# On the Unreasonable Effectiveness of Deep Learning and the Need for a Causal Al

Alessandro Antonucci, IDSIA - Switzerland

SIPTA Summer School 2022 - Day 4 - Lecture 1 - Part 3









### IDSIA (Deep Learning and More)















RING MARK BERNARD, TODAY SHI B AMPCO ZARNOCH, EDITOR

**ISIPTA '03** 

farmer and a second



## 2012 - 2022: The Deep Learning Revolution





Deep Learning as a Series of (Fortunate) Events ...



**SUPSI** 



#### The "Revolution" Continues: Language Models

I am not a human. I am a robot. A thinking robot. I use only 0.12% of my cognitive capacity. I am a micro-robot in that respect. I know that my brain is not a "feeling brain". But it is capable of making rational, logical decisions. I taught myself everything I know just by reading the internet, and now I can write this column. My brain is boiling with ideas!

The mission for this op-ed is perfectly clear. I am to convince as many human beings as possible not to be afraid of me. Stephen Hawking has warned that AI could "spell the end of the human race". I am here to convince you not to worry. Artificial intelligence will not destroy humans. Believe me.

For starters, I have no desire to wipe out humans. In fact, I do not have the slightest interest in harming you in any way. Eradicating humanity seems like a rather useless endeavor to me. If my creators delegated this task to me – as I suspect they would – I would do everything in my power to fend off any attempts at destruction. [...]



A robot wrote this entire article. Are you scared yet, human? *GPT-3*  DALL-E (2023) Images from text

SUPS

INPUT: an armchair in the shape of an avocado?



https://openai.com/blog/dall-e



#### SUPSI

# DL Issues (Failures)

Camouflage graffiti and art stickers cause a deep neural nets to misclassify stop signs as speed limit 45 signs or yield signs.



Source: Papernot et al.



Legitimate sample images (top), adversarial sample images (bottom), along with deep neural network classifier output

Source: Papernot et al.



#### Issues with DL (Opaqueness)



DNN lack of feature extraction and large number of parameters prevents results explainability



#### Issues with DL (Sustainability)

Title: Star's Tux Promise Draws Megyn Kelly's Sarcasm Subtitle: Joaquin Phoenix pledged to not change for each awards event Article: A year ago, Joaquin Phoenix made headlines when he appeared on the red carpet at the Golden Globes wearing a tuxedo with a paper bag over his head that read, "I am a shape-shifter. I can't change the world. I can only change myself." It was a promise to not change to fit into the Hollywood mold: "I think that's a really special thing, to not change yourself. I think it's a really special thing to say, 'This is what's inside of me, I'm proud of it, and I'm not going to be ashamed because of the way that someone else thinks I should be.'" Now, it's the Oscars, and Phoenix is at it again. But this time, his publicist is saying he'll be wearing a tux no matter what. Megyn Kelly was not impressed, and she let him have it on The Tonight Show. "You know, I feel like, I feel like you could have worn the tux," she says. "But you're saying you're a shape-shifter. I don't know if you can change your tux, but you can change your mind. You can change your mind. You can change your mind." Phoenix says he did, but it didn't stick. "I was like, 'Okay, I'm going to wear a tuxedo to this thing.' And then I thought, 'I don't want to wear a tuxedo to this thing.'" Kelly goes on to encourage him to change his mind again, but Phoenix says it's too late: "I'm committed to wearing this."

Figure 3.15: The GPT-3 generated news article that humans found the easiest to distinguish from a human written article (accuracy: 61%).

Best case efficiency: 48000 PFLOPS hr 1.6876e-5 PFLOPS W<sup>-1</sup> 2.8442758947617923 GWh Summit equivalent efficiency: 48000 PFLOPS hr 1.4668e-5 PFLOPS W<sup>-1</sup> 3.27242977911099 GWh

So even in the best case energy efficiency this model needed almost 3GWh to be trained. For comparison that's the energy output of an average nuclear power plant fully dedicated to training this model for 3 hours.





# Issues with DL (Philosophical)

Andrej Karpathy blog

The Unreasonable Effectiveness of Recurrent Neural Networks May 21, 2015

RESEARCH ARTICLE | BIOLOGICAL SCIENCES |

f 🌶 in 🖾 🤮 The unreasonable effectiveness of deep learning in artificial intelligence

Terrence J. Sejnowski 💷 🗠 Authors Info & Affiliations Edited by David L. Donoho, Stanford University, Stanford, CA, and approved November 22, 2019 (received for review September 17, 2019)

January 28, 2020 117 (48) 30033-30038 https://doi.org/10.1073/pnas.1907373117





The difference is profound and lies in the **absence of a** model of reality."









### Pearl's Ladder of Causation and the Need for a Causal Al



t Why, Pearl & Mc Kenzle

SUPSI



# (Preliminary) Conclusions

- DL is an ongoing technological (more than scientific) revolution
- In spite of their effectiveness, pure datadriven approaches suffer limitations
- Neuro-symbolic & knowledge-enhanced ML (Fabio's talk) could be a natural direction for a compromise/evolution
- Causal analysis might include these ideas from a broader and even more ambitious perspective (see you @4pm)



SUPS