



Cabinet Office
Digital

Using AI to Review Records in the Cabinet Office

27 June 2024

David Canning MBE, Head of Digital Knowledge and Information Management
Dr Kelcey Swain, Development Operations Lead

Enabling and Transforming the Cabinet Office

○ Inclusive ○ Quality ○ Transparent ○



The Digital Paradigm

Issues affecting access to digital records in archives include:

- Volume
- Format
- Distribution
- File Plans!

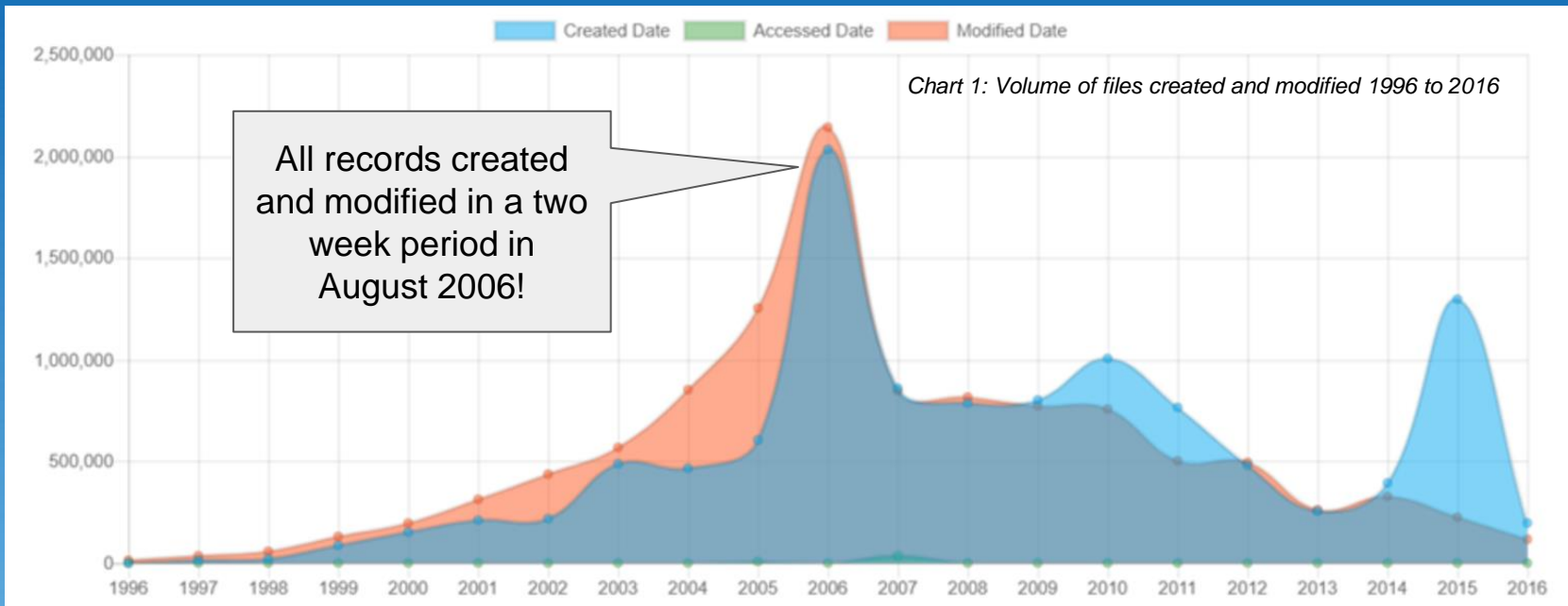
It's a digital problem, so we need digital solution



Enabling and Transforming the Cabinet Office



Live systems corrupt data!



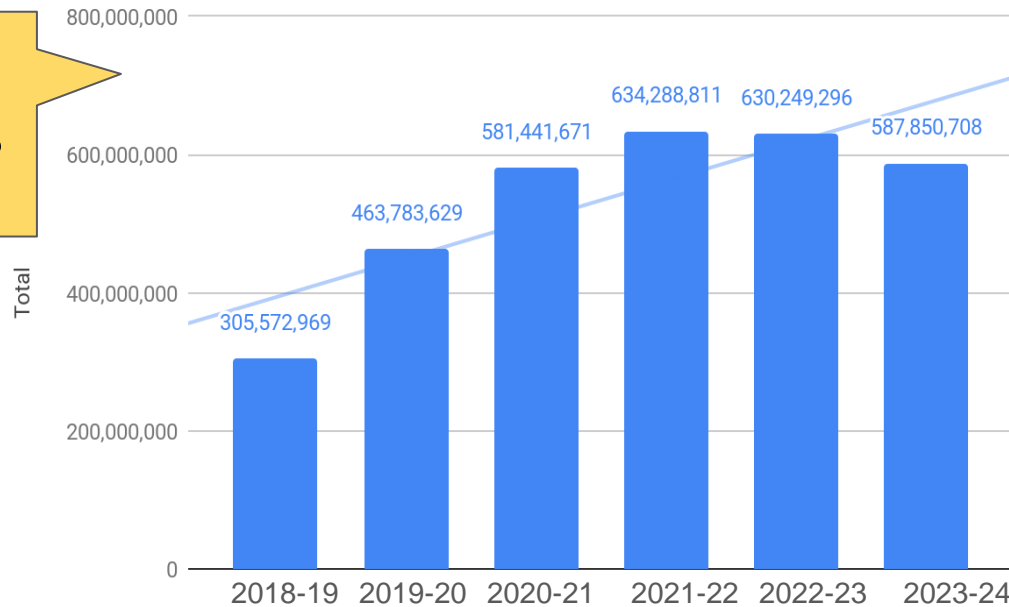
Enabling and Transforming the Cabinet Office



Data will just keep growing

How much of
this is R.O.T?

Total volume of digital objects (all information stores) 2018-19 to 2023-24



“The neuropsychological capacity of the human brain to process and record information may constitute the dominant limiting factor for the overall growth of globally stored information, with real-world economic constraints having only a negligible influence.”

Gros, Kaczor, Markovic
Goethe University,
Frankfurt, 2011

Enabling and Transforming the Cabinet Office



Digital Disposal Methodology

All information - 11.8million files

A manual review removed a further 2.8 million files through a top level review of content.

Classification Analysis - 3.3 million files removed

We removed 3.3 million files consisting of unwanted formats

Aggressive reduction - 4.6 million files removed

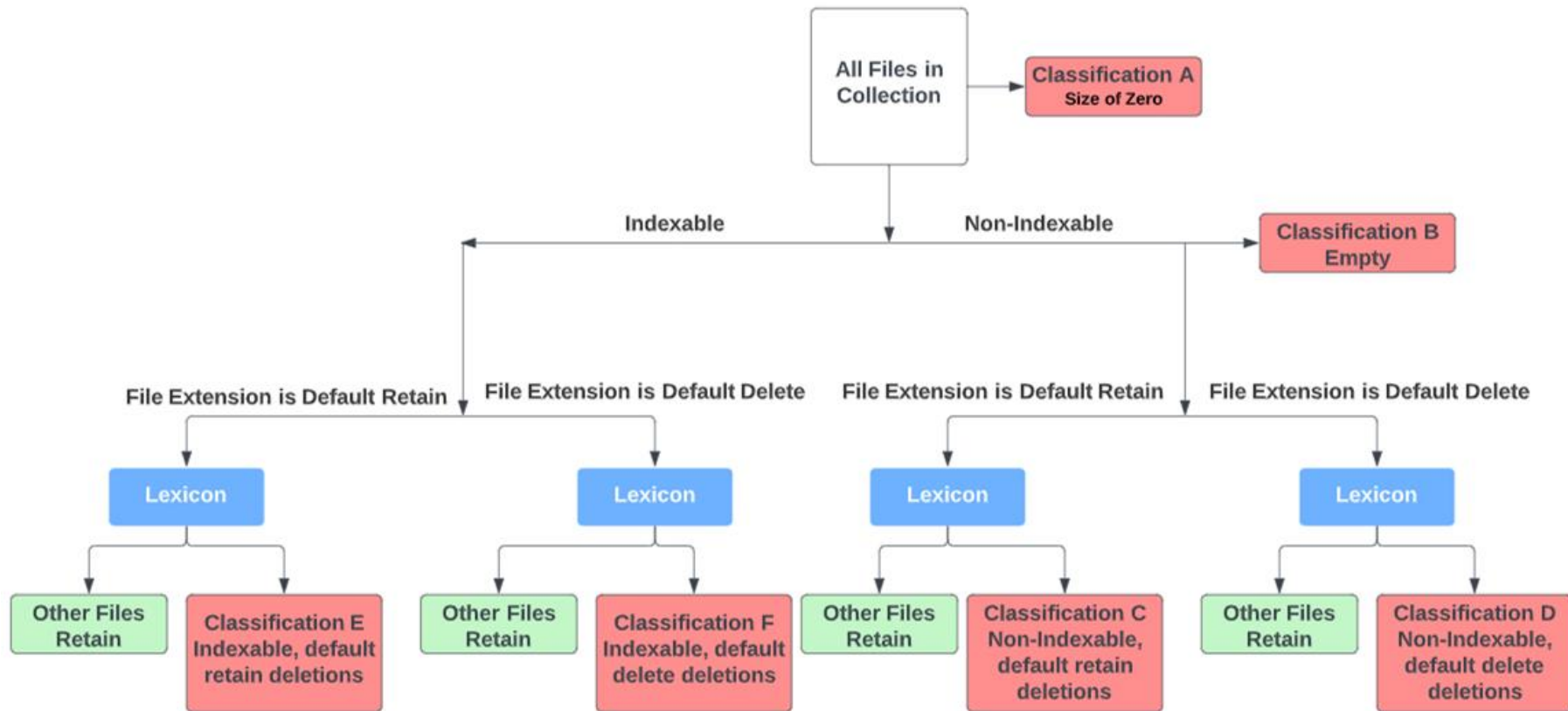
Automated processes identified 1.8 million files for removal

Human Review

3.9 million files to archive for preservation

Enabling and Transforming the Cabinet Office

Our Algorithmic Model for Document Review





Lexicon terms were categorised according to their long-term value

Value

Submission
Minister
Bilateral
Richard Wilson

ROT

Annual leave
Blah blah
Christmas card
Stuff

ROT, even when in the presence of Value terms

Auto-response
Template
Timesheet
000



Weighting the algorithm - bias

Bias towards **retention** because there is an assumption that the document will contain valuable information based on its format

Bias towards **deletion** because there is an assumption that the document will not contain valuable information based on its format

High - (Frequency of terms) - Low

High - (Frequency of terms) - Low

Low - (Weighting of terms) - High

Default Value Extension Types	Value	ROT	ROT
Value	Retain	Retain	Delete
ROT	Retain	Delete	Delete
ROT	Delete	Delete	Delete

Default ROT Extension Types	Value	ROT	ROT
Value	Retain	Delete	Delete
ROT	Delete	Delete	Delete
ROT	Delete	Delete	Delete

Files falling into this category may need to be human reviewed

Enabling and Transforming the Cabinet Office



Weighting the algorithm - Relevance

Term frequency (TF)

- The more times that a search term appears in the field we are searching in a document, the more relevant that document is.

Inverse document frequency (IDF)

- The more documents that contain a search term in the field that we are searching, the less important that term is.

Field length

- If a document contains a search term in a field that is very short (i.e. has few words), it is more likely relevant than a document that contains a search term in a field that is very long (i.e. has many words).



Outcomes and Benefits

Cost avoidance/ efficiency

- Legacy: £2.2m
- Backlog: £1.7m
- Future: £190k
pa

1xFTE pa required to
review 200,000
documents

Consistency & accuracy

- Predictable
- Auditable
- Consistent and
Reliable

Speed

- 59 years' work
completed in
one year



Next Steps Development

Machine Learning Techniques

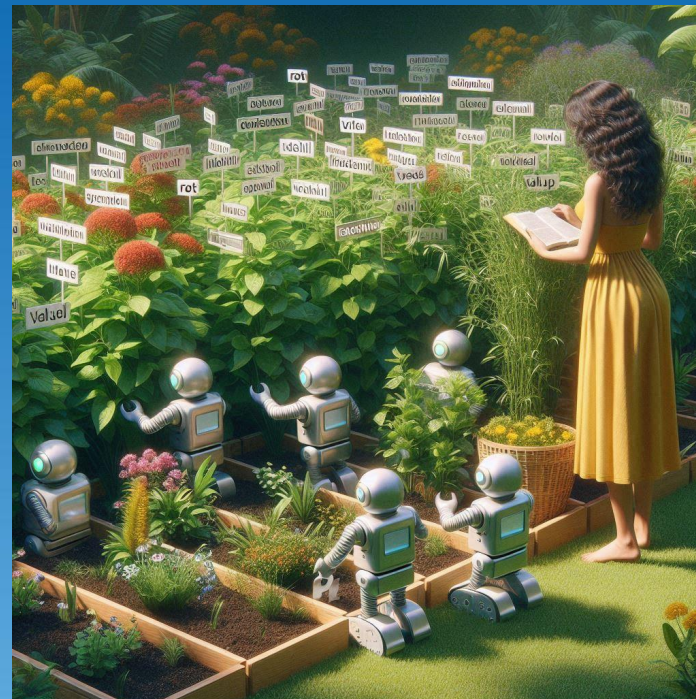
- Naive Bayes
- Markovian Discrimination

Mitigation

- Human-in-the-loop

Future-proof

- Retaining a trace of ROT



Enabling and Transforming the Cabinet Office