

# “2<sup>nd</sup>” Year PhD Seminar Presentation

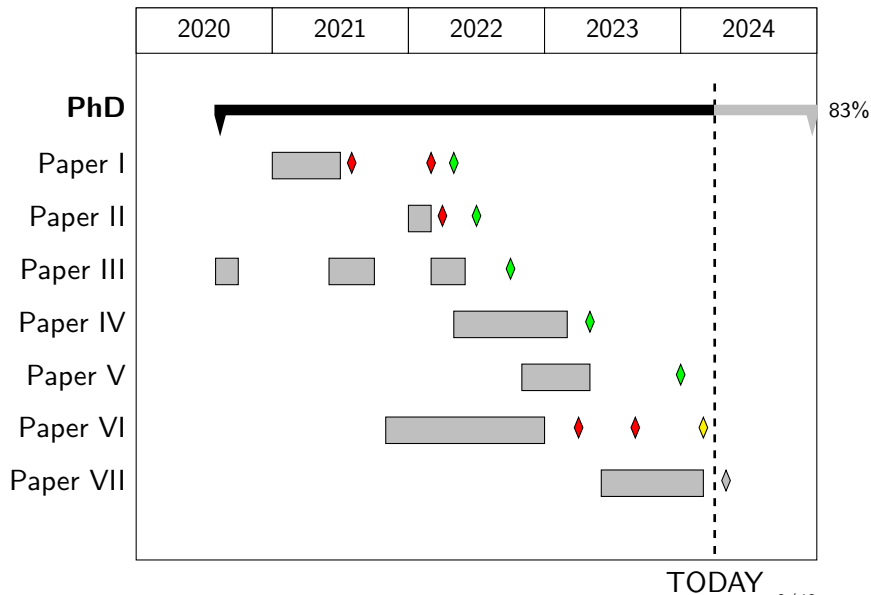
Ananth Mahadevan

February 23, 2024

# Quick Recap

- **Masters:** Aalto University
- **Started:** Aug 2020 (contract) and Jan 2021 (study right)
- **Supervisor:** Michael Mathioudakis
- **Research Group:** Algorithmic Data Science (ADS)

# PhD Progress



# Papers

Number	One Word Title	Venue	Include in Thesis?
I	Unlearning	MAKE 2022	✓
II	JANE	Entropy 2022	✓
III	Sketching	CIKM 2022	🙋
IV	ReceptionReader	JOHD 2023	🙋
V	Mandeville	DES 2023	🙋
VI	Retraining	KBS 2024	✓
VII	TextReuse	VLDB 2024	✓

# Research Projects

## Multiple Projects:

- ① Maintaining ML models
  - Paper I: Unlearning
  - Paper VI: Retraining
- ② Analyzing Historical Documents
  - Paper IV: ReceptionReader
  - Paper V: Mandeville
  - Paper VII: TextReuse
- ③ Scaling and Evaluating Algorithms
  - Paper II: JANE
  - Paper III: Sketching
  - Paper VIII ?: Diverse Sampling

# Maintaining ML models

## Research Question

How to update a trained ML model when the data changes?

- ① Machine Unlearning
  - Training data is deleted/removed
  - Update model parameters to forget information
- ② Cost-Aware Retraining
  - Streams drift over time
  - Data and Queries are present
  - Retraining consumes resources
  - When is it worth retraining?

# Analyzing Historical Documents

Shaftesbury, Anthony Ashley Cooper (1671-1713)  
1708 - A letter concerning enthusiasm

concerning ENTHUSIASM. 27.

Nation expresses it) it is necessary a People shou'd have a *Publick Leading* in Religion. For *to deny the Magistrate a Worship, or take away a National Church, is as mere Enthusiasm as the Notion which sets up Persecution.* For why shou'd there not be publick Walks, as well as private Gardens? Why not publick Librarys, as well as private Education and Home-Tutors? But to prescribe bounds to Fancy and Speculation, to regulate Mens Apprehensions and religious Beliefs or Fears, to suppress by Violence the natural Passion of Enthusiasm, or to endeavour to ascertain it, or reduce it to one Species, or bring it under any one Modification, is in truth no better Sense, nor deserves a better Character, than what the Comedian declares of the like Project in the Affair of Love——

D 2

*Nihil*

Astell, Mary (1668-1731)  
1709 - Bart'lemy fair

AN ENQUIRY AFTER WIT.

171

well'd against *Enthusiasm*, as the Title imports? I wish it may prove so, for the *poor Gentleman's* sake; for if it be otherwise, the Wit of his *Letter* will be no equivalent for the Profaneness: Let us see then his account of *Enthusiasm*. By p. 27. one wou'd take it to be nothing else but *Persecution*, from whatever hand it comes, whether the Magistrate persecutes the People, or they the Magistrate: For, says the Letter, *to deny the Magistrate a Worship, or take away a National Church, is as mere Enthusiasm as the Notion which sets up Persecution.* Whence it follows, that whatever Privilege the *good People* may claim, *Magistrates* must not pretend to be out of the hearing of a Deity. Their *P. 58.* Understandings ought to be superior to those *old and contrary Storys*, which (as it seems) puzzle little Wits. They shou'd *Know* that there is a GOD; for if they Disbelieve or Doubt, it will be the most *Ridiculous Formality in the World* to Worship one. *Magistrates*, in a word, must neither be Atheists nor Deists, for these *good People* have no manner of occasion for a *Worship*, or a *National Church*, 'twou'd be highly Ridiculous in them to pretend to any: Whereas, according to the *Letter*, a Church and a *Worship* are so absolutely necessary to the *Magistrate*, that he can't be without them, if he means to guard against *Enthusiasm*.

AGAIN

# Analyzing Historical Documents

However, OCR texts are very noisy

## Document 1 string

```
to deny the Ma-  
tifl:ate a Worflip, or take away a  
hational Church, is as mere En-  
Ihufiafin as the Notion which sets  
uip Persecution
```

## Document 2 string

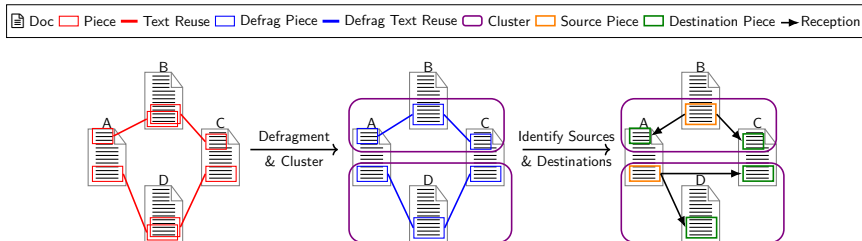
```
, to deny the Ala-  
inrate a ITorftip, or take awway a National  
'uircb, is as mere Entnztfiafm as the Notion  
.bic fJets tup Persecution. W
```

How?

- Use BLAST to do fuzzy alignment



# Pre-Processing Pipeline



- Clean up BLAST hits for downstream tasks
- Implemented in Apache Spark
- Scales up to **6.31 billion** pairs of reuses

# Related Publications

- ① Paper IV: Rosson et al. [2023]
  - *Reception reader: Exploring text reuse in early modern british publications*
  - Front-end user interface for browsing reuses
- ② Paper V: Ryan et al. [2023]
  - *A Comparative text similarity analysis of the works of Bernard Mandeville*
  - Study using the data and interfaces from Rosson et al. [2023]
- ③ Paper VI: Mahadevan et al. [2024]
  - *Optimizing a Data Science System for Text Reuse Analysis*
  - Studies design choices for the optimal performance of the system

# Scalability and Robustness

# References I

- Ananth Mahadevan, Michael Mathioudakis, Eetu Mäkelä, and Mikko Tolonen. Optimizing a data science system for text reuse analysis, 2024.
- David Rosson, Eetu Mäkelä, Ville Vaara, Ananth Mahadevan, Yann Ryan, and Mikko Tolonen. Reception reader: Exploring text reuse in early modern british publications. *Journal of Open Humanities Data*, Apr 2023. doi: 10.5334/johd.101.
- Yann Ryan, Ananth Mahadevan, and Mikko Tolonen. A comparative text similarity analysis of the works of bernard mandeville. *Digital Enlightenment Studies*, 1:28–58, 12 2023. ISSN 3029-0953. doi: 10.61147/des.6. URL <https://digitalenlightenmentstudies.org/article/id/6/>.