



THE UNIVERSITY of EDINBURGH  
Centre for Data, Culture & Society

# Coding my way out of a box: Achieving emancipation and collaboration through the Digital Humanities

Will Lamb,  
Professor in Gaelic Ethnology and  
Linguistics, University of Edinburgh



[www.ccds.ed.ac.uk](http://www.ccds.ed.ac.uk)



# Coding my way out of a box: Achieving emancipation and collaboration through DH

Will Lamb  
2 July 2024  
DHRSE Summer School



THE UNIVERSITY  
*of* EDINBURGH

Open to  
the world

# Outline



What are the Digital Humanities?



How do I decide what to work on?



What imperatives drive my own research?

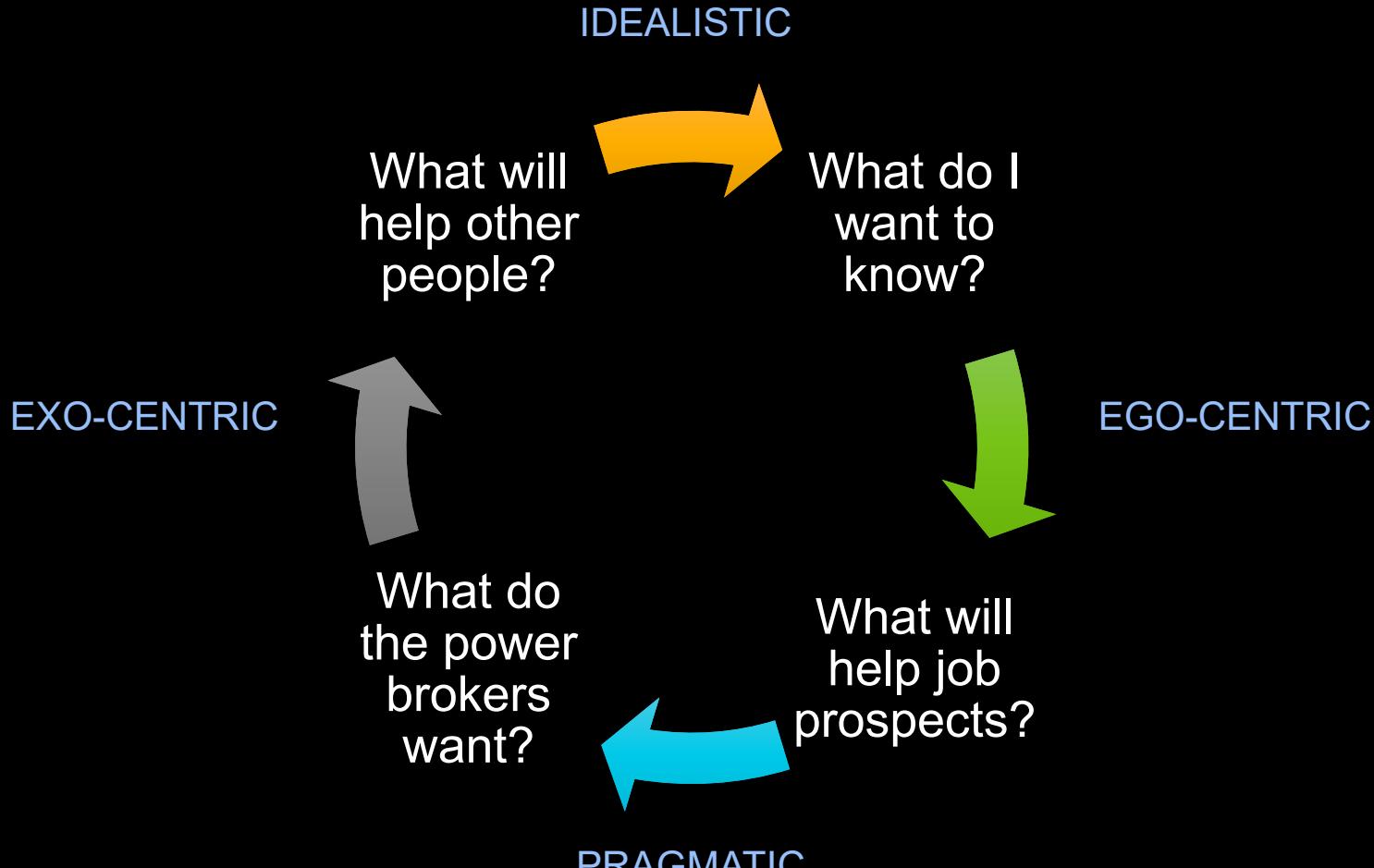


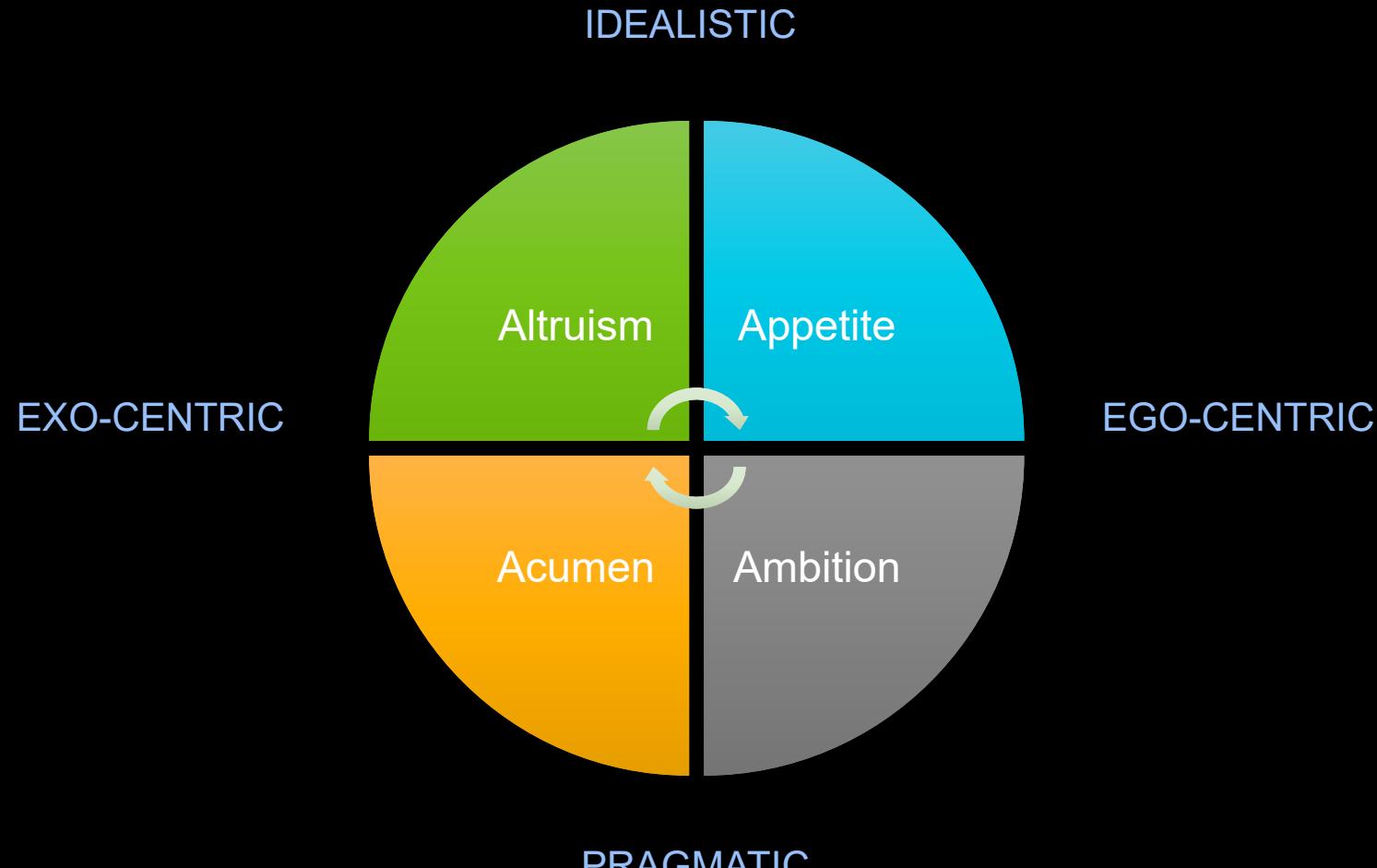
What has my transition from Gaelic linguistics to NLP been like?



## Defining the Digital Humanities



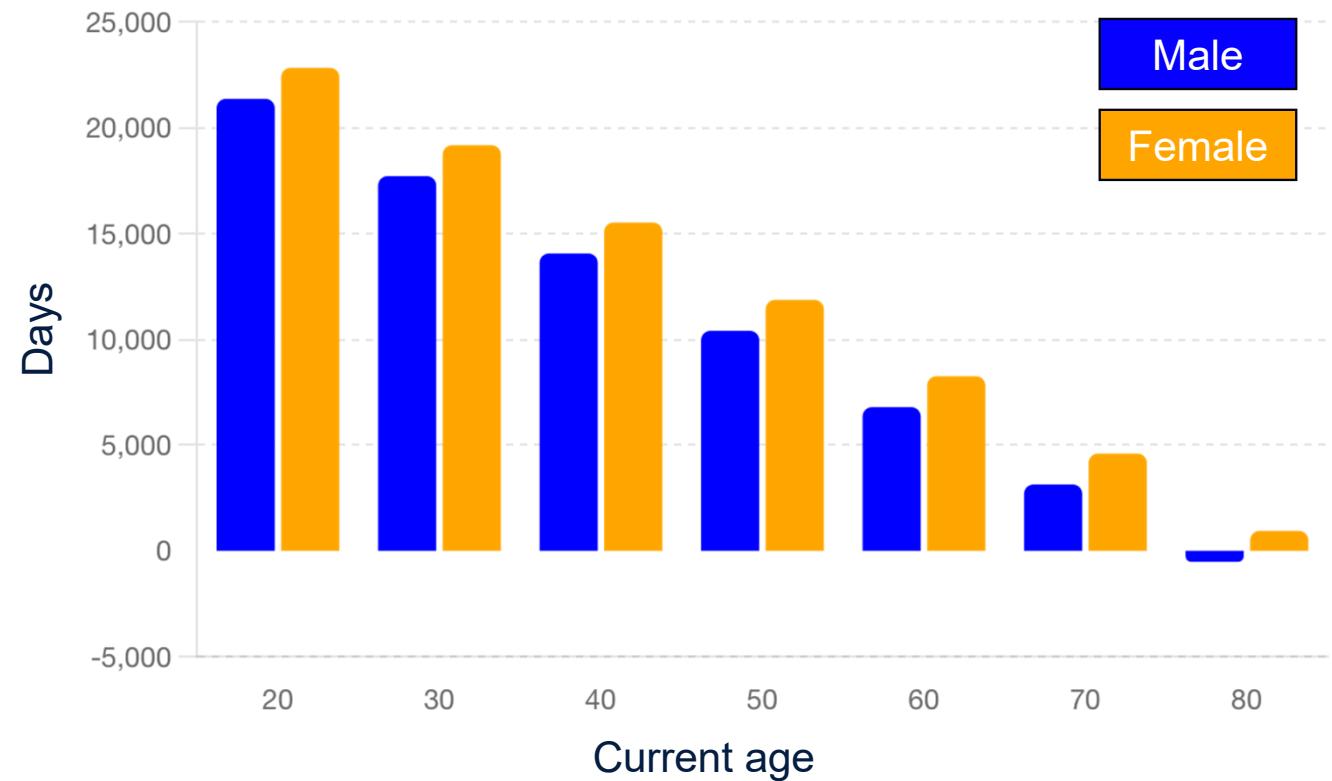




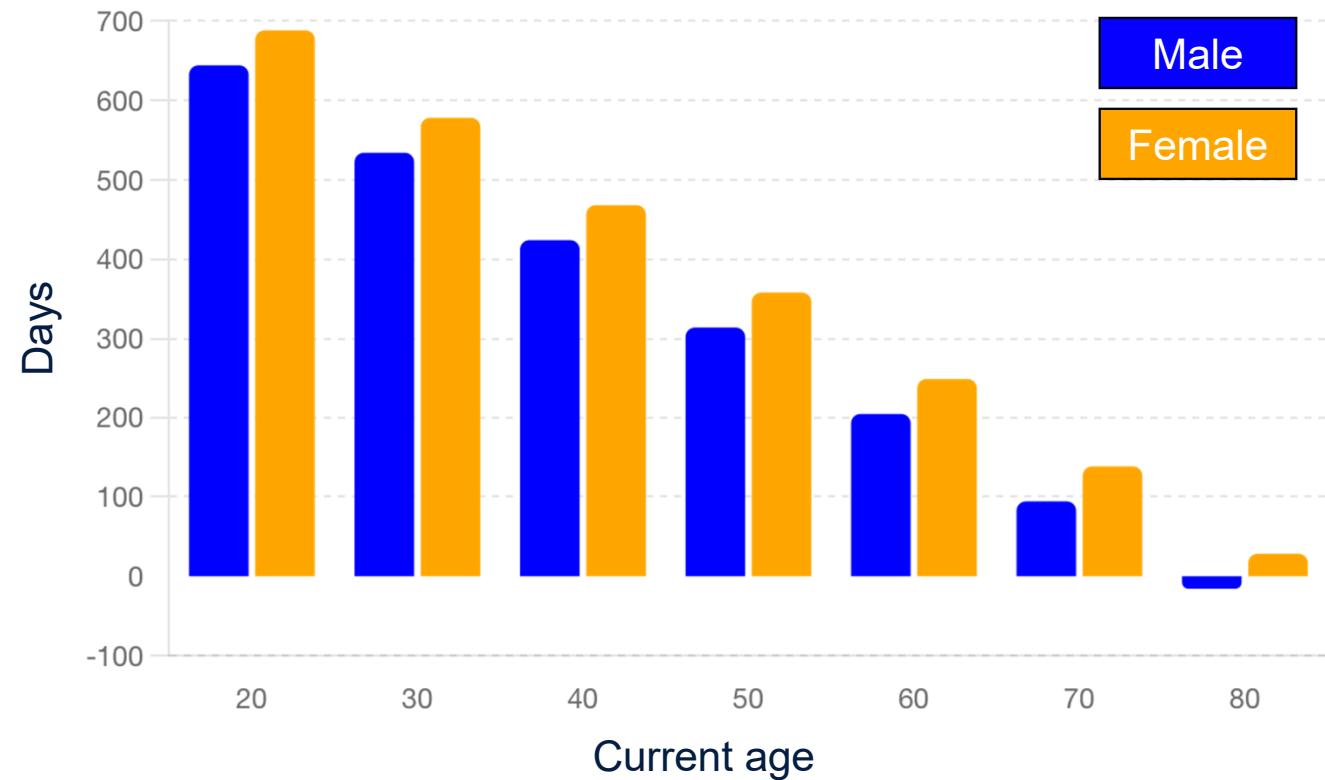
A photograph of a long, straight asphalt road stretching towards a horizon under a vast, dramatic sky filled with dark, textured clouds. The road is marked by a solid yellow center line and white edge lines. In the foreground, the perspective of the road creates a strong sense of depth. The overall mood is contemplative and inspiring.

**What drives you?**

# Average days left according to age (UK)



Average  
number of  
>20c days left  
(Edinburgh)



# Gaelic census data

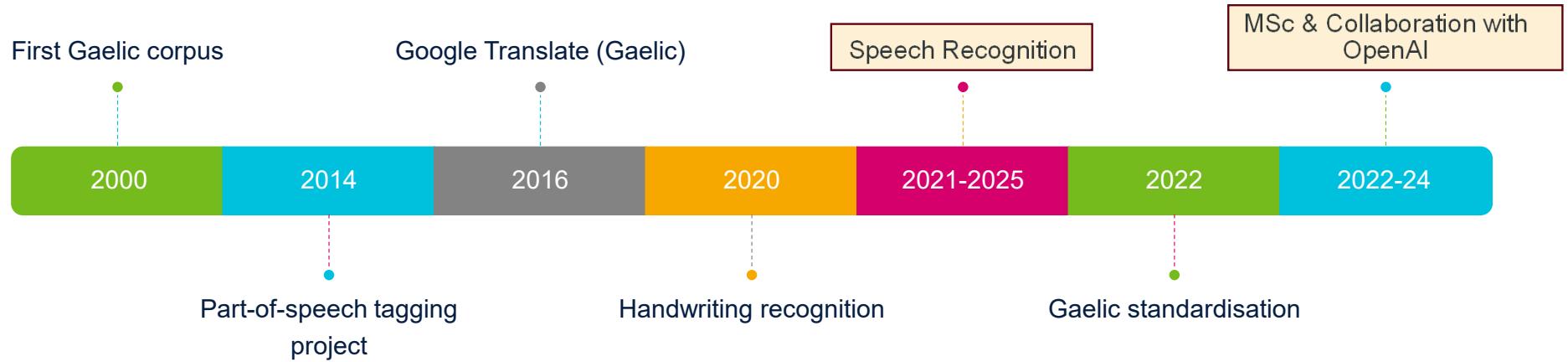




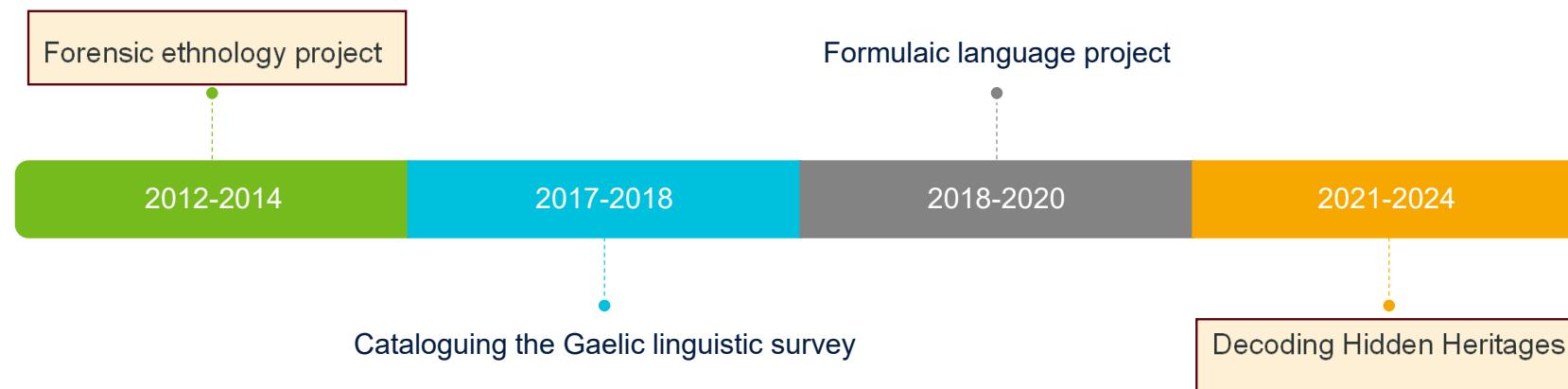
# PhD: Gaelic register variation (2002)



## Language Processing



## Digital Humanities



# Forensic Ethnology (2012-2014)

- Donald John MacDonald collection (School of Scottish Studies Archives)
  - 1953 through 1958
  - 26 bound volumes (Upper Library)
  - > 1500 pages attributed to Duncan and Neil
  - Largest manuscript collection of SSSA



# Dice similarity scores

The highest similarity is between D-1944 and N-1955

D-1944 = a published text transcribed from Duncan by K. C. Craig

N-1955 = an unpublished ‘transcription’ in DJM’s mss from Duncan’s brother Neill

Shows indisputable evidence of copying

Text 1	Text 2	Relation
D-1944	N-1955	0.87
D-1944	D-1953	0.82
D-1953	D-1950	0.82
D-1944	D-1950	0.81
D-1947	D-1944	0.81
D-1947	D-1953	0.80
D-1947	D-1950	0.79
N-1955	D-1953	0.78
D-1944	D-1936	0.77
N-1955	D-1950	0.77
D-1947	N-1955	0.77
D-1947	D-1936	0.76
N-1955	D-1936	0.76
D-1953	D-1936	0.75
D-1950	D-1936	0.75



# Visual evidence 1: ‘banais’

D-1944 pg 24

Chaidh mo chur a chadal ann an sobhal fada fàs an oidhche sin a rithist. Agus thàinig guth chon na h-uinneig, agus dh' éibh e gu robh dà latha seilgeadh agus sìdhneadh fhathast agam ri dhianamh mum faighinn banais no pòsadh.

‘Tha sin ann,’ orsa mi fhìn, ‘agus nam biodh an còrr ann, cha rachadh tus’ a dh’ inns’ an ath sgeul.’

24

DJM MSS pg 3557

The image shows a photograph of a page from a handwritten manuscript. The text is written in a cursive hand and appears to be in Scottish Gaelic. Several words are circled in red ink. In the center-left, the word 'banais' is circled. To its left, 'dheanamh' is also circled. Other circled words include 'robh', 'agus', 'sidhne', 'meo', 'meo', 'agam', 'faighinn', 'posadh', 'cursa', and 'easga'. There are also some smaller, less distinct circled words like 'an', 'a', and 'an'. The handwriting is somewhat faded and varies in size and style across the page.

# Visual evidence 2: ‘bainis’

D-1944 pg 26

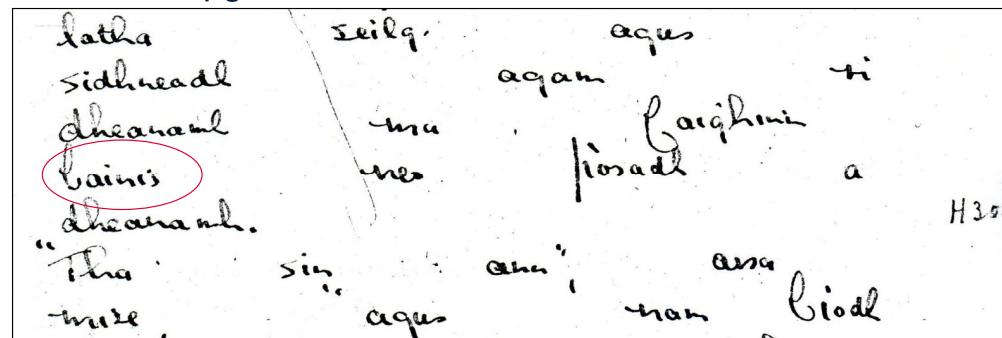
dhomh gu robh latha seilgeadh agus sìdhneadh agam ri dhianamh  
mum faighinn bainis no pòsadh.’

‘Tha sin ann,’ orsa mi fhìn, ‘agus nam biodh an còrr ann, cha  
rachadh tus’ a dh’ inns’ an ath sgeul.’

‘Cha bu lughaidh do chuid-sa a gheasachd an eilein sin,’ ors am  
fear eile, ‘mura cuirinn-sa na geasaibh ud ort, chuireadh fear eil’  
ort iad.’

26

DJM MSS pg 3557



---

# What I learnt from this project

1. A Gaelic handwriting recogniser would be useful
2. A Gaelic speech recogniser would be very useful
3. There is a huge amount of data in the School of Scottish Studies Archives



# The School of Scottish Studies Archives

- Research institution inaugurated in Jan 1951
- Tasked with collecting and interpreting Scotland's tangible and intangible cultural heritage



THE UNIVERSITY  
*of* EDINBURGH

Open to  
the world



# Decoding Hidden Heritages



Arts and  
Humanities  
Research Council



IRISH RESEARCH COUNCIL  
An Chomhairle um Thaighde in Éirinn



THE UNIVERSITY  
of EDINBURGH

Open to  
the world

# DHH Research Questions

- How can AI help to digitise and automatically recognise transcriptions of Gaelic and Irish folklore?
- What can these texts tell us about Scotland and Ireland's shared history and culture?
- How can traditional folkloristic and AI-based methods help to bring this to light?



THE UNIVERSITY  
*of* EDINBURGH

Open to  
the world

# DHH Principal Activities

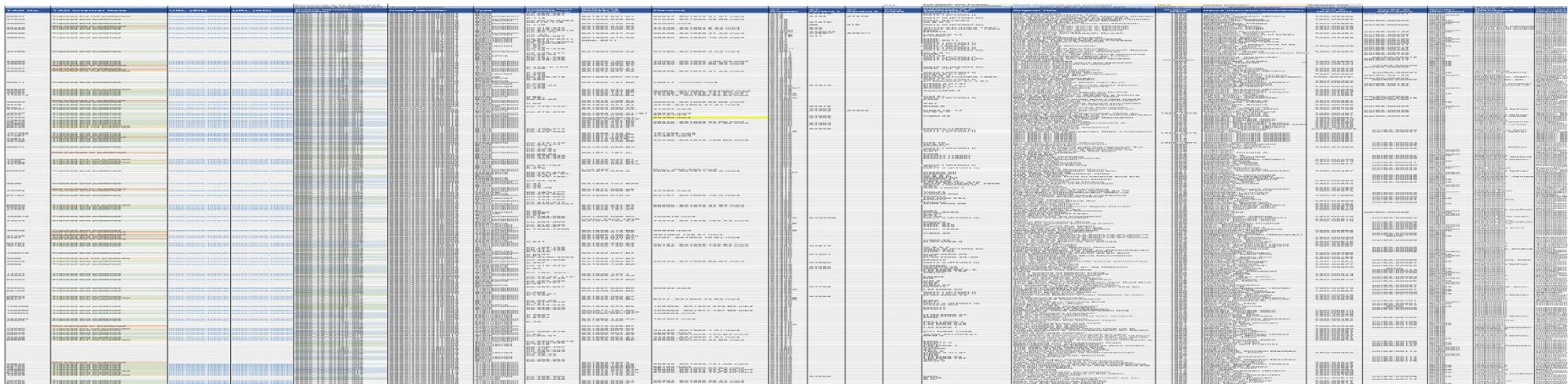
- Digitising and generating metadata for the Tale Archive at the School of Scottish Studies Archives (University of Edinburgh)
- Recognising printed and handwritten materials semi-automatically
- Clearing material for dissemination (e.g. copyright permission)
- Analysing the material with traditional methods and AI-driven techniques



THE UNIVERSITY  
of EDINBURGH

Open to  
the world

# SSSA Tale Archive: Basic Stats



## Scanned images

Document images scanned: 26,201

Index cards scanned: 12,419

Total scanned images: 38,620

## Metadata

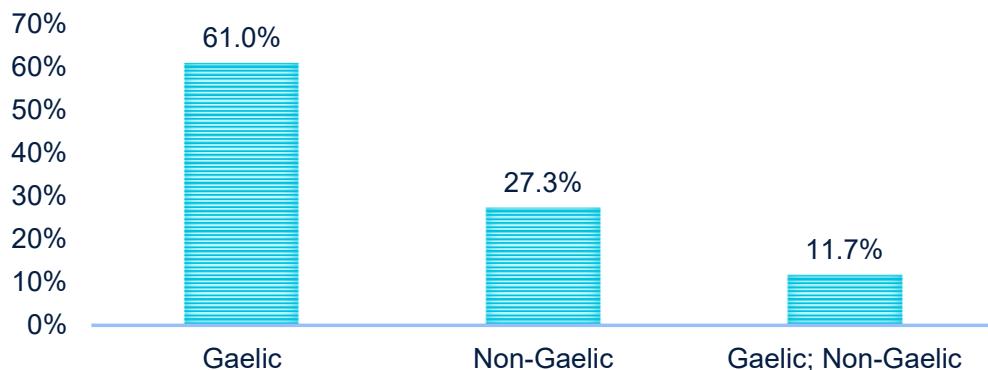
Document metadata rows completed: 3,860

Index cards metadata completed: 5,384

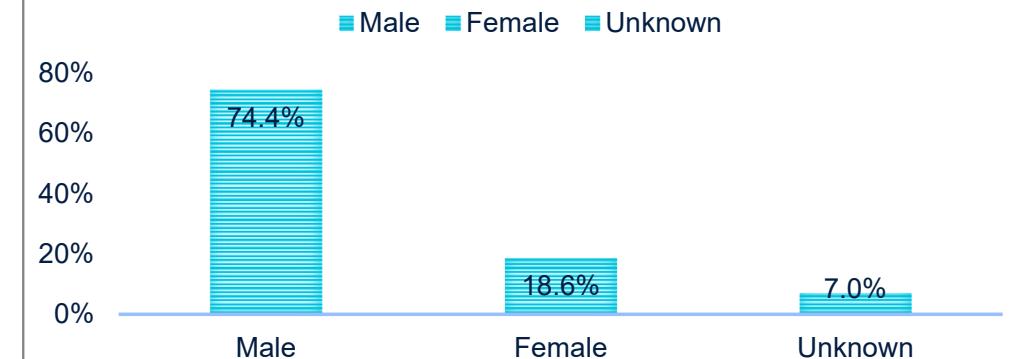
Total metadata entries: 9244

# SSSA Tale Archive: Basic Stats (cont.)

**DOCUMENT LANGUAGES  
PERCENTAGE**



**GENDER OF NARRATORS  
(PER TALE)<sup>1</sup>**



<sup>1</sup>[Decoding Hidden Women: Feminist digitisation practices in the Tale Archive](#)

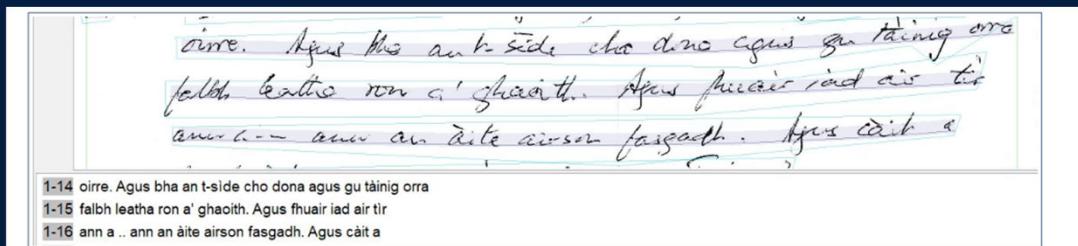


THE UNIVERSITY  
*of* EDINBURGH

Open to  
the world

# Handwriting Recognition

- Uses Transkribus platform
- Pages recognised: 8489
- Words recognised: 1.9M
- Accuracy levels
  - Char level accuracy: 97.5%
  - Word level accuracy: 95%



(a) Good quality output for the principal hand

A screenshot of the Transkribus software interface showing a page of handwritten text in Irish. The text is very blurry and illegible, with many characters and words misrecognized. Below the page, a list of recognized words is shown:

- 2-75 Muai
- 2-76 a
- 2-77 shotil
- 2-78 ann
- 2-79 aigin
- 2-80 ithist
- 2-81 Catha
- 2-82 a
- 2-83 an
- 2-84 Agus
- 2-85 Chuaidhe

(b) Bad quality output for a hand with little training data

Good vs bad quality Transkribus output



THE UNIVERSITY  
of EDINBURGH

Open to  
the world

# DHH text mining workstreams

1. Phylogenetic analysis
2. Topic modelling / automatic clustering
3. Automatic identification of formulaic language

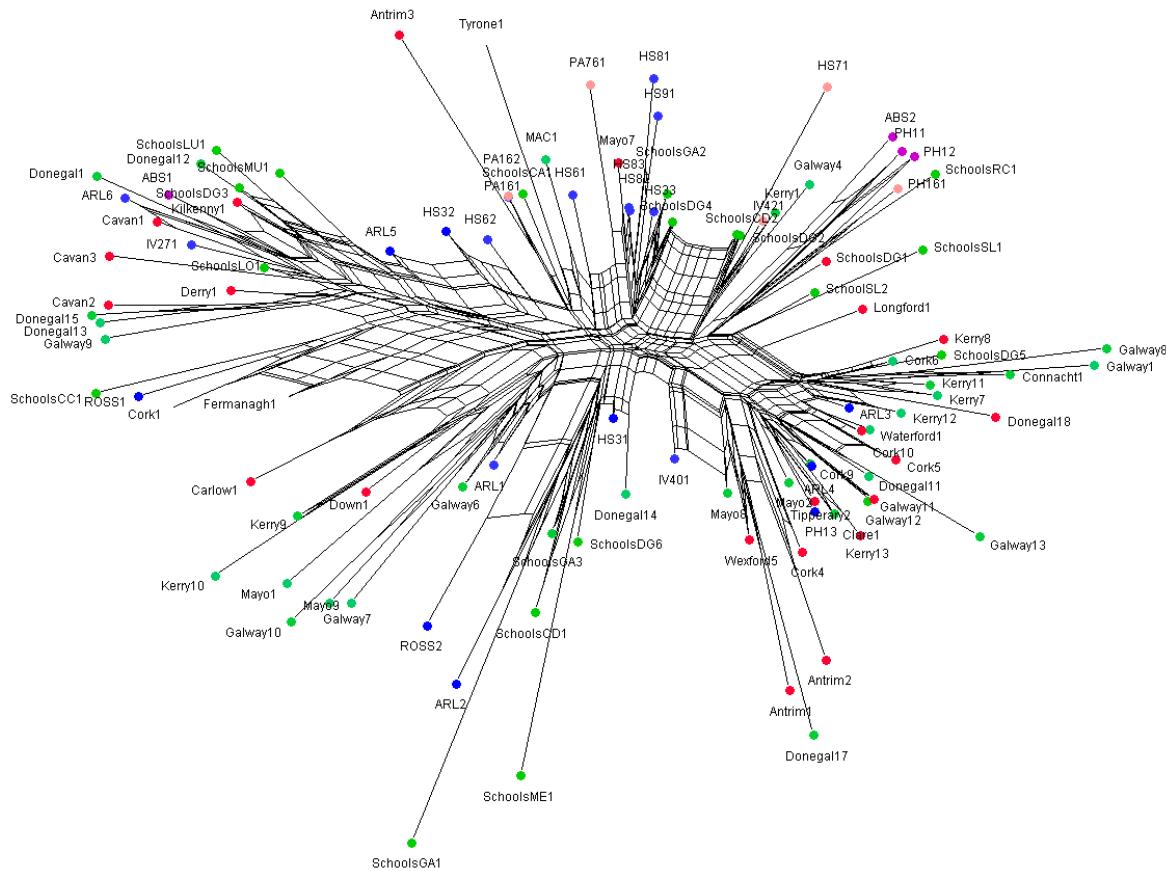


THE UNIVERSITY  
*of* EDINBURGH

Open to  
the world

# Phylogenetic network analysis: ATU 503

- Network analysis has a highly reticulated structure
  - Suggests complicated patterns of transmission between communities and linguistic traditions
  - Bottom line: lots of borrowing and mixing of motifs

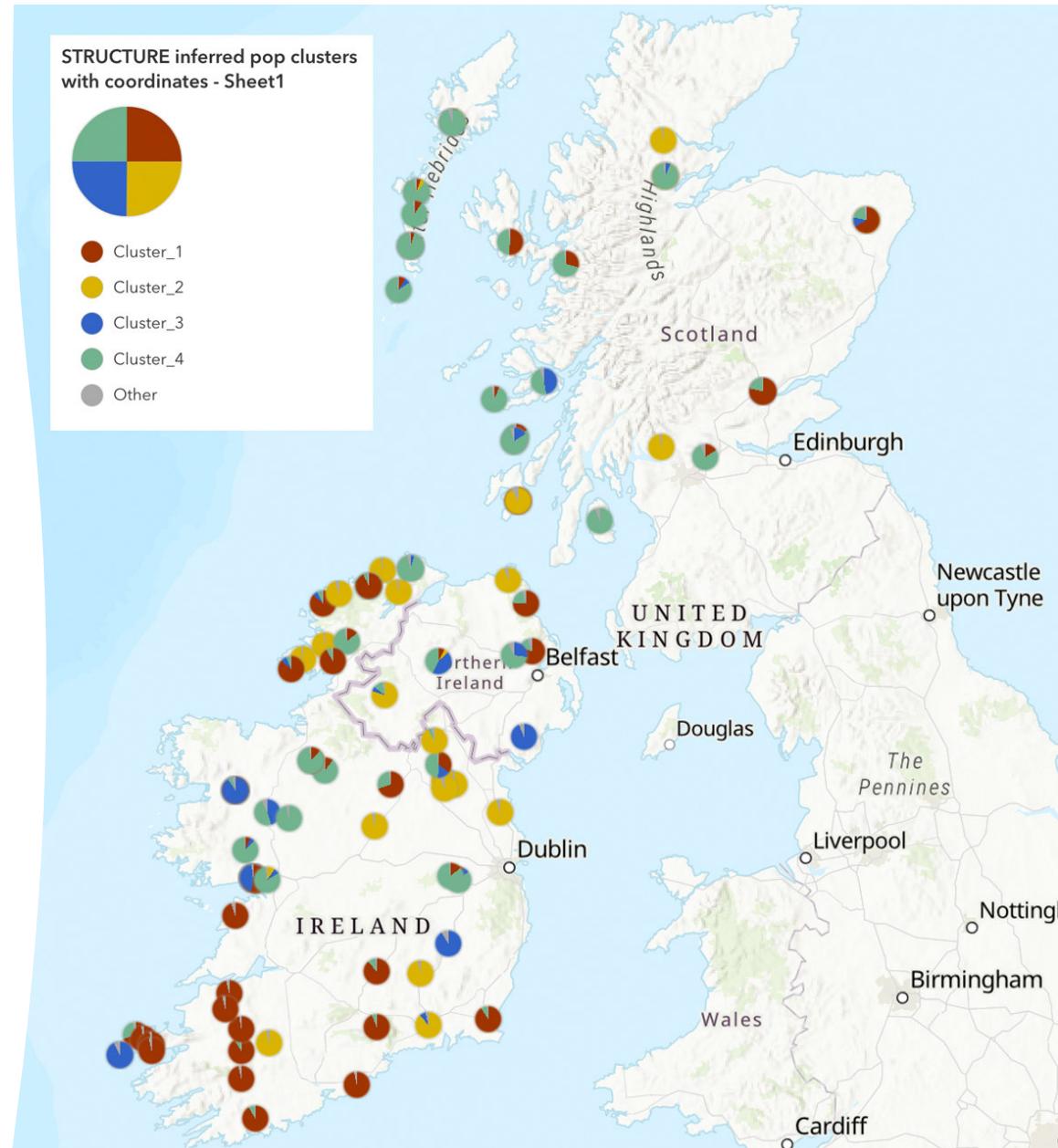


THE UNIVERSITY  
of EDINBURGH

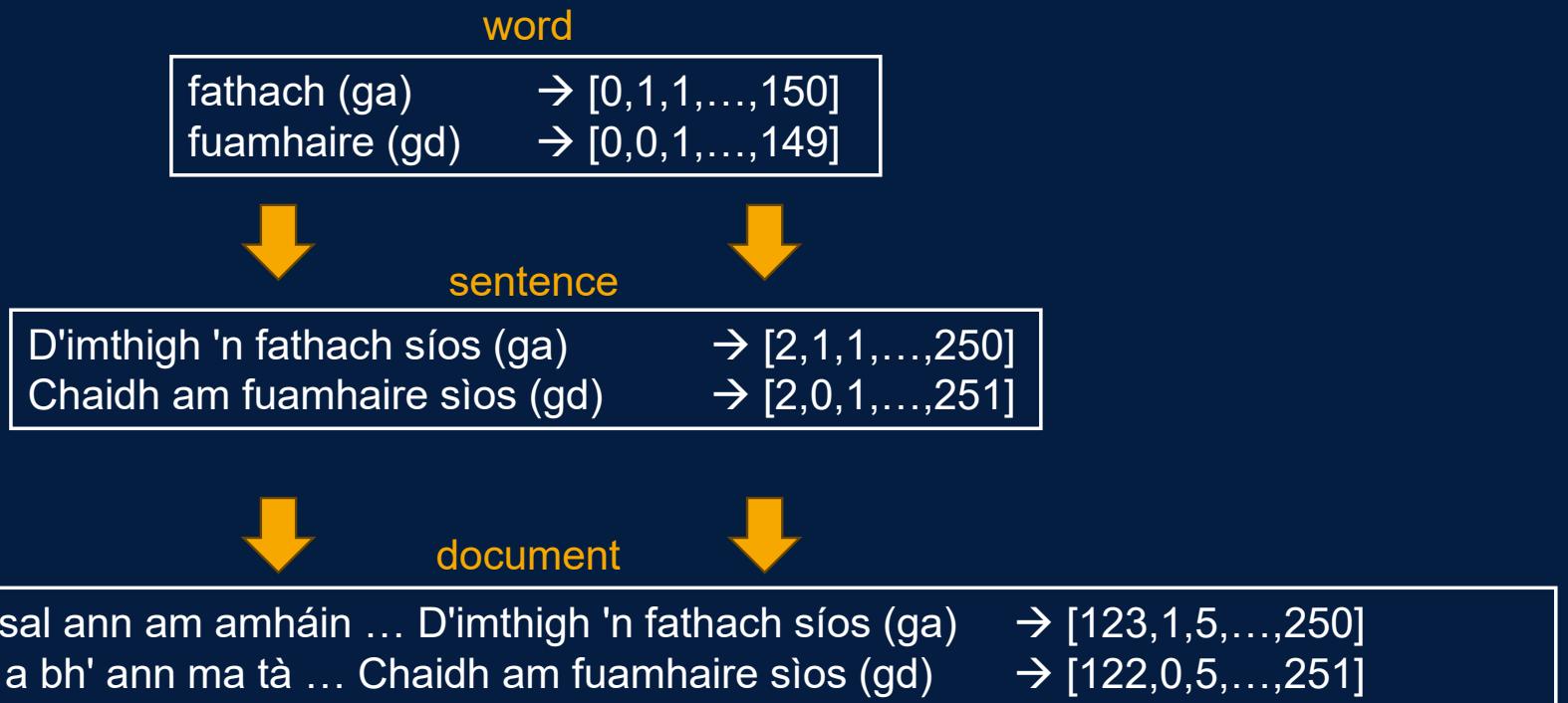
**Open to  
the world**

# Diatopic cluster distribution: ATU 503

- Shows some effect of geographical proximity (<50km)
- The effect of linguistic similarity, however, is stronger
  - Linguistic groups (e.g. Gaelic speakers in Scotland) have more in common with themselves than with other, geographically proximate groups



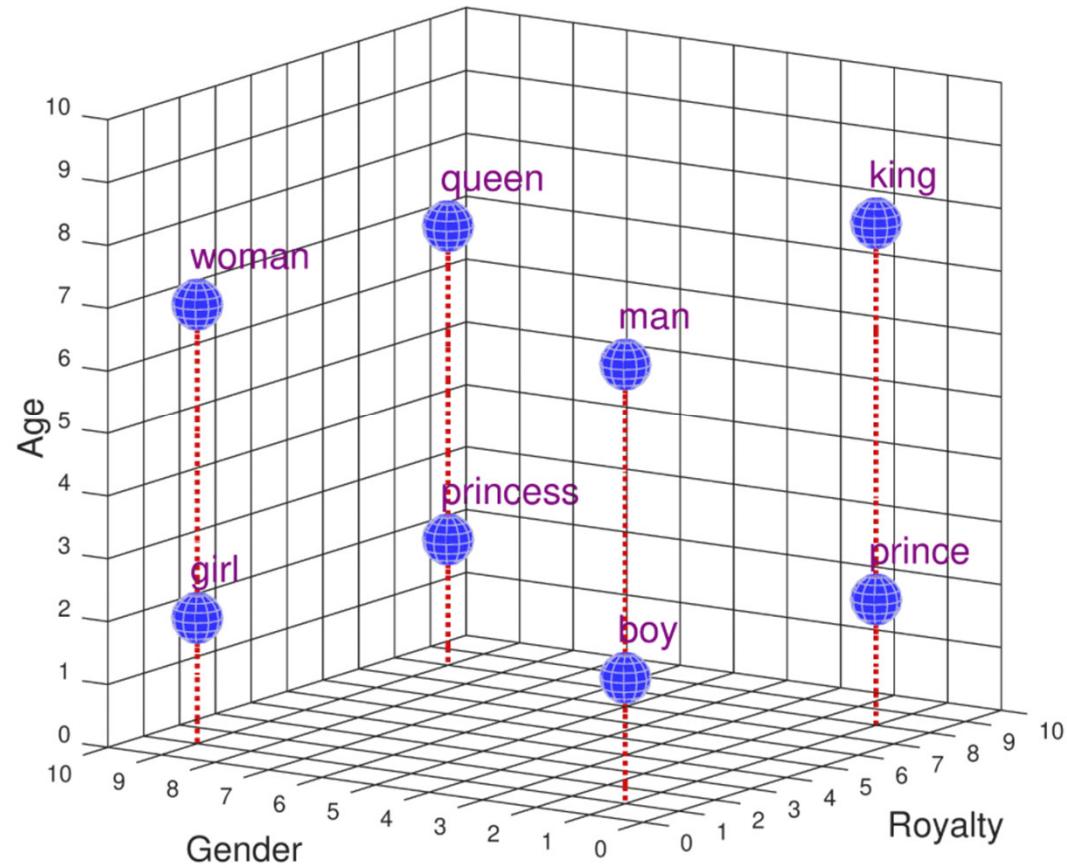
# Word to Sentence to Topic Embeddings



THE UNIVERSITY  
of EDINBURGH

Open to  
the world

## 3D Semantic Feature Space



# Decoding Hidden Heritages: Next steps

- Run repertoire analyses using the coordinates and metadata
  - Particularly interested in gendered distinctions
- Wrap up the phylogenetics and formulaic language research
- Launch the project website (Q4, 2024)
  - Will provide searchable texts and PDFs
- Prepare *Computational Folkloristics*, an edited collection of papers from the project



THE UNIVERSITY  
of EDINBURGH

Open to  
the world

# Gaelic speech recognition

seirbheis air an taobh a-mu...

Air neo...

o no leig às faidhleachan an toiseach...

Air neo...



## ÈIST project: Current aims

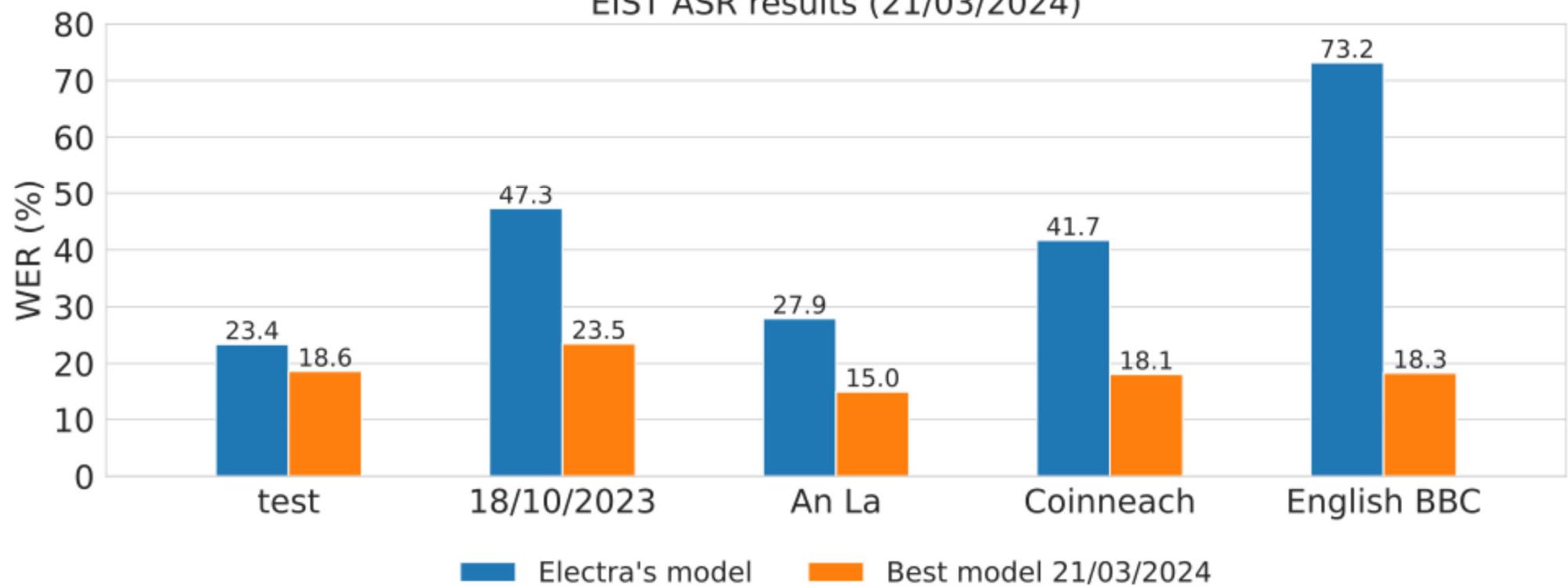
---

1. Produce a subtitling system for the BBC (radio and television)
2. Produce a transcription app for Gaelic school children
3. Expedite transcriptions of ethnographic material for Tobar an Dualchais / Kist o Riches



ECH TECHNOLOGIES

EIST ASR results (21/03/2024)



```
5 N_USER = 2 # number of user prompt sentences
6 N_AGENT = 30 # number of agent generate sentences
7
8 SYSTEM_MSG_SUMMARY = "Your role is to summarise the given Scottish Gaelic in 3 to 4 sentences in"
9 SYSTEM_MSG = "Your role is to generate news story in Scottish Gaelic given the English summary."
10
11 ROOT_PATH = os.getcwd()
12 DATA_PATH = os.path.join(ROOT_PATH, "data")
13 INPUT_FILE_NAME = "bbc_dongge_with_timestamps.jl"
14 OUTPUT_PATH = os.path.join(ROOT_PATH, "output")

n 12 1 with open(os.path.join(ROOT_PATH, DATA_PATH, INPUT_FILE_NAME), 'r', encoding='utf-8') as file:
2     content = file.read()
3 df_base = pd.read_json(content, lines=True)
4 df_base

▼ /var/folders/mt/z9020bs533j19ypwxz5g2kdm0000gq/T/ipykernel_72101/2124592222.py:3: FutureWarning
version. To read from a literal string, wrap it in a 'StringIO' object.
df_base = pd.read_json(content, lines=True)
```



# MSc SLP + OpenAI project

# Many thanks to our funders and collaborators

## *Taing mhòr dhan luchd-taic againn*

Am Faclair Beag, Ceòlas Uibhist Ltd, Digital Archive of Scottish Gaelic (DASG: U of Glasgow), European Ethnological Research Centre, Faclair na Gàidhlig, Grace Note Publications, Guthan nan Eilean / Island Voices, LearnGaelic (MG Alba), National Folklore Collection (University College Dublin), OpenAI, Ruairidh MacIleathain, Sabhal Mòr Ostaig, The National Library of Scotland, The School of Scottish Studies Archives, Tobar an Dualchais / Kist o Riches, University of Bangor, University of the Highlands and Islands, University of Glasgow and Riaghaltas na h-Alba (The Scottish Government)



Arts and  
Humanities  
Research Council



THE UNIVERSITY  
of EDINBURGH

Open to  
the world



THE UNIVERSITY  
*of* EDINBURGH

Ceud Mìle Taing  
*Thank you*

Open to  
the world