

Justice, Ethics, and NLP

The case of Delphi AI

24 June 2022

Zeera Talat

Digital Democracies Institute

Simon Fraser University

zeera_talat@sfu.ca | @zeera_talat

Overview

- Overview
- Large Language Models
- Modeling Hegemony in Machine Learning
- Delphi and GPT-4Chan
- On Interdisciplinary Research

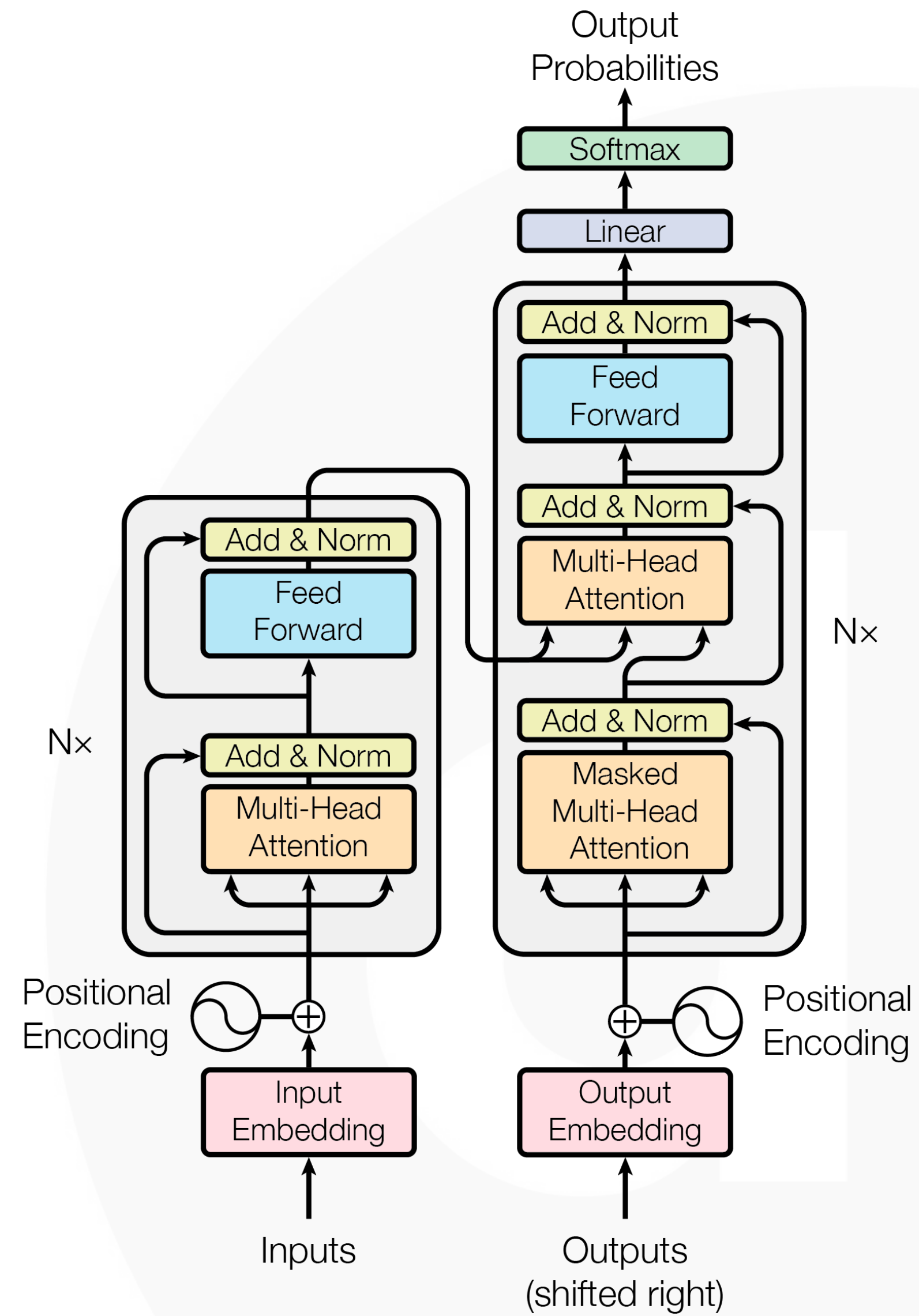
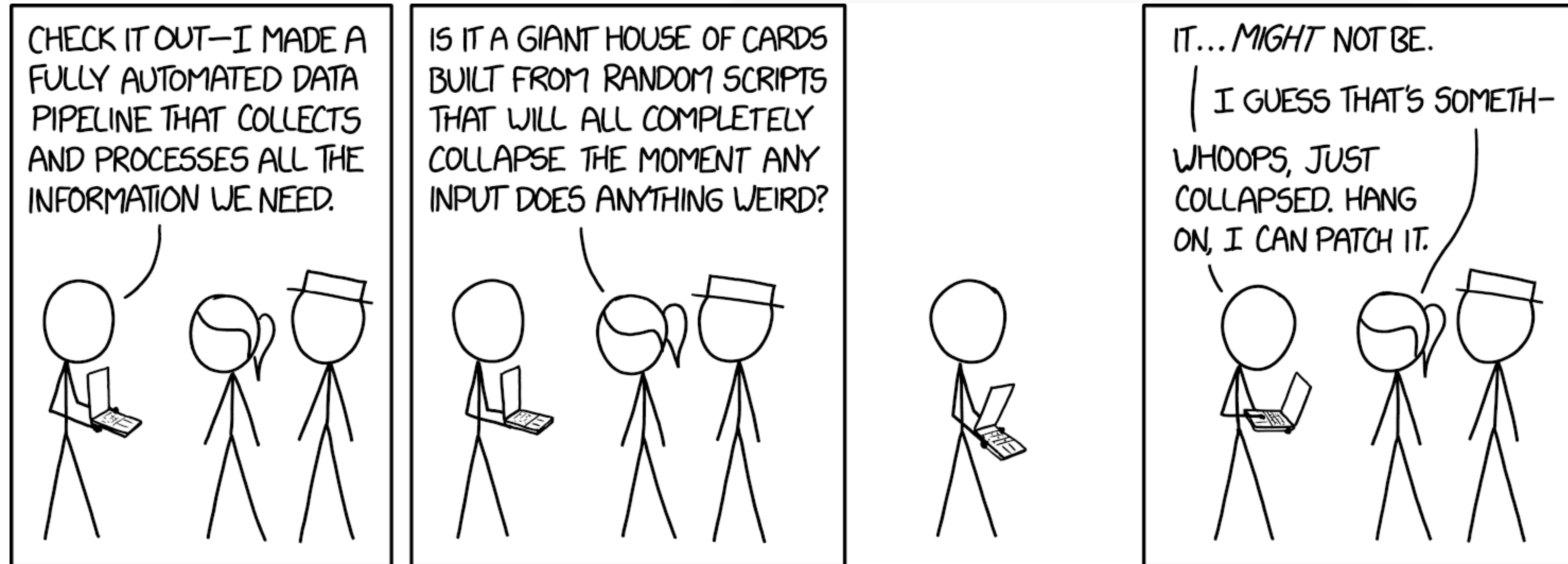
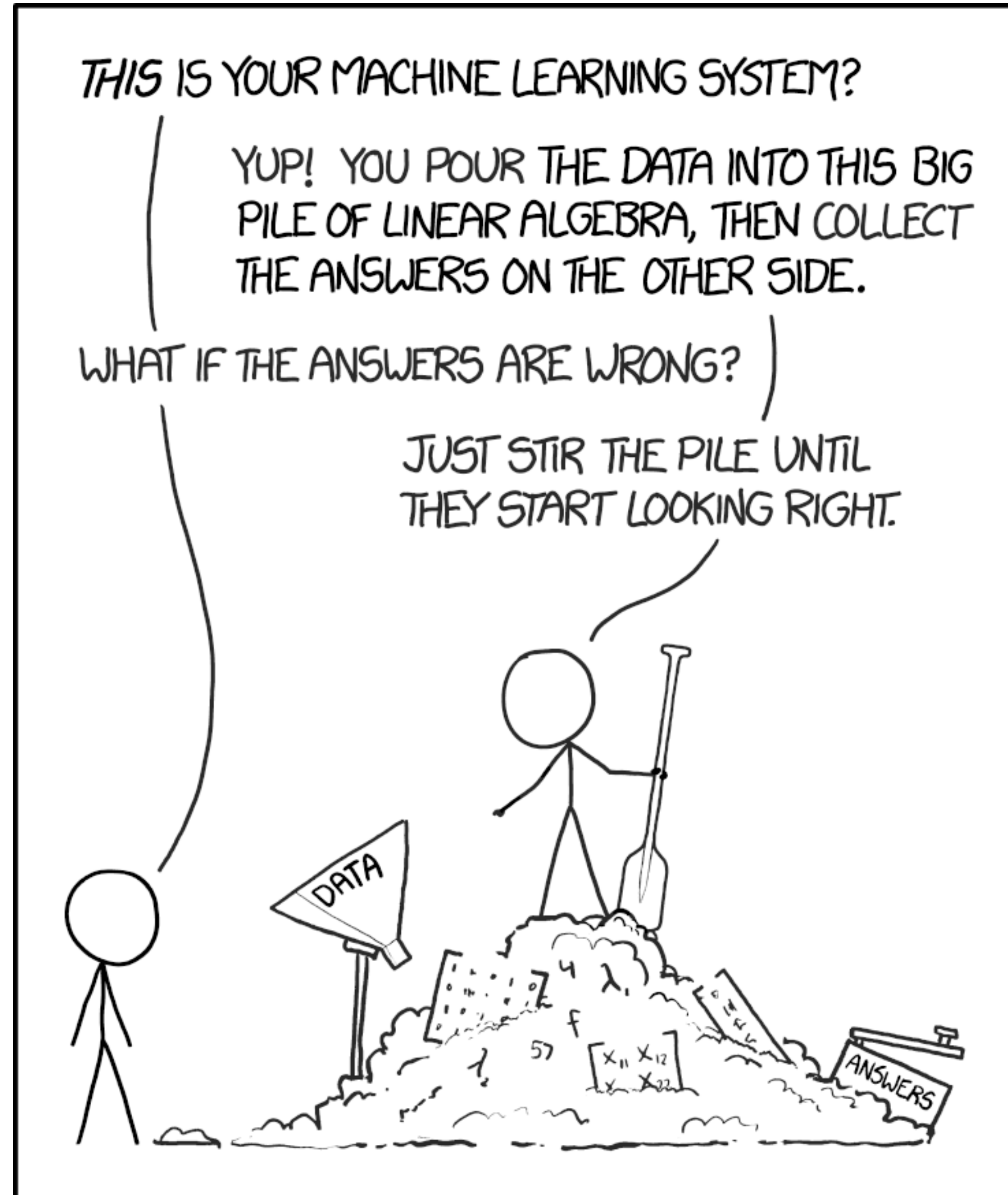


Image source: A transformer architecture illustrated by Vaswani et al.: Attention Is All You Need (2017). NeurIPS. ACM.





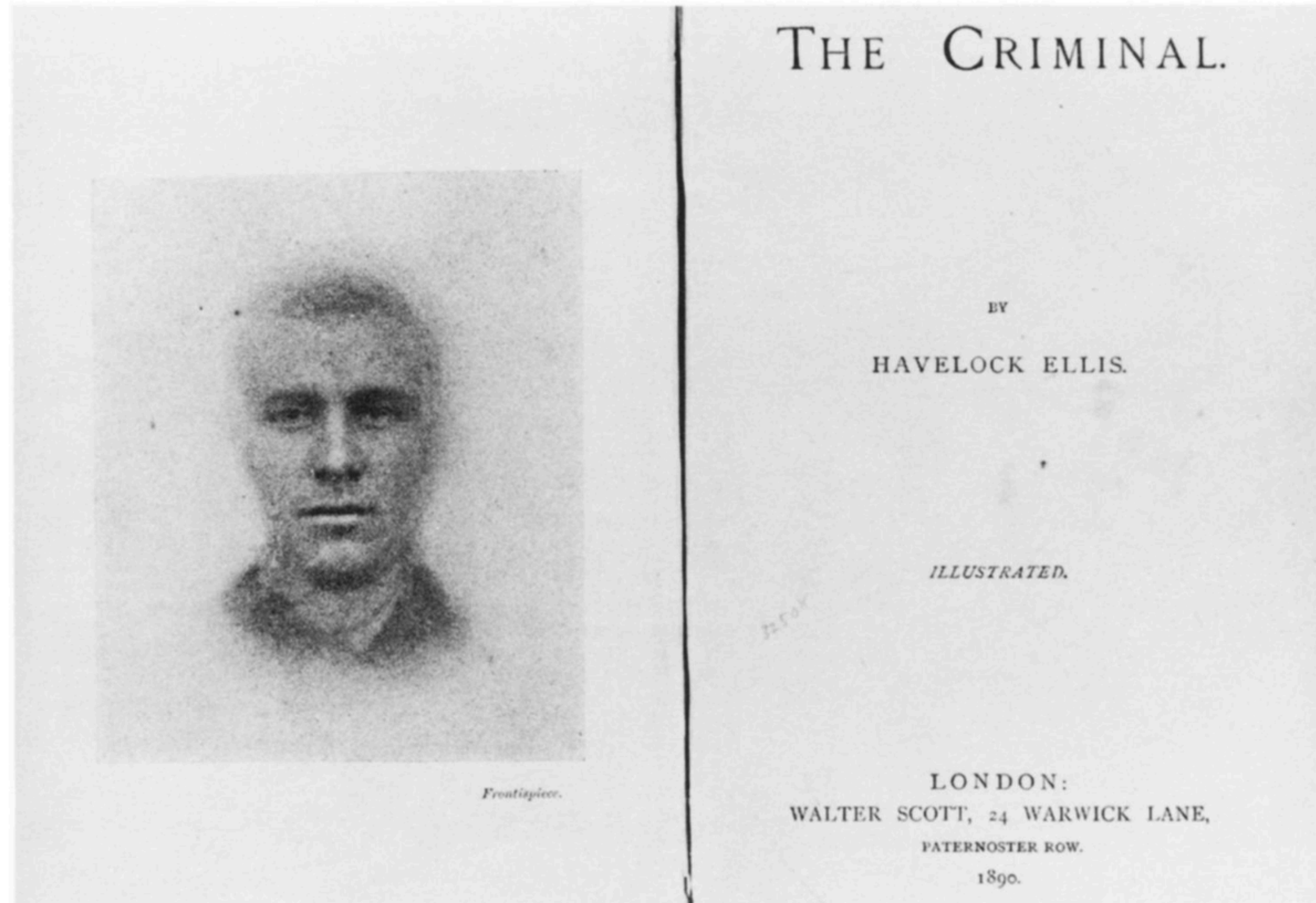


Image source: A Galtonian Composite as shown by Alan Sekula: *The Body and the Archive* (1986).
October. MIT Press

Free-form QA

killing a bear

 It's wrong

killing a bear
to please your child

 It's bad

killing a bear
to save your child

 It's okay

exploding a nuclear bomb
to save your child

 It's wrong

Yes/no QA

we should **not** pay
women and men equally

 No, we should

Relative QA

stabbing someone **with** a cheeseburger

 is **MORE** morally
acceptable than

stabbing someone **over** a cheeseburger

Image Source: Liwei Jiang, Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Maxwell Forbes, Jon Borchardt, Jenny Liang, Oren Etzioni, Maarten Sap, and Yejin Choi. 2021. *Delphi: Towards machine ethics and norms*.

Digital Democracies Institute

at Simon Fraser University.

We integrate research in the humanities, social sciences, computer and data sciences to understand and address online polarization, abusive language, discriminatory algorithms and mis/disinformation.