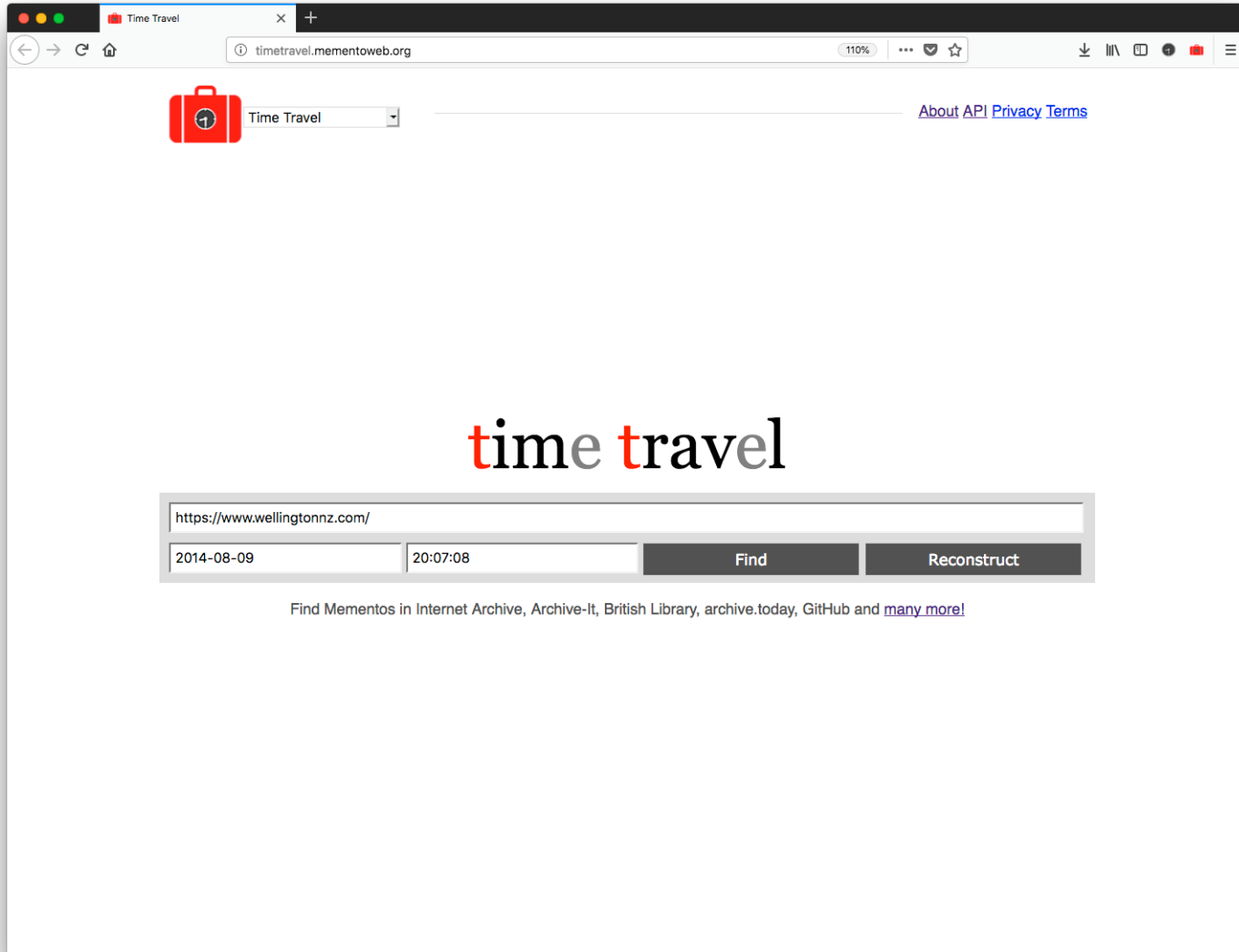# Smart Routing

## of Requests

**Martin Klein[1]**
Lyudmila Balakireva[1]
Harihar Shankar[1]
James Powell[1]
Herbert Van de Sompel[2]

[1]Research Library
Los Alamos National Laboratory

[2]Data Archiving and Networked Services
The Netherlands

# Memento



http://timetravel.mementoweb.org/

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Memento



http://timetravel.mementoweb.org/list/20140809200708/https://www.wellingtonnz.com/

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Memento



https://arquivo.pt/wayback/20141207132322/http://www.wellingtonnz.com/

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# How does this work?
## Memento Aggregator (very simplistic view)

URI_R
Datetime

URI-G

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Memento Aggregator

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Memento Aggregator

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# LANL Memento Aggregator - Problem

- As the number of archives grows, sending requests to each archive for every incoming request is not feasible

  - Response times

  - Memento infrastructure load

  - Load on distributed archives

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# What if…

- We could predict, by merely looking at a URI-R, whether or not to issue a request to a specific archive?
    - A binary classifier per archive

- We could train the classifiers using cached data?

- That would be pretty neat, indeed:
    - Retrain classifiers as web archive collections evolve
    - Not dependent on external data
    - Querying classifiers probably way faster (msec) than polling archives (sec)

# We can! Published @ JCDL 2016

- ML models based on simple URI features
  - Character count, n-grams, domain
- Common ML algorithms used per archive
  - Logistic Regression, Multinomial Bayes, SVM
- Optimized for
  - Prediction time, not training time
  - Reduction of false positive rate

**Results:**
- Requests per URI-R: **3.96 vs 11**
- Response time:
  **2.16s vs 3.08s**
- Recall:
  **84.7%**

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

**Los Alamos**
NATIONAL LABORATORY

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# In Production…

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

URI_R
Datetime

URI-G

30-day
Cache

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

Cache Hit

URI_R
Datetime

URI-G

Cache

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

Cache Hit

URI_R
Datetime

URI-G

Cache

Cache Miss

List of
predicted
archives
with
Mementos

Machine Learning
Process

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Populating the Cache

# Questions to Ask

- How effective is the cache?

  - What is the hit/miss ratio? Does it vary for different Memento services?

  - Is the cache freshness period appropriate?

- How effective is the ML process?

  - What is the false negative and false positive rate?

  - Do we need to retrain the models? How often?

# Evaluation

- Memento Aggregator currently covers
  - 23 web archives
  - 17 with native memento support
  - 6 with by-proxy memento support

- Analysis of log files
  - recorded between July 4[th] 2017 and October 17[th] 2018
  - > 11m requests in total
  - Approx. 2.6m requests against machine learning process
    - Results in 2.6m lookups to populate cache
      - Used as "truth" to assess ML prediction

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Cache Hit/Miss Rate



Memento Cache Performance

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Cache Hit/Miss Rate



Memento Cache Performance

Mostly driven by    machines    humans    machines    humans

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# False Negatives by Number of Archives



False Negatives By Number of Archives

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# False Negatives by Archive



False Negatives By Archives

| Archive | Count |
|---|---|
| York Univ | 322037 |
| UK Parliament | 1822 |
| UK Nat Arc | 1 |
| Stanford Arc | 8604 |
| PRONI | 1433 |
| perma.cc | 7560 |
| Nat Rec Scotland | 0 |
| Nat Lib Ireland | 5962 |
| LoC | 36336 |
| Icelandic Arc | 0 |
| Croatian Arc | 0 |
| Canadian Arc | 1036 |
| British Lib | 5639 |
| BibAlex.org | 42470 |
| Bayern Arc | 4061 |
| Arquivo.pt | 31400 |
| archive.is | 295512 |
| Archive-It | 60351 |

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# False Positives by Number of Archives



False Positives By Number of Archives

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# False Positives by Archive



False Positives By Archives

| Archive | Count |
|---|---|
| York Univ | 0 |
| UK Parliament | 0 |
| UK Nat Arc | 158990 |
| Stanford Arc | 6224 |
| PRONI | 6976 |
| perma.cc | 397672 |
| Nat Rec Scotland | 0 |
| Nat Lib Ireland | 87365 |
| LoC | 2915 |
| Icelandic Arc | 0 |
| Croatian Arc | 201699 |
| Canadian Arc | 10291 |
| British Lib | 365004 |
| BibAlex.org | 677999 |
| Bayern Arc | 7852 |
| Arquivo.pt | 246579 |
| archive.is | 376835 |
| Archive-It | 22692 |

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Changes in Archive Holdings

Change in Archive Holdings By Number of Archives

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Archives Added



Addition in Archive Holdings by Archive Name

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

**Los Alamos**
NATIONAL LABORATORY

# Archives Removed

Reduction in Archive Holdings by Archive Name



| Archive | Count |
|---|---|
| York Univ | 15878 |
| UK Parliament | 2930 |
| perma.cc | 2432 |
| Nat Lib Ireland | 5410 |
| LoC | 2341 |
| Internet Arc | 27374 |
| Bayern Arc | 1102 |
| BibAlex.org | 27065 |
| Archive-It | 2542 |
| archive.is | 7997 |

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

Los Alamos
NATIONAL LABORATORY

# Takeaways

- Memento Aggregator cache is very effective
  - Especially for human-driven services
- Machine learning process saves!
  - Requests & time while at acceptable recall level
  - FPR: 0.33 (std dev: 0.16)
- Re-training seems necessary, frequency TBD

Optimization

- ML model trained on archival holdings, not usage logs/cache
  - Beneficial for new archives
- Neural network classifier, based on simple URI features, show promising results

**Smart Routing of Memento Requests**
@mart1nkle1n
IIPC WAC 2018, 11/15/2018, Wellington, NZ

# Smart Routing
# of Memento Requests

**Martin Klein**[1]
Lyudmila Balakireva[1]
Harihar Shankar[1]
James Powell[1]
Herbert Van de Sompel[2]

[1]Research Library
Los Alamos National Laboratory

[2]Data Archiving and Networked Services
The Netherlands