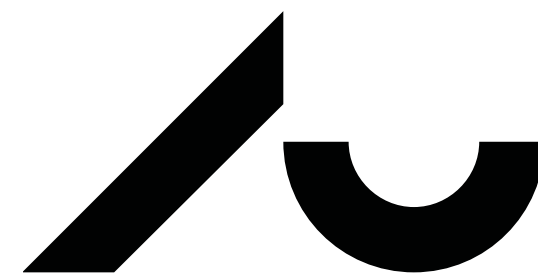# Danish Foundation Models

## Validerede sprogmodeller til dansk
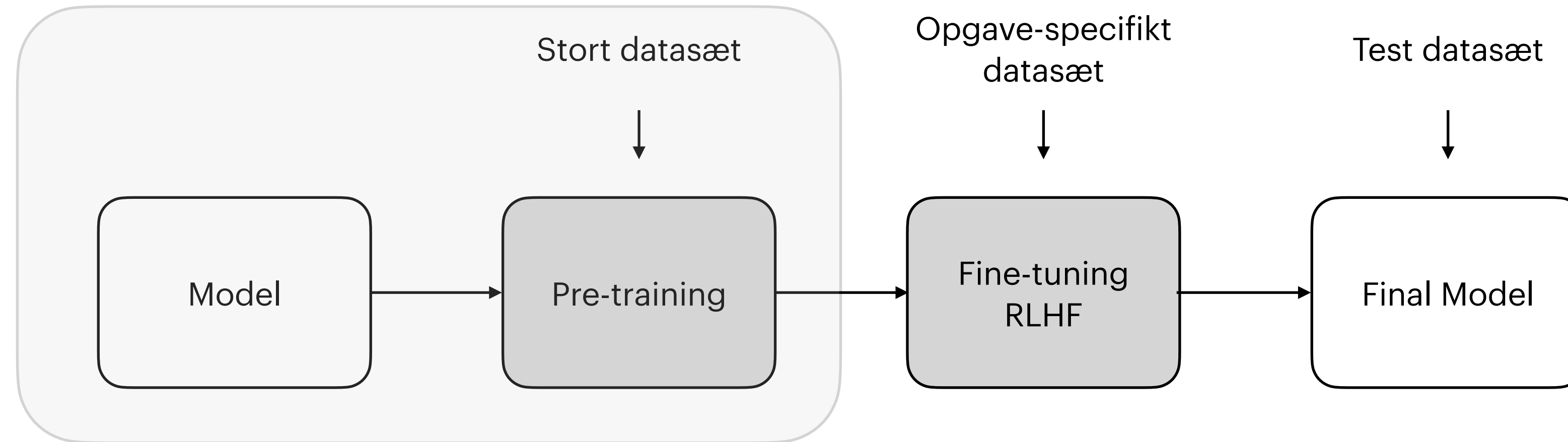
Kenneth Enevoldsen & Lasse Hansen | November 20, 2023
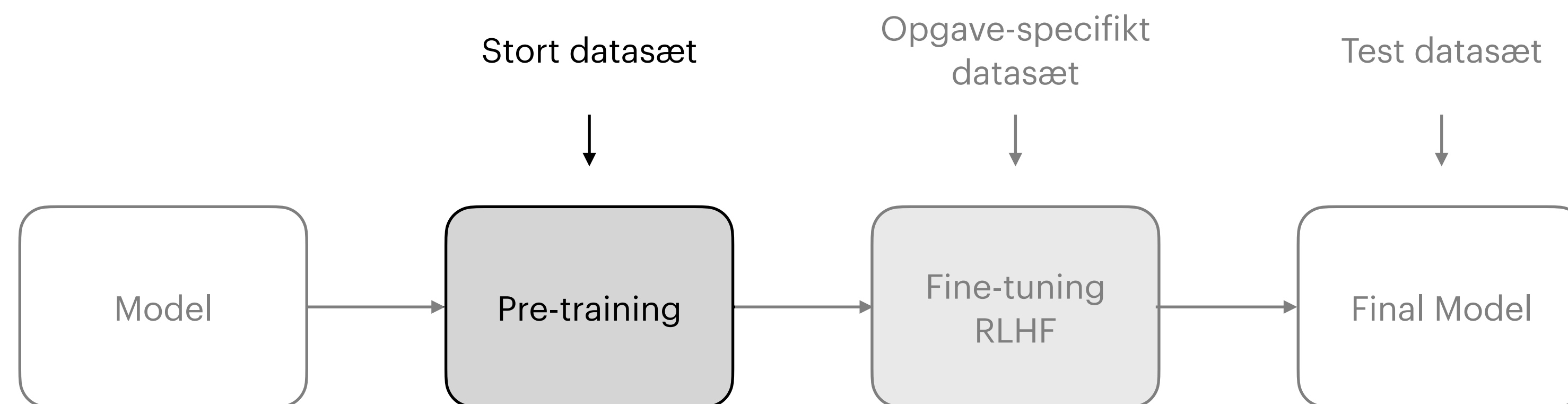
# Introduktion

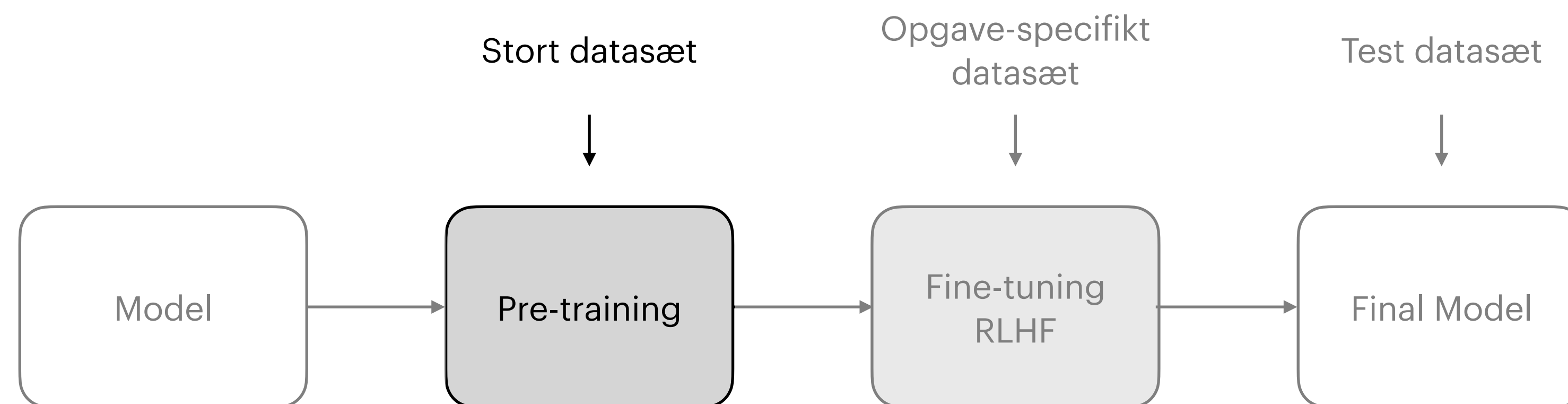# Hvad er en foundation model?

Stort datasæt

Opgave-specifikt datasæt

Test datasæt

```
Model → Pre-training → Fine-tuning RLHF → Final Model
```

CENTER FOR
HUMANITIES
COMPUTING

Stort datasæt

Opgave-specifikt datasæt

Test datasæt

Model → Pre-training → Fine-tuning RLHF → Final Model

Maskering  The quick [MASK] fox jumps over the [MASK] dog

Prædiktion  The quick brown  fox jumps over the  lazy  dog

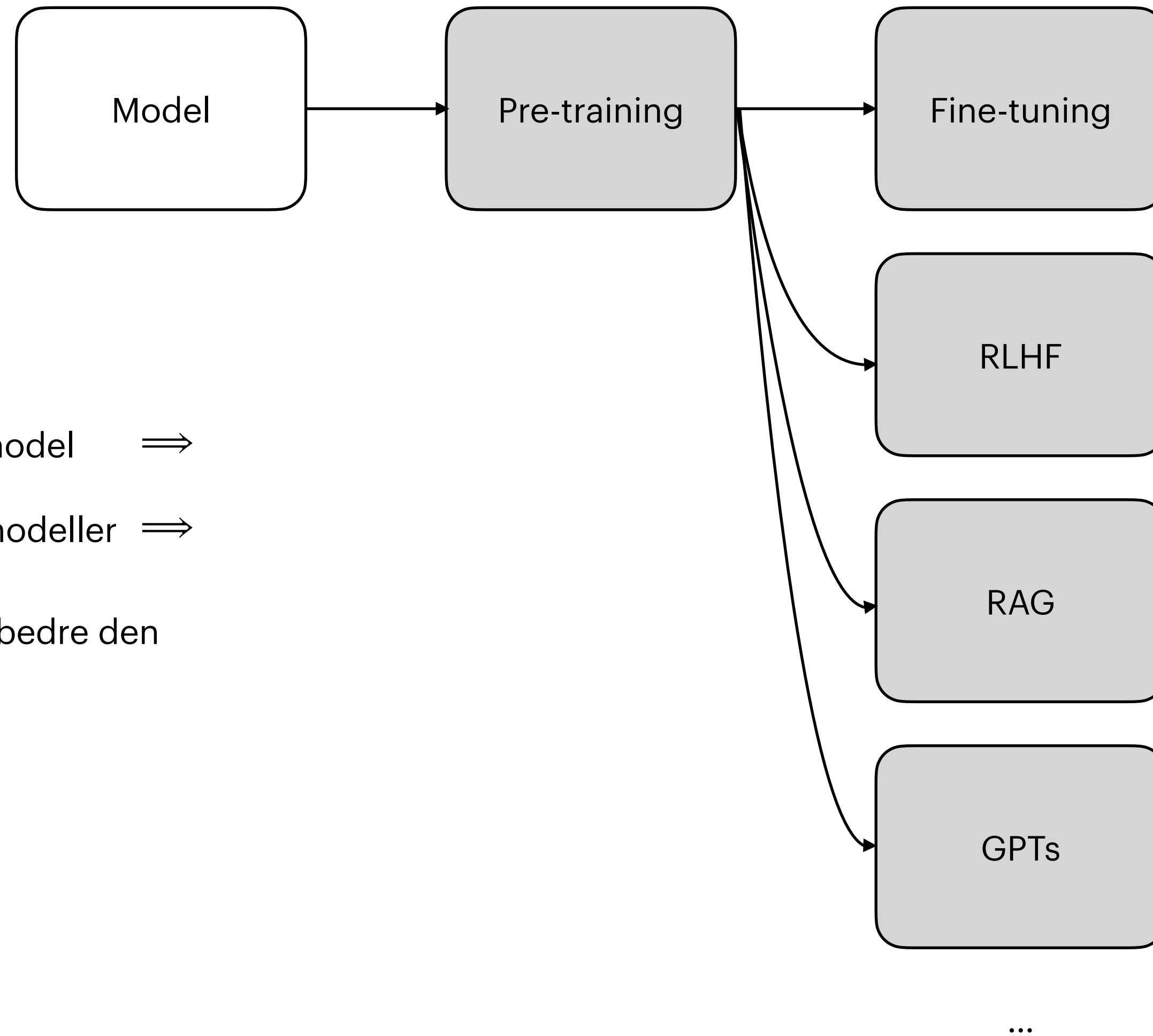Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 4171–4186. https://doi.org/10.18653/v1/N19-1423

Sources & Notes

CENTER FOR HUMANITIES COMPUTING

Stort datasæt

Opgave-specifikt datasæt

Test datasæt

Model → Pre-training → Fine-tuning RLHF → Final Model

Kontekst          Margrethe 2. er dronning af ___

Prædiktion        Danmark

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI blog, 1(8), 9.

Sources & Notes

CENTER FOR HUMANITIES COMPUTING

# Motivation

# Hvorfor danske sprogmodeller?

## Tokenisation

Lav genbrugelighed

```
antiradikaliseringskonference
['anti', '—radi', '—kal', '—isering', '—skon', '—ference']
```

Høj genbrugelighed

```
Ferskvandsfisk
['fersk', '—vand', '—sf', '—isk']
```

Færre tokens
—> højere effektivitet

```
skolehjemsamtaler
['skole', '—hjem', '—sam', '—taler']
```

CENTER FOR HUMANITIES COMPUTING

# Hvorfor danske sprogmodeller?

## Tokenisation

Lav genbrugelighed

```
antiradikaliseringskonference
['anti', '—radi', '—kal', '—isering', '—skon', '—ference']
['anti', '—radikal',     '—isering', '—s', '—konference']
```

Høj genbrugelighed

```
Ferskvandsfisk
['fersk', '—vand', '—sf', '—isk']
['fersk', '—vands',     '—fisk']
```

Færre tokens
—> højere effektivitet

```
skolehjemsamtaler
['skole', '—hjem', '—sam', '—taler']
['skole', '—hjem', '—samtale', '—r']
```

## Udtale



*Danish sentence: Røget ørred*



*Norwegian sentence: Røkt ørret*

The Unigram tokeniser have also been argued to be a better model in general, even for English:

Bostrom, Kaj, and Greg Durrett. "Byte Pair Encoding Is Suboptimal for Language Model Pretraining." In Findings of the Association for Computational Linguistics: EMNLP 2020, 4617–24. Online: Association for Computational Linguistics, 2020. https://doi.org/10.18653/v1/2020.findings-emnlp.414.

Trecca, F., Bleses, D., Madsen, T. O., & Christiansen, M. H. (2018). Does sound structure affect word learning? An eye-tracking study of Danish learning toddlers. Journal of Experimental Child Psychology, 167, 180-203.

Sources & Notes

CENTER FOR HUMANITIES COMPUTING

# Monolinguale modeller klarer sig godt

| Model ID | DA ▼ |
|---|---|
| 🇩🇰 chcaa/dfm-encoder-large-v1 | 66.17 ± 1.49 |
| 🇩🇰 KennethEnevoldsen/dfm-sentence-encoder-large-1 | 66.10 ± 1.40 |
| 🇩🇰 KennethEnevoldsen/dfm-sentence-encoder-large-2 | 65.22 ± 3.12 |
| 🇳🇴 ltg/norbert3-large | 64.40 ± 1.95 |
| 🇳🇴 (🇩🇰🇸🇪) NbAiLab/nb-bert-large | 64.40 ± 1.29 |
| 🇩🇰 vesteinn/DanskBERT | 63.87 ± 1.26 |
| 🌐 google/rembert | 63.41 ± 1.63 |
| ... | |
| 🇺🇸 (🌐) gpt-4-0613 (val) (few-shot) | 61.87 ± 2.77 |

# Monolingualle modeller skal løfte andre opgaver

```
sentiment.ts    write_sql.go    parse_expenses.py    addresses.rb

1  import datetime
2
3  def parse_expenses(expenses_string):
4      """Parse the list of expenses and return the list of triples (date, value, currency).
5      Ignore lines starting with #.
6      Parse the date using datetime.
7      Example expenses_string:
8          2016-01-02 -34.01 USD
9          2016-01-03 2.59 DKK
10         2016-01-03 -2.72 EUR
11     """
12     expenses = []
13     for line in expenses_string.splitlines():
14         if line.startswith("#"):
15             continue
16         date, value, currency = line.split(" ")
17         expenses.append((datetime.datetime.strptime(date, "%Y-%m-%d"),
18                         float(value),
19                         currency))
20     return expenses
```

Copilot

**You**

Hvordan søger jeg kontanthjælp?

CENTER FOR
HUMANITIES
COMPUTING

# Åbenhed skaber tillid



**Two DUI Arrests**

GREGORY LUGO
LOW RISK    1

MALLORY WILLIAMS
MEDIUM RISK    6

*Lugo crashed his Lincoln Navigator into a Toyota Camry while drunk. He was rated as a low risk of reoffending despite the fact that it was at least his fourth DUI.*

**Datasheets for Datasets**

TIMNIT GEBRU, Black in AI
JAMIE MORGENSTERN, University of Washington
BRIANA VECCHIONE, Cornell University
JENNIFER WORTMAN VAUGHAN, Microsoft Research
HANNA WALLACH, Microsoft Research
HAL DAUMÉ III, Microsoft Research; University of Maryland
KATE CRAWFORD, Microsoft Research

**Model Cards for Model Reporting**

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, Timnit Gebru

- Motivation

- Sammensætning

- Indsamlingsproces

- Præprocessering


- Intenderet brug

- Træningsdata

- Analyser, etiske overvejelser

- Forbehold og anbefalinger

Sources & Notes

Mattu, J. A., Jeff Larson,Lauren Kirchner,Surya. (2016.). Machine Bias. ProPublica

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Iii, H. D., & Crawford, K. (2021). Datasheets for datasets. Communications of the ACM, 64(12), 86-92.

Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019, January). Model cards for model reporting. In Proceedings of the conference on fairness, accountability, and transparency (pp. 220-229).

CENTER FOR HUMANITIES COMPUTING

# Patient Record

Desses

## Age

**Age**
**Name**
20.7

As
Sage

Lisex
Sex 4

Dstemons

Bam. 4 36 Danes

Poe 20.8 501

Gartalalicles Sielnadan Chesates

Makeyricam

Geanatie Nariqsess Methda
85 20.28312

Tom Koss Plehan Shas Yolan

Denneflers Saire FRE TE TEXT

# Manglende validering

**Domæne**

**Sprog**

# Projektet

# Målsætninger

1. At **udvikle** og **vedligeholde** state-of-the-art sprogmodeller til dansk inden for både tekst og tale

2. At **validere** foundation-modeller på dansk og danske brugsscenarier

3. At opretholde en høj standard for **dokumentation** af modeller samt data

4. At **open-source** ikke kun modellerne, men også komponenter, der kræves for reproducerbarhed, såsom datarensning, træningskode og valideringskode

Sources & Notes

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

CENTER FOR HUMANITIES COMPUTING

# Status

# The State of Foundation models for Danish

| | Model weights | Code Available | Model card | Data sheet | Language | Validated for Danish |
|---|---|---|---|---|---|---|
| **Text** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-encoder-large-v1** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓ |
| nb-bert-large | ✓ | ✓ | ✗ | ✗ | 🇳🇴 (🇩🇰 🇸🇪) | ✓ |
| XLM-Roberta | ✓ | ✓ | ✗ | ✗ | 🌐 | ✓ |
| *Generative models* | | | | | | |
| GPT-4 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✗* |
| DanskGPT | ✗ | ✗ | ✗ | ✗ | 🇩🇰 | ✗* |
| DanT5 | ✓ | ✗ | ✗ | ✗ | 🇩🇰 | ✗ |
| Llama-v2 | ✓ | ✗ | ✓ | ✗ | 🇺🇸 | ✗* |
| *Embeddings* | | | | | | |
| text-embedding-ada-2 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✓* |
| MiniLM-L12-v2[1] | ✓ | ✓ | ✗ | ✓ | 🌐 | ✓ |
| **Speech** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-xls-r-300m** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓[†] |
| wav2vec2-base-da | ✓ | ✓ | ✗ | ✗ | 🇩🇰 | ✓[†] |

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

Sources & Notes

CENTER FOR HUMANITIES COMPUTING

# The State of Foundation models for Danish

| | Model weights | Code Available | Model card | Data sheet | Language | Validated for Danish |
|---|---|---|---|---|---|---|
| **Text** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-encoder-large-v1** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓ |
| nb-bert-large | ✓ | ✓ | ✗ | ✗ | 🇳🇴 (🇩🇰 🇸🇪) | ✓ |
| XLM-Roberta | ✓ | ✓ | ✗ | ✗ | 🌐 | ✓ |
| *Generative models* | | | | | | |
| GPT-4 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✗* |
| DanskGPT | ✗ | ✗ | ✗ | ✗ | 🇩🇰 | ✗* |
| DanT5 | ✓ | ✗ | ✗ | ✗ | 🇩🇰 | ✗ |
| Llama-v2 | ✓ | ✗ | ✓ | ✗ | 🇺🇸 | ✗* |
| *Embeddings* | | | | | | |
| text-embedding-ada-2 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✓* |
| MiniLM-L12-v2[1] | ✓ | ✓ | ✗ | ✓ | 🌐 | ✓ |
| **Speech** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-xls-r-300m** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓† |
| wav2vec2-base-da | ✓ | ✓ | ✗ | ✗ | 🇩🇰 | ✓† |

CENTER FOR HUMANITIES COMPUTING

# The State of Foundation models for Danish

| | Model weights | Code Available | Model card | Data sheet | Language | Validated for Danish |
|---|---|---|---|---|---|---|
| **Text** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-encoder-large-v1** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓ |
| nb-bert-large | ✓ | ✓ | ✗ | ✗ | 🇳🇴 (🇩🇰 🇸🇪) | ✓ |
| XLM-Roberta | ✓ | ✓ | ✗ | ✗ | 🌐 | ✓ |
| *Generative models* | | | | | | |
| GPT-4 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✗* |
| DanskGPT | ✗ | ✗ | ✗ | ✗ | 🇩🇰 | ✗* |
| DanT5 | ✓ | ✗ | ✗ | ✗ | 🇩🇰 | ✗ |
| Llama-v2 | ✓ | ✗ | ✓ | ✗ | 🇺🇸 | ✗* |
| *Embeddings* | | | | | | |
| text-embedding-ada-2 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✓* |
| MiniLM-L12-v2[1] | ✓ | ✓ | ✗ | ✓ | 🌐 | ✓ |
| **Speech** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-xls-r-300m** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓† |
| **wav2vec2-base-da** | ✓ | ✓ | ✗ | ✗ | 🇩🇰 | ✓† |

CENTER FOR HUMANITIES COMPUTING

# The State of Foundation models for Danish

| | Model weights | Code Available | Model card | Data sheet | Language | Validated for Danish |
|---|---|---|---|---|---|---|
| **Text** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-encoder-large-v1** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓ | ⭐ |
| nb-bert-large | ✓ | ✓ | ✗ | ✗ | 🇳🇴 (🇩🇰 🇸🇪) | ✓ |
| XLM-Roberta | ✓ | ✓ | ✗ | ✗ | 🌐 | ✓ |
| *Generative models* | | | | | | |
| GPT-4 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✗* |
| DanskGPT | ✗ | ✗ | ✗ | ✗ | 🇩🇰 | ✗* |
| DanT5 | ✓ | ✗ | ✗ | ✗ | 🇩🇰 | ✗ |
| Llama-v2 | ✓ | ✗ | ✓ | ✗ | 🇺🇸 | ✗* |
| *Embeddings* | | | | | | |
| text-embedding-ada-2 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✓* |
| MiniLM-L12-v2[1] | ✓ | ✓ | ✗ | ✓ | 🌐 | ✓ |
| **Speech** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-xls-r-300m** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓† | ⭐ |
| wav2vec2-base-da | ✓ | ✓ | ✗ | ✗ | 🇩🇰 | ✓† |

CENTER FOR HUMANITIES COMPUTING

# Our current dataset

| Name | Description | Size | Open access | Novel corpus |
|------|-------------|------|-------------|--------------|
| **Text** | | | | |
| DAGW | Danish Gigaword | 1B tokens | ✓ | ✗ |
| reddit-da | Danish Reddit | <.1B tokens | ✓ | ✗ |
| HopeTwitter | Danish Tweets | 0.48B tokens | ✗ | ✓ |
| DaNews | Danish newspapers | 0.5B tokens | ✗ | ✓ |
| Netarkivet Text | Danish internet | >100B tokens | ✗ | ✓ |
| **Speech** | | | | |
| DaRadio | Danish talk radio | 140.000 hours | ✗ | ✓ |
| DaTV | Danish subtitled TV | 900 | ✗ | ✓ |

- Nye aftaler med:

  - Infomedia

  - Det kongelige bibliotek

- Arbejder vi på samarbejde med:

  - sundhed.dk, lex.dk, borger.dk, rigsarkivet, etc.

    - Giver bl.a. høj-kvalitets data til validering af eksisterende sprogmodeller

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

Sources & Notes

CENTER FOR HUMANITIES COMPUTING

# Datasikkerhed og Governance

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

# Næste skridt

# Er danske modeller blot navlepilleri?

**KRONIKEN** 14. JUN. 2023 KL. 15.40
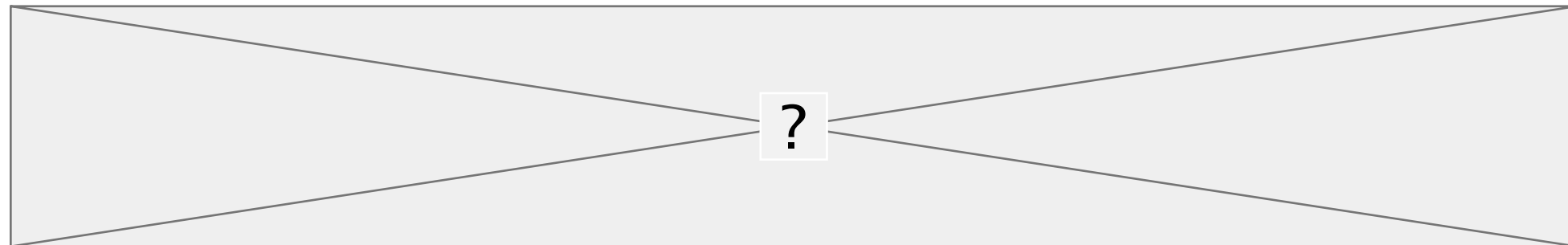
**ANDERS SØGAARD SUNE LEHMANN REBECCA ADLER-NISSEN OLE WINTHER MICHAEL BANG PETERSEN**

Anders Søgaard er professor, Københavns Universitet. Sune Lehmann er professor, DTU og Københavns Universitet. Rebecca Adler-Nissen er professor, Københavns Universitet. Ole Winther er professor, DTU og Københavns Universitet. Michael Bang Petersen er professor, Aarhus Universitet.

**TEKNOLOGI**

## Staten drømmer om sin helt egen chatbot – men manden bag Danmarks største it-virksomhed siger 'nej tak' til opgaven

Sources & Notes

https://www.dr.dk/nyheder/politik/professor-danmark-boer-udvikle-sin-egen-kunstige-intelligens

https://www.dr.dk/nyheder/viden/teknologi/staten-droemmer-om-sin-helt-egen-chatbot-men-manden-bag-danmarks-stoerste-it

https://politiken.dk/debat/kroniken/art9374522/Skab-et-offentligt-alternativ-til-techgiganterne

CENTER FOR HUMANITIES COMPUTING

# Er danske modeller blot navlepilleri?

- Store sprogmodeller er **fleksible** og **modulære**

  - **Kombinere** modeller

  - Modeller som **kulturelle vidensdatabaser**

  - …

- En potential fremtid

Sources
& Notes

Li, M., Gururangan, S., Dettmers, T., Lewis, M., Althoff, T., Smith, N. A., & Zettlemoyer, L. (2022). Branch-Train-Merge: Embarrassingly Parallel Training of Expert Language Models (arXiv:2208.03306). arXiv. https://doi.org/10.48550/arXiv.2208.03306

Feng, S., Shi, W., Bai, Y., Balachandran, V., He, T., & Tsvetkov, Y. (2023). Knowledge Card: Filling LLMs' Knowledge Gaps with Plug-in Specialized Language Models (arXiv:2305.09955; Version 2). arXiv. https://doi.org/10.48550/arXiv.2305.09955

CENTER FOR
HUMANITIES
COMPUTING

# The State of Foundation models for Danish

| | Model weights | Code Available | Model card | Data sheet | Language | Validated for Danish |
|---|---|---|---|---|---|---|
| **Text** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-encoder-large-v1** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓ |
| nb-bert-large | ✓ | ✓ | ✗ | ✗ | 🇳🇴 (🇩🇰 🇸🇪) | ✓ |
| XLM-Roberta | ✓ | ✓ | ✗ | ✗ | 🌐 | ✓ |
| *Generative models* | | | | | | |
| GPT-4 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✗* |
| DanskGPT | ✗ | ✗ | ✗ | ✗ | 🇩🇰 | ✗* |
| DanT5 | ✓ | ✗ | ✗ | ✗ | 🇩🇰 | ✗ |
| Llama-v2 | ✓ | ✗ | ✓ | ✗ | 🇺🇸 | ✗* |
| *Embeddings* | | | | | | |
| text-embedding-ada-2 | ✗ | ✗ | ✗ | ✗ | 🇺🇸 (🌐) | ✓* |
| MiniLM-L12-v2[1] | ✓ | ✓ | ✗ | ✓ | 🌐 | ✓ |
| **Speech** | | | | | | |
| *Structured learning* | | | | | | |
| **dfm-xls-r-300m** (ours) | ✓ | ✓ | ✓ | ✓ | 🇩🇰 | ✓[†] |
| wav2vec2-base-da | ✓ | ✓ | ✗ | ✗ | 🇩🇰 | ✓[†] |

Sources & Notes

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

CENTER FOR HUMANITIES COMPUTING

# Næste skridt

- Modeller

  - 7B generative modeller —> skalering til større modeller

- Datasamarbejder

  - Sikre højkvalitetsdata tæt på anvendelse

- Evalueringssamarbejder

  - Sikre evaluering på danske anvendelsesområder

CENTER FOR
HUMANITIES
COMPUTING

# Afslutning

# Særlig tak til det nuværende hold

Kenneth Enevoldsen[*1,2]        Lasse Hansen[*2,1]        Dan S. Nielsen[3]

Rasmus A. F. Egebæk[4]        Søren V. Holm[4]        Martin C. Nielsen[4]

Martin Bernstorff[2, 1]        Rasmus Larsen[3]        Peter B. Jørgensen[3]

Malte Højmark-Bertelsen[5]        Peter B. Vahlstrup[1]        Per Møldrup-Dalum[1]

Kristoffer Nielbo[1]

[1]Center for Humanities Computing, Aarhus University, Denmark
[2]Department of Clinical Medicine, Aarhus University, Denmark
[3]The Alexandra Institute, Copenhagen, Denmark
[4]Alvenir, Copenhagen, Denmark
[5]Beyond Work
kenneth.enevoldsen@cas.au.dk
lasse.hansen@clin.au.dk

CENTER FOR
HUMANITIES
COMPUTING

# Lær mere

**Videnskabelig Artikel**

**Hjemmeside**

**GitHub**

**Kontakt os**

Sources & Notes

Enevoldsen, K., Hansen, L., Nielsen, D. S., Egebæk, R. A. F., Holm, S. V., Nielsen, M. C., Bernstorff, M., Larsen, R., Jørgensen, P. B., Højmark-Bertelsen, M., Vahlstrup, P. B., Møldrup-Dalum, P., & Nielbo, K. (2023). Danish Foundation Models (arXiv:2311.07264). arXiv. http://arxiv.org/abs/2311.07264

CENTER FOR HUMANITIES COMPUTING