# LinkGate

● ● ●

Web Archive Graph Visualization

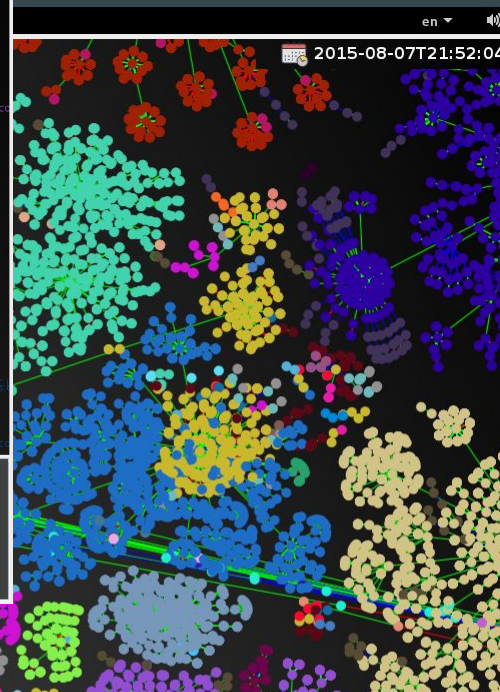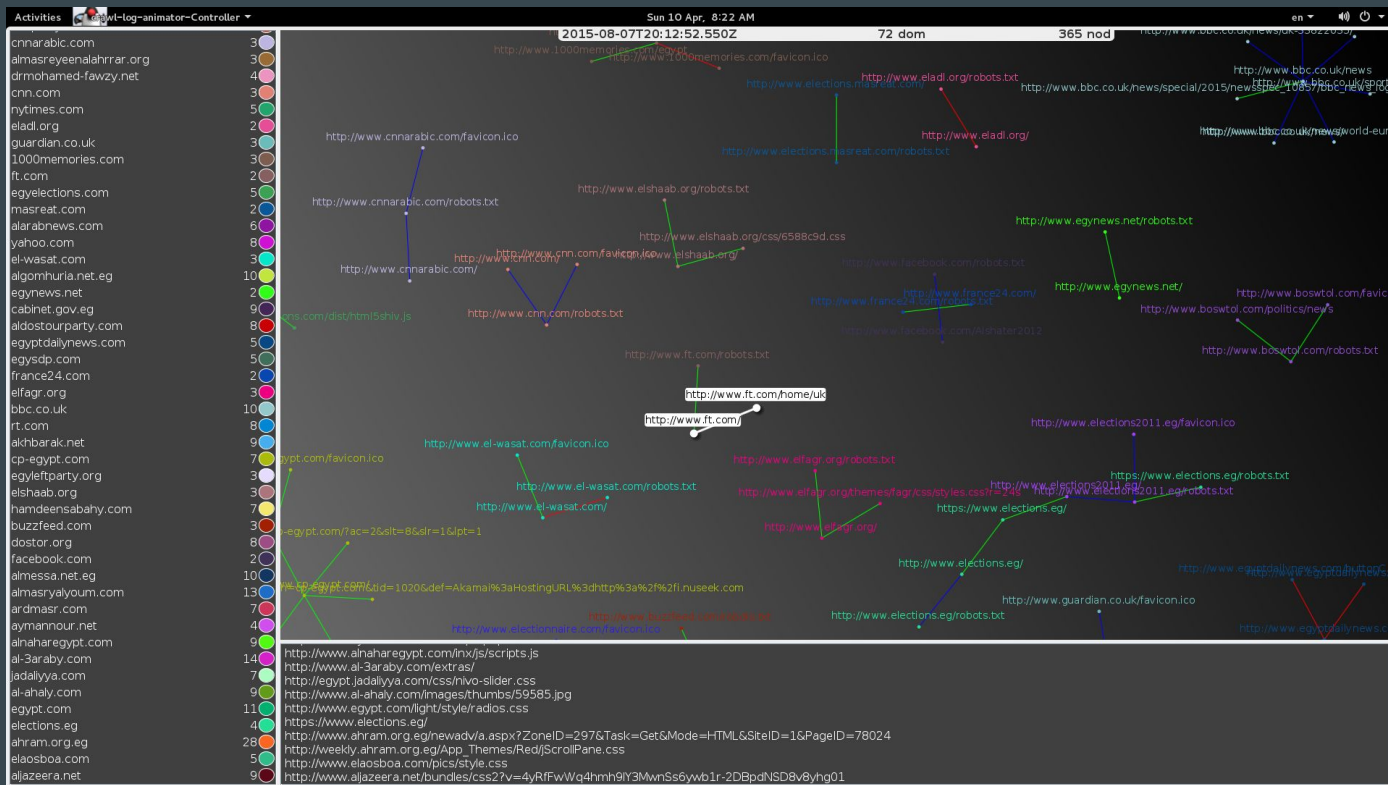# Rationale, partnership, and IIPC funding

- Visualization is an essential tool for understanding data and conducting research
  - Hyperlinks are the key concept behind the web, making the web a big graph
  - Tools exist for graph visualization, e.g., Gephi, this project addresses scalability and interoperability
- Bibliotheca Alexandrina and National Library of New Zealand worked together to develop core functionality and compile inventory of research use cases
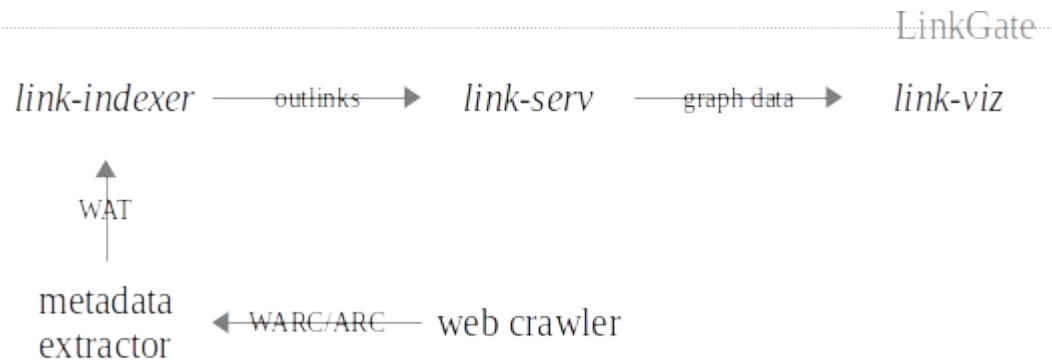- Funded by the IIPC during 2020

# Crawl Log Animator

# Overview

- 3 components:
  - *link-serv*
  - *link-indexer*
  - *link-viz*
- link-serv: scalable graph data service
- link-indexer: get link data from web archive storage and insert into link-serv
- link-viz: visualize and explore graph data from link-serv inside a web browser
- Inventory of research use cases to guide futur development

LinkGate

*link-indexer* ——outlinks——▶ *link-serv* ——graph data——▶ *link-viz*

▲
WAT
|

metadata extractor ◀——WARC/ARC—— web crawler

*link-viz*: Web Archive Graph Visualization Frontend

*link-indexer*: Linked Data Collection Tool

# What is *link-indexer*?

- Data flow begins with *link-indexer*
- From input data (e.g., WAT files), for each record, extract the following at a minimum:
  - Identifier/URI
  - Timestamp
  - Outlinks
- Produces graph data to insert into *link-serv*
- Written in Python
- Uses *webarchive-commons*, *urlcanon*, and *warcio* from the web archiving tool ecosystem

# *link-indexer* features

- Input format handling implemented as separate modules
- Generate WAT on the fly
- Post data to API endpoint or dry-run for troubleshooting
- Improve performance via batch processing
- Configurable network tolerance
- Script-friendly logging output
- Configurable input data handling behavior
- Configurable error handling behavior
- Load options from a configuration file

# What's next?

- Continue large-scale testing
- Enhance logging, including providing extra details, writing to a remote logging service
- Support more input formats
- More command-line options
- Extract additional metadata

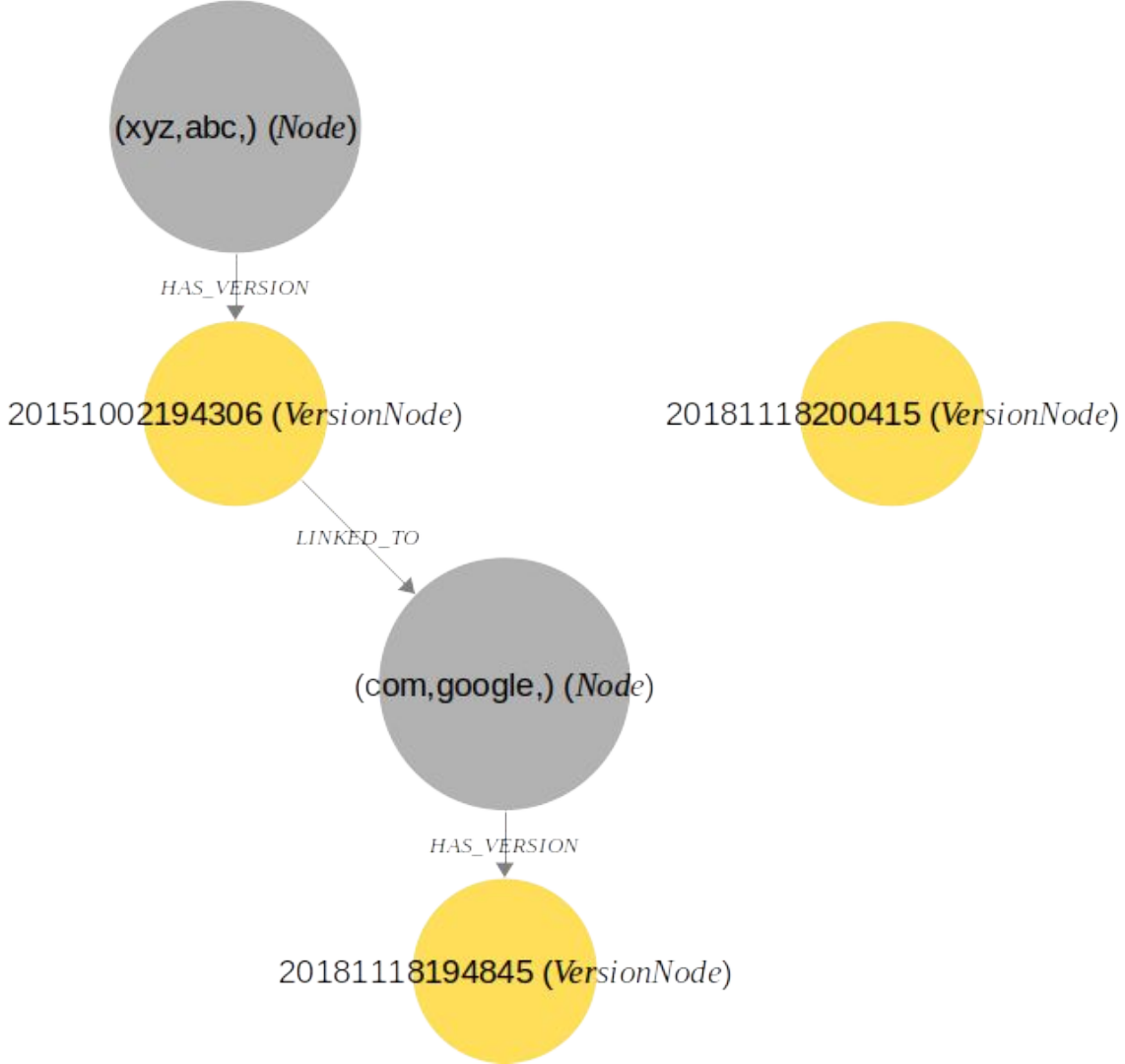*link-serv*: Temporal Graph Data Service

# What is *link-serv*?

- A service to provide a RESTful API for inserting temporal graph data extracted from a web archive into a central data store
- This service is used for retrieving back that data for rendering and navigation
- Design goals:
  - Data store scalability
  - Use publicly licensed technology
  - Data schema for temporal (i.e., versioned) graph data
  - RESTful API and Gephi compatibility

# Graph data stores

- A technical survey for data stores has been conducted
- Virtuoso, 4Store, OrientDB, Neo4j, ArangoDB, and others were considered
- Neo4j and ArangoDB are chosen as the best fit for the application
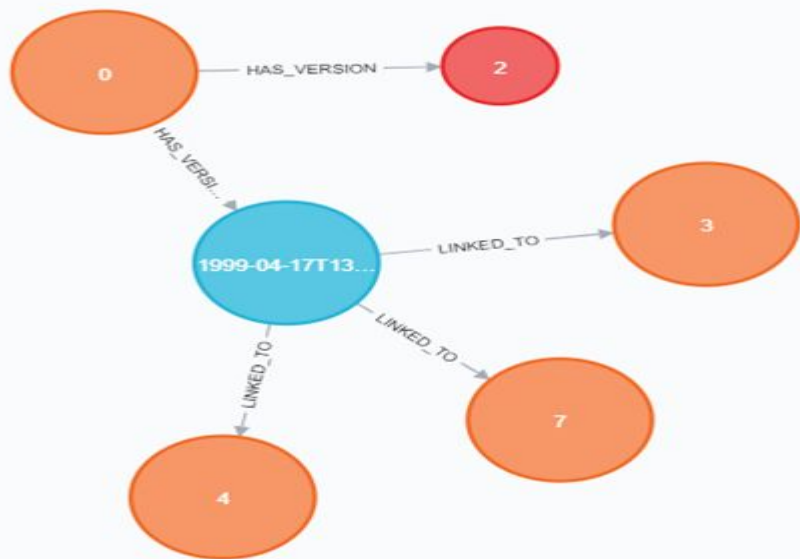- Both provide high performant read and write operations

# Neo4j data model

# Neo4j data model

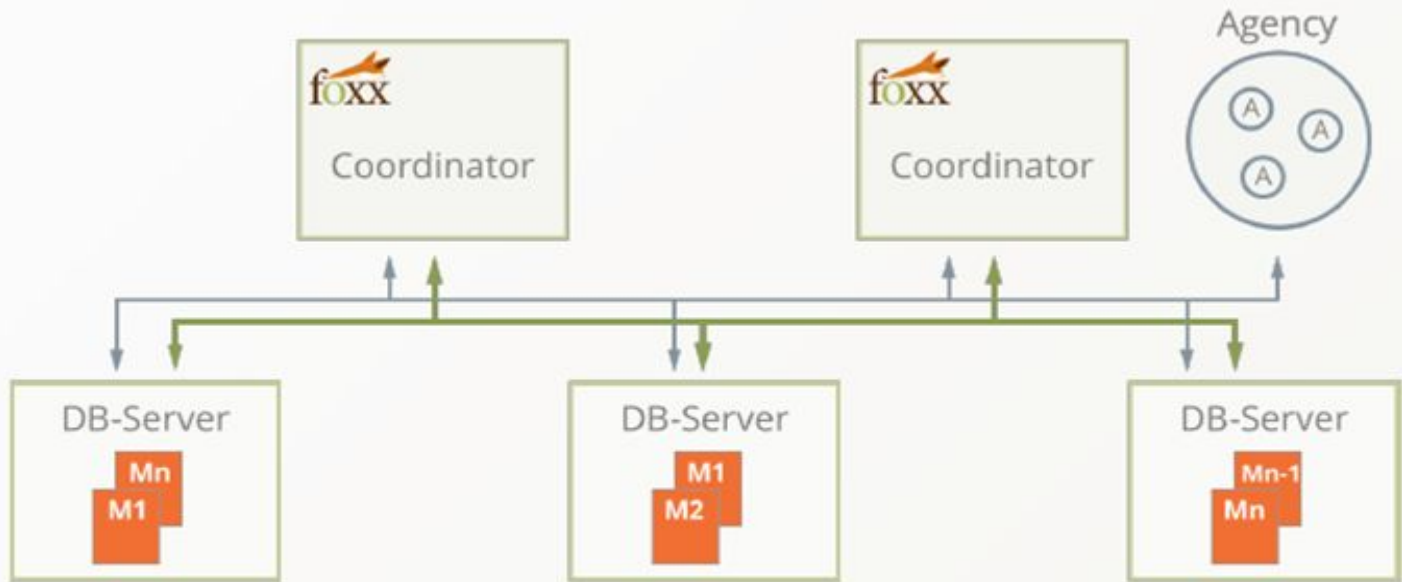- This mainly depends on the concept of "edgecuts"

Shard 1



Shard 2

# ArangoDB data model

## nodes

_key
_id
identifier
timestamp
label

---

**System Index:**
Type: primary index
Attribute: _key
Unique: true

**User-defined Index:**
Type: persistent index
Attributes: identifier,
timestamp
Unique: true

## linked_to

_key
_id
_from
_to

---

**System Index:**
Type: primary index
Attribute: _key
Unique: true

**System Index:**
Type: edge index
Attributes: _from, _to
Unique: false

**User-defined Index:**
Type: persistent index
Attributes: _from, _to
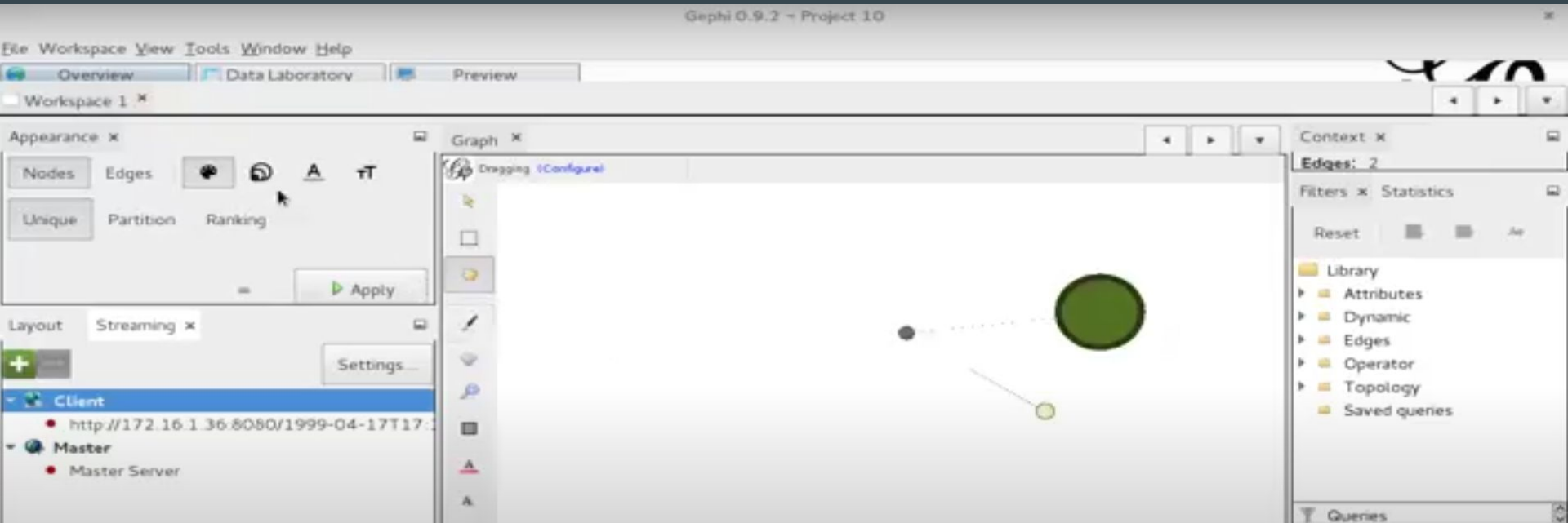Unique: true

# ArangoDB clustering

# API

- *link-serv* is implemented as a web service using Java and the Spring framework
- API for exposing functionality provided by the data model
- Currently implemented operations:
  - updateGraph
  - getGraph
  - getVersionCountsYearly
  - getVersionCountsMonthly
  - getVersionCountsDaily
  - getVersions
  - getLatestVersion

# Gephi compatibility

- *link-serv* is compatible with the API used by the Gephi streaming plugin
- Gephi can be used to render data from *link-serv*

# What's next?

- Design and implementation of replication mechanism according to the expected and real workload
- Large scale testing is a challenge for both data store solutions
- Enhance logging
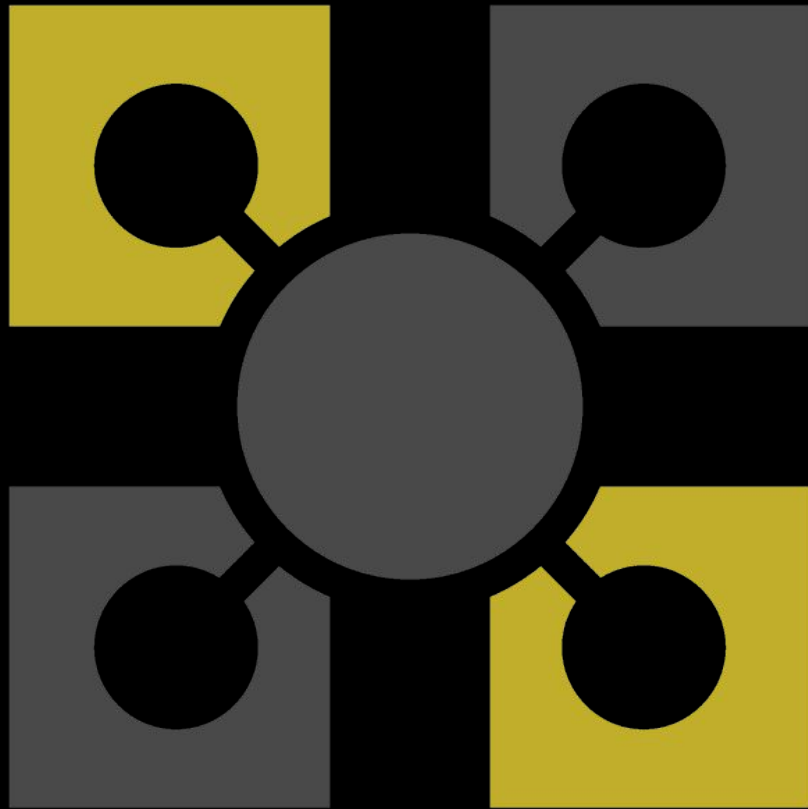- Support alternative data stores

Research Use Cases for Web Archive Graph Visualization

# Use cases to guide future development

- Tracking the promulgation of content through a web archive
- Providing tailored viewshafts into web archives
- Tagging and grouping web archive content with attributes
- Visualizing images and texts
- Preprocessing links before loading into visualization software
- Creating visualizations of crawl log data
- Creating curated web archive views for classrooms, or different audiences

Inventory of use cases: https://github.com/arcalex/linkgate/wiki/Use-cases

LinkGate:
https://linkgate.bibalex.org

Stay tuned for updates:
https://netpreserveblog.wordpress.com/tag/linkgate/

Get in touch:
linkgate@iipc.simplelists.com

On GitHub:
https://github.com/arcalex/linkgate
https://github.com/arcalex/link-serv
https://github.com/arcalex/link-indexer
https://github.com/arcalex/link-viz