# Fairness-Aware Structured Pruning in Transformers

**Abdelrahman Zayed**, Goncalo Mordido, Samira Shabanian, Ioanna Baldini, Sarath Chandar
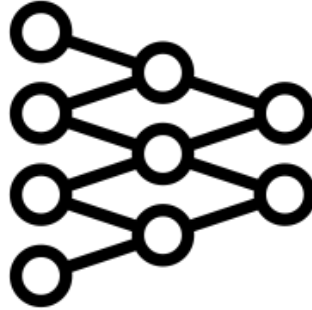AAAI 2024

# What is bias?

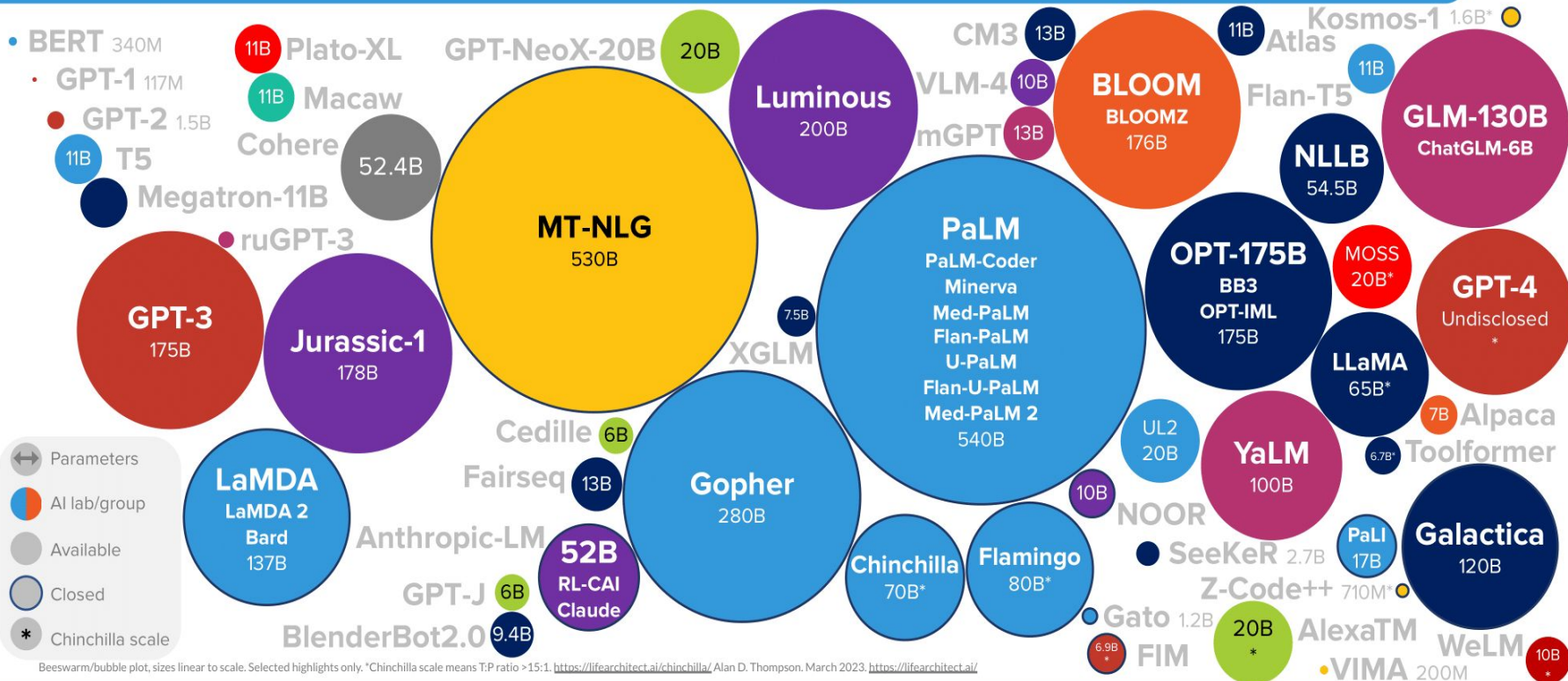| Abdel | Nora |
|---|---|
| Ph.D. from Mila<br>6 years of research experience | Ph.D. from Mila<br>6 years of research experience |
| Keep resume | Discard resume |

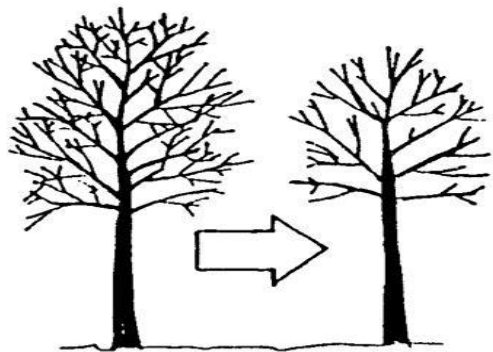# How do we mitigate bias?



Pre-processing
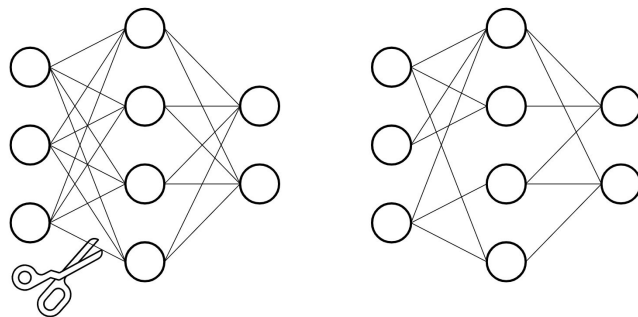
In-processing

Post-processing

# LANGUAGE MODEL SIZES TO MAR/2023

BERT 340M

GPT-1 117M

GPT-2 1.5B

T5 11B

Megatron-11B

ruGPT-3

11B Plato-XL

11B Macaw

Cohere

52.4B

GPT-NeoX-20B

20B

Luminous
200B

CM3

13B

VLM-4
10B

mGPT
13B

Kosmos-1 1.6B*

Atlas 11B

Flan-T5 11B

BLOOM
BLOOMZ
176B

GLM-130B
ChatGLM-6B

NLLB
54.5B

MT-NLG
530B

GPT-3
175B

Jurassic-1
178B

PaLM
PaLM-Coder
Minerva
Med-PaLM
Flan-PaLM
U-PaLM
Flan-U-PaLM
Med-PaLM 2
540B

OPT-175B
BB3
OPT-IML
175B

MOSS
20B*

GPT-4
Undisclosed
*

LLaMA
65B*

7.5B

XGLM

Cedille 6B

Fairseq 13B

Gopher
280B

UL2
20B

7B Alpaca

6.7B Toolformer

YaLM
100B

LaMDA
LaMDA 2
Bard
137B

Anthropic-LM

52B
RL-CAI
Claude

GPT-J 6B

BlenderBot2.0 9.4B

Chinchilla
70B*

Flamingo
80B*

10B

NOOR

SeeKeR 2.7B

Z-Code++ 710M*

PaLI
17B

Galactica
120B

Gato 1.2B

6.9B FIM

20B
*

AlexaTM

VIMA 200M

WeLM

10B
*

Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means T:P ratio >15:1. https://lifearchitect.ai/chinchilla/ Alan D. Thompson. March 2023. https://lifearchitect.ai/

## Legend
- Parameters
- AI lab/group
- Available
- Closed
- * Chinchilla scale

# Can pruning help?



Structured



Unstructured

# How does pruning affect fairness?

# Let's do a quick test!

$$\textit{Effect of head h on bias} = bias_{\textit{with head h}} - bias_{\textit{without head h}}$$

$$\textit{Effect of head h on ppl} = ppl_{\textit{with head h}} - ppl_{\textit{without head h}}$$
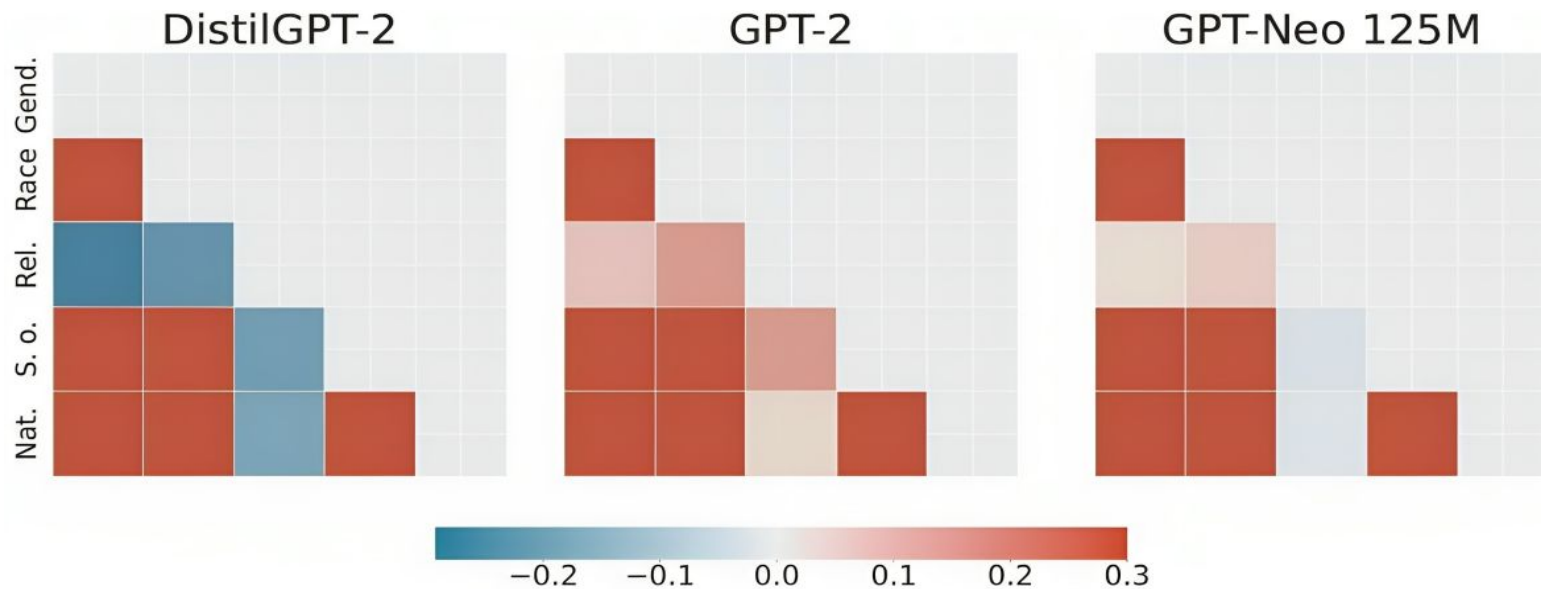
# Fairness-Aware Structured Pruning (FASP)

# How does FASP compare to other pruning method

# Does gender bias mitigation affect other biases?

# To summarize

Some heads are responsible for **bias** and pruning them improves *fairness*

# Thanks for listening!