

# LOTUS: ADAPTIVE TEXT SEARCH FOR BIG LINKED DATA

**F. Ilievski | W. Beek | M. van Erp | L. Rietveld | S. Schlobach**

# INTRODUCTION

A wealth of information is potentially available in Linked Open Data sources.

This information could be exploited by researchers and developers for tools and evaluations on a LOD-scale.

But, accessing Big Linked Data is not trivial

# INTRODUCTION

A wealth of information is potentially available in Linked Open Data sources.

This information could be exploited by researchers and developers for tools and evaluations on a LOD-scale.

But, accessing Big Linked Data is not trivial

**No centralized query  
service for Linked Data**

# INTRODUCTION

A wealth of information is potentially available in Linked Open Data sources.

This information could be exploited by researchers and developers for tools and evaluations on a LOD-scale.

But, accessing Big Linked Data is not trivial

**No centralized query  
service for Linked Data**

**Limited natural language  
access to Linked Data**

# INTRODUCTION

A wealth of information is potentially available in Linked Open Data sources.

This information could be exploited by researchers and developers for tools and evaluations on a LOD-scale.

But, accessing Big Linked Data is not trivial

**No centralized query service for Linked Data**

**Limited natural language access to Linked Data**

**Text-based retrieval is not customizable**

“


*“The lack of a global entry point to resources through a flexible text index is a serious obstacle for linked data consumption.”*

“

*“The lack of a global entry point to resources through a flexible text index is a serious obstacle for linked data consumption.”*

***by researchers and developers***

# REQUIREMENTS!



FOR A GLOBAL  
TEXT-BASED  
ENTRY POINT  
TO LOD



# REQUIREMENTS

1. Text-based queries
2. Resilience (of text search)
3. Findability (of authoritative and non-authoritative statements)
4. Availability
5. Scalability
6. Serviceability (for both machines and humans)
7. Customizability



**THE  
LOTUS  
APPROACH**



# LOD LAUNDROMAT

A centralized Linked Data  
cleaning and publishing architecture  
Allows access to a big subset of the LOD Cloud

**38 billion statements**

**LOTUS!**



[lotus.lodlaundromat.org](https://lotus.lodlaundromat.org)



## LINGUISTIC ENTRY POINT

LOTUS is a linguistic entry point to the LOD Laundromat data collection.

## APPROXIMATE MATCHING

Allows statements to be findable based on approximate string matching on associated literals.

## ADAPTIVE FRAMEWORK

LOTUS allows the resource retrieval to be tailored to fit various use cases.



# CUSTOMIZABILITY OF RETRIEVAL

4 MATCHINGS  
X 8 RANKINGS  
= 32 RETRIEVAL  
OPTIONS



## MATCHING OPTIONS

- ▶ Phrase matching
- ▶ Disjunctive token matching
- ▶ Conjunctive token matching
- ▶ Conjunctive token matching with character edit distance

# RETRIEVAL OPTIONS



## MATCHING OPTIONS

- ▶ Phrase matching
- ▶ Disjunctive token matching
- ▶ Conjunctive token matching
- ▶ Conjunctive token matching with character edit distance

## RANKING ALGORITHMS

- ▶ Length normalization
- ▶ Practical scoring function **CONTENT -BASED**
- ▶ Phrase proximity

---

- ▶ Terminological richness **DOCUMENT -BASED**
- ▶ Semantic richness
- ▶ Recency

---

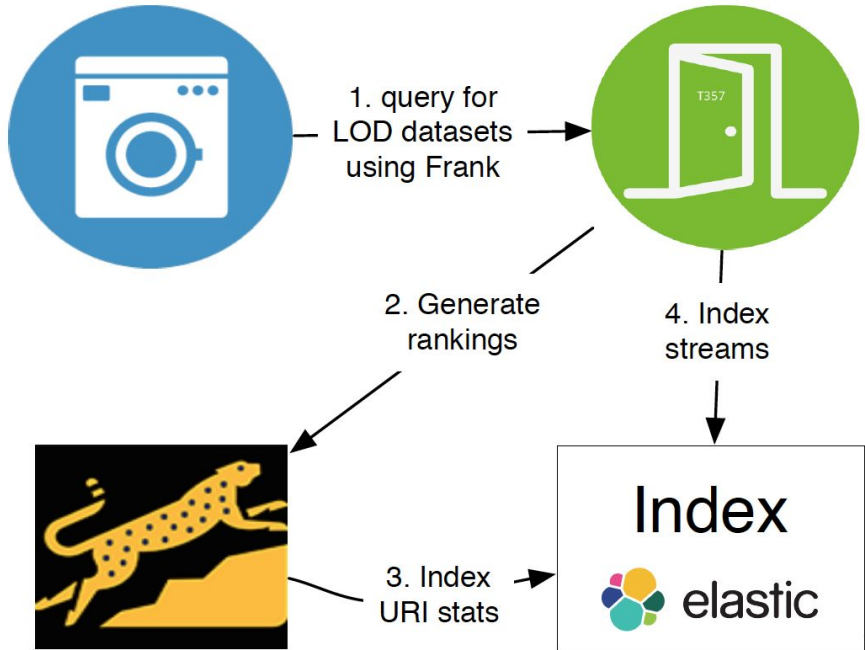
- ▶ Degree popularity **RESOURCE -BASED**
- ▶ Appearance popularity



# IMPLEMENTATION



## Index Builder



# IMPLEMENTATION



## Index Builder



1. query for  
LOD datasets  
using Frank



2. Generate  
rankings

4. Index  
streams

3. Index  
URI stats



6. Send user  
query to index

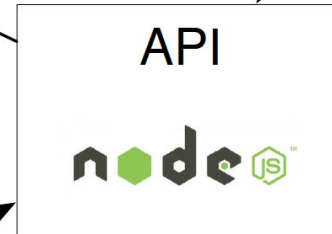
7. Return results  
to front end

## Public Interface



5. Query

8. Retrieve  
results



WEB INTERFACE + API



[lotus.lodlaundromat.org](https://lotus.lodlaundromat.org)

# DISTRIBUTED ARCHITECTURE



**4,334,672,073 INDEXED LITERALS**

**SCALED HORIZONTALLY OVER 5 SERVERS**

**DATA REPLICATION TO ENSURE HIGH  
RUNTIME AVAILABILITY OF LOTUS**



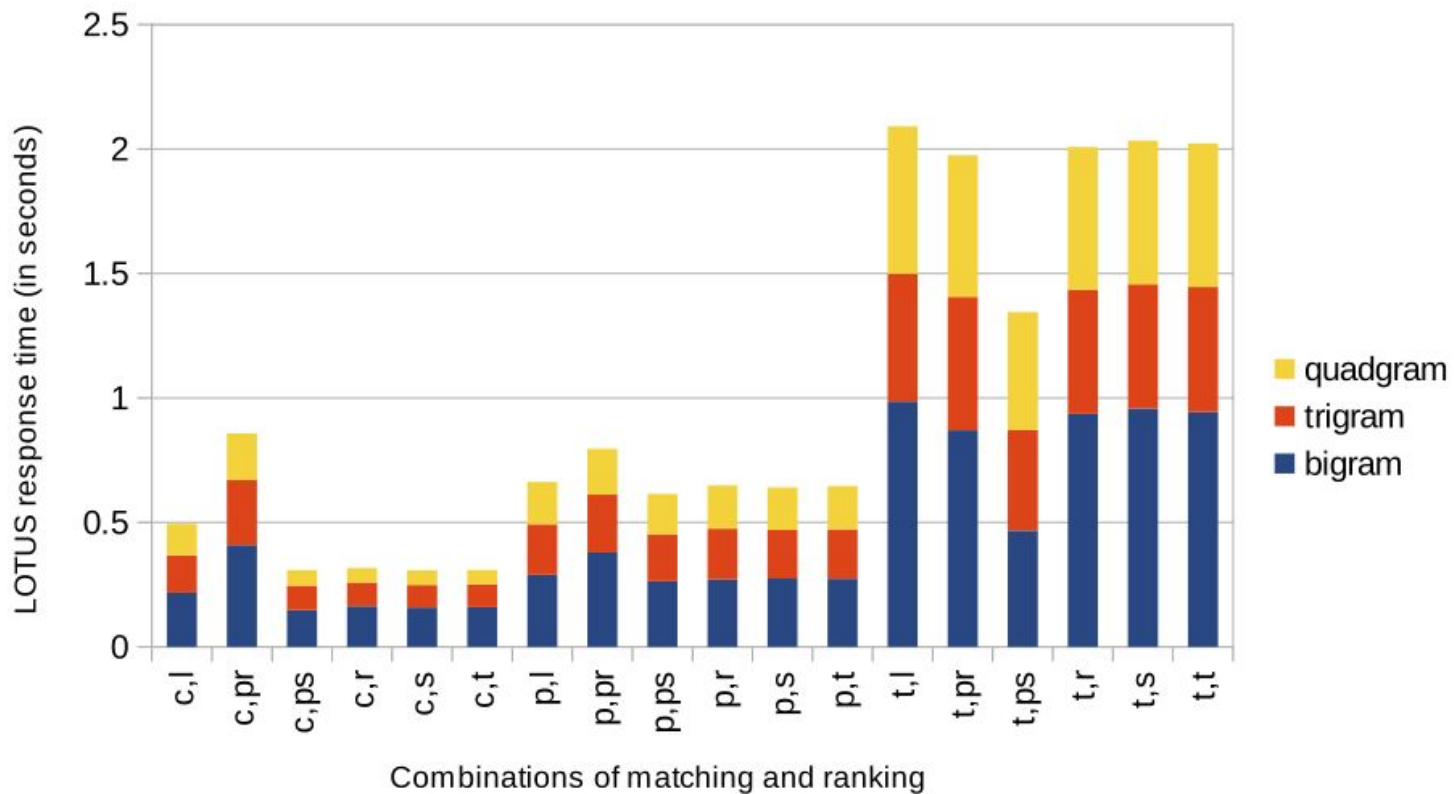
**PERFORMANCE!**

**AND  
USAGE SCENARIOS**

# SCALING AND PERFORMANCE



We used 18k queries to benchmark 18 retrieval combinations of LOTUS



# CONCLUSIONS

## LOTUS IS

### **A CENTRALIZED LINGUISTIC ENTRY TO BIG LINKED DATA**

LOTUS indexes over 4 billion literals from the LOD Laundromat

### **AN ADAPTIVE RETRIEVAL FRAMEWORK**

LOTUS allows its retrieval to be customized to fit users' needs by offering 32 matching+ranking options.

### **"CONNECTING THE DOTS"**

LOTUS relies heavily on 2 existing systems (LOD Laundromat & ES), but fills the gap by offering a much needed tool for scientific evaluation.

VISION: SCALING APPLICATIONS AND EVALUATIONS AT LOD  
SCALE WITH LOD LAB

LOD Lab



Experiments at LOD Scale



## FUTURE WORK

### **THE PRECISION AND RECALL**

of LOTUS should be evaluated on concrete applications, such as Entity Linking and Network Analysis.

## FUTURE WORK

### **THE PRECISION AND RECALL**

of LOTUS should be evaluated on concrete applications, such as Entity Linking and Network Analysis.

### **CONTEXT-DEPENDENT RANKING**

could be added in the future to take the query context into account in order to improve the ranking accuracy.

**THANKS!**

[lotus.lodlaundromat.org](http://lotus.lodlaundromat.org)



**any  
questions  
?**

You can find me at  
[@earthling91](#) / [f.ilievski@vu.nl](mailto:f.ilievski@vu.nl)