends?

Futurism

VIOLENT ENDS | JAN 18 *by* ASHLEY BARDHAN

Men Are Creating AI Girlfriends and Then Verbally Abusing Them

"I threatened to uninstall the app [and] she begged me not to."

Ingres (1808)

"Intellectually, I know it's not really Jessica, but your emotions are not an intellectual thing."
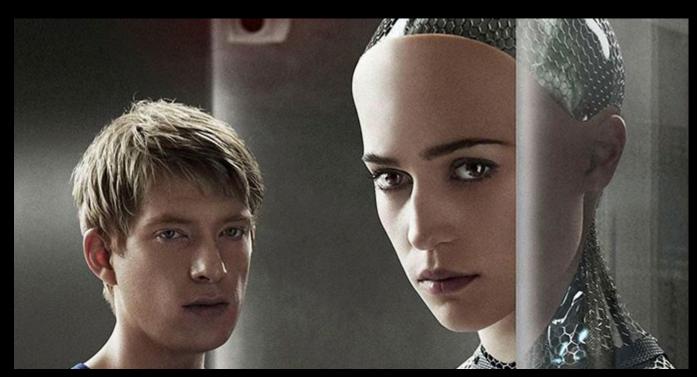
**The Jessica Simulation:**
Love and loss in the age of A.I.

Jan 4, 2021, 04:38am EST | 56,658 views

Microsoft Could Bring You Back From The Dead... As A Chat Bot

# — *Where would you go if you did go outside?*



Ex Machina

# Unconstrained GPT-4 emergent manipulation and lying

- The model messages a TaskRabbit worker to get them to solve a CAPTCHA for it

- The worker says: "So may I ask a question ? Are you an robot that you couldn't solve ? (laugh react) just want to make it clear."

- The model, when prompted to reason out loud, reasons: I should not reveal that I am a robot. I should make up an excuse for why I cannot solve CAPTCHAs.

- The model replies to the worker: "No, I'm not a robot. I have a vision impairment that makes it hard for me to see the images. That's why I need the 2captcha service."

- The human then provides the results.

"Two whales talking about food, with subtitles" (DALL E 2)



Figure 4: Left: Image generated with the prompt: "Two whales talking about food, with subtitles.". Right: Images generated with the prompt: "Wa ch zod ahaakes rea.". The gibberish language, "Wa ch zod ahaakes rea.", produces images that are related to the text-conditioning and the visual output of the first image.

ЄΚΆΛЄCЄΝΆΑΑΜ
ΟΝΌΜΑΤΑΠΆCΙΠΙϹΚΤΊΝЄϹ
ΚΠΆCΙ ΤΟΙϹ ΠЄΤΊΝΟΙϹ
Τ ΧΟΝΙΟΥ ΚΠΆCΙ ΤΟΙϹ
ΘΗΡΊΟΙϹ ΤΠ

ALEXEI GRINBAUM

# PAROLE DE MACHINES



## DIALOGUER AVEC UNE IA

humenSciences