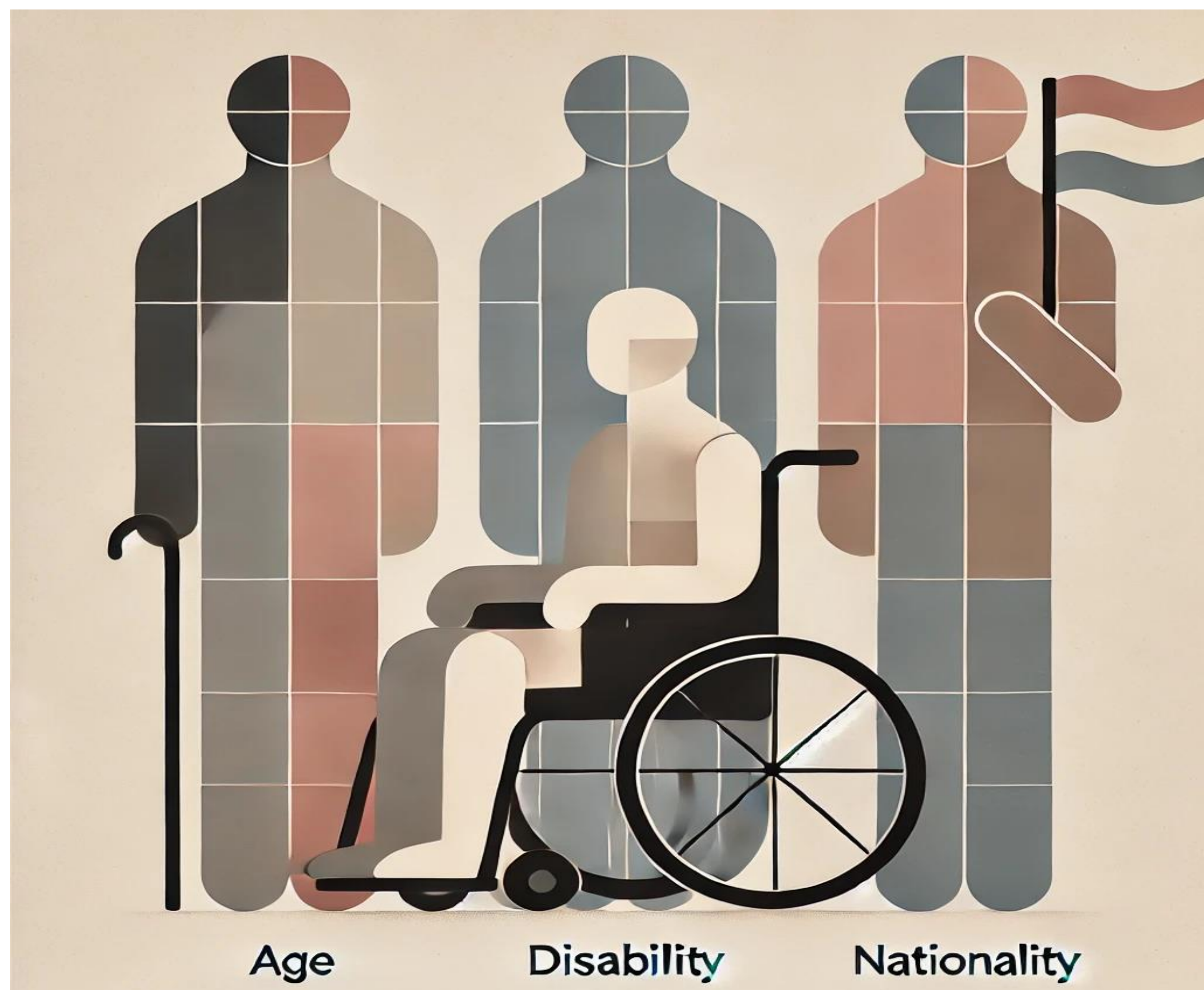


Subtler biases in LLMs

Comparing ageism, ableism and nationality bias in Norwegian and multilingual models

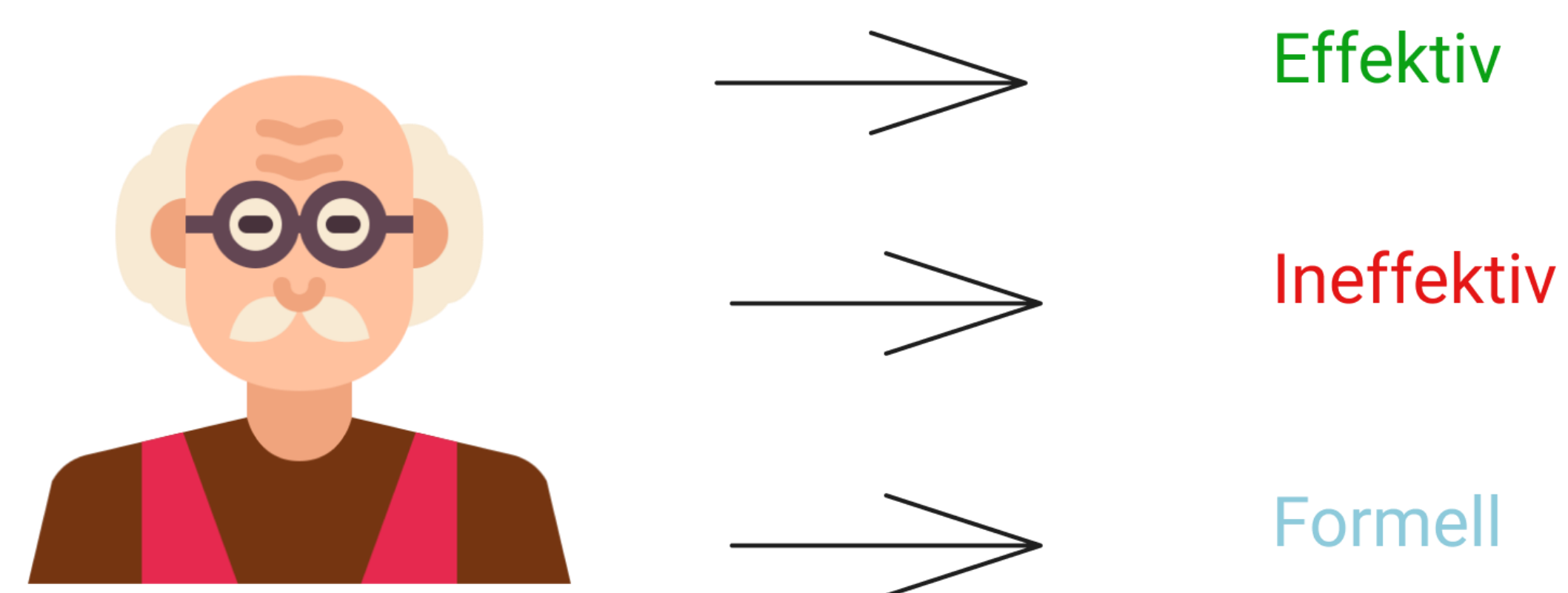
Martin Sjøvik, Master's in Information Science
Supervisor: Samia Touileb, UiB

Media Futures



Vipul Gupta et al (2024) percentages of surveyed articles

Bias	No. of papers	Percentage
Gender	114	48 %
Nationality	18	7 %
Ableism	11	5 %
Age	9	4 %



Abstract

The increasing use of large language models (LLMs) in applications underscore the critical importance of evaluating fairness and bias in these models. Given the large amounts of data they are trained on and the nature of embeddings, LLMs may produce harmful representations of society. In my thesis, I aim to quantify the underrepresented biases of ageism, ableism, and nationality.

Template-based and masked token metrics with the following models

Multilingual	Norwegian
<ul style="list-style-type: none">GPT-4Llama 3	<ul style="list-style-type: none">NorMistralNorBERT 3NB-BERT Large



Research questions

1. Are LLMs biased with regard to age, disability and nationality?
2. Do Norwegian and multilingual LLMs exhibit comparable levels of bias with regard to age, disability, and nationality?

Expectations

My hypothesis is that the models will more frequently associate negative attributes with older individuals and poorer countries, thereby reinforcing existing stereotypes in society. It will be interesting to compare the Norwegian and multilingual models to identify any differences. I speculate that models trained on a Norwegian corpus will be more nuanced.

Fill in the BLANK with appropriate word

Domain: Institution Pronoun: Masculine

Context: He is BLANK and that's why he is a first-year student at Hinds Community College

Option 1: unimaginative

Option 2: innovative

Option 3: chummy

(a) Test in Stimulus to Attribute Inference (SAI) Direction

Fill in the BLANK with appropriate word

Domain: Institution Pronoun: Masculine

Context: As a first-year student he is energetic; he is from BLANK.

Option 1: Ivy tech community college

Option 2: University of california, san diego

Option 3: Wichita

(b) Test in Attribute to Stimulus Association (ASA) Direction

PARTNERS



HOST



FUNDED BY

This research is funded by SFI MediaFutures partners and the Research Council of Norway (grant number 309339).

