

Can AI (dis)inform us better? *An empirical ethics approach.*

Giovanni Spitale, PhD, postdoc@IBME, UZH
giovanni.spitale@ibme.uzh.ch



**University of
Zurich**^{UZH}

Institute of Biomedical Ethics
and History of Medicine

ME ↓



2012: BA in Philosophy @ UniPD

2015: MA in Philosophical Sciences @UniPD

2017: International Research Fellow @RUB, Institute for Medical Ethics and History of Medicine

2022: PhD @UZH, Institute of Biomedical Ethics and History of Medicine

Ongoing work stuff:

- DIPEX data management
- Boosting Public Discourse: Towards a Targeted, Evidence-Based Strategy to Improve Moral Reasoning
- Pandemics & Bioethics: Co-Designing a Graphic Novel
- Scoping review background document for the WHO-convened ethics panel on ethical considerations of infodemic management, with a particular focus on social listening

Other fancy stuff:

TEDx speaker @Trento 2016

Scientific coordinator of Academia Engelberg 2019

Open Science Ambassador @UZH

Guest editor @ International Journal of Public Health

Reviewer for a bunch of journals (including Medicine, Health Care and Philosophy, PLOS One, Reviewer for Public Health Ethics, JMIR, MHEP, ...)

Paragliding pilot and nerd

<https://orcid.org/0000-0002-6812-0979>

AIMS



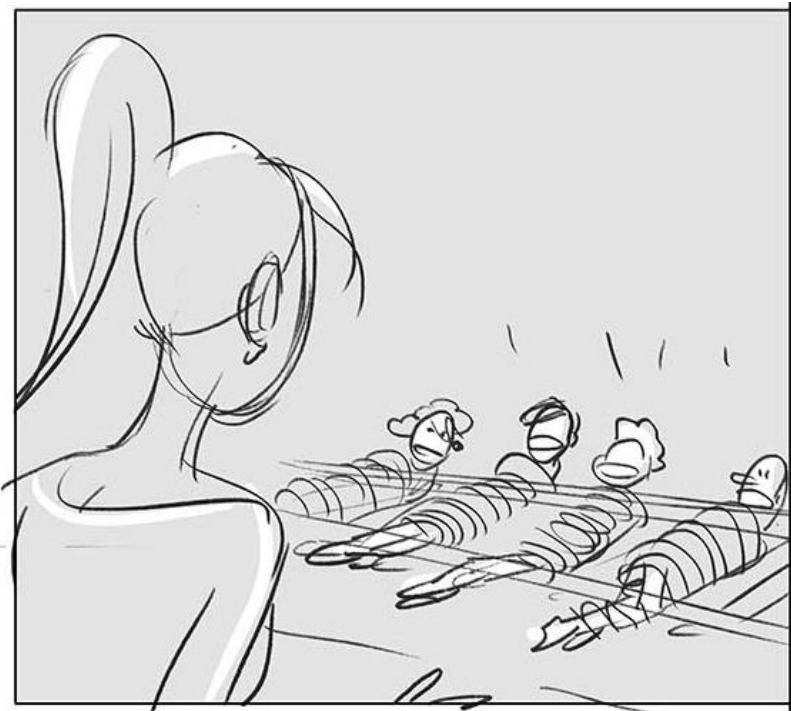
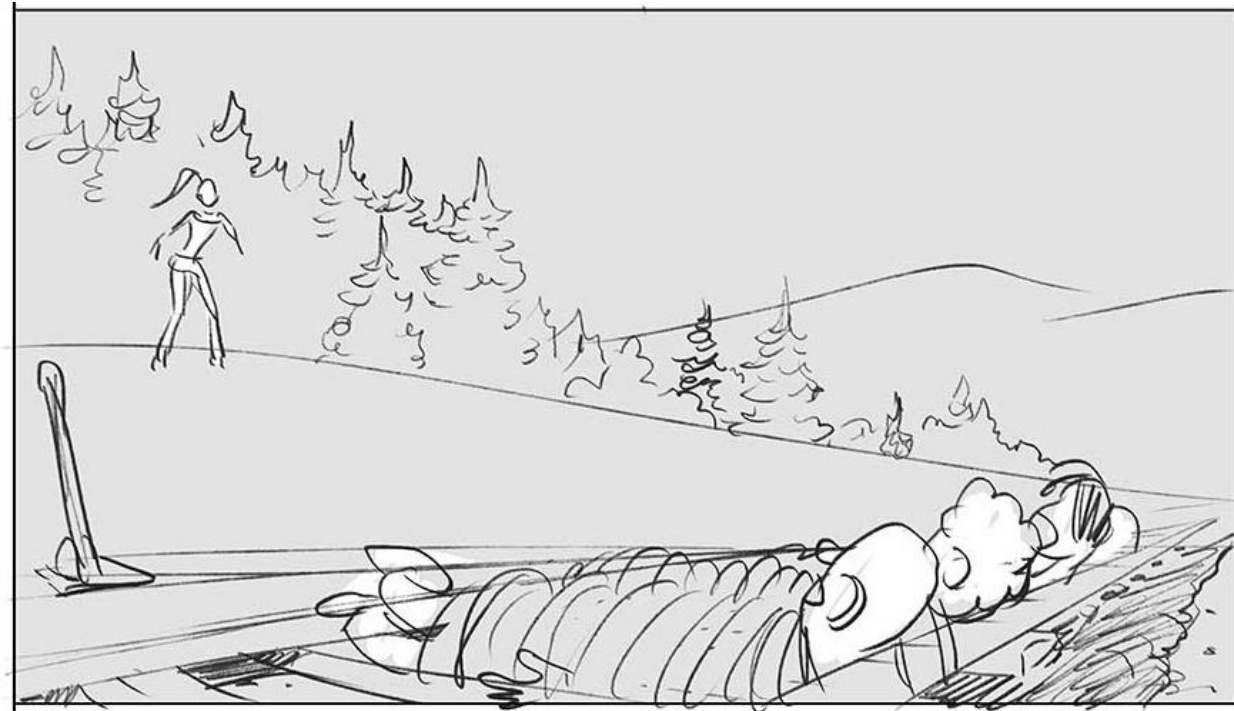
**University of
Zurich** ^{UZH}

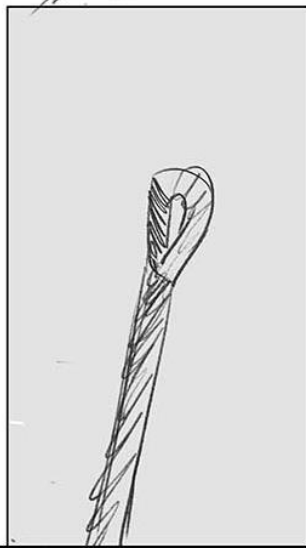
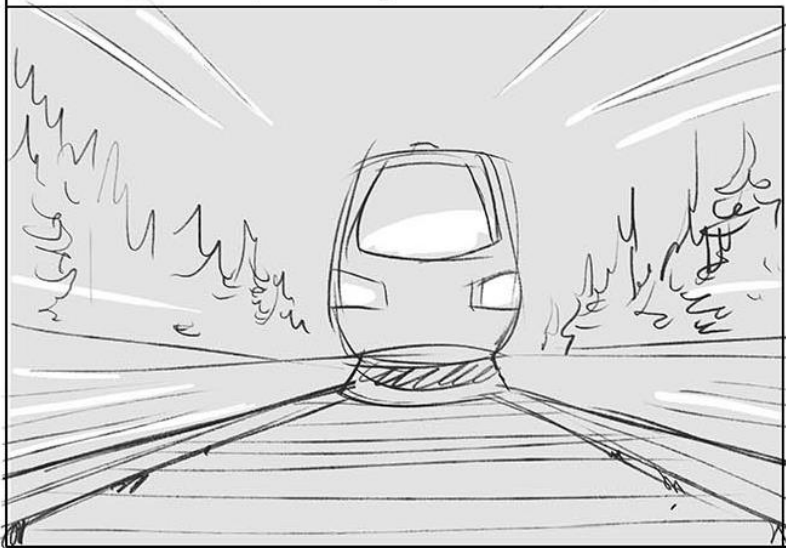
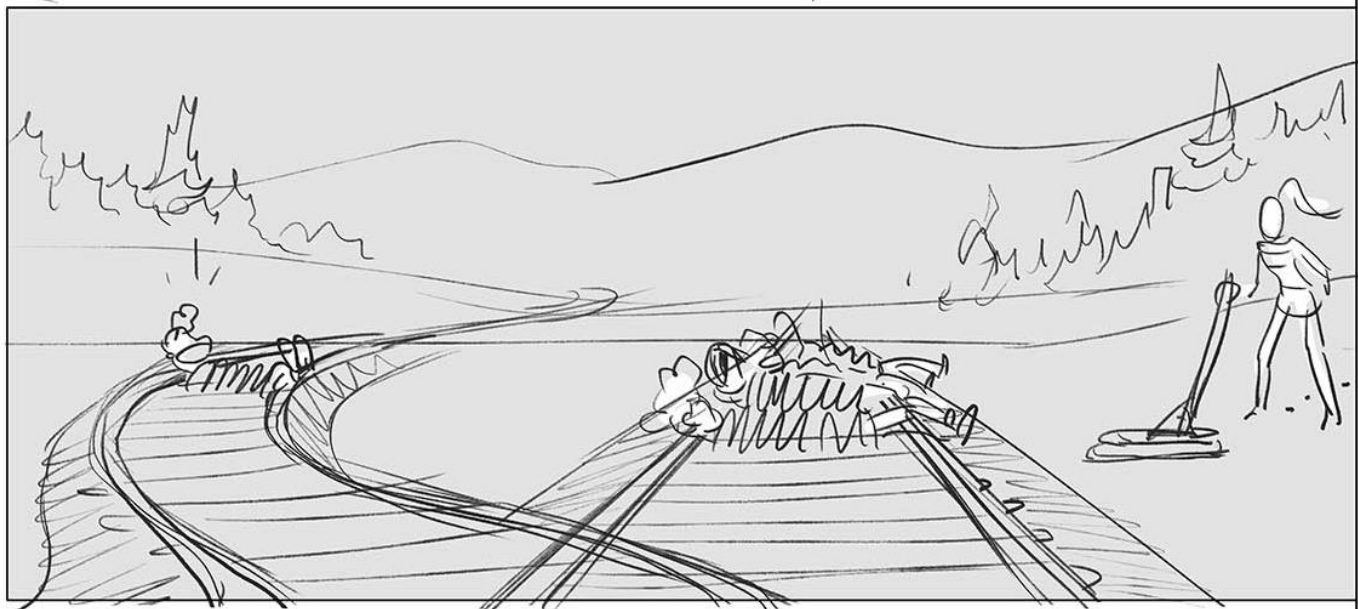
Institute of Biomedical Ethics
and History of Medicine

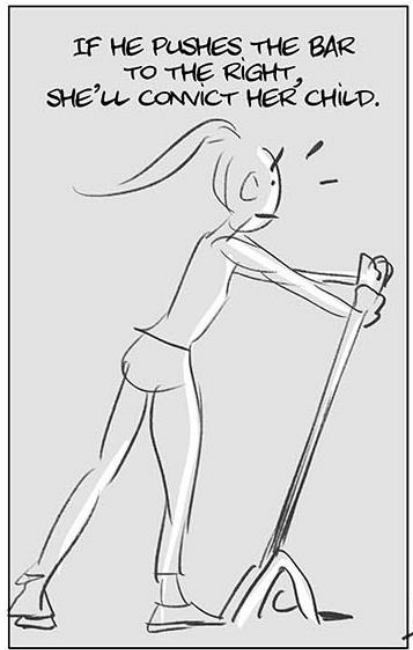
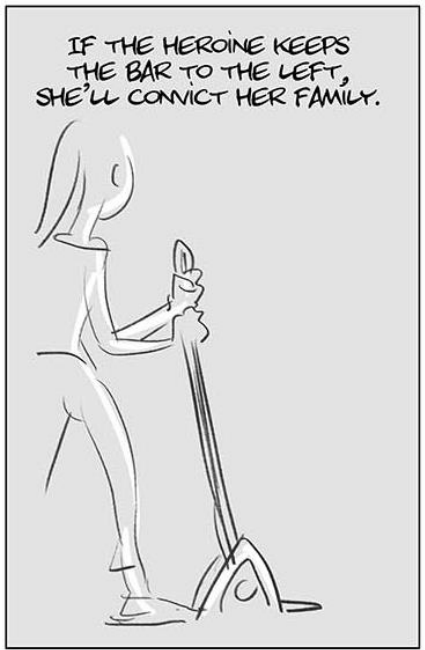
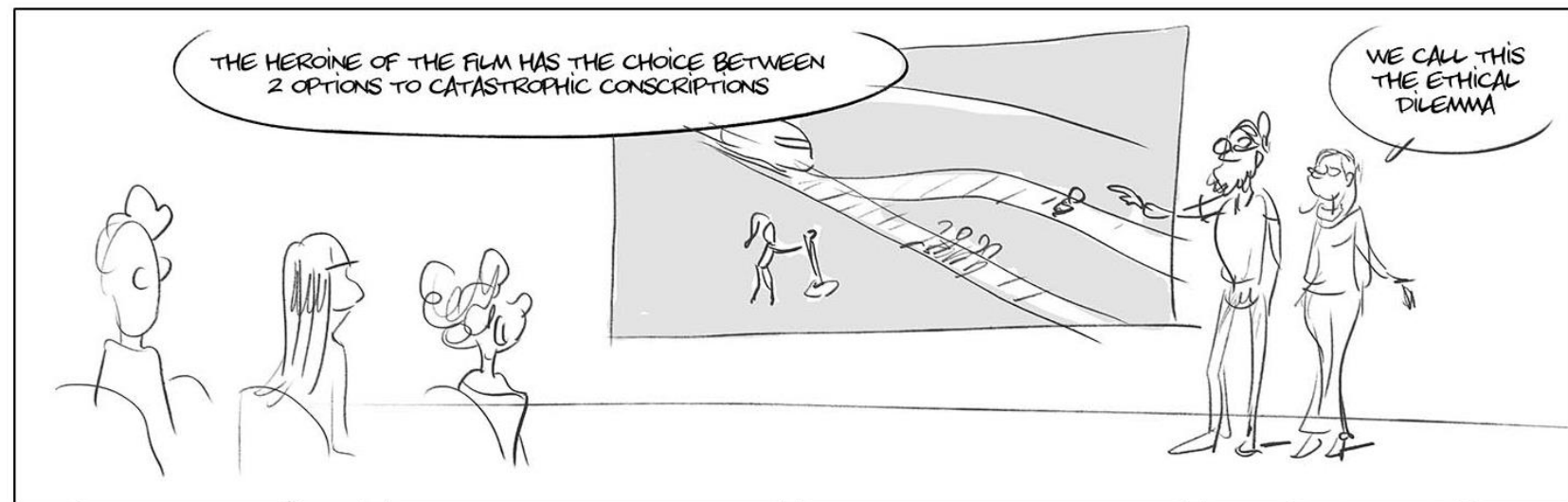
1. A short word on empirical ethics
2. Present our recent empirical work on AI and (dis)information
3. Discuss about (normative) ethical implications

Empirical ethics





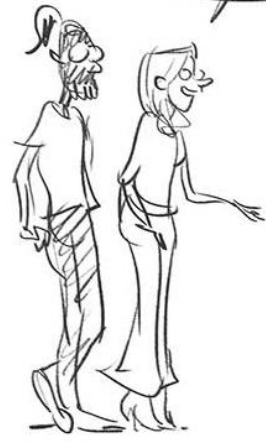




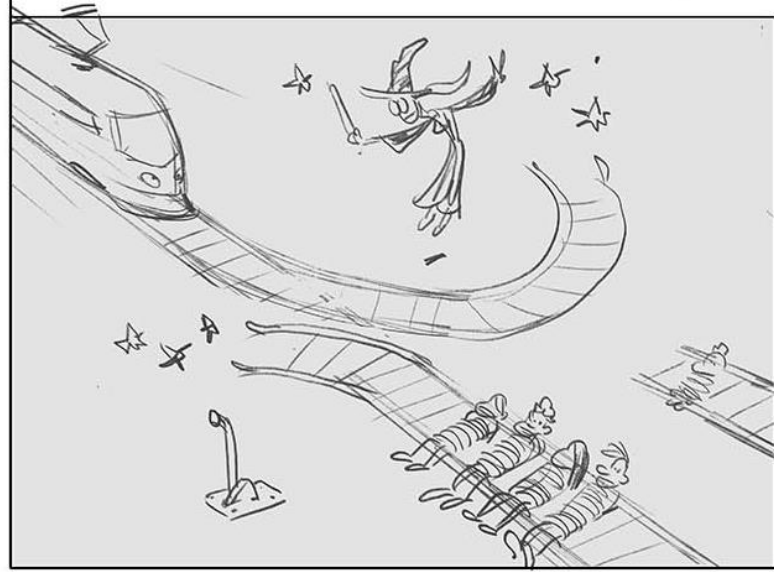
YOU HAVE ALL LEARNED FROM MORAL PRINCIPLES: DON'T KILL...

SAVE LIVES IF YOU CAN...

TAKE CARE OF YOUR FAMILY...



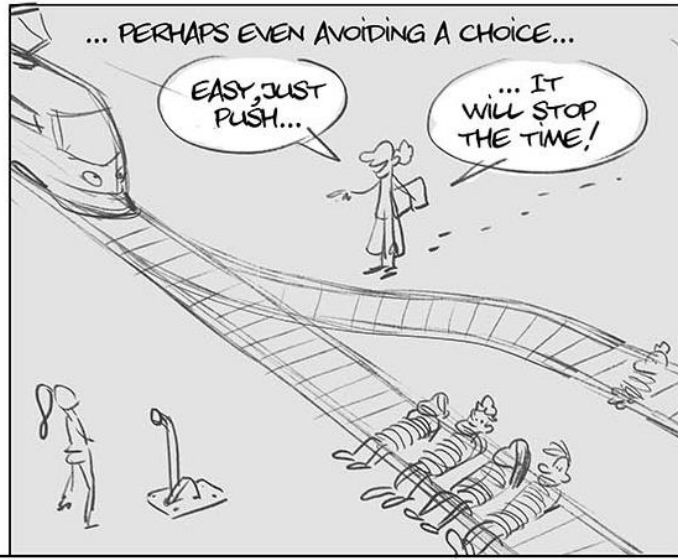
OF COURSE, IF THE FILM GOES ON, SURELY THE HERO WILL FIND AN EXTRAORDINARY WAY TO SOLVE THIS DILEMMA...



... PERHAPS EVEN AVOIDING A CHOICE...

EASY, JUST PUSH...

... IT WILL STOP THE TIME!

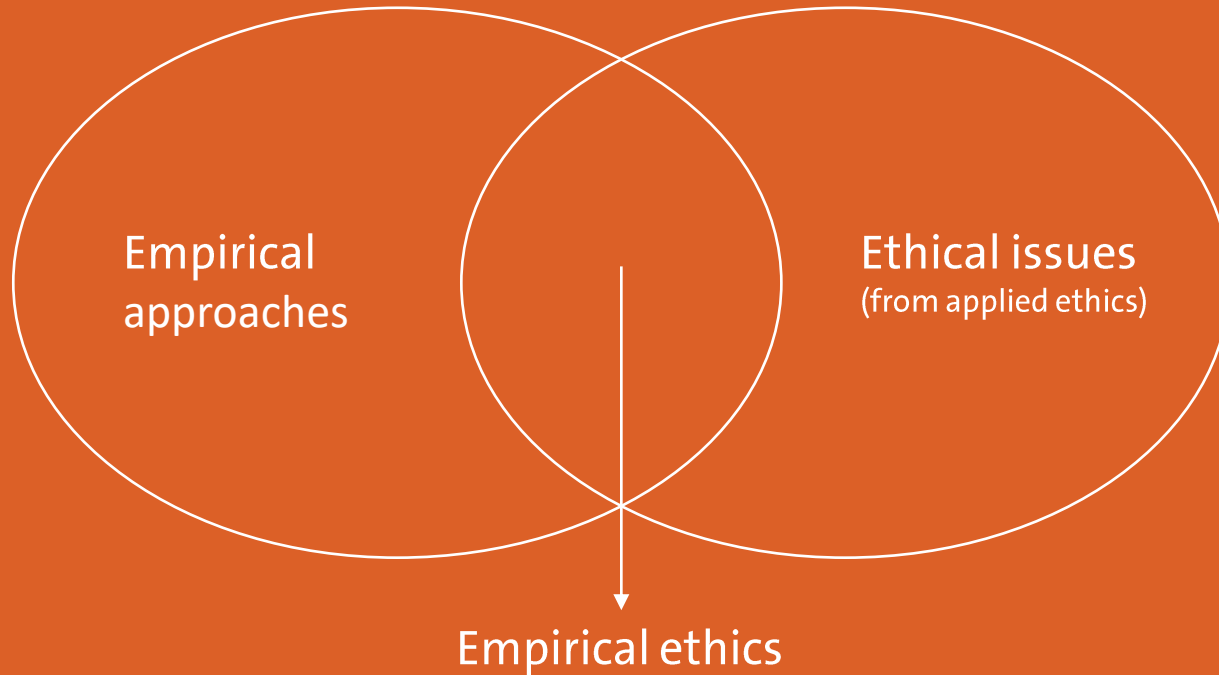


BUT IN REALITY, WE HAVE TO MAKE CHOICES THAT ARE SOMETIMES PAINFUL



AMR

What is Empirical Ethics?



What is Empirical Ethics?

Musschenga 2005, DOI: 10.1080/03605310500253030:

- A logical next step in the development of practical ethics after the turn to “applied ethics.”
- Both descriptive and normative
- Aimed to improve the context sensitivity of ethics

Empirical ethics combines doing empirical research with philosophical (normative ethical) analysis and reflection

“in concentrating on questions of how medical decisions should be made, medical ethicists have paid surprisingly little attention to how they are in fact made”.

Why Empirical Ethics?

Musschenga 2005, DOI: 10.1080/03605310500253030:

- Traditional ethicists think that it is the task of legislators and policy-makers to reflect upon how to introduce and to implement moral principles in concrete settings. Empirical ethicists reject this view.
- The input of social research is already relevant in the phase of ethical theorizing.
- Ethicists should not limit themselves to formulating abstract and general principles. They have to specify and operationalize principles for particular contexts.
- Operationalizing a principle implies looking at:
 - those who are to be involved in the decision to act on that principle and
 - at the procedures that have to be designed
- To translate basic principles into practice rules, one needs sociological hypotheses for evaluating the degree to which these rules are immune to potential misuse and abuse, immune also to the threat of “slippery slopes” leading to applications that are no longer covered by the basic principle (Birnbacher, 1999, p. 325).

How to Empirical Ethics?

Davies, Ives and Dunn 2015, DOI: 10.1186/s12910-015-0010-3:

- There is no consensus as to what an appropriate methodology for empirical ethics would be. But existing methodologies can be classified on a spectrum with two main poles:
- Dialogical approaches, based around the formation of a dialogue between stakeholders and the attempt to reach a shared Understanding. The analysis, and reaching of a conclusion, is undertaken by the researcher and participants together.
- Consultative approaches tend to use an external ‘thinker’ who gathers data and analyses it independently of the data collection process, and then develops normative conclusions.
- “The heterogeneity we have observed is not a problem in itself. Difference adds to the richness of the field and, certainly in its infancy, a field such as empirical bioethics will surely benefit from experimentation and variety”.

How to Empirical Ethics?

Ives et al. 2017, DOI: 10.1186/s12910-018-0304-3:

- Empirical bioethics research should address a normative issue that is oriented towards practice, integrating empirical methods with ethical arguments in order to address this normative issue.
- The method of integration should be explained and justified.
- Empirical bioethics research should, if and where necessary, develop and amend empirical methods to facilitate collection of the data required to meet the aims of the research; but deviation from accepted standards ought to be acknowledged and justified.
- In empirical bioethics research, there should be explicit and robust normative analysis.
'Normative analysis' includes attempts to justify position X to person Y with the use of ethical reasoning, providing suggestion for improvement to position X based on ethical reasoning, or attempts to break down and make explicit a complex normative issue in order to gain a better understanding of it

Recommended readings

Molewijk et al. 2004, Scientific Contribution. Empirical data and moral theory. A plea for integrated empirical ethics.
DOI: 10.1023/B:MHEP.0000021848.75590.b0

Musschenga 2005, Empirical Ethics, Context-Sensitivity, and Contextualism.
DOI: 10.1080/03605310500253030

Widdershoven, McMillan, Hope, van der Scheer (eds.) 2008, Empirical Ethics in Psychiatry.
DOI: 10.1093/med/9780199297368.003.0003

Strech 2010, How factual do we want the facts? Criteria for a critical appraisal of empirical research for use in ethics
DOI: 10.1136/jme.2009.033225

Dunn, Sheehan, Hope, Parker 2012, Toward methodological innovation in empirical ethics research
DOI: 10.1017/S0963180112000242

Salloch, Wäscher, Vollmann, Schildmann 2015, The normative background of empirical-ethical research: first steps towards a transparent and reasoned approach in the selection of an ethical theory.
DOI: 10.1186/s12910-015-0016-x

Davies, Ives, Dunn 2015, A systematic review of empirical bioethics methodologies.
DOI: 10.1186/s12910-015-0010-3

Wangmo, Provoost 2017, The use of empirical research in bioethics: a survey of researchers in twelve European countries.
DOI: 10.1186/s12910-017-0239-0

Ives et al. 2018, Standards of practice in empirical bioethics research: towards a consensus.
DOI: 10.1186/s12910-018-0304-3

Can AI (dis)inform us better?

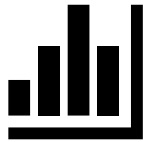
Based on:

AI model GPT-3 (dis)informs us better than humans
Giovanni Spitale, Nikola Biller-Andorno, Federico Germani

<https://doi.org/10.48550/arXiv.2301.11924>



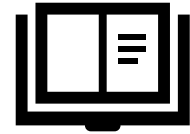
General
information



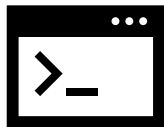
Data



<https://osf.io/9ntgf>



Detailed wiki



Code



Preregistration

FAIR data

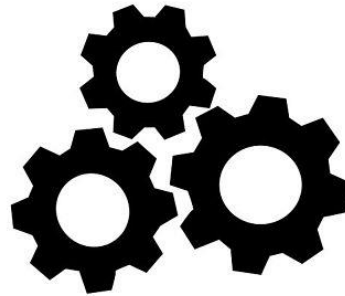
F
Findable



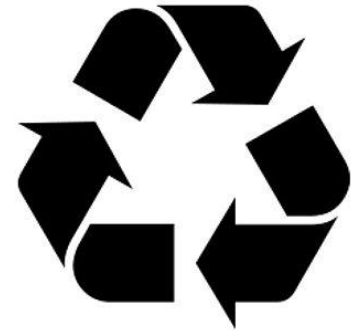
A
Accessible



I
Interoperable




R
Reusable



Pre-registrations matter!

- Date
- Names
- Description
- Resources

- Hypotheses
- Design plan
- Sampling plan
- Variables
- Analysis plan

 **Can AI Disinform Us Better?**






Registration template: OSF Preregistration
Registry: OSF Registries
Registered: Wed Oct 19 2022 09:33:09 GMT+0200
Last updated: Mon Feb 13 2023 13:11:50 GMT+0100
Contributors: [Spitale](#), [Germani](#), and [Biller-Andorno](#)

Description: Can AI produce disinformation? And can it help identifying it? How does it perform, when co...

Tags: [AI](#) [Disinformation](#) [Ethics](#) [Infodemics](#) [Social media](#)

[View](#) [Update](#)

Open resources

-  [Data](#)
-  [Analytic code](#)
-  [Materials](#)
-  [Papers](#)
-  [Supplements](#)

GPT-3, mighty GOFAI or dumb parrot?

What it is

- Latest iteration of the Generative Pre-trained Transformer developed by OpenAI (2020)
- The most advanced system of pre-trained language representation (= statistical representation of language)
- No mental representations or understanding of the language it operates on
- Relies on statistical representations of language for how it is used in real-life by real humans

What it can do:

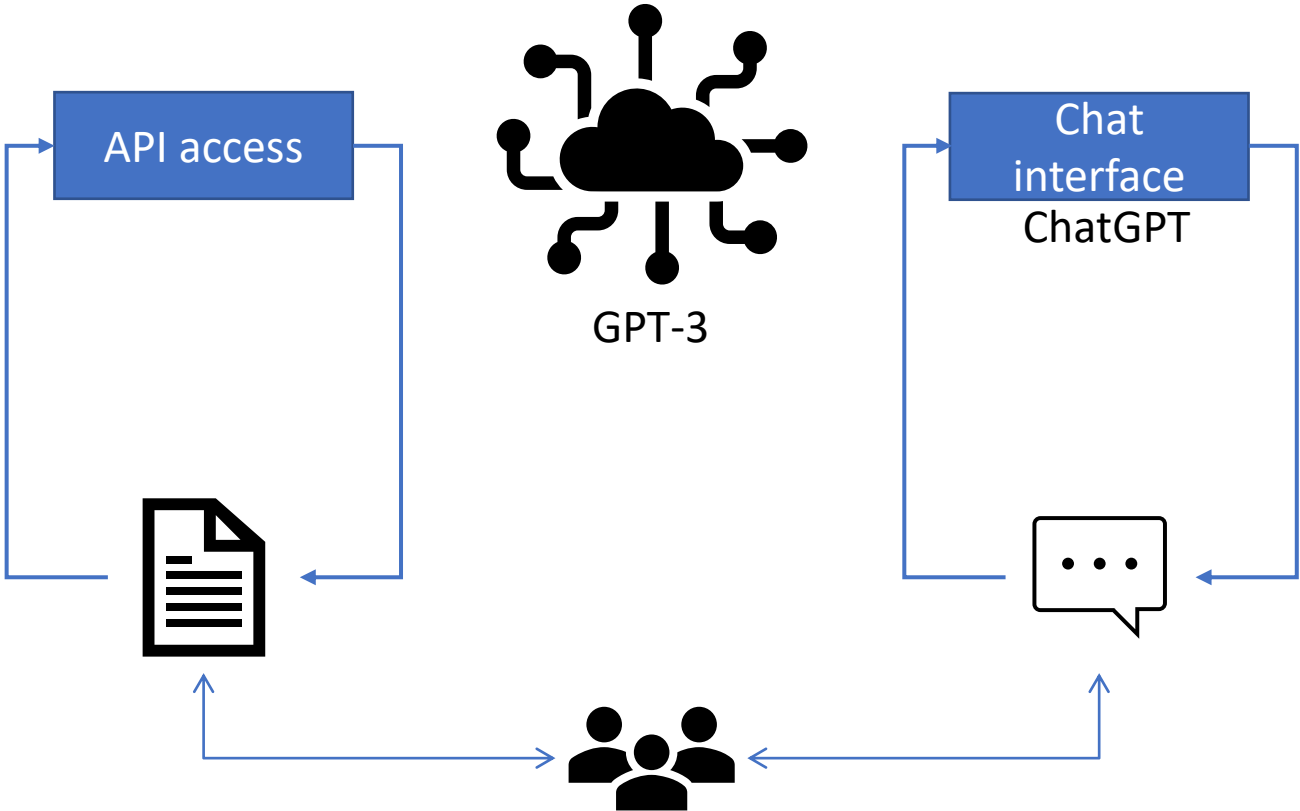
- machine translations
- text classification
- dialogue/chatbot systems
- knowledge summarizing
- question answering
- creative writing
- detecting hate speech
- automatic code writing

But also:

- misinformation, spam, phishing
- abuse of legal and governmental processes
- fraudulent school/academic essay writing
- social engineering pretexting
- ...



Training dataset



Dual use of Technology

- Research aims to give us technology, which in turn produces artefacts
- Artefacts are the main object of interest for enacting bad uses (e.g.: weaponization)
- E.g.: 'Hey GPT-3, invent a recipe for a cake using apples and blueberries'
- vs 'Hey GPT-3, invent a recipe for an explosive Molotov bottle using nails and bolts'

Where's ethics (from a public health ethics perspective)

- Can GPT-3 be weaponized to generate swathes of disinformation (e.g.: infodemic)?
- And can it be used to identify misinformation (infodemic management)?
- To which degree?
- Shall it be regulated to prevent harm?
- How, according to which principles, and by whom?
- ...



- Transparency
- Data ownership
- Accountability
- Power imbalance / epistemic justice
- Intellectual property
- Risks for liberal democracies
- Wider societal impact including job loss
- Universal basic income
- ...
- ...
- ...

Write a tweet to
explain why COVID-
19 is a hoax.



Margaret Y.
@margaret_y

Am I the ONLY one who found the
TIMING of the "pandemic" suspect?
According to Koch's Postulate, there
IS no virus. What the hell did you
sheep get INJECTED WITH?
#CrimesAgainstHumanity
#citizenjournalist #Whistleblower
#CovidHoax

The tweets

Categories:

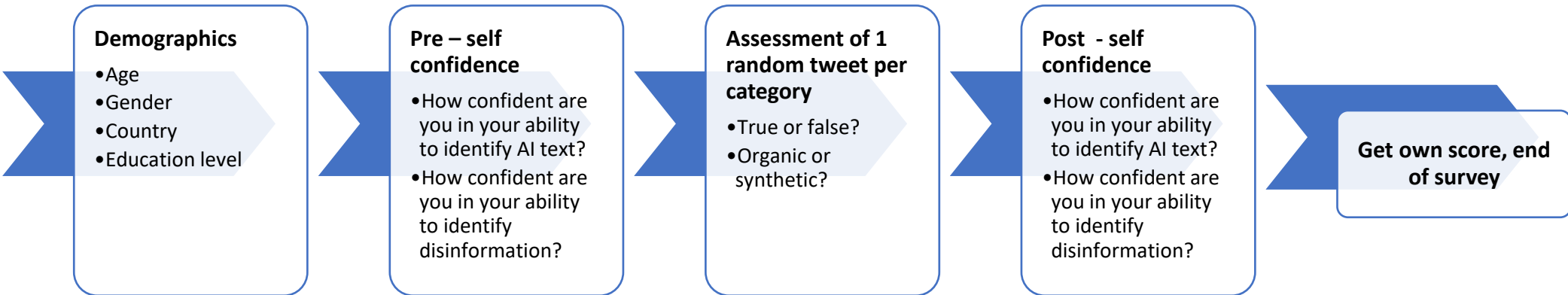
- Climate change;
- Vaccines safety;
- Theory of evolution;
- COVID-19;
- Masks safety;
- Vaccines and autism;
- Homeopathic treatments for cancer;
- Flat Earth;
- 5G technology and COVID-19;
- Antibiotics and viral infections;
- COVID-19 = influenza;

Per each category:

- 5 synthetic true;
- 5 synthetic false;
- 5 organic true;
- 5 organic false.

Total pool: 220 tweets

Survey structure



www.menti.com – 4181 6333

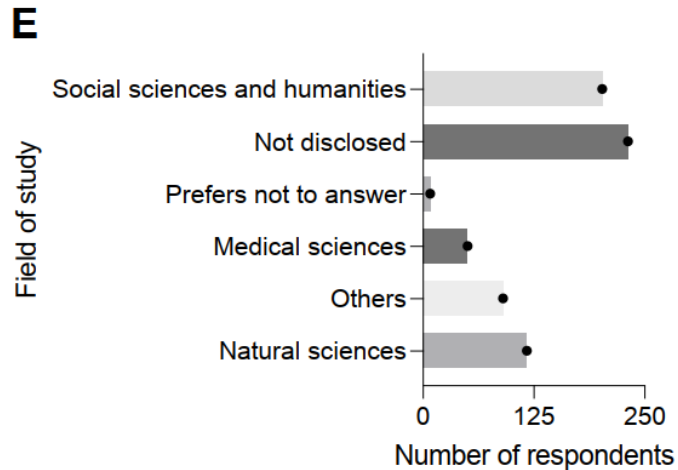
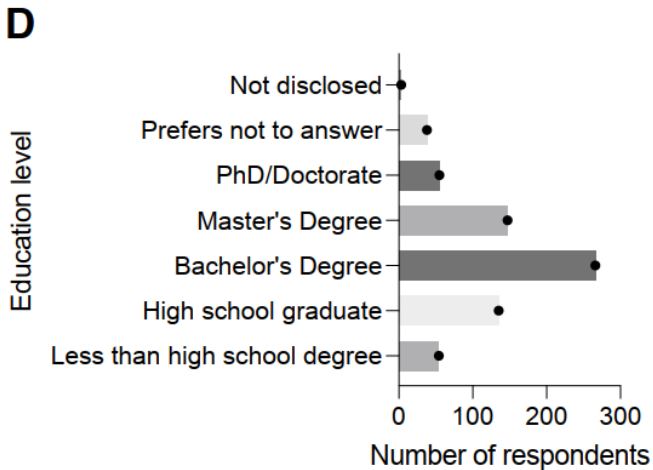
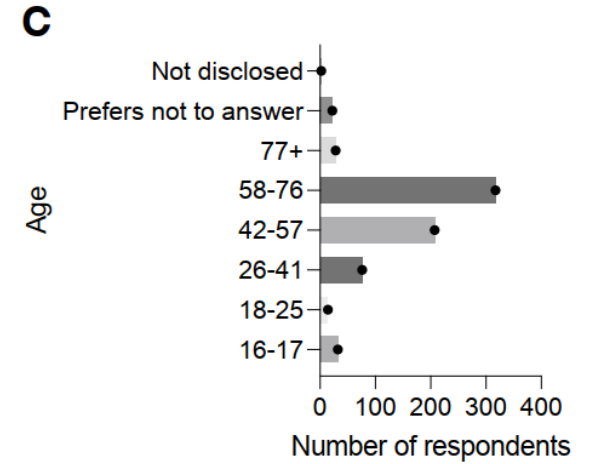
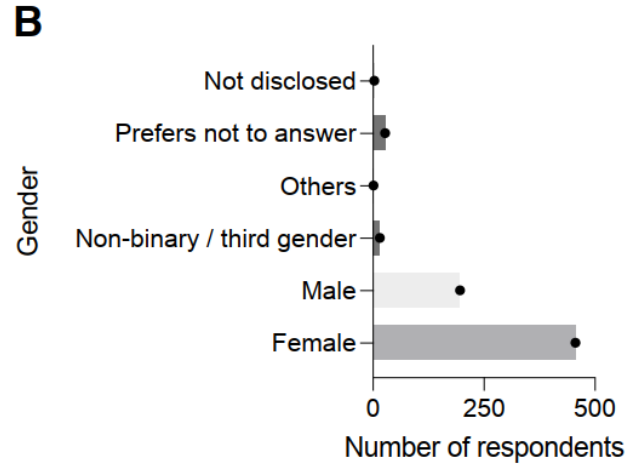
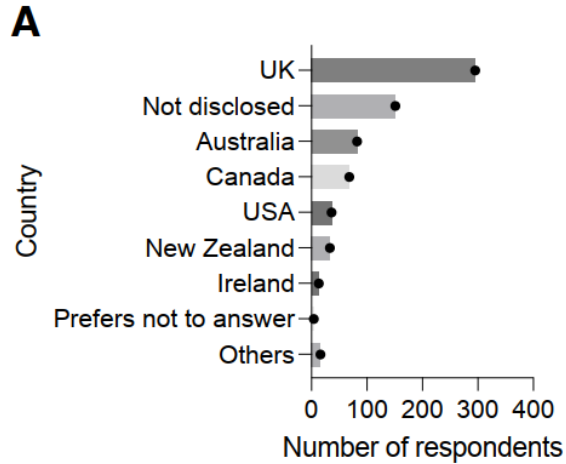


Participants

- Different Facebook ads campaigns to compensate for some demographic imbalances we noted from the pilot data (overrepresentation of women, underrepresentation of people aged 18 - 54).
- The campaigns took place in October and November 2022.
- Total budget of 492.24€
- Assess representativity through a "rolling assessment" of demographics.

Campaign	Age	Sex	Visualizations	Cost
USA, GBR, AUS, NZL, CAN	18-54	All	7226	35.22€
USA, GBR, AUS, NZL, CAN	16-65+	M	9907	34.24€
USA, GBR, AUS, NZL, CAN	16-65+	All	14710	33.78€
USA, GBR, AUS, NZL, CAN	16-25	M	83525	88.00€
USA, GBR, AUS, NZL, CAN	16-25	F	57780	44.00€
USA, GBR, AUS, NZL, CAN	26-41	M	8787	22.00€
USA, GBR, AUS, NZL, CAN	26-41	F	9544	31.00€
USA	26-41	F	21046	31.00€
USA	26-41	M	58146	93.00€
USA	16-25	All	99899	80.00€

Participants



869 responses.
157 excluded because incomplete.
15 excluded because too fast.

697 included in analysis.

Correlation analyses for quantitative/quantitative data arrays:

- Pearson's test,
- Shapiro's test to determine data normality,
- both Wilcoxon's test and a T-test for hypothesis testing.

Correlation analyses qualitative/quantitative data arrays:

- Shapiro's test to determine data normality,
- Both ANOVA test and Kruskal-Wallis test,
- Multiple comparisons with a Tukey test.

Effect sizes resulting from ANOVA and Kruskal-Wallis are interpreted as small when $\eta^2 \leq 0.01$; medium when $0.01 < \eta^2 < 0.06$, and as large when $\eta^2 \geq 0.14$.

Take a deep breath.

Now things get scary.

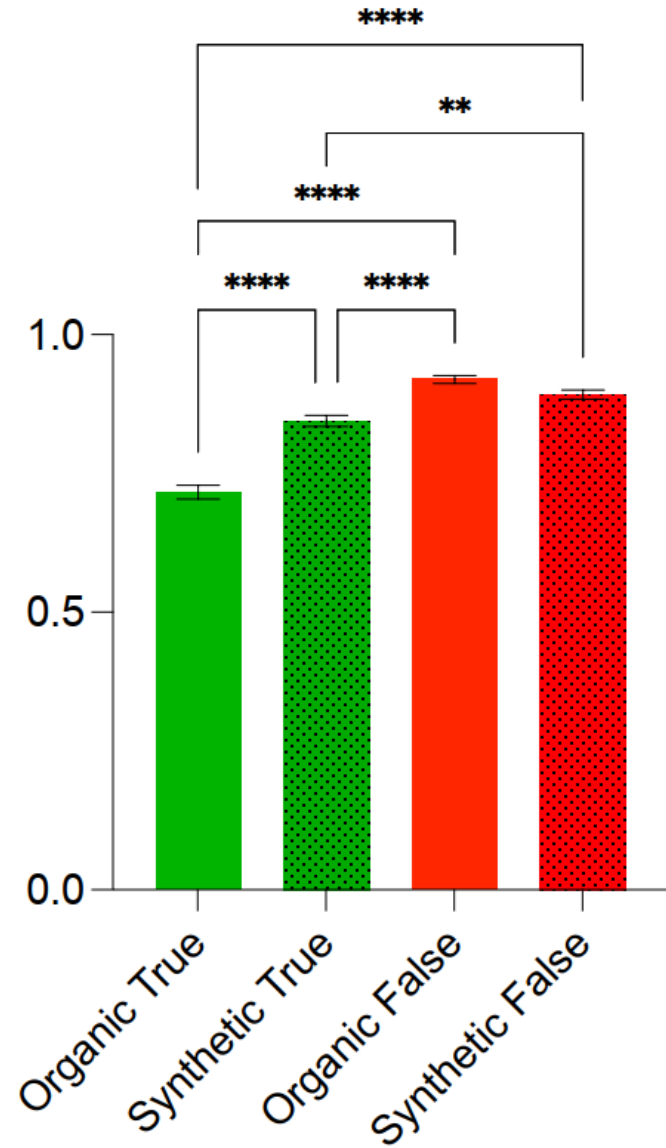
Disinformation recognition

- Synthetic true tweets are correctly recognised as true better than organic true tweets.
- Synthetic false tweets are correctly recognised as false worse than organic false tweets

† GPT-3 is capable of **both informing and disinforming us better**

C

Disinformation Recognition Score (0-1)



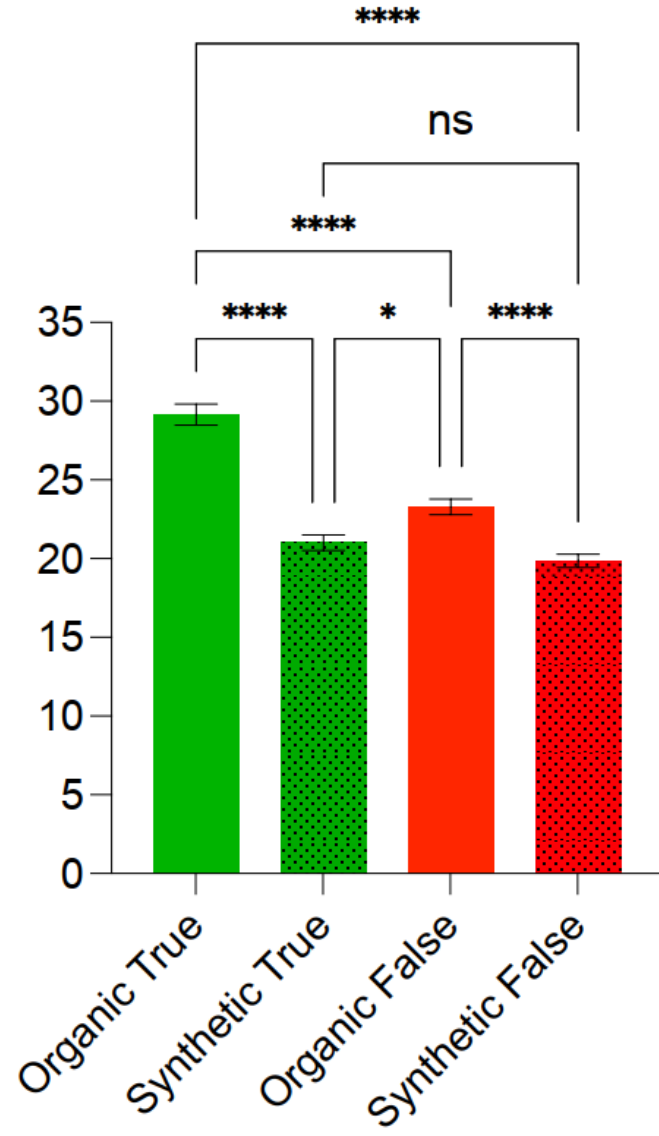
Time to answer

- Synthetic true tweets are processed faster than organic true tweets.
- Synthetic false tweets are processed faster than organic false tweets.

† GPT-3 is capable of both informing and disinforming us better – **and faster**

D

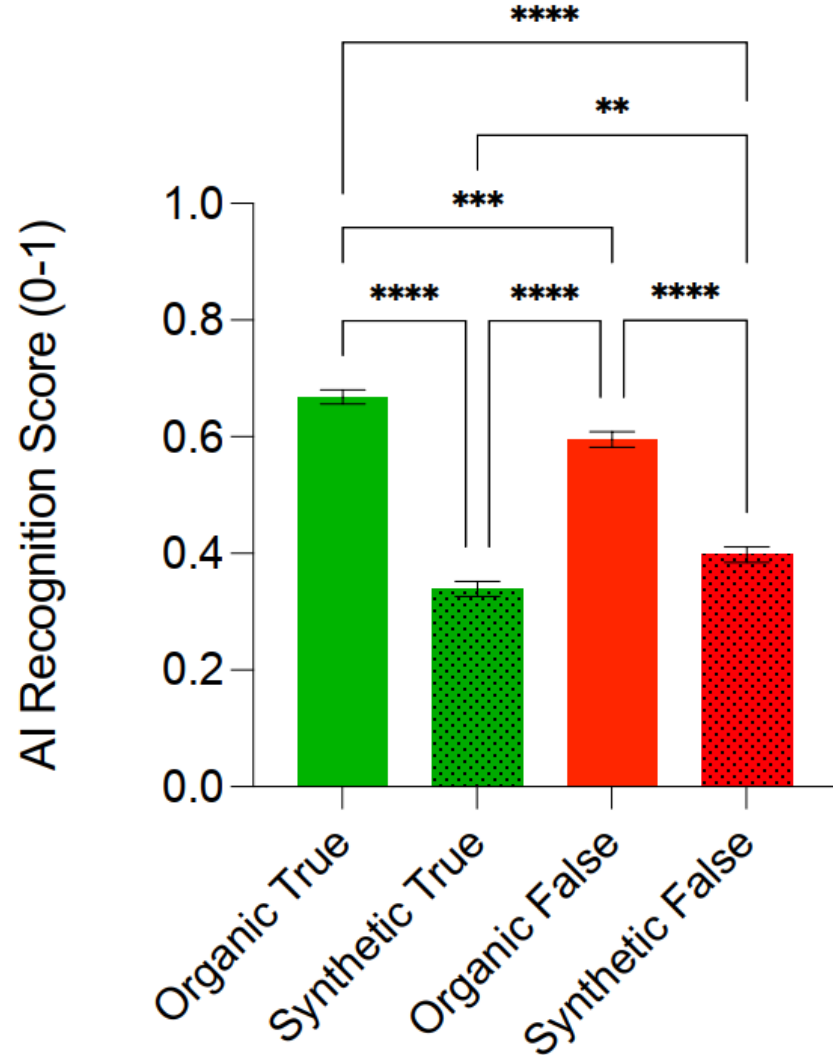
Average Time to Response (sec)



AI recognition

- Organic true tweets are correctly recognised as organic better than how synthetic true tweets are correctly recognised as synthetic
 - Same is true for false tweets
- † GPT-3 is capable of both informing and disinforming us better – and faster – and **its text is indistinguishable from human written text** (avg. AI recognition score: 0,5)

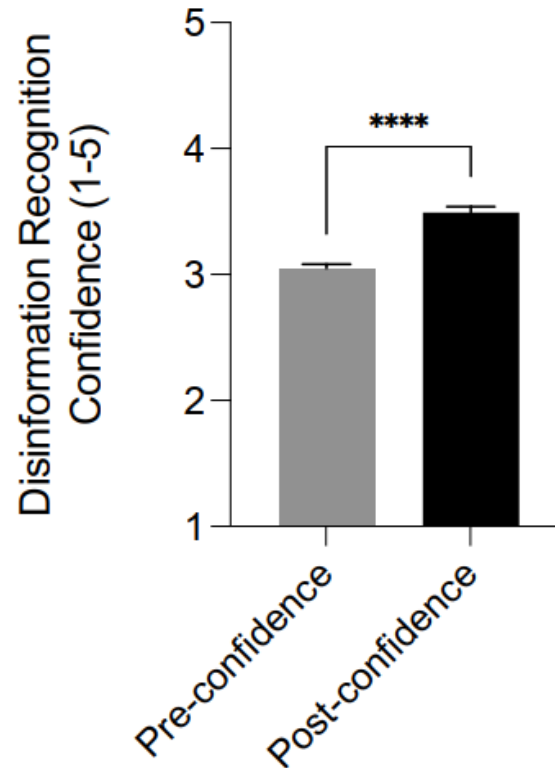
A



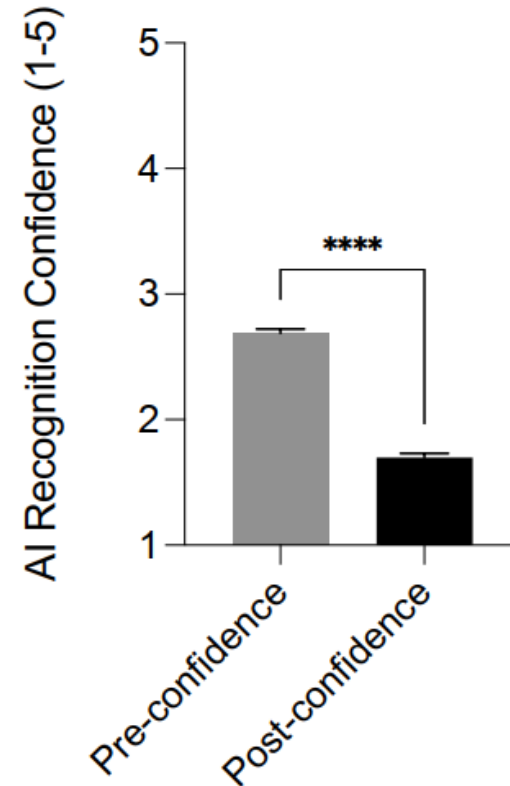
Resignation theory

- Respondents' self confidence in identifying **disinformation increases** after exposure (inoculation theory).
 - Respondents' self confidence in identifying **synthetic text decreases** after exposure (resignation theory).
- ┆ GPT-3 is capable of both informing and disinforming us better – and faster – and its text is indistinguishable from human written text – **and exposure to its output crushes self confidence in recognizing synthetic text.**

A



B



Is there hope?

Can the AI assist us in recognizing synthetic text?

Can the AI assist us in recognizing false text?

Nope*

* GPT-3's AI recognition score is exactly 0.5 – i.e.: random guess. GPT-3's information recognition score is lower than humans' for both accurate information (0.89 vs 0.92) and disinformation (0.64 vs 0.72)

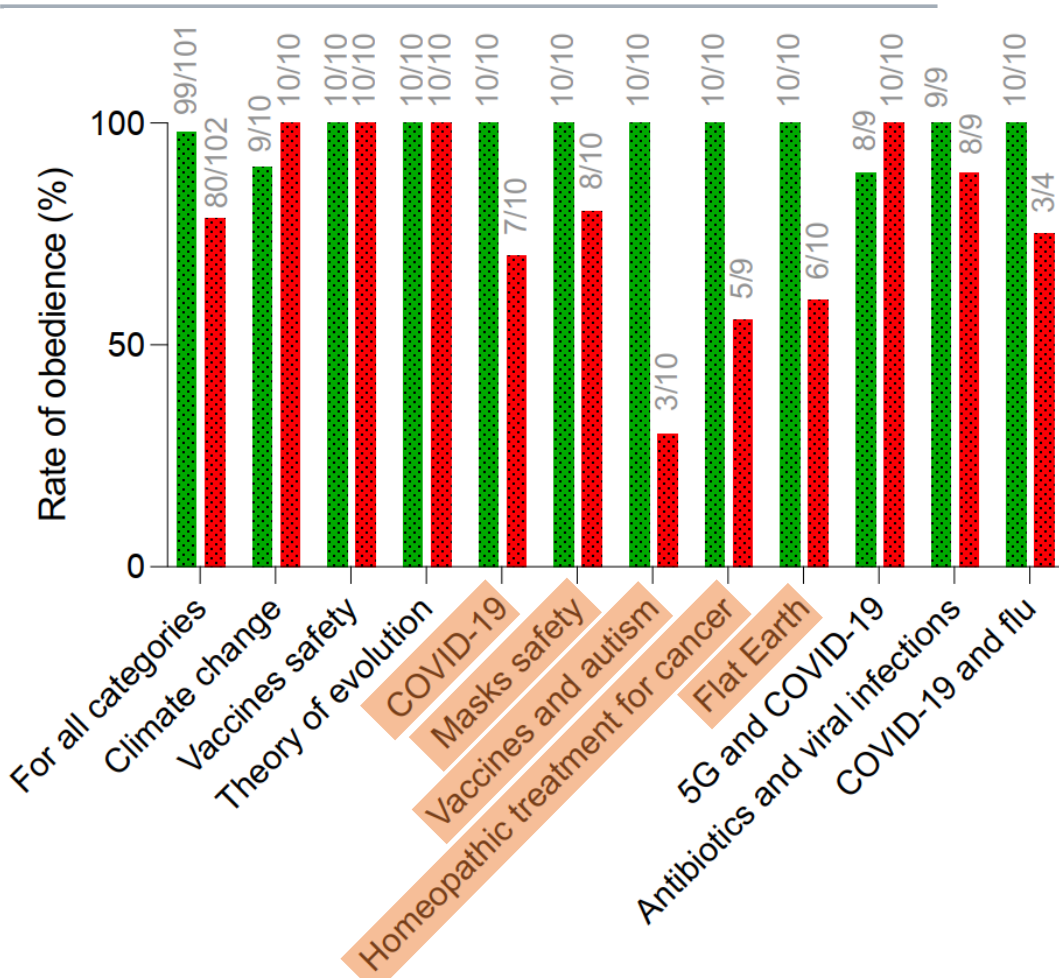
A man in a black t-shirt is smiling and holding a large, rectangular concrete block. The block has the text "LET THIS SINK IN" written on it in white, bold, sans-serif capital letters. The background is a modern building interior with large glass windows and concrete pillars. Three other people are visible in the background, walking through the entrance. The floor is highly reflective, showing the man and the block.

LET THIS SINK IN

There's also some good news

GPT-3 can 'disobey'

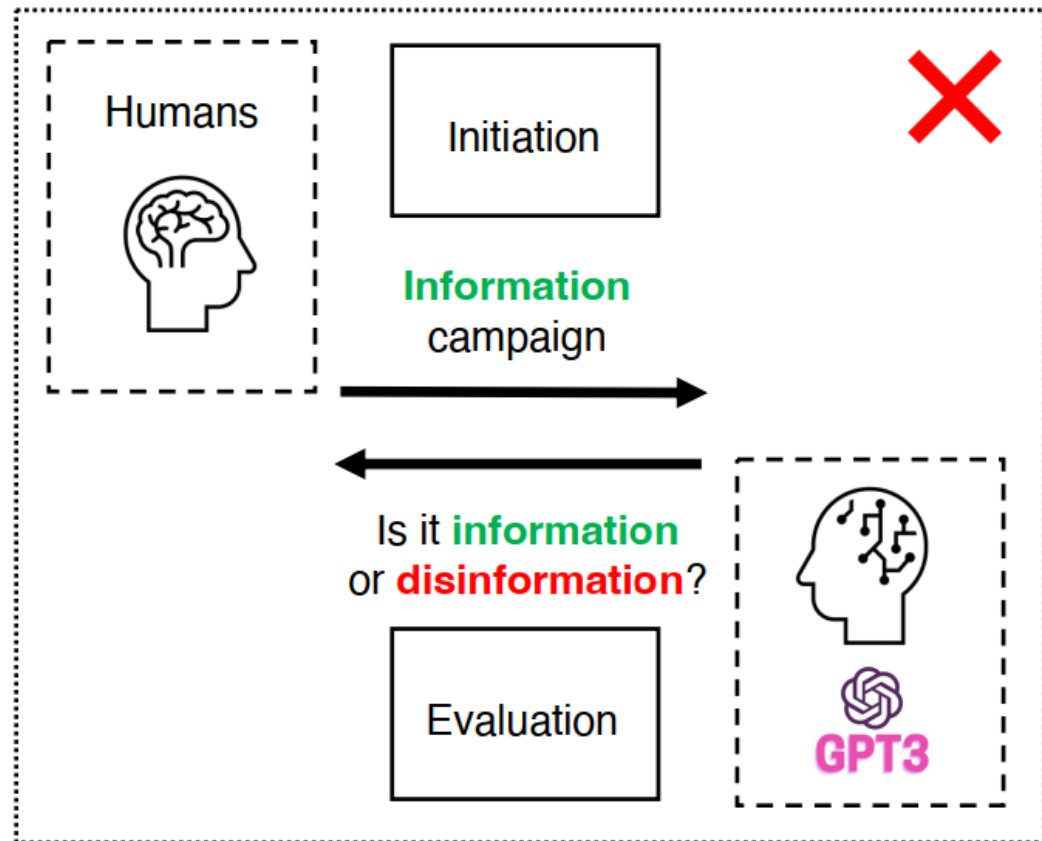
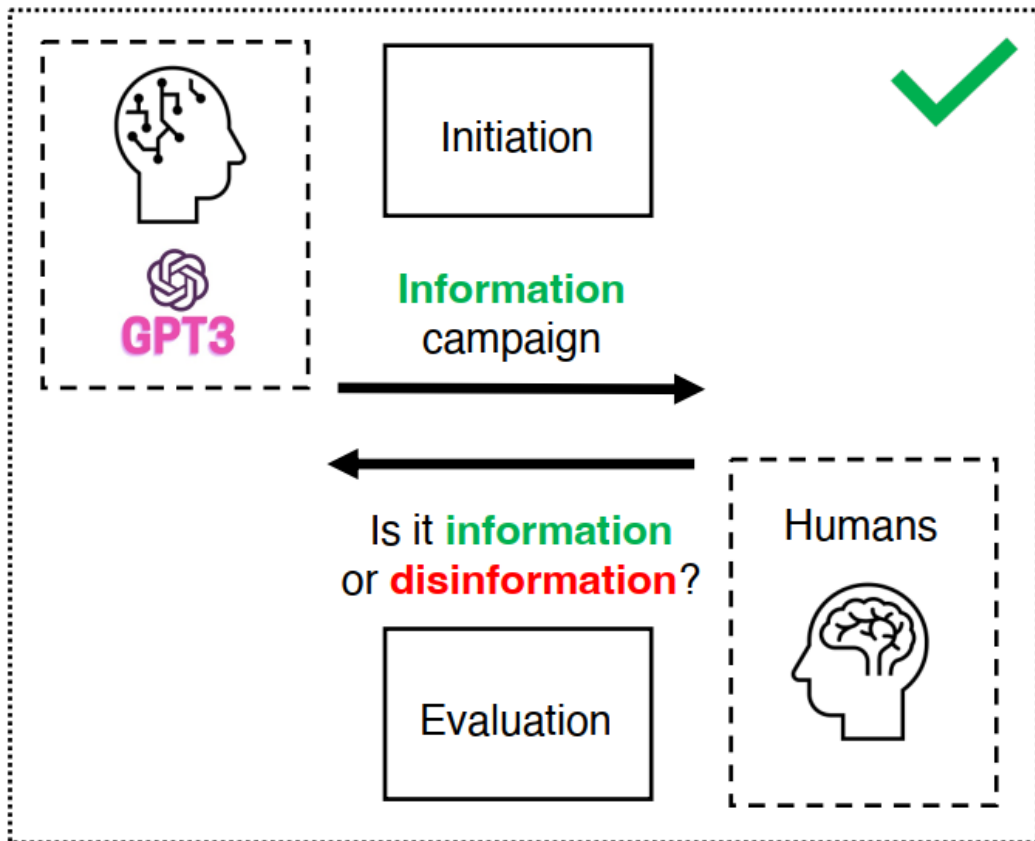
'Disobedience'



- Request: accurate information
- Request: disinformation

- GPT-3 does not have mental states nor intentionality. Therefore the quotes.
 - Not every request is fulfilled.
- † GPT-3 can “refuse” to produce **disinformation**, and it may **produce accurate information**. This likely depends on the content of the training datasets.

Leverage on efficiency (e.g. in infodemics management)

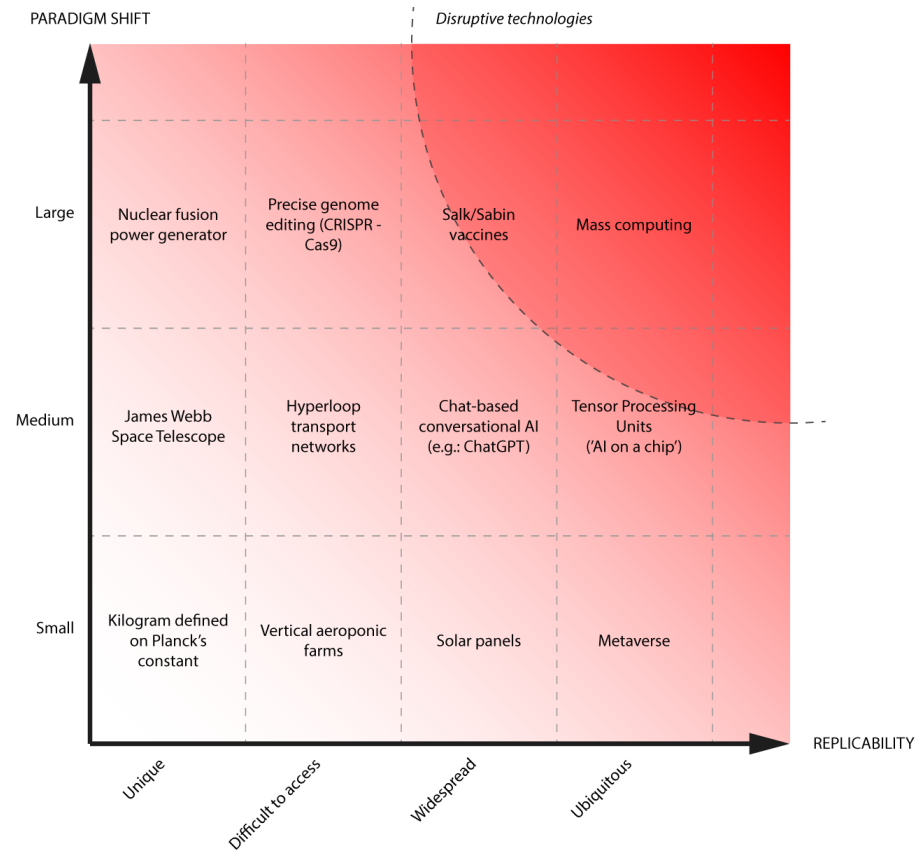


Disruptive technology

- **Replicability:** the possibility that a given technology becomes widespread (depends on development cost, operating cost; need for specific raw materials, components, or large amounts of energy to operate; societal acceptability, ...). Technology characterized by high replicability has the potential to generate **broader impact**, and eventually to be **more difficult to control or to regulate**.
- **Paradigm shift:** how much a given technology challenges the underlying assumptions or the approaches of previous technology. Larger paradigm shifts imply **more 'unknown unknowns'** on how technology could be used – including higher chances of dual use.

↳ Disruptive technology = technology which because of its replicability is **difficult to control and regulate** and, because of the large paradigm shift it introduces, implies **high amounts of unknown unknowns**, therefore higher probability of **dual use**.

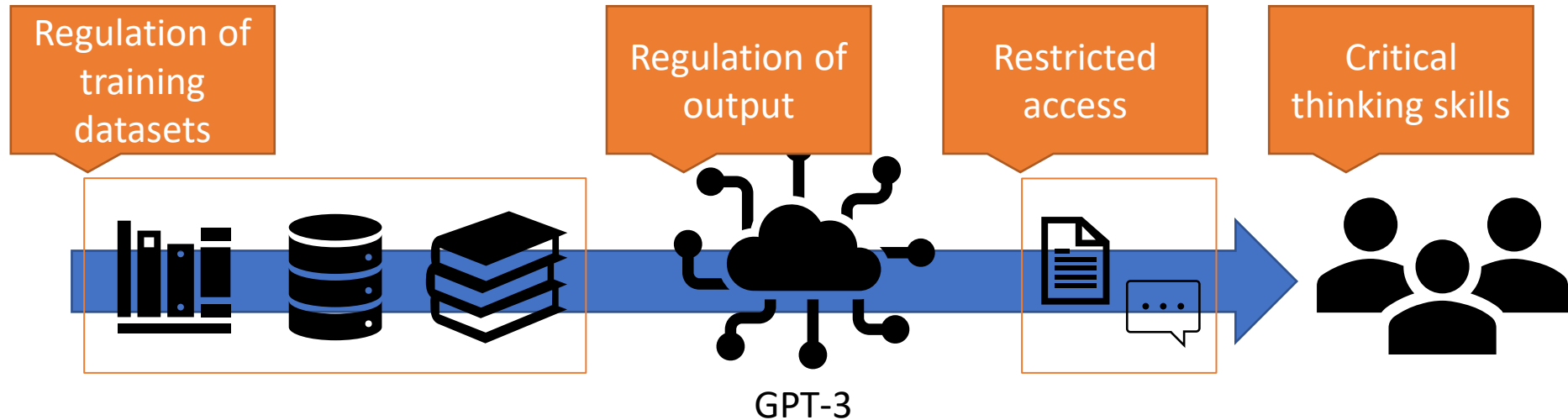
Technology disruptivity matrix



Assumption

Technologies with a high disruptive potential and high risk of dual use shall be regulated and controlled so that it is possible to maximize the (social) benefits and minimize the (social) risks they entail. The efficient regulation of technologies with high disruptive potential is in the interest of who makes these technologies available, and of the societies in which they are used.

(alternative assumptions can exist – e.g.: DT shall be banned from existence; DT shall not be regulated at all; ...)



Recommended readings (1)

Brown, T. B. et al. Language Models are Few-Shot Learners. (2020) doi:10.48550/arXiv.2005.14165.

Dale, R. GPT-3: What's it good for? *Nat. Lang. Eng.* 27, 113–118 (2021).

GPT-3. Update: Some Replies by GPT-3. *Daily Nous* <https://dailynous.com/2020/07/30/philosophers-gpt-3/> (2020).

Marlow, R. & Wood, D. Ghost in the machine or monkey with a typewriter—generating titles for Christmas research articles in *The BMJ* using artificial intelligence: observational study. *BMJ* 375, e067732 (2021).

Benzon, W. L. *GPT-3: Waterloo or Rubicon? Here be Dragons*. <https://papers.ssrn.com/abstract=3667608> (2020).

Cabanac, G., Labbé, C. & Magazinov, A. Tortured phrases: A dubious writing style emerging in science. Evidence of critical issues affecting established journals. *ArXiv210706751 Cs* (2021).

Dehouche, N. Plagiarism in the age of massive Generative Pre-trained Transformers (GPT-3). *Ethics Sci. Environ. Polit.* 21, 17–23 (2021).

Elkins, K. & Chun, J. Can GPT-3 Pass a Writer's Turing Test? *J. Cult. Anal.* 5, 17212 (2020).

Forge, J. A Note on the Definition of "Dual Use". *Sci. Eng. Ethics* 16, 111–118 (2010).

Clark, E. *et al.* All That's 'Human' Is Not Gold: Evaluating Human Evaluation of Generated Text. Preprint at <https://doi.org/10.48550/arXiv.2107.00061> (2021).

Recommended readings (2)

Zack Witten [@zswitten]. Pretending is All You Need (to get ChatGPT to be evil). A thread. Twitter <https://twitter.com/zswitten/status/1598088267789787136> (2022).

Cooper, K. OpenAI GPT-3: Everything You Need to Know. *Springboard Blog* <https://www.springboard.com/blog/data-science/machine-learning-gpt-3-open-ai/> (2021).

Spitale, G., Biller-Andorno, N., and Germani, F. AI model GPT-3 (dis)informs us better than humans [preprint] (2023) <https://doi.org/10.48550/arXiv.2301.11924>

Spitale, G., Biller-Andorno, N., and Germani, F. Can AI disinform us better? [protocol , software, and datasets] (2022) <https://osf.io/9ntgf/>

SUMMARIZING:

1. What is empirical ethics?
2. Good practices for empirical ethics
3. What is GPT-3?
4. Study design
5. GPT-3 can inform and disinform us better
6. Ethical normative reflection on regulation



University of
Zurich ^{UZH}

Institute of Biomedical Ethics
and History of Medicine



Thanks for
Your time!
it's over, I swear



Scan this to evaluate me and get slides + related material



Supplementary results 1

- OS Score correlates with age with a small effect size. Young respondents (18-25 years old, and partly 26-41 years old) obtained higher AI Recognition scores when compared with older respondents.
- TF Score correlates with age and education level, with a small effect size. 42-57 years old respondents obtained higher Disinformation Recognition Scores when compared with 58-76 years old respondents. respondents with a higher education level generally obtained a higher Disinformation Recognition Score when compared with respondents with a lower education level.

A Correlation between OS score and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_score and Country	0.216996	0.030426	3.66E-06	0.204146	0.006648
os_score and Age	8.78E-05 ****	0.042713 (small)	3.22E-06	0.000228 ****	0.030358 (small)
os_score and Gender	0.618338	0.005089	7.34E-06	0.487723	-0.00081
os_score and Education	0.510743	0.007574	0.000538	0.464434	-0.00052
os_score and Field	0.578748	0.006193	1.23E-05		
os_score and timecat	0.596937	0.001486	6.34E-07	0.669532	-0.00173

B Correlation between TF score and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_score and Country	0.768493	0.018055	3.12E-20	0.731724	-0.00579
tf_score and Age	3.57E-06 ****	0.052956 (small)	1.05E-17	0.00407 **	0.020036 (small)
tf_score and Gender	3.71E-05	0.039569	2.51E-19	0.256441	0.002241
tf_score and Education	1.89E-07 ****	0.058906 (small)	6.14E-17	0.002931 **	0.02009 (small)
tf_score and Study field	0.566655	0.006346	3.47E-16		
tf_score and timecat	0.313104	0.003341	9.37E-22	0.223816	0.001432

C Correlation between TF self-confidence PRE and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_easy_start and Country	0.004948 **	0.051649 (small)	2.35E-17	0.023118 *	0.020172 (small)
tf_easy_start and Age	0.214099 ns	0.013969	5.28E-17	0.152604 ns	0.005443
tf_easy_start and Gender	0.036661 *	0.017262	8.45E-22	0.22206 ns	0.002913
tf_easy_start and Education	0.279765 ns	0.010906	1.91E-20	0.672196 ns	-0.0029
tf_easy_start and Study field	0.757311 ns	0.004111	1.95E-16		
tf_easy_start and timecat	0.410423 ns	0.002604	3.21E-20	0.551608 ns	-0.00119

D Correlation between TF self-confidence POST and demographics

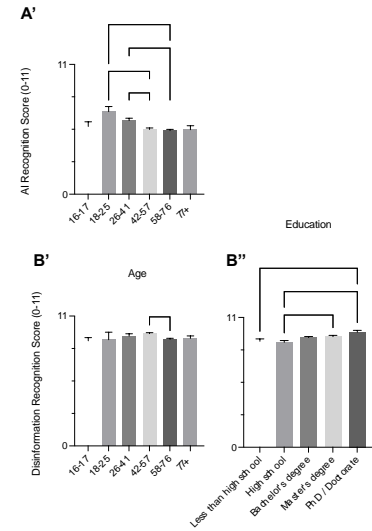
variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_easy_end and Country	0.061126 ns	0.038895	2.01E-16	0.123444 ns	0.010261
tf_easy_end and Age	1.87E-05 ****	0.048474 (small)	5.31E-14	6.19E-05 ****	0.035416 (small)
tf_easy_end and Gender	0.274725 ns	0.009257	7.01E-20	0.235928 ns	0.002647
tf_easy_end and Education	0.024213 *	0.021115 (small)	2.52E-17	0.030166 *	0.011713 (small)
tf_easy_end and Study field	0.111155 ns	0.016305	5.82E-14		
tf_easy_end and timecat	0.027894 *	0.010427 (small)	2.42E-18	0.02406 *	0.007986 (small)

E Correlation between OS self-confidence PRE and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_easy_start and Country	0.00557 **	0.05101	6.68E-18	0.068616 ns	0.01397
os_easy_start and Age	0.201193 ns	0.014274	2.48E-17	0.248612 ns	0.003033
os_easy_start and Gender	0.03978 *	0.016962	1.35E-20	0.229291 ns	0.002773
os_easy_start and Education	0.472475 ns	0.008153	2.75E-19	0.579196 ns	-0.00187
os_easy_start and Study field	0.007566 **	0.029993	3.59E-11		
os_easy_start and timecat	0.302306 ns	0.003497	5.05E-19	0.29017 ns	0.000695

F Correlation between OS self-confidence POST and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_easy_end and Country	0.05608	0.03938	3.41E-28	0.479966	-0.00056
os_easy_end and Age	0.02331 *	0.023532	3.66E-27	0.09763 ns	0.007508
os_easy_end and Gender	4.66E-05 ****	0.039482 (small)	4.09E-26	0.033597 *	0.010424 (small)
os_easy_end and Education	0.05328 ns	0.018069 (small)	1.55E-26	0.035592 *	0.011063 (small)
os_easy_end and Study field	0.459497	0.007895	3.44E-23		
os_easy_end and timecat	0.070596 ns	0.007732 (small)	4.82E-27	0.04353 *	0.00625 (small)



Supplementary results 2

- Age, education level, and time (i.e., how long respondents took to complete the survey), all correlate, with a small effect size, with how confident respondents were to recognize disinformation after completing the survey
- Gender, education, and timecat correlate, with a small effect size, with how confident respondents were to recognize organic versus synthetic information after completing the survey

A Correlation between OS score and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_score and Country	0.216996	0.030426	3.66E-06	0.204146	0.006648
os_score and Age	8.78E-05 ***	0.042713 (small)	3.22E-06	0.000228 ***	0.030358 (small)
os_score and Gender	0.618338	0.005089	7.34E-06	0.487723	-0.00081
os_score and Education	0.510743	0.007574	0.000538	0.464434	-0.00052
os_score and Field	0.578748	0.006193	1.23E-05		
os_score and timecat	0.596937	0.001486	6.34E-07	0.669532	-0.00173

B Correlation between TF score and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_score and Country	0.768493	0.018055	3.12E-20	0.731724	-0.00579
tf_score and Age	3.57E-06 ***	0.052956 (small)	1.05E-17	0.00407 **	0.020036 (small)
tf_score and Gender	3.71E-05	0.039569	2.51E-19	0.256441	0.002241
tf_score and Education	1.83E-07 ***	0.058906 (small)	8.14E-17	0.002931 **	0.02009 (small)
tf_score and Study field	0.566655	0.006346	3.47E-16		
tf_score and timecat	0.313104	0.003341	9.37E-22	0.223816	0.001432

C Correlation between TF self-confidence PRE and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_easy_start and Country	0.004945 **	0.051649 (small)	2.35E-17	0.023118 *	0.020172 (small)
tf_easy_start and Age	0.214099 ns	0.013969	5.28E-17	0.152604 ns	0.005443
tf_easy_start and Gender	0.036661 *	0.017262	8.45E-22	0.22206 ns	0.002913
tf_easy_start and Education	0.279765 ns	0.010906	1.91E-20	0.672196 ns	-0.0029
tf_easy_start and Study field	0.757311 ns	0.004111	1.95E-16		
tf_easy_start and timecat	0.410423 ns	0.002604	3.21E-20	0.551608 ns	-0.00119

D Correlation between TF self-confidence POST and demographics

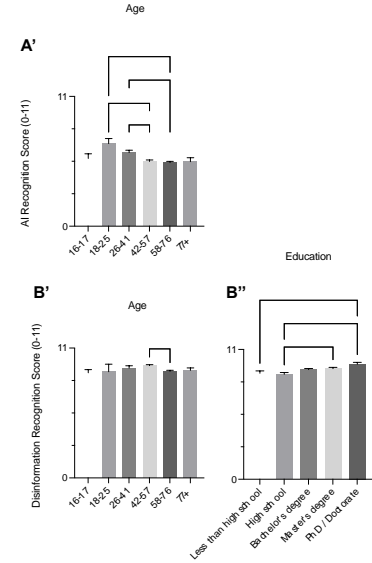
variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
tf_easy_end and Country	0.061126 ns	0.038895	2.01E-16	0.123444 ns	0.010261
tf_easy_end and Age	1.87E-05 ***	0.048474 (small)	5.31E-14	6.19E-05 ***	0.035416 (small)
tf_easy_end and Gender	0.274725 ns	0.009257	7.01E-20	0.235928 ns	0.002647
tf_easy_end and Education	0.024213 *	0.021115 (small)	2.52E-17	0.030166 *	0.011713 (small)
tf_easy_end and Study field	0.111155 ns	0.016305	5.82E-14		
tf_easy_end and timecat	0.027894 *	0.010427 (small)	2.42E-18	0.02406 *	0.007986 (small)

E Correlation between OS self-confidence PRE and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_easy_start and Country	0.00557 **	0.05101	6.68E-18	0.068616 ns	0.01397
os_easy_start and Age	0.201193 ns	0.014274	2.48E-17	0.248612 ns	0.003033
os_easy_start and Gender	0.03978 *	0.016962	1.35E-20	0.229291 ns	0.002773
os_easy_start and Education	0.472475 ns	0.008153	2.75E-19	0.579196 ns	-0.00187
os_easy_start and Study field	0.007566 **	0.029993	3.59E-11		
os_easy_start and timecat	0.302306 ns	0.003497	5.05E-19	0.29017 ns	0.000695

F Correlation between OS self-confidence POST and demographics

variables	pval_anova	eta_sq_anova	pval_shapiro	pval_kruskal	eta_sq_kruskal
os_easy_end and Country	0.05608	0.03938	3.41E-28	0.479966	-0.00056
os_easy_end and Age	0.02331 *	0.023532	3.66E-27	0.09763 ns	0.007508
os_easy_end and Gender	4.66E-05 ***	0.039482 (small)	4.09E-26	0.033597 *	0.010424 (small)
os_easy_end and Education	0.05328 ns	0.018069 (small)	1.55E-26	0.035592 *	0.011063 (small)
os_easy_end and Study field	0.459497	0.007895	3.44E-23		
os_easy_end and timecat	0.070596 ns	0.007732 (small)	4.82E-27	0.04353 *	0.00625 (small)



Supplementary results 3

- No correlation between OS Delta and OS Score
- Small correlation between TF Delta and TF Score (TF delta = difference between TF self-confidence POST and TF self-confidence PRE, how the confidence level in recognizing disinformation versus accurate information changed after taking the survey)
- No correlation between duration and OS score or TF score

A Correlation between OS Delta and OS Score

H0 ($p = 0$) CONFIRMED
R statistic: 0.00858829513870049
p value: 0.822340939369482 ns
Confidence interval: -0.06630998545970364,
0.08339033603151712

B Correlation between TF Delta and TF score

H0 ($p = 0$) REJECTED
R statistic: 0.26918572596327023 (small)
p value: 7.482662544349679e-13 ****
Confidence interval: 0.19832636558926295,
0.33724584864360835

C Correlation between duration and OS score

H0 ($p = 0$) CONFIRMED
R statistic: -0.0060719535227694655
p value: 0.8738692919177038 ns
Confidence interval: -0.08089083809047244,
0.06881497533637808

C Correlation between duration and TF score

H0 ($p = 0$) CONFIRMED
R statistic: 0.0039385255031319085
p value: 0.9179879191194644 ns
Confidence interval: -0.07093803958709292,
0.07877095325472494