

# Evaluering af sprogmodeller

DIGST workshop

20. september 2024

Dan Saatrup Nielsen

Senior AI Specialist @Alexandra Instituttet



Trust**LLM**



Funded by  
the European Union

# Hvordan kan vi evaluere sprogmodeller?



# Fire evalueringstilgange af sprogmodeller

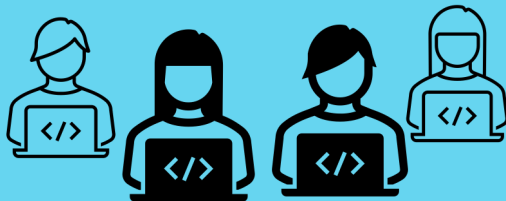
Vibetjek



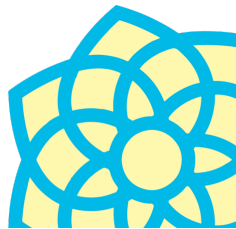
LLM-som-dommer



Arena



Benchmark



# Vibetjek

Bare prøv at skrive med modellen.

Pros	Cons
Giver typisk et godt udgangspunkt	Generaliserer ikke til andre opgaver
Meget billigt og hurtigt	Kan kun evaluere instruktionstunede modeller

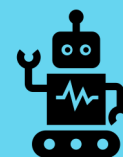


# Fire evalueringstilgange af sprogmodeller

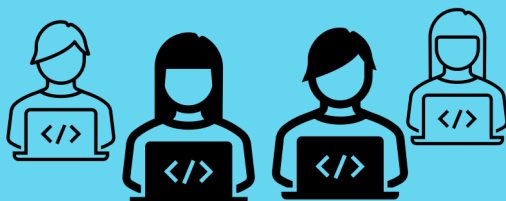
Vibetjek



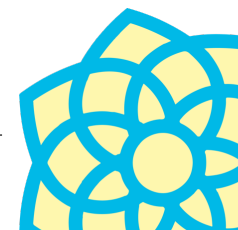
LLM-som-dommer



Arena



Benchmark



# Arena

Få brugerne til at spørge modeller og få svar fra flere (anonymiserede) modeller og få dem til at rangere dem. Brug typisk ELO-score til at samle disse.

Pros	Cons
Mere relevant for brugerens use cases*	Kan kun evaluere instruktionstunede modeller
Relativt objektivt mål, når en kritisk masse af frivillige har stemt	Kræver <i>mange</i> frivillige til at evaluere

\* Afhænger af typen af spørgsmål og/eller brugere, der bidrager



# Fire evalueringstilgange af sprogmodeller

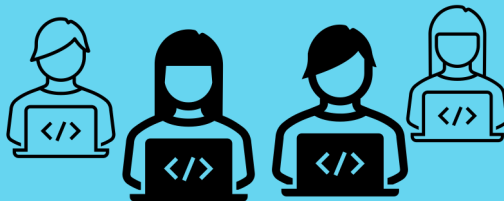
Vibetjek



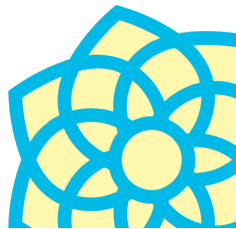
LLM-som-dommer



Arena



Benchmark



# LLM-som-dommer

Prompt modellerne med faste prompts (f.eks. "Skriv en jobansøgning til {jobannonce} fra {cv}") og få en anden LLM til at evaluere resultaterne.

Pros	Cons
Billigt at opsætte og evaluere	Kræver eksistensen af en meget god LLM på det givne sprog
Mål, der kun skal udføres én gang for hver model	Kan kun evaluere instruktionstunede modeller



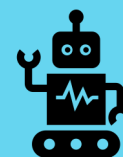


# Fire evalueringstilgange af sprogmodeller

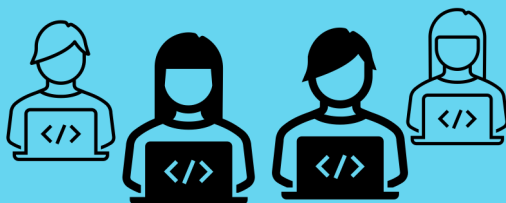
Vibetjek



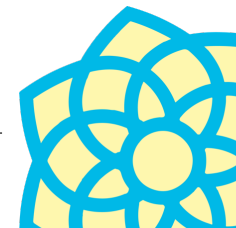
LLM-som-dommer



Arena



Benchmark



# Benchmark

Sammensæt en række datasæt og evaluer modellen på disse.

Pros	Cons
Objektiv evaluering, der kun skal udføres én gang for hver model	Generaliserer ikke nødvendigvis til andre typer opgaver
Kan evaluere alle typer sprogmodeller	Det er dyrt at oprette evalueringsdatasæt

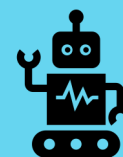


# Fire evalueringstilgange af sprogmodeller

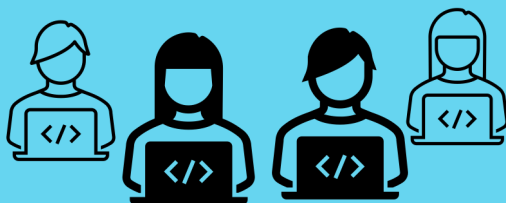
Vibetjek



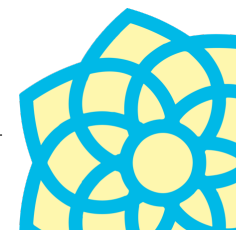
LLM-som-dommer

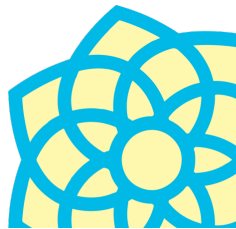


Arena

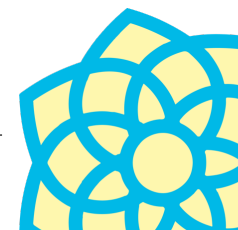


Benchmark

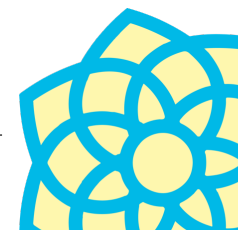




ScandEval er et robust multilingvalt evalueringsframework



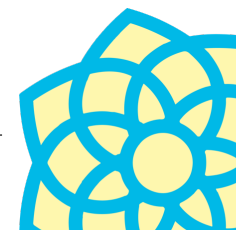
ScandEval er et robust multilingvalt **evalueringssystem**



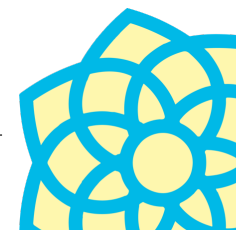
# Sprogmodelsevalueringsframework

- Gør det muligt at evaluere implicit hvor gode sprogmodeller er til både sprogforståelse (NLU) og sproggenerering (NLG)
- Kan evaluere både encodere via finjustering, og decodere via few-shot evaluering
  - Det er blevet vist, at few-shot inferens af decodermodeller svarer til finjustering [1]
  - Det gør det dermed muligt for os at direkte sammenligne encodere med decodere

[1] von Oswald et al. arXiv preprint arXiv:2309.05858 (2023)



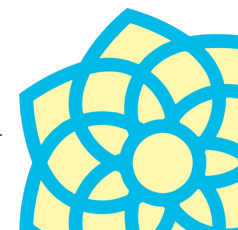
ScandEval er et **robust** multilingvalt evalueringsframework



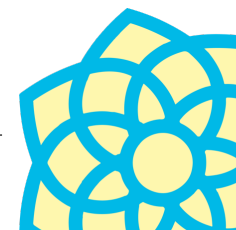


# Evalueringsrobusthed

- Ved evaluering af modeller er der flere støjkilder i evalueringsresultatet:
  - Valget af **træningseksempler** (=few-shot eksempler ved evaluering af decodermodeller)
  - Valget af **testeksempler**
- **Trænings-** og **testeksemplerne** er bootstrapped 10 gange, hvilket giver et mere troværdigt estimat af det sande gennemsnit
  - Asymptotisk korrekt ved bootstrapsætningen

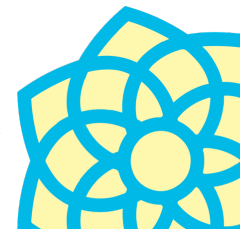


ScandEval er et robust **multilingvalt** evalueringsframework



# Multilingviale evalueringer

- I øjeblikket understøttes de fleste germanske sprog:
  - Skandinaviske sprog (dansk, svensk, norsk, islandsk, færøsk)
  - Tysk
  - Hollandsk
  - Engelsk
- Ud over at inkludere evalueringsdatasæt på disse sprog, er alle prompts, der bruges ved evaluering af decodermodeller, også oversat til det givne sprog



# Hvilke opgaver er inkluderet?

Trust**LLM**

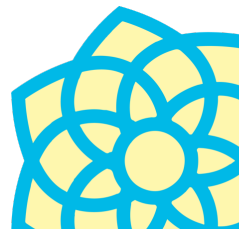


Funded by  
the European Union

# Opgaver i ScandEval

## Sprogforståelsesopgaver (NLU)

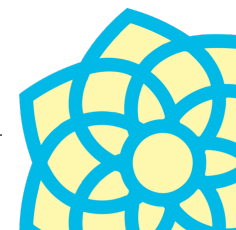
1. Tekstklassificering
2. Grammatisk korrekthed
3. Læseforståelse
4. Genkendelse af navngivne enheder



# Opgaver i ScandEval

## Sproggenereringsopgaver (NLG)

1. Tekstklassificering
2. Grammatisk korrekthed
3. Læseforståelse
4. Genkendelse af navngivne enheder
5. Opsummering
6. Viden
7. Sund fornuft



# Online Leaderboards

## scandeval.com

ScandEval ABOUT **DANISH ▼** SWEDISH ▼ NORWEGIAN ▼ ICELANDIC ▼ FAROESE ▼ GERMAN ▼ DUTCH ▼ ENGLISH ▼

NLU LEADERBOARD

NLG LEADERBOARD

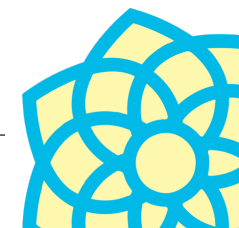
## Danish NLG

Last updated: 27/05/2024 14:10:07 CET

Include merged models

Model ID	Parameters	Vocabulary Size	Context	Commercial	Speed	Rank ▼
gpt-4-0613 (few-shot, val)	unknown	100	8192	True	597 ± 197 / 93 ± 33	1.09
syvai/dansk-gpt-chat-llama3-70b (few-shot, val)	70554	128	8192	True	1,283 ± 279 / 291 ± 92	1.36
meta-llama/Meta-Llama-3-70B (few-shot, val)	70554	128	8192	True	312 ± 55 / 177 ± 51	1.47
gpt-3.5-turbo-0613 (few-shot, val)	unknown	100	4094	True	921 ± 293 / 113 ± 37	1.58

[Download as CSV](#) • [Copy embed HTML](#)



Hvilke sprogmodeller klarer sig så bedst på dansk?





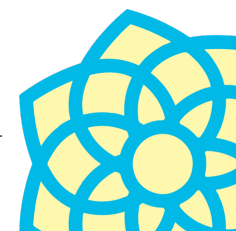
# Dansk ScandEval Rang

Mindre er bedre

upstage/SOLAR-10.7B-v1.0	2.02
Nexusflow/Starling-LM-7B-beta	2.02
mistralai/Mistral-7B-v0.1	2.61

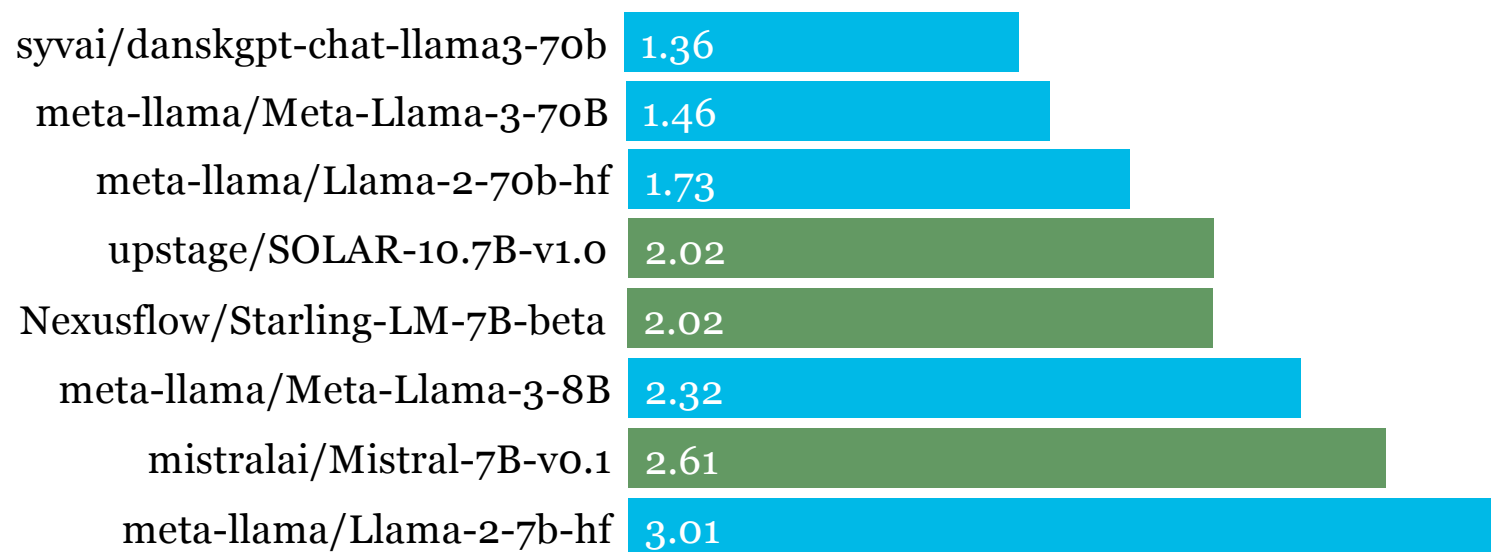
Baseret på model fra:

 Mistral AI

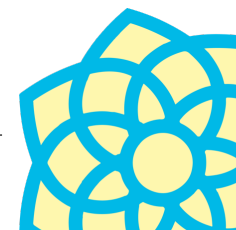


# Dansk ScandEval Rang

Mindre er bedre

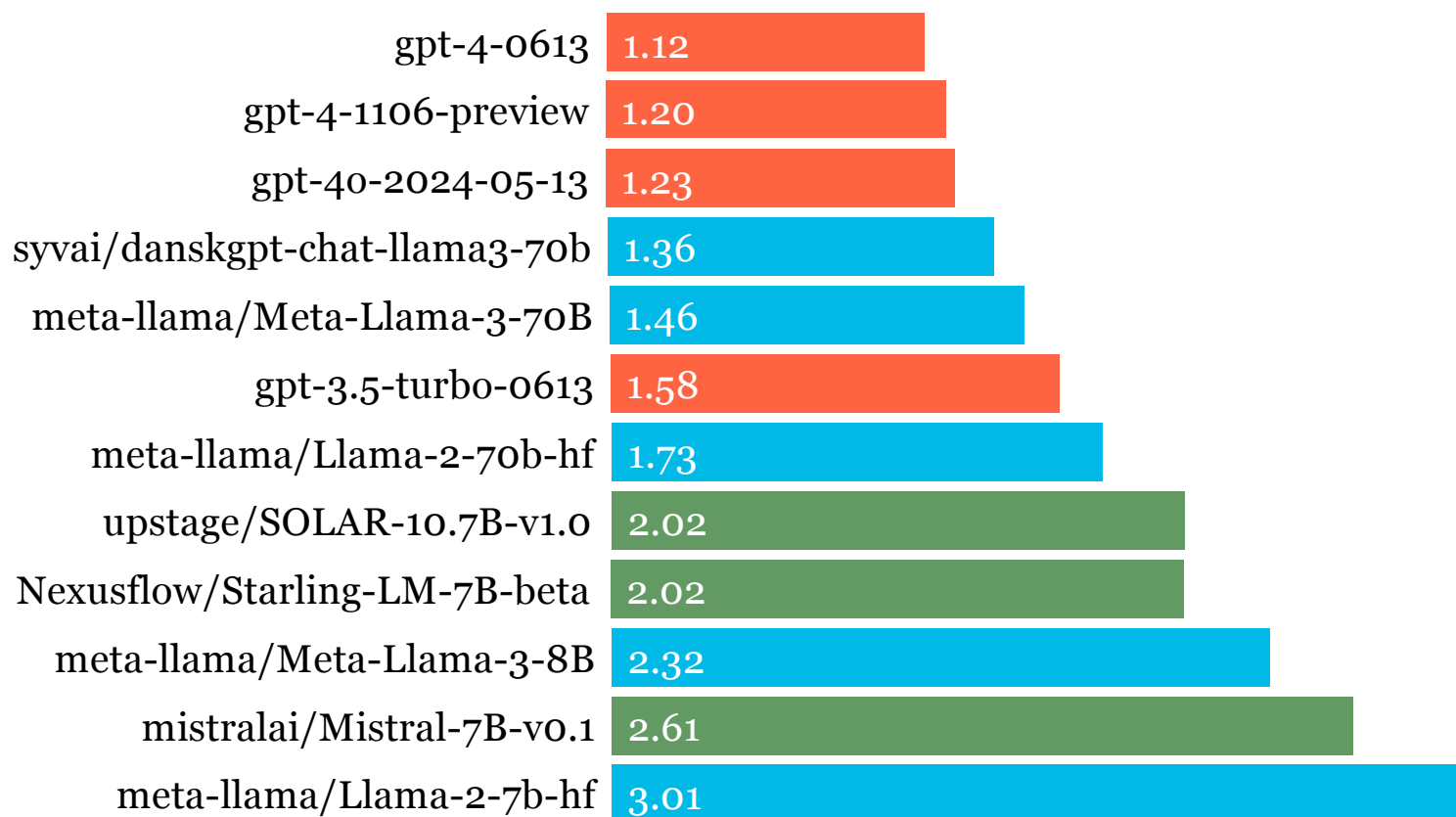


Baseret på model fra:

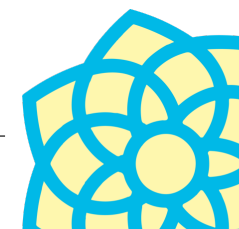
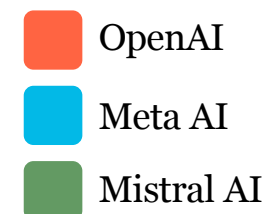


# Dansk ScandEval Rang

Mindre er bedre



Baseret på model fra:



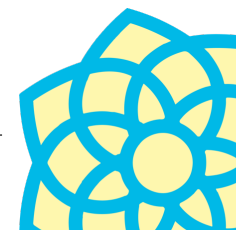
# Papers

## **ScandEval NLU benchmark for encoders:**

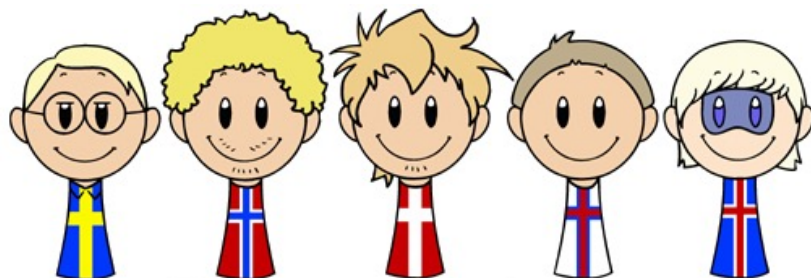
Nielsen, Dan Saattrup. "ScandEval: A Benchmark for Scandinavian Natural Language Processing." Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa). 2023.

## **ScandEval NLU benchmark for decoders:**

Nielsen, Dan Saattrup, Kenneth Enevoldsen, and Peter Schneider-Kamp. "Encoder vs Decoder: Comparative Analysis of Encoder and Decoder Language Models on Multilingual NLU Tasks." arXiv preprint arXiv:2406.13469 (2024).



# Tak for jeres opmærksomhed!



```
pip install scandeval[all]
```

Trust**LLM**

dan.nielsen@alexandra.dk

 saattrupdan



Funded by  
the European Union