CASE 1: HOPEING TO FAIR WELL

Kristoffer Nielbo kln@cas.au.dk





OUTLINE

1. During COVID-19

HOPE FAIR ambition Media monitoring

2. Post COVID-19

The unstructured challenge Models cards and data sheets





"Open Science is a complex matter, where the demands for openness and transparency are difficult to comply with."





OUTLINE

1. During COVID-19

HOPE FAIR ambition Media monitoring

2. Post COVID-19

The unstructured challenge Models cards and data sheets





How Democracies Cope with COVID19

How do democracies react and cope as the COVID-19 crisis unfolds and with what effects?

- trajectory of the COVID-19 pandemic
- decisions of governments and international organisations
- decisions of media and social media landscapes
- citizens' behavior and well-being

Open Science and FAIR ambitions from day one







Dana Schutz's 'Building the Boat While Sailing', source: The New Yorker, 2012





ENTER FOR

OMPUTING

Types of data, source: Coveo, 2022





Public facing information

CENTER FOR HUMANITIES COMPUTING

Three modes of data-driven public communication

- News-like website: hope-project.dk
- Interactive dashboard: hope-project.dk/dashboard/
- HOPE reports: hope-project.dk/reports



Oplevet læring fra corona-krisen



SLIDE 7 IN 16

All data were stored and shared internally in a GDPR-compliant VRE and a shared data processing agreement

Open (structured) data

- Data pushed in regular intervals and available upon request
- All survey materials and data deposited in a publicly accessible database (OSF and Zenodo)
- FAIR ambition using a minimal DCAT v2 template





Open source

Proto FAIR software approach:

- GitHub publicly accessible repo with version control
- MIT license added license
- Zenodo registered in community registry
- CITATION.cff citation of software
- online sustainability evaluation quality checklist













CENTER FOR HUMANITIES COMPUTING

OUTLINE

1. During COVID-19

HOPE FAIR ambition Media monitoring

2. Post COVID-19

The unstructured challenge Models cards and data sheets







The vast majority of HOPE's data are (potentially) sensitive and proprietary and cannot be shared in their original form. Can we find a way to convert those data to wel





Transformers and Transfer Learning

1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp. patterns in language. By the end of the training process

2 - Supervised training on a specific task with a labeled dataset Supervised Learning Step



BERT broke several records for how well models can handle language tasks. Soon after the release of the paper describing the model, the team also open-sourced the model code and made it available for download versions of the model that were already pre-trained on massive datasets.





Danish Foundation Models

- 1. Develop and maintain **state-of-the-art** models for Danish
- 2. which are **well-validated** across a wide range of tasks

3. Furthermore, we wish to **ensure good documentation**, which allows users to assess the model for their use-case critically

4. **Open-source**, both model and source code





National Improvements

	Text		Audio	
	Danish	English	Danish	English
Architectures	BERT, (ELECTRA), (T5)	BERT, DeBERTaV3, GPT	(Wav2vec2)	Wav2Vec, HuBERT, WavLM
Largest pre-training corpus	1-1.6B Tokens	>10 0B tokens	1.300 hours	60.000 hours
Quality	Duplicate removal Removed cookies (Varied domains)	Near duplicate removal Quality fittered Varied domains Repetitious text removal	Single domain (audiobooks)	Audiobooks YouTube Podcasts
Number of parameters	14M - 100 M (770M)	8M - 6008	135M	135M - 1B
Compute resource	20 days on 4 A100	Multiple months on >1024 TPUs w. Multiple of experimentation	7 days 4 titan GPUs	1 month on 8 A100 GPUs

Danish was lacking behind - 2021 the Royal Library of Norway releases a model that outperforms Danish models by 4-5 percentage points.



Not including multilingual models, ex. Whisper. Parenthesis indicates that we have the model, but it is small or trained on limited data.







Mainland Scandinavian NLU Benchmark: ScandEval





THANKS

kln@cas.au.dk

chc.au.dk

SLIDES

knielbo.github.io/files/kln-hope-os.pdf

ACKNOWLEDGEMENTS

This research was supported the "HOPE - How Democracies Cope with COVID-19"project funded by The Carlsberg Foundation with grant CF20-0044, NeiC's Nordic Digital Humanities Laboratory project, and DeiC Type-1 HPC with project DeiC-AU1-L-000001.



