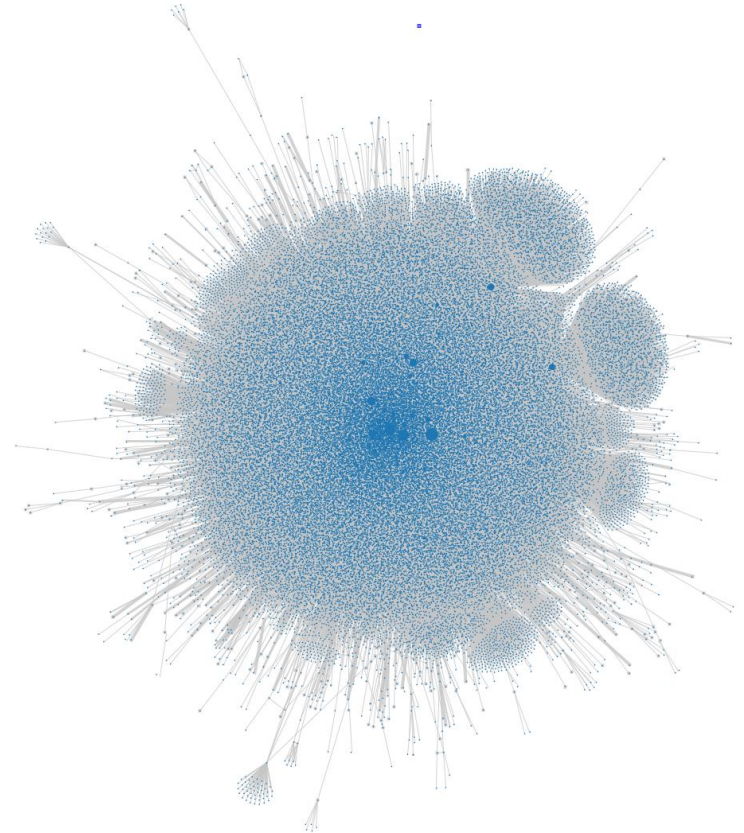# Ten years of the UK Web Archive:
## What have we saved?

Andy Jackson (@anjacks0n)

UK Web Archive Technical Lead

# The UK Web Archive

- Three collections:
    - Open Archive (since 2004)
    - Legal Deposit Archive (since 2013)
    - JISC Historical Archive (1996-2013)

- Statistics:
    - Over eight billion resources
    - Over 160TB compressed data

- Goals:
    - Preserve UK web history
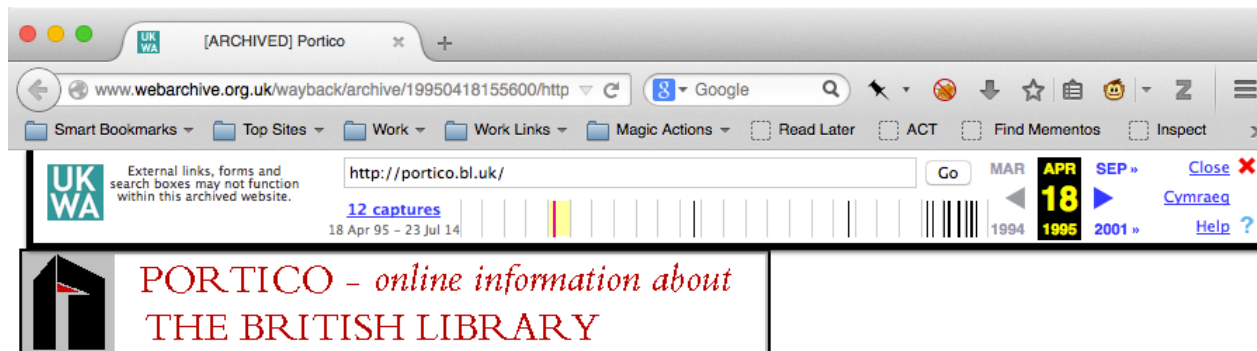    - Support access
    - Enable research

# Understanding Our Collections

# Resource-Level Access

# Curated Collections

# Full-Text Discovery & Trend Analysis

# Secondary Datasets

- JISC UK Web Domain Dataset (1996-2013):
  – Format Profile
  – Geo-Index
  – Host-Level Links
  – Crawled URL Index
  – WATs (not released yet)

- UK Open (Selective) Web Archive:
  – Website Classification Dataset

- Available as CC0 downloads:
  – http://data.webarchive.org.uk/opendata/

# Links From 1996

# Format & Feature Analysis

<applet>

<font>

<blink>

<script>

# Putting Our Archives In Context

- Looking inward is not enough:
  - To understand the value of our collection, we need to look beyond our walls and put it in context.

- Was it worth archiving?
  - How much of our collection is still on the live web?
  - How bad is reference rot in the UK domain?

# Open UKWA Crawl History



**Open UKWA Archived URLs**

# Sampling The URLs

- Use a random sample 1,000 URLs per year:
  - If the host name does not resolve, or is unreachable:
    - **GONE**
  - If the server responds with an error:
    - **ERROR**
  - If the server response leads to *404 Not Found*:
    - **MISSING**
  - If the server response leads to a valid resource:
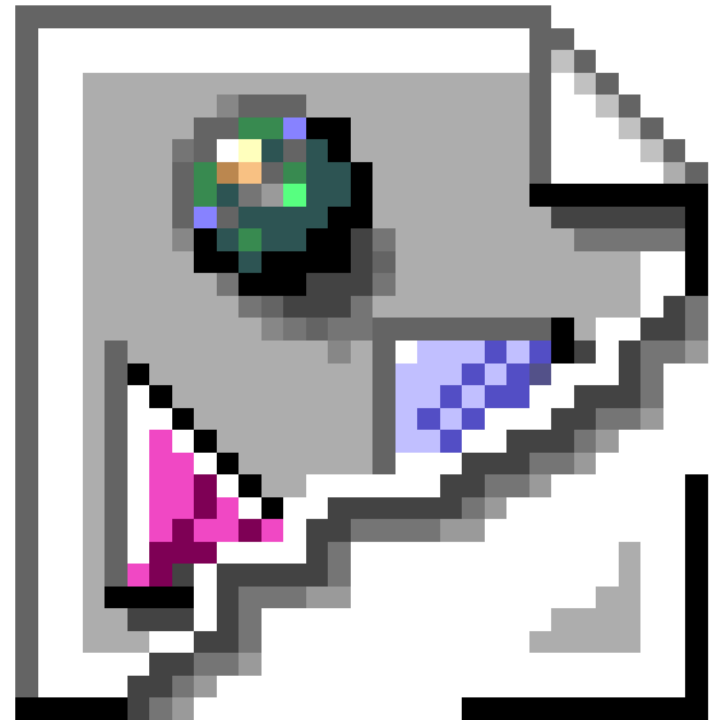    - **MOVED** (if via redirects)
    - **OK** (otherwise)

- n.b. 'soft 404s' are surprisingly rare (< 1%)

# Where Are They Now?

# NICE Example

# Extract The Text

CG121 Lung cancer: full guideline appendix 11 Sign In | Register Home News Get involved About NICE Find guidance NICE Pathways Quality standards Into practice QOF Conditions and diseases Blood and immune system ? Cancer ? Cardiovascular ? Central nervous system ? Digestive system ? Ear and nose ? Endocrine, nutritional and metabolic ? Eye ? Gynaecology, pregnancy and birth ? Infectious diseases ? Injuries, accidents and wounds ? Mental health and behavioural conditions ? Mouth and dental ? Musculoskeletal ? Respiratory ? Skin ? Urogenital Public health Accidents and injuries ? Alcohol ? Behaviour change ? Cancer ? Cardiovascular disease ? Child health ? Child social care ? Chronic illness ? Diabetes ? Drugs ? Environmental health ? Infectious diseases ? Maternal health ? Mental health ? Non-communicable diseases ? Obesity and diet ? Occupational health ? Older people ? Physical activity ? Sexual health ? Smoking and tobacco ? Transport ? Vaccine preventable diseases ? Working with and involving communities Treatments, Procedures and Devices Bones and joint surgery ? Cardiovascular surgery ? Cardiovascular system drug treatments ?

# Generate Fingerprints

- We use the 'ssdeep' fuzzy hash algorithm to generate a fingerprint for the extracted text
  - Compare fingerprints instead of content

- Earlier This Year:
  - aDJjTi6KVkfrehQfnSSXWYjqyBmiF8H9

- From The 2013 Archive:
  - aDJjTi6KWkfrehQfN+SSXWZjbO4kiF+H2LZcn

- Similarity Result: 50%

# NICE Example (Archived in 2013)

# NICE Example (this year)

# Page Footer Problem
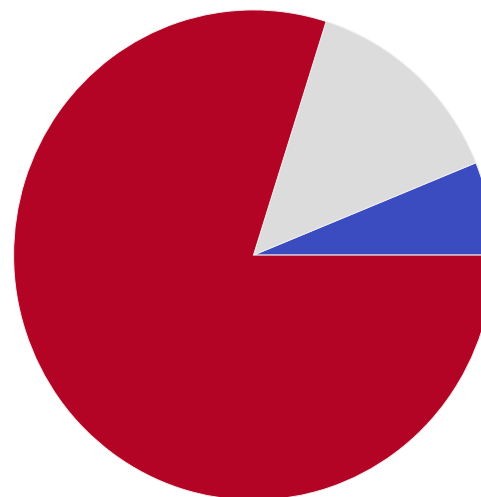
# Not Really OK

# OK versus MOVED



OK (2013)
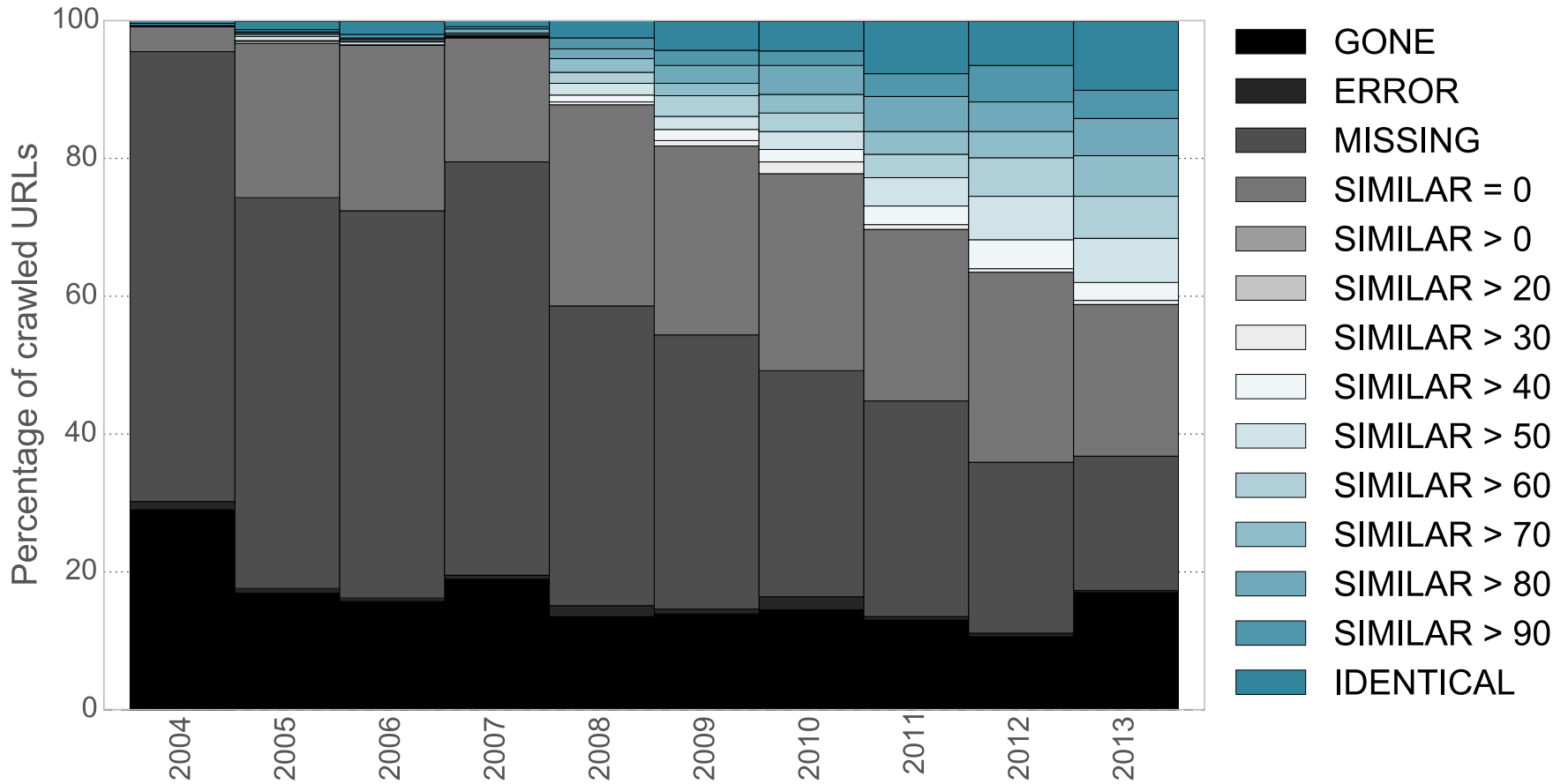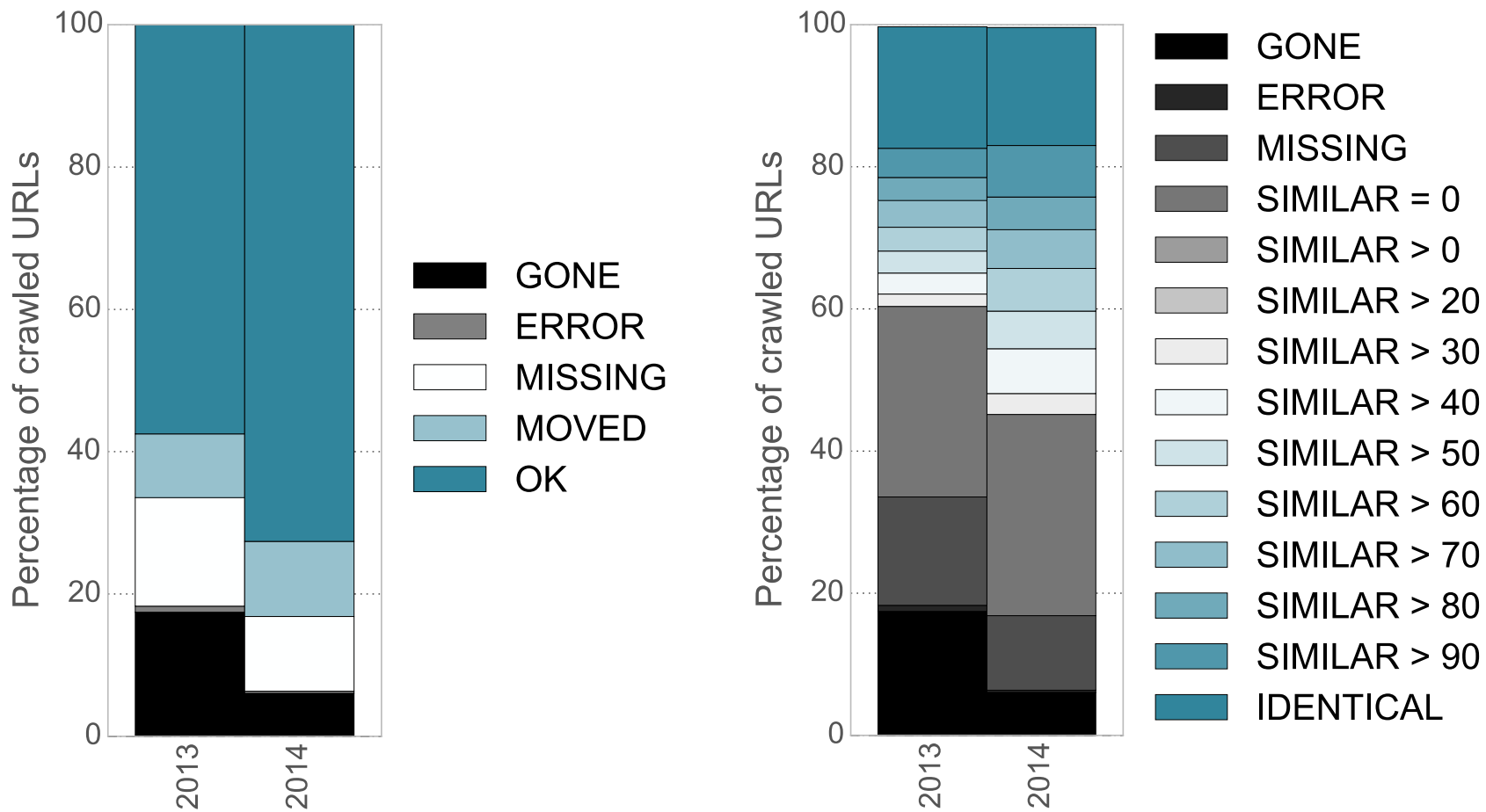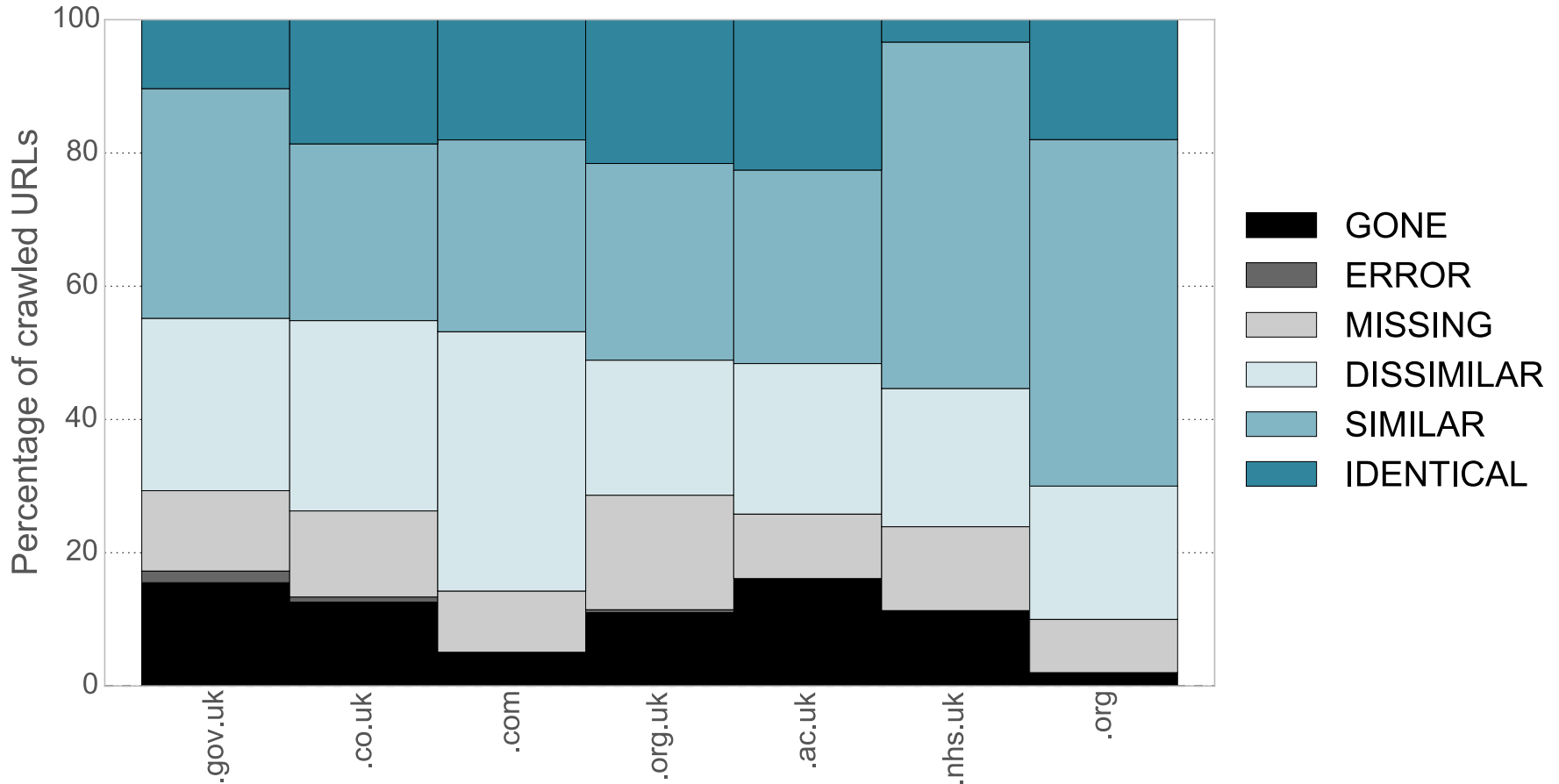
MOVED (2013)

IDENTICAL
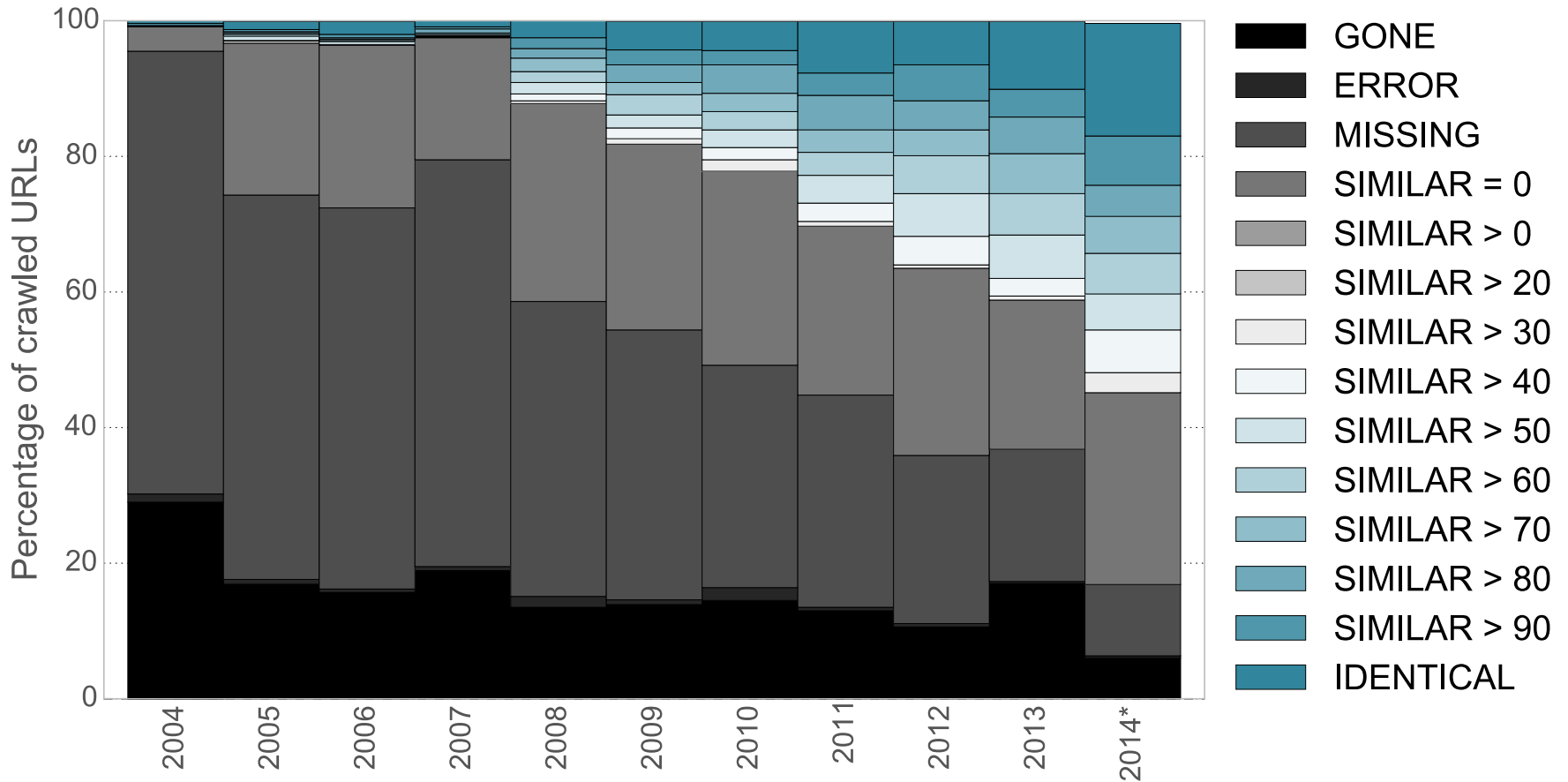SIMILAR
DISSIMILAR

# The URLs Ain't Cool

# Results For The Legal Deposit Collection

# Legal Deposit 2013-2014 By Domain Type

# What We've Saved (2004-2014)

# Summary

- Link rot & content drift dominate:
  - 50% of resources unrecognisable or gone after 1 year
  - 60% after 2 years, 65% after 3 years (islands of stability)
  - Noticeably higher rot rate than results for legal/academic web

- Simple similarity measure provides insight, although:
  - Only sensitive to text changes
  - Overly sensitive to header/footer changes

- Future work:
  - Look for old content at new URLs via hash similarity
  - Compare archival holdings via Memento

# Thank you!

## Questions?

Getting in touch:

Twitter: @ukwebarchive
Email: web-archivist@bl.uk
UK Web Archive:
http://www.webarchive.org.uk