Link Discovery Tutorial Introduction

Mohamed Ahmed Sherif⁽¹⁾, Irini Fundulaki⁽²⁾, Axel-Cyrille Ngonga Ngomo⁽¹⁾

(1) Institute for Applied Informatics, Germany (2) FORTH, Greece



May 28th, 2017 Portoroz, Slovenia.

Iohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 23, 2017 1 / 24

Introduction

HOBBIT AKSW

Disclaimer

- No pursuit of completeness [Nen+15]
- Focus on
 - Basic ideas and principles
 - Principles
 - Evaluation
 - Open questions and challenges

(InfAI & FORTH



LD TUTORIAL:INTRO

May 23, 2017 2 / 24

Introduction

AKSW

Disclaimer

• No pursuit of completeness [Nen+15]

- Focus on
 - Basic ideas and principles
 - Principles
 - Evaluation
 - Open questions and challenges

Guidelines

- Question? Just ask!
- Comment? Go for it!
- Be kind :)

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 23, 2017 2 / 24

Let's Go!

🚓 HOBBIT AKSW



Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 23, 2017 3 / 24

Linked Data Principles

BOBBIT AKSW



May 23, 2017 4 / 24

Why Link Discovery?

BIT AKSW



Cross-Ontology QA

HOBBIT AKSW

Example

Give me the name and description of all drugs that cure their side-effect. [SNA13]

TUTORIAL:INTRO

- Need information from
 - Drugbank (Drug description)
 - Sider (Side-effects)
 - DBpedia (Description)
- Gathering information via SPARQL query using links

(INFAI & FORTH)



May 23, 2017 6 / 24

Cross-Ontology QA

HOBBIT AKSW

Example

Give me the name and description of all drugs that cure their side-effect.

```
SELECT ?drug ?name ?desc WHERE
{
    ?drug a drugbank:Drug .
    ?drug rdfs:label ?name .
    ?drug drugbank:cures ?disease .
    ?drug owl:sameAs ?drug2 .
    ?drug owl:sameAs ?drug3 .
    ?drug2 sider:hasSideEffect ?effect .
    ?effect owl:sameAs ?disease .
    ?drug3 dbo:hasWikiPage ?desc .
}
```



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D TUTORIAL:INTRO

May 23, 2017 <u>7 / 24</u>

Cross-Ontology QA

HOBBIT AKSW



Give me flats near kindergartens in Kobe. [Leh+12]

FORTH

```
SELECT ?flat WHERE
{
    ?flat a deqa:Flat .
    ?flat deqa:near ?school .
    ?school a lgdo:School .
    ?school lgdo:city lgdo:Kobe .
}
```



LD TUTORIAL:INTRO

May 23, 2017 8 / 24

Data Integration

BOBBIT AKSW

May 23, 2017

• Federated Queries on Patient Data [Kha+14]



Federated Queries

S HOBBIT AKSW

Example (FedBench CD2)

Return Barack Obama's party membership and news pages. [Sal+15]

SELECT ?party ?page WHERE
{
 dbr:Barack_Obama dbo:party ?party .
 ?x nytimes:topicPage ?page .
 ?x owl:sameAs dbr:Barack_Obama .
}



HOBBIT AKSW

Definition (Link Discovery, informal)

- \bullet Given two knowledge bases S and T, find links of type ${\cal R}$ between S and T
- Here, declarative link discovery

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 23, 2017 11 / 24

HOBBIT AKSW

Definition (Link Discovery, informal)

- \bullet Given two knowledge bases S and T, find links of type ${\cal R}$ between S and T
- Here, declarative link discovery

Definition (Declarative Link Discovery, formal, similarities)

- \bullet Given sets S and T of resources and relation ${\cal R}$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Common approach: Find $M' = \{(s,t) \in S \times T : \sigma(s,t) \ge \theta\}$

'utorial:Intro

May 23, 2017 11 / 24

AKSW

Definition (Link Discovery, informal)

- \bullet Given two knowledge bases S and T, find links of type ${\cal R}$ between S and T
- Here, declarative link discovery

Definition (Declarative Link Discovery, formal, similarities)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Common approach: Find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

Definition (Declarative Link Discovery, formal, distances)

- \bullet Given sets S and T of resources and relation ${\cal R}$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Common approach: Find $M' = \{(s,t) \in S \times T : \delta(s,t) \le \tau\}$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 23, 2017 11 / 24

HOBBIT AKSW

- Most common: $\mathcal{R} = \texttt{owl:sameAs}$
- Also known as deduplication [Nen+15]



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 23, 2017 12 / 24

HOBBIT AKSW

- \bullet Goal: Address all possible relations ${\cal R}$
- Declarative Link Discovery: Similarity/distance defined using property values (incl. property chains)

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 23, 2017 13 / 24

HOBBIT AKSW

- \bullet Goal: Address all possible relations ${\cal R}$
- Declarative Link Discovery: Similarity/distance defined using property values (incl. property chains)
- Example: $\mathcal{R} = : \texttt{sameModel}$



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

Tutorial:Intro

May 23, 2017 <u>13 / 24</u>

Why is it difficult?

🚓 HOBBIT 🗛 📕

• Time complexity

- Large number of triples (e.g., LinkedTCGA with 20.4 billion triples [Sal+14])
- Quadratic a-priori runtime

(InfAI & FORTH

• 69 days for mapping cities from DBpedia to Geonames

LD TUTORIAL:INTRO

• Solutions usually in-memory (insufficient heap space)



MAY

Why is it difficult?

BBIT AKSW



- Large number of triples (e.g., LinkedTCGA with 20.4 billion triples [Sal+14])
- Quadratic a-priori runtime
- 69 days for mapping cities from DBpedia to Geonames
- Solutions usually in-memory (insufficient heap space)



Accuracy

- Combination of several attributes required for high precision
- Tedious discovery of most adequate mapping
- Dataset-dependent similarity functions

(INFAI & FORTH

Structure

SHOBBIT AKSW

Time complexity

- LIMES algorithm [NA11]
- MultiBlock [IJB11]
- HR³ [Ngo12]
- AEGLE [GSN16]
- Summary and Challenges

Accuracy

- RAVEN [Ngo+11]
- EAGLE [NL12]
- COALA [NLC13]
- Summary and Challenges

Benchmarking

- Benchmarking [NGF16]
- Synthetic Benchmarks [Sav+15]
- Real Benchmarks [Mor+11]
- Summary and Challenges

Hands-On Session

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

That's all Folks!

HOBBIT AKSW







Axel Ngonga AKSW Research Group Institute for Applied Informatics ngonga@informatik.uni-leipzig.de

Irini Fundulaki ICS FORTH fundul@ics.forth.gr

Mohamed Ahmed Sherif AKSW Research Group Institute for Applied Informatics sherif@informatik.uni-leipzig.de

LD TUTORIAL:INTRO

Iay 23, 2017 16 / 24

A cknowledgment

(INFAI & FORTH

HOBBIT AKSW



This work was supported by grants from the EU H2020 Framework Programme provided for the project HOBBIT (GA no. 688227).

LD TUTORIAL:INTRO

References I

HOBBIT AKSW

Kleanthi Georgala, Mohamed Ahmed Sherif, and Axel-Cyrille Ngonga Ngomo. "An Efficient Approach for the Generation of Allen Relations". In: ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands - Including Prestigious Applications of Artificial Intelligence (PAIS 2016). Ed. by Gal A. Kaminka et al. Vol. 285. Frontiers in Artificial Intelligence and Applications. IOS Press, 2016, pp. 948–956. ISBN: 978-1-61499-671-2. DOI: 10.3233/978-1-61499-672-9-948. URL: http://dx.doi.org/10.3233/978-1-61499-672-9-948.

Robert Isele, Anja Jentzsch, and Christian Bizer. "Efficient Multidimensional Blocking for Link Discovery without losing Recall". In: Proceedings of the 14th International Workshop on the Web and Databases 2011, WebDB 2011, Athens, Greece, June 12, 2011. Ed. by Amélie Marian and Vasilis Vassalos. 2011. URL: http://webdb2011.rutgers.edu/papers/Paper%2039/silk.pdf.

D TUTORIAL:INTRO

May 23, 2017 18 / 24

References II

BOBBIT AKSW

Yasar Khan et al. "SAFE: Policy Aware SPARQL Query Federation Over RDF Data Cubes". In: Proceedings of the 7th International Workshop on Semantic Web Applications and Tools for Life Sciences, Berlin, Germany, December 9-11, 2014. Ed. by Adrian Paschke et al. Vol. 1320. CEUR Workshop Proceedings. CEUR-WS.org, 2014. URL: http://ceur-ws.org/Vol-1320/Preface_SWAT4LS2014.pdf.

Jens Lehmann et al. "deqa: Deep Web Extraction for Question Answering". In: The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference, Boston, MA, USA, November 11-15, 2012, Proceedings, Part II. Ed. by Philippe Cudré-Mauroux et al. Vol. 7650. Lecture Notes in Computer Science. Springer, 2012, pp. 131–147. ISBN: 978-3-642-35172-3. DOI: 10.1007/978-3-642-35173-0_9. URL: http://dx.doi.org/10.1007/978-3-642-35173-0_9.

References III

Ahmed Sherif et al

Mohamed

(INFAI & FORTH)

🔅 HOBBIT AKSW

Mohamed Morsey et al. "DBpedia SPARQL Benchmark - Performance Assessment with Real Queries on Real Data". In: The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I. Ed. by Lora Aroyo et al. Vol. 7031. Lecture Notes in Computer Science. Springer, 2011, pp. 454–469. ISBN: 978-3-642-25072-9. DOI: 10.1007/978-3-642-25073-6_29. URL: http://dx.doi.org/10.1007/978-3-642-25073-6_29.
Axel-Cyrille Ngonga Ngomo and Sören Auer. "LIMES - A Time-Efficient Approach for Large-Scale Link Discovery on the Web of Data". In: IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011. Ed. by Toby Walsh. IJCAI/AAAI, 2011, pp. 2312–2317. ISBN: 978-1-57735-516-8. DOI: 10.5591/978-1-57735-516-8/IJCAI11-385. URL: http://dx.doi.org/10.5591/978-1-57735-516-8/IJCAI11-385.
Markus Nentwig et al. "A survey of current Link Discovery frameworks". In: Semantic Web Preprint (2015), pp. 1–18.

LD TUTORIAL:INTRO

May 23, 2017 20 / 24

References IV

🚵 HOBBIT AKSW

Axel-Cyrille Ngonga Ngomo, Alejandra Garcıa-Rojas, and Irini Fundulaki. "HOBBIT: Holistic Benchmarking of Big Linked Data". In: ERCIM News 2016.105 (2016). URL: http://ercim-news.ercim.eu/en105/r-i/hobbit-holistic-benchmarking-of- big-linked-data.
Axel-Cyrille Ngonga Ngomo et al. "RAVEN - active learning of link specifications". In: Proceedings of the 6th International Workshop on Ontology Matching, Bonn, Germany, October 24, 2011. Ed. by Pavel Shvaiko et al. Vol. 814. CEUR Workshop Proceedings. CEUR-WS.org, 2011. URL: http://ceur-ws.org/Vol-814/om2011_Tpaper3.pdf.
Axel-Cyrille Ngonga Ngomo. "Link Discovery with Guaranteed Reduction Ratio in Affine Spaces with Minkowski Measures". In: The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference, Boston, MA, USA, November 11-15, 2012, Proceedings, Part I. Ed. by Philippe Cudré-Mauroux et al. Vol. 7649. Lecture Notes in Computer Science. Springer, 2012, pp. 378–393. ISBN: 978-3-642-35175-4. DOI: 10.1007/978-3-642-35176-1_24. URL: http://dx.doi.org/10.1007/978-3-642-35176-1_24.

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 23, 2017 21 / 24

References V

HOBBIT AKSW

Axel-Cyrille Ngonga Ngomo and Klaus Lyko. "EAGLE: Efficient Active Learning of Link Specifications Using Genetic Programming". In: The Semantic Web: Research and Applications - 9th Extended Semantic Web Conference, ESWC 2012, Heraklion, Crete, Greece, May 27-31, 2012. Proceedings. Ed. by Elena Simperl et al. Vol. 7295. Lecture Notes in Computer Science. Springer, 2012, pp. 149–163. ISBN: 978-3-642-30283-1. DOI: 10.1007/978-3-642-30284-8_17. URL: http://dx.doi.org/10.1007/978-3-642-30284-8_17.

Axel-Cyrille Ngonga Ngomo, Klaus Lyko, and Victor Christen. "COALA - Correlation-Aware Active Learning of Link Specifications". In: The Semantic Web: Semantics and Big Data, 10th International Conference, ESWC 2013, Montpellier, France, May 26-30, 2013.
Proceedings. Ed. by Philipp Cimiano et al. Vol. 7882. Lecture Notes in Computer Science. Springer, 2013, pp. 442–456. ISBN: 978-3-642-38287-1. DOI: 10.1007/978-3-642-38288-8_30. URL: http://dx.doi.org/10.1007/978-3-642-38288-8_30.

LD TUTORIAL:INTRO

fay 23, 2017 22 / 24

References VI

Mohamed Ahmed

(INFAI & FORTH

🚵 HOBBIT AKSW

Muhammad Saleem et al. "TopFed: TCGA Tailored Federated Query Processing and Linking to LOD". In: J. Biomedical Semantics 5 (2014), p. 47. DOI: 10.1186/2041-1480-5-47. URL: http://dx.doi.org/10.1186/2041-1480-5-47.
Muhammad Saleem et al. "A fine-grained evaluation of SPARQL endpoint federation systems". In: Semantic Web 7.5 (2015), pp. 493–518. DOI: 10.3233/SW-150186. URL: http://dx.doi.org/10.3233/SW-150186.
Tzanina Saveta et al. "LANCE: Piercing to the Heart of Instance Matching Tools". In: The Semantic Web - ISWC 2015 - 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part I. Ed. by Marcelo Arenas et al. Vol. 9366. Lecture Notes in Computer Science. Springer, 2015, pp. 375–391. ISBN: 978-3-319-25006-9. DOI: 10.1007/978-3-319-25007-6_22. URL: http://dx.doi.org/10.1007/978-3-319-25007-6_22.

LD TUTORIAL:INTRO

May 23, 2017 23 / 24

References VII

HOBBIT AKSW

Saeedeh Shekarpour, Axel-Cyrille Ngonga Ngomo, and Sören Auer. "Question answering on interlinked data". In: 22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013. Ed. by Daniel Schwabe et al. International World Wide Web Conferences Steering Committee / ACM, 2013, pp. 1145–1156. ISBN: 978-1-4503-2035-1. URL: http://dl.acm.org/citation.cfm?id=2488488.

D Tutorial:Intro

May 23, 2017 24 / 24

Link Discovery Tutorial Part I: Efficiency

Mohamed Ahmed Sherif⁽¹⁾, Irini Fundulaki⁽²⁾, Axel-Cyrille Ngonga Ngomo⁽¹⁾

(1) Institute for Applied Informatics, Germany (2) FORTH, Greece



May 28th, 2017 Portoroz, Slovenia.

LD TUTORIAL: EFFICIENC

FORTH

May 28, 2017 1 / 117

Table of Contents

HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link Discovery
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

Mohamed

LD TUTORIAL: EFFICIENC

May 28, 2017 2 / 117

Table of Contents

🚵 HOBBIT AKSW

1 Introduction
2 LIMES
3 MultiBlock
4 Reduction-Ratio-Optimal Link Disco
5 AEGLE
6 RADON
7 GNOME
8 GPUs and Hadoop
9 Summary and Conclusion

(INFAI & FORTH

Mohamed

LD TUTORIAL: EFFICIENCY

May 28, 2017 3 / 117

Introduction

Link Discovery



Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

OHAMED AHMED SHERIF ET AL. (INFAI & FORTH)

D TUTORIAL: EFFICIENCY

May 28, 2017 4 / 117

Introduction

Link Discovery

AKSW

Definition (Declarative Link Discovery)

- \bullet Given sets S and T of resources and relation ${\cal R}$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

Problem

- Naïve complexity $\in O(S \times T)$, i.e., $O(n^2)$
- $\bullet\,$ Example: \approx 70 days for cities in DBpedia and LinkedGeoData

Introduction

Link Discovery

BOBBIT AKSW

Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

Problem

- Naïve complexity $\in O(S \times T)$, i.e., $O(n^2)$
- $\bullet\,$ Example: \approx 70 days for cities in DBpedia and LinkedGeoData

Solutions

- Reduced number of comparisons $C(A) \ge |M'|$ [NA11; IJB11; Ngo12; Ngo13]
- Use planning [Ngo14]
- Reduce the complexity of linking to $< n^2$ [GSN16]
- Use features of hardware environment [Ngo+13; NH16]

Iohamed Ahmed Sherif et al. (InfAI & FORTH)

Table of Contents

HOBBIT AKSW

Intr	odu	ctior

2 LIMES

- 3 MultiBlock
- 4 Reduction-Ratio-Optimal Link Discovery

(INFAI & FORTH

- **5** AEGLE
- 6 RADON
- **GNOME**
- 8 GPUs and Hadoop

Mohamed Ahmed Sherif et

9 Summary and Conclusion

LD TUTORIAL: EFFICIEN

May 28, 2017 5 / 117


Intuition

🚓 HOBBIT 🗛 KSW

Definition (Declarative Link Discovery)

- \bullet Given sets S and T of resources and relation ${\cal R}$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s,t) \in S \times T : \delta(s,t) \le \tau\}$
- Some δ are distances (in the mathematical sense) [NA11]
- Examples
 - Levenshtein distance
 - Minkowski distance

Tutorial: Efficiency

May 28, 2017 6 / 117

LIMES

Intuition

AKSW

Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s,t) \in S \times T : \delta(s,t) \le \tau\}$
- Some δ are distances (in the mathematical sense) [NA11]
- Examples
 - Levenshtein distance
 - Minkowski distance

Intuition

- Distances abide by triangle inequality
- $\delta(x,z) \delta(z,y) \le \delta(x,y) \le \delta(x,z) + \delta(z,y)$
- Use exemplars for pessimistic approximation
- Reduce number of computations

Ahmed Sherif et al. (InfAI & FORTH)

torial: Efficiency

May 28, 2017 6 / 117





 $\bullet \quad \text{Start with random } e_1 \in \mathcal{T}$

$$e_{n+1} = \operatorname*{arg\,max}_{x \in T} \sum_{i=1}^{n} \delta(x, e_i)$$

- Assign each t to e(t) with $e(t) = \underset{e_i}{\arg\min \delta(t, e_i)}$
- Approximate $\delta(s, t)$ by using $\delta(s, t) \geq \delta(s, e(t)) \delta(e(t), t)$



LIMESApproach

HOBBIT AKSW



• Start with random $e_1 \in T$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 8 / 117

LIMESApproach

🚓 HOBBIT (AKSW)



• Start with random $e_1 \in T$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 9 / 117



HOBBIT AKSW



• Start with random $e_1 \in T$ • $e_{n+1} = \underset{x \in T}{\operatorname{arg max}} \sum_{i=1}^{n} \delta(x, e_i)$

FORTH

D TUTORIAL: EFFICIENCY

May 28, 2017 <u>10 / 117</u>

LIMESApproach

HOBBIT AKSW



 $\bullet \quad \text{Start with random } e_1 \in \mathcal{T}$

$$e_{n+1} = \operatorname*{arg\,max}_{x \in T} \sum_{i=1}^{n} \delta(x, e_i)$$

• Assign each t to e(t) with $e(t) = \underset{e_i}{\arg\min \delta(t, e_i)}$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 <u>11 / 1</u>17

LIMESApproach

HOBBIT AKSW

• Start with random $e_1 \in T$

$$e_{n+1} = \operatorname*{arg\,max}_{x \in T} \sum_{i=1}^{n} \delta(x, e_i)$$

- Solution Assign each t to e(t) with $e(t) = \underset{e_i}{\arg\min \delta(t, e_i)}$
- Approximate $\delta(s, t)$ by using $\delta(s, t) \geq \delta(s, e(t)) \delta(e(t), t)$

(INFAI & FORTH

utorial: Efficie





HOBBIT AKSW

- Measure = Levenshtein
- Threshold = 0.9

Mohamed Ahmed Sherii

- KB base size = 10^3
- x-axis is number of exemplars
- y-axis is number of comparisons

(INFAI & FORTH)



LD TUTORIAL: EFFICIENCY

May 28, 2017 13 / 117

LIMES Evaluation

HOBBIT AKSW

- Measure = Levenshtein
- Threshold = 0.9
- KB base size = 10^4
- x-axis is number of exemplars
- y-axis is number of comparisons
- Optimal number of examplars $\approx \sqrt{|T|}$

(INFAI & FORTH)



LD TUTORIAL: EFFICIENCY

LIMES Conclusion

HOBBIT AKSW

- Time-efficent approach
- Reduces number of comparisons through approximation

(INFAI & FORTH

• No guarantee pertaining to reduction ratio

Mohamed Ahmed Sherif et



LD TUTORIAL: EFFICIENCY

May 28, 2017 15 / 117

Table of Contents

HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link L
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

Mohamed Ahme

LD TUTORIAL: EFFICIENCY

May 28, 2017 16 / 117

Intuition

HOBBIT AKSW

Idea

- Create multidimensional index of data
- Partition the space so that $\sigma(x, y) < \theta \Rightarrow index(x) \neq index(y)$
- Decrease number of computations by only comparing pairs (x, y) with index(x) = index(y)

Mohamed Ahmed Sherif et al. (InfAI & FORTH

D TUTORIAL: EFFICIENCY

May 28, 2017 17 / 117

Intuition

S HOBBIT AKSW

Idea

- Create multidimensional index of data
- Partition the space so that $\sigma(x, y) < \theta \Rightarrow index(x) \neq index(y)$
- Decrease number of computations by only comparing pairs (x, y) with index(x) = index(y)



Approach





- Similarity functions are **aggregations** of atomic measures *sim*, e.g., $\sigma = MIN(levenshtein, jaccard)$
- Each atomic measure must define a blocking function $block:(S\cup T) imes [0,1] o \mathcal{P}(\mathbb{N}^n)$
- Each aggregation must define a similarity, a block and a threshold aggregation

) Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

Example



- $sim = \frac{levenshtein(s,t)}{\max(|s|,|t|)}$ (normalized levenshtein)
- block : $\Sigma^q \to \mathbb{N}$
 - Maps each q-gram to exactly one block (e.g., Goedelisation)
 - Only $c = \max(|s|, |t|)(1 \theta) \cdot q + 1$ q-grams to be indexed per string
- $index(s, \theta) = \{block(qgrams(s_{[0...c]}))\}$
- String mapped to several blocks
- Comparison only within blocks

TUTORIAL: EFFICIENCY

May 28, 2017 19 / 117

Example

BOBBIT AKSW

- $sim = \frac{levenshtein(s,t)}{\max(|s|,|t|)}$ (normalized levenshtein)
- block : $\Sigma^q \to \mathbb{N}$
 - Maps each q-gram to exactly one block (e.g., Goedelisation)
 - Only $c = \max(|s|, |t|)(1 \theta) \cdot q + 1$ q-grams to be indexed per string
- $index(s, \theta) = \{block(qgrams(s_{[0...c]}))\}$
- String mapped to several blocks
- Comparison only within blocks
- Comparison of drugs in DBpedia and Drugbank

Setting	Comparisons	Runtime	Links
Full comparison	22,242,292	430s	1,403
Blocking (100 blocks)	906,314	44s	1,349
Blocking (1,000 blocks)	322,573	14s	1,287
MultiBlock	122,630	бs	1,403

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

ORIAL: EFFICIENCY

AY 28, 2017 19 / 117

Table of Contents

HOBBIT AKSW

	Inter	- Jai	ati	
	111110		CU	O'II
-				

- 2 LIMES
- 3 MultiBlock
- Reduction-Ratio-Optimal Link Discovery

(INFAI & FORTH

- **5** AEGLE
- 6 RADON
- **G**NOME
- 8 GPUs and Hadoop

Mohamed Ahmed Sherif

9 Summary and Conclusion

LD TUTORIAL: EFFICIENC

May 28, 2017 20 / 117





Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s,t) \in S \times T : \delta(s,t) \le \theta\}$

MOHAMED AHMED SHERIF ET AL. (INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 21 / 117



BOBBIT AKSW

Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s,t) \in S \times T : \delta(s,t) \le \theta\}$
- Intuition: Reduce the number of comparisons $\mathcal{C}(\mathcal{A}) \geq |\mathcal{M}'|$
- Maximize reduction ratio: $RR(A) = 1 \frac{C(A)}{|S||T|}$



Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \delta(s, t) \le \theta\}$
- Intuition: Reduce the number of comparisons $C(\mathcal{A}) \geq |\mathcal{M}'|$
- Maximize reduction ratio: $RR(A) = 1 \frac{C(A)}{|S||T|}$

Question

 \mathcal{HR}^3

Link Discovery

- Can we devise lossless approaches with guaranteed RR?
- Advantages
 - Space management
 - Runtime prediction
 - Resource scheduling

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D TUTORIAL: EFFICIENCY

May 28, 2017 21 / 117





• Best achievable reduction ratio: $RR_{max} = 1 - \frac{|M'|}{|S||T|}$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 22 / 117



HOBBIT AKSW

- Best achievable reduction ratio: $RR_{max} = 1 \frac{|M'|}{|S||T|}$
- Approach $\mathcal{H}(\alpha)$ fulfills RR guarantee criterion, iff:

FORTH

 $\forall r < RR_{\max}, \exists \alpha : RR(\mathcal{H}(\alpha)) \geq r$

May 28, 2017 22 / 117



- Best achievable reduction ratio: $RR_{max} = 1 \frac{|M'|}{|S||T|}$
- Approach $\mathcal{H}(\alpha)$ fulfills RR guarantee criterion, iff:

 $\forall r < \mathit{RR}_{\mathsf{max}}, \exists \alpha : \mathit{RR}(\mathcal{H}(\alpha)) \geq r$

• Here, we use relative reduction ratio (*RRR*):

FORTH

$$RRR(\mathcal{A}) = rac{RR_{\max}}{RR(\mathcal{A})}$$

fay 28, 2017 22 / 117





Formal Goal

Devise $\mathcal{H}(\alpha)$: $\forall r > 1, \exists \alpha : RRR(\mathcal{H}(\alpha)) \leq r$

FORTH

LD TUTORIAL

May 28, 2017 23 / 117



HOBBIT AKSW



Mohamed Ahmed Sherif et al.

(INFAI & FORTH





May 28, 2017 24 / 117



HOBBIT AKSW

HYPPO

- $\delta(s,t) \le \theta$ describes a hypersphere
- Approximate hypersphere by using a hypercube
 - Easy to compute
 - No loss of recall (blocking)



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 25 / 117





• Set width of single hypercube to $\Delta=\theta/\alpha$

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 26 / 117





- Set width of single hypercube to $\Delta = \theta / \alpha$ Tile $\Omega = S \cup T$ into the adjacent cubes C
 - - Coordinates: $(c_1, \ldots, c_n) \in \mathbb{N}^n$
 - Contains points $\omega \in \Omega : \forall i \in \{1 \dots n\}, c_i \Delta \leq \omega_i < (c_i + 1) \Delta$

FORTH

Efficie





- Set width of single hypercube to $\Delta = \theta / \alpha$ Tile $\Omega = S \cup T$ into the adjacent cubes C
 - - Coordinates: $(c_1, \ldots, c_n) \in \mathbb{N}^n$
 - Contains points $\omega \in \Omega$: $\forall i \in \{1 \dots n\}, c_i \Delta \leq \omega_i < (c_i + 1)\Delta$



(INFAI & FORTH

LD TUTORIAL: EFFICIEN

HYPPO

BOBBIT AKSW

• Combine $(2\alpha + 1)^n$ hypercubes around $C(\omega)$ to approximate hypersphere



LD TUTORIAL:

Efficie

• $RRR(HYPPO(\alpha)) = \frac{(2\alpha+1)^n}{\alpha^n S(n)}$

Mohamed Ahmed Sherif et al. (InfAI & FORTH

• $\lim_{\alpha \to \infty} RRR(HYPPO(\alpha)) = \frac{2^n}{S(n)}$

IAY 28, 2017 27 / 117



Mohamed Ahmed Sheri



• RRR(HYPPO) for p = 2, n = 2, 3, 4 and $2 \le \alpha \le 50$

(INFAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 28 / 117



Mohamed Ahmed Sherif et al.



- RRR(HYPPO) for p = 2, n = 2, 3, 4 and $2 \le \alpha \le 50$
- $\lim_{\alpha \to \infty} RR(HYPPO(\alpha)) = \frac{4}{\pi} \approx 1.27 \ (n = 2)$
- $\lim_{\alpha \to \infty} RR(HYPPO(\alpha)) = \frac{6}{\pi} \approx 1.91 \ (n = 3)$
- $\lim_{\alpha \to \infty} RRR(HYPPO(\alpha)) = \frac{32}{\pi^2} \approx 3.24 \ (n = 4)$

(INFAI & FORTH

LD TUTORIAL: EFFICIENC



Aay 28, 2017 28 / 117





$$index(C,\omega) = \begin{cases} 0 \text{ if } \exists i : |c_i - c(\omega)_i| \le 1, 1 \le i \le n, \\ \sum_{i=1}^n (|c_i - c(\omega)_i| - 1)^p \text{ else,} \end{cases}$$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 29 / 117





- Compare $\mathcal{C}(\omega)$ with \mathcal{C} iff $\mathit{index}(\mathcal{C},\omega) \leq lpha^p$
- $\alpha = 4, p = 2$

18	13	10	0	0	0	10	13	18
13	8	5	0	0	0	5	8	13
10	5	2	0	0	0	2	5	10
0	0	0	0	0	0	0	0	0
0	0	0	0	0•	0	0	0	0
0	0	0	0	0	0	0	0	0
10		2	0	0	0	2	5	10
13	8	S	0	0	0	5	8	13
18	13	10	0	0	0	10	13	18

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 30 / 117





Claims

• No loss of recall

Mohamed Ahmed Sherif et

• $\lim_{\alpha \to \infty} RRR(\mathcal{HR}^3(\alpha)) = 1$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 31 / 117




Lemma

 $index(C,s) = x \Rightarrow \forall t \in C \ \delta^p(s,t) > x\Delta^p$

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 32 / 117





Lemma

 $index(C,s) = x \Rightarrow \forall t \in C \ \delta^p(s,t) > x\Delta^p$

• p = 2, n = 2, index(C, s) = 5



Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 32 / 117



BOBBIT AKSW

Lemma

 $\mathit{index}(\mathit{C}, \mathit{s}) = x \Rightarrow \forall t \in \mathit{C} \ \delta^{p}(\mathit{s}, t) > x\Delta^{p}$

Proof.

• index
$$(C,s) = x \Rightarrow \sum_{i=1}^{n} (|c_i - c_i(s)| - 1)^p = x$$

• Out of definition of cube index follows $|s_i - t_i| > (|c_i - c_i(s)| - 1)\Delta$

LD TUTORIAL:

Efficienc

• Thus,
$$\sum_{i=1}^{n} |s_i - t_i|^p > \sum_{i=1}^{n} (|c_i - c_i(s)| - 1)^p \Delta^p$$

• Therewith,
$$\delta^p(s,t) > x\Delta^p$$

May 28, 2017 33 / 117





Lemma

 $\forall s \in S : index(C,s) > lpha^p$ implies that all $t \in C$ are non-matches

TUTORIAL:

Efficies

Proof.

Follows directly from Lemma 1: $index(C, s) > \alpha^p \Rightarrow \forall t \in C, \delta^p(s, t) > \Delta^p \alpha^p = \theta^p$

(INFAI & FORTH

May 28, 2017 34 / 117





Claims

Mohamed

- No loss of recall 🗸
- $\lim_{\alpha \to \infty} RRR(\mathcal{HR}^3(\alpha)) = 1$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 35 / 117





Lemma

 $\forall \alpha > 1 \ RRR(\mathcal{HR}^{3}(2\alpha)) < RRR(\mathcal{HR}^{3}(\alpha))$

• *p* = 2, *α* = 4



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

May 28, 2017 36 /



Lemma

Proof (idea)

 \mathcal{HR}^3

 $\forall \alpha > 1 \ RRR(HR^{3}(2\alpha)) < RRR(HR^{3}(\alpha))$

• $p = 2, \alpha = 8$



AY 28, 2017 37 / 117





Mohamed Ahmed Sherif et al. (InfAI & FORT)

LD TUTORIAL: EFFICIENCY

May 28, 2017 39 / 117



🚓 HOBBIT AKSW

Theorem

 $\lim_{\alpha \to \infty} RRR(\mathcal{HR}^3(\alpha)) = 1$

Proof.

•
$$\alpha \to \infty \Rightarrow \Delta \to 0$$

Mohamed Ahmed Sherif

• Thus,
$$C(s) = \{s\}, C = \{t\}$$

• Index function
$$:\sum_{i=1}^{n} \Delta^{p} (|c_{i}(s) - c_{i}| - 1)^{p} \leq \Delta^{p} \alpha^{p}$$

(INFAI & FORTH

LD TUTORIAL: EFFICIENC

•
$$\Delta \to 0 \Rightarrow \sum_{i=1}^n |s_i - t_i|^p \le \theta^p$$

May 28, 2017 40 / 117





Claims

• No loss of recall 🗸

Mohamed Ahmed Sherif et

• $\lim_{\alpha \to \infty} RRR(\mathcal{HR}^3(\alpha)) = 1 \checkmark$

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 41 / 117



Evaluation

 \mathcal{HR}^3

 $\bullet~\mbox{Compare}~\ensuremath{\mathcal{HR}^3}$ with LIMES 0.5's HYPPO and SILK 2.5.1

(InfAI & FORTH

- Experimental Setup:
 - Deduplicating DBpedia places by minimum elevation, elevation and maximum elevation $(\theta = 49m, 99m)$.
 - Geonames and LinkedGeoData by longitude and latitude ($heta=1^\circ,9^\circ$)
- Windows 7 Enterprise machine 64-bit computer with a 2.8GHz i7 processor with 8GB RAM.



 \mathcal{HR}^3 Evaluation (Comparisons)

• Experiment 2: Deduplicating DBpedia places, $\theta = 99m$ • 0.64×10^6 less comparisons





 \mathcal{HR}^3 Experiments (Comparisons)

• Experiment 4: Linking Geonames and LinkedGedData, $\theta=9^\circ$ • 4.3×10^6 less comparisons





 \mathcal{HR}^3 Experiments (RRR)

Mohamed Ahmei

• Experiment 3,4: Geonames and LGD, $\theta=1,9^\circ$





46

 $\bullet\,$ Experiment 3,4: Geonames and LGD, $\theta=1,9^\circ$

 \mathcal{HR}^3

Experiments (Runtime)





Experiments (Runtime)

 \mathcal{HR}^3

Mohamed

- Experiment 1, 2: DBpedia, $\theta = 49,99m$
- Experiment 3, 4: Geonames and LGD, $heta=1,9^\circ$



May 28, 2017 47 / 117



- Main result: new category of algorithms for link discovery
- Outperforms the state of the art (runtime, comparisons)
- Future Work

 \mathcal{HR}^3

Conclusion

- $\bullet~\mbox{Combine}~\ensuremath{\mathcal{HR}^3}$ with multi-indexing approach
- Devise resource management approach

(INFAI & FORTH

• Develop other algorithms (esp. for strings) with the same/similar theoretical guarantees

Table of Contents

🚵 HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link Di
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

LD TUTORIAL: EFFICIENC

Mohamed

May 28, 2017 49 / 11







:E1	rdfs:label "Engine failure"@en
:E1	rdf:type :Error
:E1	:beginDate :"2016-04-22T11:39:35"
:E1	:endDate :"2015-04-22T11:39:37"

:E2 rdfs:label "Car accident"@en :E2 rdf:type :Accident :E2 :beginDate :"2015-06-28T11:45:22" :E2 :endDate :"2015-06-28T11:45:24"

- Need to create links between events, e.g., :startsEvent
- Need to deal with volume and velocity

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 50 / 117



HOBBIT AKSW

Definition (Event)

- Events can be modeled as time intervals: v = (b(v), e(v))
 - b(v) is the beginning time (:beginDate)
 - e(v) is the end time (:endDate)

(INFAI & FORTH)

• b(v) < e(v)

LD TUTORIAL: EFFICIENCY

May 28, 2017 51 / 117

AEGLE

Allen's Interval Algebra

Mohamed Ahmed Sherif et al

🚓 HOBBIT AKSW

May 28, 2017

52 / 117

Relation	Notation	Inverse	Illustration	
			X	
X before Y	bf(X, Y)	bfi(X, Y)	Y	
			X	
X meets Y	m(X, Y)	mi(X, Y)	Y	
			X	
X finishes Y	f(X, Y)	fi(X, Y)	Y	
			X	
X starts Y	st(X, Y)	sti(X, Y)	Y	
			X	
X during Y	d(X, Y)	di(X, Y)	Y	
			X	
X equal Y	eq(X, Y)	eq(X, Y)	Y	
			X	
X overlaps with Y	ov(X, Y)	ovi(X, Y)	Y	
F FT AL (INFAL & FORTH)	I.D.	TUTORIAL · FEFICIENCY	Ma	v 9







(INFAI & FORTH)

AEGLE: <u>A</u>llen's int<u>E</u>rval al<u>G</u>ebra for <u>L</u>ink discov<u>E</u>ry

- Efficient computation of temporal relations between events
- Allen's Interval Algebra: distinct, exhaustive, and qualitative relations between time intervals
- Intuition
 - Expressing 13 Allen relations using 8 atomic relations
 - Time is ordered: Find matching entities using two sorted lists

fay 28, 2017 53 / 117

AEGLE Express st(s, t) using atomic relations

(INFAI & FORTH





LD TUTORIAL: EFFICIEN

May 28, 2017 54 / 117

AEGLE Express st(s, t) using atomic relations





(INFAI & FORTH

TUTORIAL

Efficien



AEGLE

Express st(s, t) using atomic relations





AEGLE Express st(s, t) using atomic relations

BOBBIT AKSW





- Compute 8 atomic Boolean relations between begin and end points
 - BeginBegin (BB) for b(s), b(t):

(INFAI & FORTH

• $BB^{1}(s,t) \Leftrightarrow (b(s) < b(t))$ • $BB^{0}(s,t) \Leftrightarrow (b(s) = b(t))$ • $BB^{-1}(s,t) \Leftrightarrow (b(s) > b(t)) \Leftrightarrow \neg (BB^{1}(s,t) \lor BB^{0}(s,t))$

LD TUTORIAL: EFFICIENC

- BeginEnd(BE) for b(s), e(t)
- EndBegin(EB) for e(s), b(t)
- EndEnd(EE) for e(s), e(t)



May 28, 2017 56 / 117





AEGLE	
Algorithm for st, sti	





(INFAI & FORTH

Mohamed Ahmed Sher

LD TUTORIAL: EFFICIENCY

May 28, 2017 57 / 117



LD TUTORIAL: EFFICIENC

May 28, 2017

57

Mohamed Ahmed Sheri

(INFAI & FORTH





LD TUTORIAL: EFFICIE

201'

(InfAI & FORTH



FORTH

🚓 HOBBIT AKSW






Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 58 / 117



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 58 / 117



(INFAI & FORTH

HOBBIT AKSW

2017



LD TUTORIAL: EFFICIEN



BOBBIT AKSW

201



FORTH



AEGLE

Experimental Set-Up

Mohamed Ahmed Sherii

(INFAI & FORTH)

BOBBIT AKSW

• Datasets: S = T

Log Type	Dataset name	Size	Unique $b(s)$	Unique <i>e</i> (<i>s</i>)
	3KMachines	3,154	960	960
Machinery	30KMachines	28,869	960	960
	300KMachines	288,690	960	960
	3KQueries	3,888	3,636	3,638
Query	30KQueries	30,635	3,070	3,070
	300KQueries	303,991	184	184

LD TUTORIAL: EFFICIENCY

AEGLE

Experimental Set-Up

BOBBIT AKSW

• Datasets: S = T

Log Type	Log Type Dataset name		Unique <i>b</i> (<i>s</i>)	Unique <i>e</i> (<i>s</i>)
	3KMachines	3,154	960	960
Machinery	30KMachines	28,869	960	960
	300KMachines	288,690	960	960
	3KQueries	3,888	3,636	3,638
Query	30KQueries	30,635	3,070	3,070
	300KQueries	303,991	184	184

• State-of-the-art:

- $\bullet~\mathrm{SILK}$ extended to deal with spatio-temporal data
- Baseline for eq using brute-force

Iohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 59 / 117

AEGLE Experimental Set-Up

HOBBIT AKSW

• Datasets: S = T

Log Type	Dataset name	Size	Unique <i>b</i> (<i>s</i>)	Unique $e(s)$
	3KMachines	3,154	960	960
Machinery	30KMachines	28,869	960	960
	300KMachines	288,690	960	960
	3KQueries	3,888	3,636	3,638
Query	30KQueries	30,635	3,070	3,070
	300KQueries	303,991	184	184

• State-of-the-art:

- $\bullet~\mathrm{SILK}$ extended to deal with spatio-temporal data
- Baseline for *eq* using brute-force
- Evaluation measures:
 - atomic runtime of each of the atomic relations
 - relation runtime required to compute each Allen's relation
 - total runtime required to compute all 13 relations

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

IAY 28, 2017 59 / 117

AEGLEResults



25 atomic runtime 20 **3K Queries** 15 30K Queries 300K Queries 10 5. 0 BE1 EB0 EB1 EE0 BB1 BE0 EE1 BB0 atomic relations LD TUTORIAL

Query Datasets

30-

• Q₁: Does the reduction of Allen relations to 8 atomic relations influence the overall runtime of the approach? AEGLEResults



120 100 atomic runtime 80 60 **3K Events** 30K Events 300K Events 40 20 -0 -BB1 BE0 BE1 EB0 EB1 EE0 EE1 BB0 atomic relations LD

Machine Datasets

• Q₁: Does the reduction of Allen relations to 8 atomic relations influence the overall runtime of the approach?



• Q_2 : How does AEGLE perform when compared with the state of the art in terms of time efficiency?

	Datasat Nama	Total Runtime			
Log Type	Dataset Mame	Aegle	Aegle $*$	Silk	
	3KMachines	11.26	5.51	294.00	
Machine	30KMachines	1,016.21	437.79	29,846.00	
	300KMachines	189,442.16	78,416.61	NA	
	3KQueries	26.94	17.91	541.00	
Query	30KQueries	988.78	463.27	33,502.00	
	300KQueries	211,996.88	86,884.98	NA	

May

HOBBIT AKSW

			Machine			Query	
Relation	Approach	3KMachines	30KMachines	300KMachines	3KQueries	30KQueries	300KQueries
m	Aegle	0.02	0.19	3.42	0.02	0.21	3.89
111	Silk	23.00	2,219.00	NA	41.00	2,466.00	NA
	Aegle	0.05	0.79	49.84	0.05	0.45	348.51
eq	Silk	23.00	2,250.00	NA	41.00	2,473.00	NA
	baseline	2.05	171.10	23,436.30	3.15	196.09	31,452.54
	Aegle	3.16	222.27	38,226.32	11.97	257.59	42,121.68
OVI	Silk	22.00	2,189.00	NA	42.00	2,503.00	NA

Conclusion

- AEGLE = efficient temporal linking
 - Reduction of 13 Allen Interval relations to 8 atomic relations
 - Efficiency: simple sorting with complexity $O(n \log n)$
 - Scalable LD

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 63 / 117

Table of Contents

Mohamed

🚵 HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link Disc
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 64 / 117

HOBBIT AKSW

• Large Geo-spatial datasets

RADON Motivation

- LinkedGeoData contains > 20+ billion triples
- NUTS contains up to 1,500 points per resources
- Only 7.1% of the links between resources connect geo-spatial entities (Ngonga Ngomo, 2013)
- Geo-spatial resources available on the LOD
 - Described using polygons
 - Large in number
 - Demands the computation of topological relations
- Naïve computation of topological relations is impracticable for geo-spatial resources

(INFAI & FORTH





- Standard to describe the topological relations in 2D space.
- DE-9IM is to based on the intersection matrix:

(INFAI & FORTH)

$$DE9IM(a, b) \begin{bmatrix} dim(I(g_1) \cap I(g_2)) & dim(I(g_1) \cap B(g_2)) & dim(I(g_1) \cap E(g_2)) \\ dim(B(g_1) \cap I(g_2)) & dim(B(g_1) \cap B(g_2)) & dim(B(g_1) \cap E(g_2)) \\ dim(E(g_1) \cap I(g_2)) & dim(E(g_1) \cap B(g_2)) & dim(E(g_1) \cap E(g_2)) \end{bmatrix}$$

HOBBIT AKSW

- There must be at least one shared point for a relation to be hold
- Except for the disjoint relation \Rightarrow inverse of the intersects relation
- Accelerating the computation of whether two geometries share at least one point, accelerates the computation of any topological relation

JTORIAL: EFFICIENCY

RADON Basic Idea

HOBBIT AKSW





• RADON implements improved indexing approach based on

Minimum bounding boxes (MBB)

(INFAI & FORTH)

Space tiling

Mohamed Ahmed Sherif

LD TUTORIAL: EFFICIENCY

May 28, 2017 67 / 117

RADON

I. Swapping Strategy

HOBBIT AKSW

- Large geometries that span over a large number of hypercubes \Rightarrow large spatial index when used as S
- *Estimated Total Hypervolume* (*ETH*) of a set of geometries *X*

$$ETH(X) = |X| \prod_{i=1}^{d} \frac{1}{|X|} \sum_{x \in X} \left(\max_{p \in x} \{ \kappa_i(p) \} - \min_{p \in x} \{ \kappa_i(p) \} \right)$$

- If ETH(S) > ETH(T), swaps S and T and computes the reverse relation r' instead of r
- e.g. if r is covered and ETH(S) > ETH(T), then swaps
 S and T and computes coveredBy

(INFAI & FORTH



Since ETH(NUTS) > ETH(CLC), then S = CLC and T = NUTS

ay 28, 2017 68 / 117





- Insert all geometries $s \in S$ into index I(s)
 - Computes MBB(s)
 - ⁽²⁾ Maps each s to all hypercubes over MBB(s) spans
- Same procedure for all $t \in T$ but only index geometries t that are potentially in hypercubes already contained in I(S)



Radon

III. Link Generation

AKSW

- \bullet Discards unnecessary computations using the $\rm Test MBB$ procedure
- $\bullet~{\rm TESTMBB}$ optimizes the subset of DE-9IM relations for relations where
 - one geometry has interior or boundary points in the exterior of the other geometry
 - e.g. equals, covers and within
- \bullet For other relations, $\rm TestMBB$ returns true

(INFAI & FORTH)

- $\bullet~\mbox{If TestMBB}$ returns false
 - No need to compute the expensive computation of the topological relation

TUTORIAL: EFFICIENC



- TESTMBB(within, blue) = false
- TESTMBB(within, green) = true

HOBBIT AKSW

- Topological relations
 - Subset of the 7 topological relations
 - i.e. within, touches, overlaps, intersects, equals, crosses and covers
- Hardware

RADON

Evaluation Setup

- 64-core 2.3 GHz, OpenJDK 64-Bit Server
- 20 GB RAM with timeout limit of 2 hours
- State of the art
 - Silk
 - STRABON
- Datasets
 - NUTS
 - ORINE Land Cover (CLC)



D TUTORIAL: EFFICIENCY

🚓 HOBBIT 🗛 📕

• RADON **vs**. SILK

Evaluation – Linear Speedup

- 44 subsets of the CLC vs. the full NUTS
- 7 basic topological relations
- 308 experiments
- Single core

RADON

- RADON achieves a linear speedup relative to the dataset sizes
- Up to 450 times faster for the within relation

(INFAI & FORTH



Y 28, 2017 72 / 117

🚓 HOBBIT AKSW

Evaluation – Topological Relations Computations

RADON

- Same setting as in previous experiments
- $\bullet \ \mathrm{Radon}$ runs significantly less computations of the relations
- RADON carries out only 3 and 4 computations for the equals and within relations respectively
- On average, 449 times less computations per relation

(INFAI & FORTH



D TUTORIAL: EFFICIENCY



Evaluation – Runtime

RADON

- Same setting as in previous experiments
- On average, RADON is faster then
 - SILK by 65.62 times
 - $\bullet~\mathrm{Strabon}$ by 11.99 times
- STRABON outperforms RADON on the intersects relation
 - STRABON uses an **R-tree-over-GiST** spatial index over the stored geometries in the underlying **PostGIS** database
 - R-tree-over-GiST is highly optimized for the retrieval of spatially connected objects

(INFAI & FORTH



LD TUTORIAL: EFFICIENC

Iay 28, 2017 74 / 117

HOBBIT AKSW

Evaluation – Speedup Quantification (Parallel Implementation)

- Merge all 44 sub-datasets of CLC (CLC_M)
- CLC_M contains 2, 209, 538 resource

RADON

- CLC_M as both source and target datasets
- Simple round robin load balancing policy
- On average, within the 2 hours time limit
 - RADON finishes in 20.83 minutes
 - $\bullet~{\rm SILK}$ finalizes 1.16% of the tasks
 - SILK would need 4.36 days with 8 threads (linear extrapolation)
 - $\bullet~\mathrm{RADON}$ is 834.69 times faster than Silk

Relation	# Thr.	Radon	Silk	Speedup
	1	24.11	36,500 (0.33%)	1,513.58
	2	13.15	21,667 (0.55%)	1,647.58
equais	4	6.81	11,750 (1.02%)	1,725.77
	8	3.79	6,286 (1.91%)	1,658.78
	1	93.17	37,500 (0.32%)	402.50
	2	49.03	20,667 (0.58%)	421.53
Intersects	4	25.11	12,000 (1.00%)	477.81
	8	13.04	6,300 (1.90%)	483.24
	1	36.47	35,000 (0.34%)	959.74
	2	18.26	20,667 (0.58%)	1,131.86
within	4	9.44	11,765 (1.02%)	1,246.34
	8	5.92	6,202 (1.93%)	1,048.34
	1	35.62	36,000 (0.33%)	1,010.75
	2	18.51	21,029 (0.57%)	1,136.10
covers	4	10.23	12,000 (1.00%)	1,172.50
	8	5.33	6,300 (1.90%)	1,182.13
	1	94.50	35,500 (0.34%)	375.68
	2	47.71	22,196 (0.54%)	465.18
touches	4	25.09	12,121 (0.99%)	483.08
	8	13.30	6,381 (1.88%)	479.75

LD TUTORIAL: EFFICIENCY

HOBBIT AKSW

Evaluation – Speedup Quantification (Parallel Implementation)

- STRABON does not finish any of the experiments within the 2-hours time limit
- No progress feedback from STRABON
- STRABON performance estimation

RADON

- STRABON vs. sub-CLC dataset 10 times < CLC_m
- Optimistic (STRABON scales linearly), average speedup of 24

(INFAL & FORTH

• Realistic (STRABON scales in $O(n^2)$), average speedup of 241

Relation	#Thr.	Radon	Silk	Speedup	
	1	24.11	36,500 (0.33%)	1,513.58	
	2	13.15	21,667 (0.55%)	1,647.58	
equais	4	6.81	11,750 (1.02%)	1,725.77	
	8	3.79	6,286 (1.91%)	1,658.78	
	1	93.17	37,500 (0.32%)	402.50	
	2	49.03	20,667 (0.58%)	421.53	
Intersects	4	25.11	12,000 (1.00%)	477.81	
	8	13.04	6,300 (1.90%)	483.24	
	1	36.47	35,000 (0.34%)	959.74	
	2	18.26	20,667 (0.58%)	1,131.86	
within	4	9.44	11,765 (1.02%)	1,246.34	
	8	5.92	6,202 (1.93%)	1,048.34	
	1	35.62	36,000 (0.33%)	1,010.75	
	2	18.51	21,029 (0.57%)	1,136.10	
covers	4	10.23	12,000 (1.00%)	1,172.50	
	8	5.33	6,300 (1.90%)	1,182.13	
touches	1	94.50	35,500 (0.34%)	375.68	
	2	47.71	22,196 (0.54%)	465.18	
	4	25.09	12,121 (0.99%)	483.08	
	8	13.30	6,381 (1.88%)	479.75	

LD TUTORIAL: EFFICIENCY

May 28, 2017 76 / 117



- Conclusion & Future Work
 - Conclusion

RADON

- Presented RADON, an approach for rapid discovery of topological relations among geo-spatial resources
- Sparse index combines space tiling and MBB approximation
- RADON is complete and correct
- Outperforms the SOTA by up to 3 orders of magnitude
- Future work
 - More sophisticated load balancing approaches
 - Other topology approximation methods
 - Topological relations in higher dimensions

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 77 / 117

Table of Contents

HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link D
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

LD TUTORIAL: EFFICIENC

Mohamed

May 28, 2017 78 / 11

Problem

🚓 HOBBIT AKSW

What if ...

- Data does not fit memory C, i.e.,
- |S| + |T| > |C|



- Common problem
- Growing size and number of datasets
 - $pprox 150 imes 10^9$ triples
 - $\approx 10 \times 10^3$ datasets
 - $\bullet~\mbox{Largest}$ dataset with $>20\times10^9 \mbox{triples}$
- Mostly in-memory solutions

Idea

🚓 HOBBIT AKSW

Insight

- Most approaches rely on divide-and-merge paradigm
- Example: HR³
 - $\sigma(s,t) \geq \theta \Leftrightarrow \delta(s,t) \leq \Delta$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 80 / 117

Idea

HOBBIT AKSW

Insight

- Most approaches rely on divide-and-merge paradigm
- Example: HR³
 - $\sigma(s,t) \ge \theta \Leftrightarrow \delta(s,t) \le \Delta$



May 28, 2017 80 / 117

Formal Model



• Define $S = \{S_1, \ldots, S_n\}$ with $S_i \subseteq S \land \bigcup_i S_i = S$

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 81 / 117

Formal Model



- Define $S = \{S_1, \ldots, S_n\}$ with $S_i \subseteq S \land \bigcup_i S_i = S$ Define $T = \{T_1, \ldots, S_m\}$ with $T_j \subseteq T \land \bigcup_j T_j = T$

Formal Model

HOBBIT AKSW

- Define $S = \{S_1, \ldots, S_n\}$ with $S_i \subseteq S \land \bigcup_i S_i = S$ Define $T = \{T_1, \ldots, S_m\}$ with $T_j \subseteq T \land \bigcup_j T_j = T$
- Find mapping function $\mu: \mathcal{S} \to 2^{\mathcal{T}}$ with
 - elements of S_i only compared with elements of sets in $\mu(S_i)$
 - union of results over all $S_i \in S$ is exactly M'.

Formal Model

HOBBIT AKSW

- Define $S = \{S_1, \ldots, S_n\}$ with $S_i \subseteq S \land \bigcup_i S_i = S$ Define $T = \{T_1, \ldots, S_m\}$ with $T_j \subseteq T \land \bigcup_j T_j = T$
- $\textbf{ § Find mapping function } \mu: \mathcal{S} \rightarrow 2^{\mathcal{T}} \text{ with }$
 - elements of S_i only compared with elements of sets in $\mu(S_i)$
 - union of results over all $S_i \in S$ is exactly M'.



Task Graph

HOBBIT AKSW

Definition

• A task E_{ij} stands for comparing S_i with $\mathcal{T}_j \in \mu(S_i)$

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 82 / 117
Task Graph

HOBBIT AKSW

Definition

• A task E_{ij} stands for comparing S_i with $T_j \in \mu(S_i)$

• Task Graph $G = (V, E, w_v, w_e)$, with

•
$$V = S \cup T$$

•
$$w_v(v) = |V|$$

•
$$w_e(e_{ij}) = |S_i||T_j|$$



 $Task \ Graph$

HOBBIT AKSW

Definition

- A task E_{ij} stands for comparing S_i with $T_j \in \mu(S_i)$
- Task Graph $G = (V, E, w_v, w_e)$, with
 - $V = S \cup T$
 - $w_v(v) = |V|$
 - $w_e(e_{ij}) = |S_i||T_j|$



Task Graph

🚓 HOBBIT 🗛 🖌

Definition

- A task E_{ij} stands for comparing S_i with $T_j \in \mu(S_i)$
- Task Graph $G = (V, E, w_v, w_e)$, with
 - $V = S \cup T$
 - $w_v(v) = |V|$
 - $w_e(e_{ij}) = |S_i||T_j|$



Mohamed Ahmed

Problem Reformulation

HOBBIT AKSW

Locality maximization

- Two-step approach:
 - Clustering: Find groups of nodes that fit in memory and

FORTH

Oscheduling: Compute sequence of groups that minimizes hard drive access



Step1: Clustering

🚓 HOBBIT AKSW

Naïve Approach

Mohamed Ahmei

FORTH

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

🚓 HOBBIT 🗛 🖌

Naïve Approach

Mohamed Ahme

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

🚓 HOBBIT 🗛 🖌

Naïve Approach

Mohamed Ahme

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

🚓 HOBBIT 🗛 🖌

Naïve Approach

Mohamed Ahme

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

🚓 HOBBIT AKSW

Naïve Approach

Mohamed

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

HOBBIT AKSW

Naïve Approach

Mohamed

- Cluster by S_i
- Example: |C| = 7



Step1: Clustering

Mohamed

HOBBIT AKSW

Greedy Approach

- Start by largest task
- Add connected largest tasks until none fits in $\ensuremath{\mathcal{C}}$
- Example: |C| = 7



Step1: Clustering

Mohamed

HOBBIT AKSW

Greedy Approach

- Start by largest task
- Add connected largest tasks until none fits in $\ensuremath{\mathcal{C}}$
- Example: |C| = 7



Step1: Clustering

Mohamed

HOBBIT AKSW

Greedy Approach

- Start by largest task
- Add connected largest tasks until none fits in $\ensuremath{\mathcal{C}}$
- Example: |C| = 7



Step1: Clustering

Mohamed

HOBBIT AKSW

Greedy Approach

- Start by largest task
- Add connected largest tasks until none fits in $\ensuremath{\mathcal{C}}$
- Example: |C| = 7



Step1: Clustering

Mohamed

HOBBIT AKSW

Greedy Approach

- Start by largest task
- Add connected largest tasks until none fits in $\ensuremath{\mathcal{C}}$
- Example: |C| = 7



Step2: Scheduling

🚓 HOBBIT 🗛 📕

Insights

- **Output of clustering**: Sequence G_1, \ldots, G_N of clusters
- Intuition: Consecutive clusters should share data
- Goal: Maximize overlap of generated sequence

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

Step2: Scheduling

HOBBIT AKSW

Insights

- **Output of clustering**: Sequence G_1, \ldots, G_N of clusters
- Intuition: Consecutive clusters should share data
- Goal: Maximize overlap of generated sequence

AI & FORTH

• Overlap
$$o(G_i, G_j) = \sum_{v \in V(G_i) \cap V(G_j)} |v|$$

•
$$o(G_4, G_1) = 4$$



Step2: Scheduling

SHOBBIT AKSW

Insights

- **Output of clustering**: Sequence G_1, \ldots, G_N of clusters
- Intuition: Consecutive clusters should share data
- Goal: Maximize overlap of generated sequence
- Overlap $o(G_i, G_j) = \sum_{v \in V(G_i) \cap V(G_j)} |v|$ $o(G_4, G_1) = 4$

• Overlap
$$o(G_1, ..., G_N) = \sum_{i=1}^{N-1} o(G_i, G_{i+1})$$

• $o(G_4, G_1, G_2, G_1) = 9$





Step2: Scheduling

🚓 HOBBIT AKSW

Best-Effort

- Select random pair of clusters
- If permutation improves overlap, then permute

fAI & FORTH

• Relies on local knowledge for scalability

LD TUTORIAL: EFFICIENCY

Step2: Scheduling

AKSW

Best-Effort

- Select random pair of clusters
- If permutation improves overlap, then permute
- Relies on local knowledge for scalability

(INFAI & FORTH

• Trick:

$$\Delta(G_i, G_j) = (o(G_{i-1}, G_j) + o(G_j, G_{i+1}) + o(G_{j-1}, G_i) + o(G_i, G_{j+1})) - (o(G_{i-1}, G_i) + o(G_i, G_{i+1}) + o(G_{j-1}, G_j) + o(G_j, G_{j+1}))$$

$$(1)$$



LD TUTORIAL: EFFICIEN



Step 2: Scheduling

🔅 HOBBIT AKSW

Greedy

- Start with random cluster
- Choose next cluster with largest overlap

(InfAI & FORTH

• Global knowledge needed

LD TUTORIAL: EFFICIENCY



Step2: Scheduling

🔅 HOBBIT AKSW

2017

May 28

Greedy

- Start with random cluster
- Choose next cluster with largest overlap

(INFAI & FORTH

• Global knowledge needed



LD TUTORIAL: EFFICIEN

Experimental Setup

BOBBIT AKSW



Datasets

- **OBP:** 1 million labels from DBpedia version 04-2015
- IGD: 0.8 million places from LinkedGeoData

• Hardware

Intel Xeon E5-2650 v3 processors (2.30GHz)

(INFAI & FORTH

- Obuntu 14.04.3 LTS
- 10GB RAM

• Measures

Mohamed Ahmed Sherif

- Total runtime
- e Hit ratio

LD TUTORIAL: EFFICIENCY

Mohamed Ahmed Sherii

Evaluation of Clustering

BIT AKSW

- Only show results of LGD
- Results on DBP lead to similar insights

(INFAI & FORTH)

Runtimes		Hit F	Ratio	
<i>C</i>	Naive	Greedy	Naive	Greedy
100	568.0	646.3	0.57	0.77
200	518.3	594.0	0.66	0.80
400	532.0	593.3	0.67	0.80
1,000	5,974.0	118,454.7	0.51	0.64
2,000	6,168.0	115,450.0	0.51	0.63
4,000	7,118.3	121,901.7	0.50	0.63

LD TUTORIAL: EFFICIENCY

May 28, 2017 90 / 117

Evaluation of Clustering

HOBBIT AKSW

		Runtimes		Hit F	Ratio
	<i>C</i>	Naive	Greedy	Naive	Greedy
	100	568.0	646.3	0.57	0.77
• Only show results of LGD	200	518.3	594.0	0.66	0.80
• Results on DBP lead to similar	400	532.0	593.3	0.67	0.80
Insignts	1,000	5,974.0	118,454.7	0.51	0.64
	2,000	6,168.0	115,450.0	0.51	0.63
	4,000	7,118.3	121,901.7	0.50	0.63

Conclusion

- Naïve approach is more efficient
- Greedy approach is more effective
- Select naïve approach for clustering

Iohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 90 / 117

Mohamed Ahmed Sheri

Evaluation of Scheduling

HOBBIT AKSW

	Runtime	es (ms)	Hit rat	io
<i>C</i>	Best-Effort	Greedy	Best-Effort	Greedy
100	571.3	1,599.3	0.56	0.68
200	565.7	1,448.3	0.66	0.85
400	581.0	1,379.3	0.67	0.88
1,000	5,666.0	814,271.7	0.51	0.86
2,000	6,268.0	810,855.0	0.51	0.86
4,000	6,675.7	814,041.7	0.50	0.86

LD TUTORIAL: EFFICIENC

- Only show results of LGD
- Results on DBP lead to similar insights

(INFAI & FORTH)

Evaluation of Scheduling

🔅 HOBBIT AKSW

		Runtimes (ms)		Hit ratio	
	<i>C</i>	Best-Effort	Greedy	Best-Effort	Greedy
	100	571.3	1,599.3	0.56	0.68
• Only show results of LGD	200	565.7	1,448.3	0.66	0.85
• Results on DBP lead to similar	400	581.0	1,379.3	0.67	0.88
insights	1,000	5,666.0	814,271.7	0.51	0.86
	2,000	6,268.0	810,855.0	0.51	0.86
	4,000	6,675.7	814,041.7	0.50	0.86

\underline{Con} clusion

- Best-effort approach more time-efficient
- Oreedy more effective
- Best-effort approach is to be used for scheduling

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

utorial: Efficiency

May 28, 2017 91 / 117

🚓 HOBBIT AKSW

GNOME Comparison with Caching Approaches

(INFAI & FORTH)

Mohamed Ahmed Sherif et al

			Runtimes (r	ns)		
C	GNOME	FIFO	F2	LFU	LRU	SLRU
1,000	5,974.0	37,161.0	42,090.3	45,906.7	54,194.3	56,904.3
2,000	6,168.0	31,977.0	39,071.3	39,872.0	45,473.0	46,795.0
4,000	7,118.3	21,337.0	40,860.0	28,028.3	26,816.7	27,200.0
			Hit ratio			
1,000	0.51	0.17	0.16	0.19	0.17	0.17
2,000	0.51	0.29	0.30	0.32	0.30	0.30
4,000	0.51	0.54	0.55	0.59	0.55	0.56

LD TUTORIAL: EFFICIENCY

Comparison with Caching Approaches

HOBBIT AKSW

			Runtimes (r	ns)		
C	GNOME	FIFO	F2	LFU	LRU	SLRU
1,000	5,974.0	37,161.0	42,090.3	45,906.7	54,194.3	56,904.3
2,000	6,168.0	31,977.0	39,071.3	39,872.0	45,473.0	46,795.0
4,000	7,118.3	21,337.0	40,860.0	28,028.3	26,816.7	27,200.0
			Hit ratio			
1,000	0.51	0.17	0.16	0.19	0.17	0.17
2,000	0.51	0.29	0.30	0.32	0.30	0.30
4,000	0.51	0.54	0.55	0.59	0.55	0.56

Conclusion

- O GNOME is more time-efficient
- Leads to higher hit rates in most cases

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 92 / 117

AKSW

	100k	200k	400k	800k
LGD	362,141.3	1,452,922.0	5,934,038.7	20,001,965.7
DBP	434,630.7	1,790,350.7	6,677,923.0	12,653,403.3

Conclusion

- Sub-quadratic growth of runtime
- Runtime grows linearly with number of mappings

fAI & FORTH

• For LGD, 360 - 370 mappings/s

) Tutorial: Efficiency

May 28, 2017 93 / 117

Table of Contents

Mohamed

HOBBIT AKSW

1	Introduction
2	LIMES
3	MultiBlock
4	Reduction-Ratio-Optimal Link Dis
5	AEGLE
6	RADON
7	GNOME
8	GPUs and Hadoop
9	Summary and Conclusion

(INFAI & FORTH

LD TUTORIAL: EFFICIENC

Reach for the Cloud? Dealing with Time Complexity



Devise better algorithms



- Blocking
- Algorithms for given metrics
 - PPJoin+, EDJoin
 - \mathcal{HR}^3

f et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

Reach for the Cloud? Dealing with Time Complexity



201



Reach for the Cloud? Parallel Implementations



• Different architectures

Mohamed Ahmed Sheri

• Memory (shared, hybrid, distributed)

(INFAI & FORTH

- Execution paths (different, same)
- Location (remote, local)



LD TUTORIAL: EFFICIENC

Reach for the Cloud? Parallel Implementations



• Different architectures

- Memory (shared, hybrid, distributed)
- Execution paths (different, same)
- Location (remote, local)



hamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: EFFICIENCY

May 28, 2017 96 / 117

Reach for the Cloud? Premises and Goals



Premises

• Given an algorithm that runs on all three architectures ...

Mohamed Ahmed Sherif et al. (InfAI & FORTH

D TUTORIAL: EFFICIENCY

May 28, 2017 97 / 117
Reach for the Cloud? Premises and Goals

HOBBIT AKSW

Premises

- Given an algorithm that runs on all three architectures ...
- Note to self: Implement one
 - Picked \mathcal{HR}^3
 - Reduction-ratio-optimal



Mohamed Ahmed Sherif et al. (InfAI & FORTH

D TUTORIAL: EFFICIENCY

May 28, 2017 97 / 117

Reach for the Cloud? Premises and Goals

HOBBIT AKSW

Premises

- Given an algorithm that runs on all three architectures ...
- Note to self: Implement one
 - Picked \mathcal{HR}^3
 - Reduction-ratio-optimal





Reach for the Cloud?

Goals

- Compare runtimes on all three parallel architectures
- Find break-even points

LD TUTORIAL: EFFICIENCY

AY 28, 2017 97 / 117



• Assume spaces with Minkowski metric and $p \ge 2$



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

May 28, 2017 98 / 117

Mohamed Ahmed Sherif et

(INFAI





2 Create grid of width $\Delta = \theta/\alpha$

May 28, 2017

Mohamed Ahmed

Sherif

HOBBIT AKSW



D TUTORIAL

Efficiei

S Link discovery condition describes a hypersphere

FORTH

(INFAI

May 28, 2017 100 / 117

Mohamed Ahmed

HOBBIT AKSW



Approximate hypersphere with hypercube

May 28, 2017 101 / 117

Mohamed Ahmed

Sherif





FORTH

(INFAI



D TUTORIAL

Efficies

May 28, 2017 102 / 117

Reach for the Cloud? \mathcal{HR}^{3} in GPUs



- Large number of simple compute cores
- Same instruction, multiple data
- Bottleneck: PCI Express Bus
 - Run discretization on CPU
 - Run indexing on GPU
 - Run comparisons on CPU



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 103 / 117

Reach for the Cloud? \mathcal{HR}^3 on the Cloud



• Naive Approach

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 104 / 117

Reach for the Cloud? \mathcal{HR}^3 on the Cloud

HOBBIT AKSW

• Naive Approach

- Rely on Map Reduce paradigm
- Run discretization and assignment to cubes in map step
- In the second second

FORTH



Efficien

LD TUTORIAL

May 28, 2017 104 / 117

Reach for the Cloud? \mathcal{HR}^3 on the Cloud

HOBBIT AKSW

• Load Balancing

Mohamed Ahmed Sherii

- Run two jobs
- Job1: Compute cube population matrix

FORTH

Job2: Distribute balanced linking tasks across mappers and reducers



LD TUTORIAL: EFFICIENCY

May 28, 2017 105 / 117

Reach for the Cloud? Evaluation Hardware

HOBBIT AKSW

- CPU (Java)
 - 32-core server running Linux 10.0.4
 - AMD Opteron 6128 clocket at 2.0GHz

(INFAI & FORTH

- GPU (C++)
 - AMD Radeon 7870 with 20 compute units, 64 parallel threads
 - Host program ran on Intel Core i7 3770 CPU with 8GB RAM and Linux 12.10
 - Ran the Java code on the same machine for scaling
- Cloud (Java)
 - 10 c1.medium nodes (2 cores, 1.7GB) for small experiments
 - 30 c1.large nodes (8 cores, 7GB) for large experiments





Reach for the Cloud? Experimental Setup

HOBBIT AKSW

- Run deduplication task
- Evaluate behavior on different number of dimensions
- Important: Scale results
 - Different hardware (2-7 times faster C++ workstation)
 - Programming language
- Evaluate scalability (DS_4)

Dataset	Source	Size	Features
DS ₁	DBPedia	25,781	min/medium/max elevation
DS ₂	DBPedia	475,000	latitude, longitude
DS ₃	Linked Geo Data	500,000	latitude, longitude
DS ₄	Linked Geo Data	6,000,000	latitude, longitude

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

al: Efficiency

May 28, 2017 107 / 117

Reach for the Cloud?



- DBpedia, 3 dimensions, 26K
- GPU scales better



AY 28, 2017 108 / 117

Reach for the Cloud? Results - DS2



- DBpedia, 2 dimensions, 475K
- GPUs scale better across different θ
- Break-even point $\approx 10^8$ results



May 28, 2017 109 / 117

Reach for the Cloud? Results - DS3



- LinkedGeoData, 2 dimensions, 500K
- Similar picture



Mohamed Ahmed Sherif et al. (InfAI & FORT

TUTORIAL: EFFICIENCY

May 28, 2017 110 / 117

Reach for the Cloud? Results - Scalability



- LinkedGeoData, 2 dimensions, 6M
- $\bullet\,$ Cloud (with load balancing) better for $\approx 10^{10}+\,$ results



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: EFFICIENCY

May 28, 2017 111 / 117

Reach for the Cloud?

Summary



- Question: When should we use which hardware for link discovery?
- Results
 - Implemented \mathcal{HR}^3 on different hardware
 - Provided the first implementation of link discovery on GPUs
 - Devised a load balancing approach for linking on the cloud
 Discovered 10⁸ and 10¹⁰ results as break-even points



Reach for the Cloud?

Summary



- Question: When should we use which hardware for link discovery?
- Insights
 - Load balancing important for using the cloud
 - GPUs: Need faster buses that PCIe (e.g., Firewire speed)
 - Accurate use of local resources sufficient for most of the current applications



 ${\rm May}\ 28,\ 2017 \qquad 113\ /\ 117$

Table of Contents



8	GPUs and Hadoop
7	GNOME
6	RADON
5	AEGLE
4	Reduction-Ratio-Optimal Link Dis
3	MultiBlock
2	LIMES
1	Introduction

(INFAI & FORTH

Mohamed

LD TUTORIAL: EFFICIENCY MAY 28, 2017 114 / 11

Summary and Conclusion

HOBBIT AKSW

- Four approaches to scalability
 - Improve reduction ratio (LIMES, MultiBlock, HYPPO, \mathcal{HR}^3 , ...)
 - Reduce runtime complexity (AEGLE, RADON)
 - Better use of hardware (GNOME)
 - Improve execution of specifications

AI & FORTH

May 28, 2017 1

Summary and Conclusion

BOBBIT AKSW

- Four approaches to scalability
 - Improve reduction ratio (LIMES, MultiBlock, HYPPO, \mathcal{HR}^3 , ...)
 - Reduce runtime complexity (AEGLE, RADON)
 - Better use of hardware (GNOME)
 - Improve execution of specifications
- Challenges include
 - Increase number of reduction-ratio-optimal approaches (HR³, ORCHID)
 - Adaptive resource scheduling
 - Self-regulating approaches
 - Distribution in modern in-memory architecture (SPARK)
 - 5 ...



A cknowledgment

(INFAI & FORTH

HOBBIT AKSW



This work was supported by grants from the EU H2020 Framework Programme provided for the project HOBBIT (GA no. 688227).

Efficienc

May 28, 2017 116 / 117

References I

HOBBIT AKSW

Kleanthi Georgala, Mohamed Ahmed Sherif, and Axel-Cyrille Ngonga Ngomo. "An Efficient Approach for the Generation of Allen Relations". In: ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands - Including Prestigious Applications of Artificial Intelligence (PAIS 2016). Ed. by Gal A. Kaminka et al. Vol. 285. Frontiers in Artificial Intelligence and Applications. IOS Press, 2016, pp. 948–956. ISBN: 978-1-61499-671-2. DOI: 10.3233/978-1-61499-672-9-948. URL: http://dx.doi.org/10.3233/978-1-61499-672-9-948.

Robert Isele, Anja Jentzsch, and Christian Bizer. "Efficient Multidimensional Blocking for Link Discovery without losing Recall". In: Proceedings of the 14th International Workshop on the Web and Databases 2011, WebDB 2011, Athens, Greece, June 12, 2011. Ed. by Amélie Marian and Vasilis Vassalos. 2011. URL: http://webdb2011.rutgers.edu/papers/Paper%2039/silk.pdf.

LD TUTORIAL: EFFICIENCY

May 28, 2017 117 / 117

References II

BOBBIT AKSW

- Axel-Cyrille Ngonga Ngomo and Sören Auer. "LIMES A Time-Efficient Approach for Large-Scale Link Discovery on the Web of Data". In: IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011. Ed. by Toby Walsh. IJCAI/AAAI, 2011, pp. 2312–2317. ISBN: 978-1-57735-516-8. DOI: 10.5591/978-1-57735-516-8/IJCAI11-385. URL: http://dx.doi.org/10.5591/978-1-57735-516-8/IJCAI11-385.
 - Axel-Cyrille Ngonga Ngomo et al. "When to Reach for the Cloud: Using Parallel Hardware for Link Discovery". In: The Semantic Web: Semantics and Big Data, 10th International Conference, ESWC 2013, Montpellier, France, May 26-30, 2013. Proceedings. Ed. by Philipp Cimiano et al. Vol. 7882. Lecture Notes in Computer Science. Springer, 2013, pp. 275–289. ISBN: 978-3-642-38287-1. DOI: 10.1007/978-3-642-38288-8_19. URL: http://dx.doi.org/10.1007/978-3-642-38288-8_19.

References III

HOBBIT AKSW

Axel-Cyrille Ngonga Ngomo. "Link Discovery with Guaranteed Reduction Ratio in Affine Spaces with Minkowski Measures". In: The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference, Boston, MA, USA, November 11-15, 2012, Proceedings, Part I. Ed. by Philippe Cudré-Mauroux et al. Vol. 7649. Lecture Notes in Computer Science. Springer, 2012, pp. 378–393. ISBN: 978-3-642-35175-4. DOI: 10.1007/978-3-642-35176-1_24. URL:
http://dx.doi.org/10.1007/978-3-642-35176-1_24.
Axel-Cyrille Ngonga Ngomo. "ORCHID - Reduction-Ratio-Optimal Computation of

Geo-spatial Distances for Link Discovery". In: The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part I. Ed. by Harith Alani et al. Vol. 8218. Lecture Notes in Computer Science. Springer, 2013, pp. 395–410. ISBN: 978-3-642-41334-6. DOI: 10.1007/978-3-642-41335-3_25. URL: http://dx.doi.org/10.1007/978-3-642-41335-3_25.

AY 28, 2017 119 / 117

References IV

HOBBIT AKSW

- Axel-Cyrille Ngonga Ngomo. "HELIOS Execution Optimization for Link Discovery". In: The Semantic Web - ISWC 2014 - 13th International Semantic Web Conference, Riva del Garda, Italy, October 19-23, 2014. Proceedings, Part I. Ed. by Peter Mika et al. Vol. 8796. Lecture Notes in Computer Science. Springer, 2014, pp. 17–32. ISBN: 978-3-319-11963-2. DOI: 10.1007/978-3-319-11964-9_2. URL: http://dx.doi.org/10.1007/978-3-319-11964-9_2.
 - Axel-Cyrille Ngonga Ngomo and Mofeed M. Hassan. "The Lazy Traveling Salesman -Memory Management for Large-Scale Link Discovery". In: The Semantic Web. Latest Advances and New Domains - 13th International Conference, ESWC 2016, Heraklion, Crete, Greece, May 29 - June 2, 2016, Proceedings. Ed. by Harald Sack et al. Vol. 9678. Lecture Notes in Computer Science. Springer, 2016, pp. 423–438. ISBN: 978-3-319-34128-6. DOI: 10.1007/978-3-319-34129-3_26. URL: http://dx.doi.org/10.1007/978-3-319-34129-3_26.

Iay 28, 2017 120 / 117

Link Discovery Tutorial Part II: Accuracy

Mohamed Ahmed Sherif⁽¹⁾, Irini Fundulaki⁽²⁾, Axel-Cyrille Ngonga Ngomo⁽¹⁾

(1) Institute for Applied Informatics, Germany (2) FORTH, Greece



May 28th, 2017 Portoroz, Slovenia.

LD TUTORIAL: ACCURA

FORTH

May 25, 2017 1 / 68

Table of Contents





Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

May 25, 2017 2 / 68

Table of Contents





Mohamed Ahmed Sherif et al. (InfAI & FORTH) LD TUTORIAL: ACCURACY

May 25, 2017 3 / 68

Introduction
<u>Link Discovery as Classification Task</u>



Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

May 25, 2017 4 / 68

Introduction Link Discovery as Classification Task

BOBBIT AKSW

Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

Definition (Classification perspective)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : C(s, t) = +1\}$
- Here, $\mathcal{C}(s,t) = +1 \leftrightarrow \sigma(s,t) \geq heta$



Introduction Link Discovery as Classification Task

Definition (Declarative Link Discovery)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : \mathcal{R}(s, t)\}$
- Here, find $M' = \{(s, t) \in S \times T : \sigma(s, t) \ge \theta\}$

Definition (Classification perspective)

- \bullet Given sets S and T of resources and relation ${\cal R}$
- Find $M = \{(s, t) \in S \times T : C(s, t) = +1\}$
- Here, $\mathcal{C}(s,t) = +1 \leftrightarrow \sigma(s,t) \geq heta$
- Classical machine learning problem [Ngo+11; NL12]
- Dedicated techniques perform better
- Unsupervised, active and unsupervised techniques possible

OHAMED AHMED SHERIF ET AL. (INFAI & FORTH)

Iay 25, 2017 4 / 68

Introduction

Challenge

BOBBIT AKSW

Challenges

- Creation of labeled training data tedious
- O Need automated means for automatic class and property matching
- O Need for efficient execution of link specifications
- Obdicated machine learning approaches necessary

AL & FORTH

May 25, 2017 5 / 68

Introduction

Challenge

🚓 HOBBIT AKSW

Challenges

- Oreation of labeled training data tedious
- Need automated means for automatic class and property matching
- O Need for efficient execution of link specifications
- Obdicated machine learning approaches necessary

Solutions

- Use active learning approach for link discovery
- Rely on hospital/resident algorithm
- See previous section
- Topic of this section

DHAMED AHMED SHERIF ET AL. (INFAI & FORTH)

LD TUTORIAL: ACCURACY

May 25, 2017 5 / 68

Table of Contents

Mohamed Ahmed Sherif





(INFAI & FORTH) LD TUTORIAL: ACCURACY MAY 25, 2017


HOBBIT AKSW

Definition (Classification perspective)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : C(s, t) = +1\}$
- Here, $\mathcal{C}(s,t) = +1 \leftrightarrow \sigma(s,t) \geq heta$

May 25, 2017 7 / 68



AKSW

Definition (Classification perspective)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : C(s, t) = +1\}$
- Here, $\mathcal{C}(s,t) = +1 \leftrightarrow \sigma(s,t) \geq heta$
- Learning classifier C involves learning
 - Two sets of restrictions that specify the sets S resp. T,
 - **2** the components $\sigma_1 \dots \sigma_n$ of a complex similarity measure σ
 - **③** a set of thresholds $\theta_1, ..., \theta_n$ for $\sigma_1, ..., \sigma_n$



HOBBIT AKSW

Definition (Classification perspective)

- Given sets S and T of resources and relation $\mathcal R$
- Find $M = \{(s, t) \in S \times T : C(s, t) = +1\}$
- Here, $\mathcal{C}(s,t) = +1 \leftrightarrow \sigma(s,t) \geq heta$
- Learning classifier C involves learning
 - $\textcircled{\sc 0}$ Two sets of restrictions that specify the sets S resp. T,
 - **2** the components $\sigma_1 \dots \sigma_n$ of a complex similarity measure σ
 - **(a)** a set of thresholds $\theta_1, ..., \theta_n$ for $\sigma_1, ..., \sigma_n$

Assumptions

- Restrictions are class restrictions
- Classifier shape is given (e.g., linear combination)

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

May 25, 2017 7 / 68



🚓 HOBBIT 🗛 📈

Class and Property Restrictions

- Define class similarity function
- Solve corresponding hospital-resident problem
- Based on extension of stable marriage problem

AI & FORTH

TUTORIAL

May 25, 2017 8 / 68



🚓 HOBBIT 🗛 📈

Class and Property Restrictions

- Define class similarity function
- Solve corresponding hospital-resident problem
- Based on extension of stable marriage problem





HOBBIT AKSW

Class and Property Restrictions

- Define class similarity function
- Solve corresponding hospital-resident problem
- Based on extension of stable marriage problem



May 25, 2017 9 / 68



HOBBIT AKSW

Class and Property Restrictions

- Define class similarity function
- Solve corresponding hospital-resident problem
- Based on extension of stable marriage problem



LD TUTORIAL:

Accu

May 25, 2017 10 / 68



HOBBIT AKSW

Class and Property Restrictions

- Define class similarity function
- Solve corresponding hospital-resident problem
- Based on extension of stable marriage problem



LD TUTORIAL:

ACCURA

May 25, 2017 11 / 68

$\overline{R}AVEN$

Approach

🚓 HOBBIT 🗛 📈

Class Restrictions

- Similarity function
 - String similarity
 - Number of shared property values amongst instances

TUTORIAL:

Accur

- . . .
- Solve corresponding hospital-resident problem

INFAI & FORTH

Approach

HOBBIT AKSW

Class Restrictions

- Similarity function
 - String similarity
 - Number of shared property values amongst instances
 - ...
- Solve corresponding hospital-resident problem

Source	Target	5	Т
Drugbank	Disesome	Targets	Genes
Sider	Diseasome	Side-Effect	Diseases
DBpedia	Dailymed	Organization	Organization
Sider	Dailymed	Drugs	Offer
Drugbank	DBpedia	Targets	Protein

- Property mapping similar
- Leads to $\sigma_1 \dots \sigma_n$

HAMED AHMED SHERIF ET AL. (INFAI & FORTH)

Approach

🚓 HOBBIT 🗛 🖌

Learning Threshold

- Active perceptron learning
- Begin with educated guess, e.g., $\theta_i = 0.9$

AI & FORTH

• Update thresholds based on most informative examples

May 25, 2017 13 / 68

$\overline{R}AVEN$

Approach

🚓 HOBBIT 🗛 🖌

Learning Threshold

- Active perceptron learning
- Begin with educated guess, e.g., $\theta_i = 0.9$
- Update thresholds based on most informative examples



• Guess initial classifier

Mohamed Ahmed Sherif et al. (InfAI & FORTH

May 25, 2017 13 / 68

Approach

HOBBIT AKSW

Learning Threshold

HAMED AHMED SHERIF ET AL

- Active perceptron learning
- Begin with educated guess, e.g., $\theta_i = 0.9$

(INFAI & FORTH

• Update thresholds based on most informative examples



• Pick most informative examples, i.e., unclassified and closest to boundary

Approach

🚓 HOBBIT 🗛 🖌

Learning Threshold

- Active perceptron learning
- Begin with educated guess, e.g., $\theta_i = 0.9$
- Update thresholds based on most informative examples



• Ask for classification from oracle

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

Approach

🚓 HOBBIT AKSW

Learning Threshold

- Active perceptron learning
- Begin with educated guess, e.g., $\theta_i = 0.9$

(INFAI &

FORTH

• Update thresholds based on most informative examples



• Update classifier

May 25, 2017 16 / 68



- Evaluation on Diseases (Diseasome to DBpedia)
- Learning rate = 0.02

Evaluation

- 10 questions/iteration
- F-measure of up to 92%

(INFAI & FORTH

TUTORIAL:

Accu



RAVENEvaluation

Mohamed Ahmed Sheri



Learning rate = 0.02
10 questions/iteration

(INFAI & FORTH

LD TUTORIAL: ACCURAC



Table of Contents

Mohamed Ahmed Sherif

(INFAI & FORTH





LD TUTORIAL: ACCURAC

May 25, 2017

EAGLE Efficient Active Learning of Link Specifications using Genetic Programming





LD TUTORIAL: ACCURAC

Mohamed Ahmed Sherif et al. (InfAI & FORTH

May 25, 2017 20 / 68

EAGLE

Efficient Active Learning of Link Specifications using Genetic Programming

BOBBIT AKSW

EAGLE

- Provides means for automatic class and property matching
- Minimizes human labeling effort through active learning
- Allow for learning generic specs (limitation of RAVEN)
- Similar approaches [NIK+12; ISE+12]





• Same formal setting as RAVEN

(INFAI & FORTH)

- Two sets of restrictions resp. that specify the sets S resp. T,
- a specification of mapping properties $(p_1, q_1), \ldots, (p_n, q_n)$ for the elements of S and T and
- a specification of a complex similarity measure σ as the combination of several atomic similarity measures $\sigma_1, \ldots, \sigma_n$ and of a set of thresholds $\theta_1, \ldots, \theta_n$ such that θ_i is the threshold for σ_i .







LD TUTORIAL: ACCURAC

(INFAI & FORTH

Ahmed Sheri

May 25, 2017 23 / 68

EAGLE

Idea & Goal

HOBBIT AKSW

EAGLE

- Idea: Specifications are trees
- Goal: Learn elements of trees through genetic operations until best LS is found



EAGLE Algorithm Step 1: Generate initial population



- Random process (property pairs, thresholds)
- Compute fitness
- Fitness = F-Measure w.r.t known data



EAGLE Algorithm Step 2: Evolve population

Mohamed



MAY

- Tournament between two individuals
- Two operators: Mutation and crossover



Eagle Algorithm

Step 2: Evolve population

Mohamed



- Tournament between two individuals
- Two operators: Mutation and crossover



May 25, 2017 26 / 68

EAGLE Algorithm

Step 2: Evolve population

Mohamed

HOBBIT AKSW

- Tournament between two individuals
- Two operators: Mutation and crossover



May 25, 2017 26 / 68

Eagle Algorithm

Step 2: Evolve population

Mohamed

HOBBIT AKSW

- Tournament between two individuals
- Two operators: Mutation and crossover



May 25, 2017 26 / 68

Eagle Algorithm

Step 2: Evolve population

Mohamed

HOBBIT AKSW

- Tournament between two individuals
- Two operators: Mutation and crossover







- Previous approaches define amount of information of link as closeness to the decision boundary
- Here, use disagreement amongst elements of population of size n

$$\delta((s,t)) = (n - |\mathcal{M}_i^t:(s,t) \in \mathcal{M}_i)|)(n - |\mathcal{M}_i^t:(s,t) \notin \mathcal{M}_i|)$$

- Function is maximal when $\frac{n}{2}$ count (s,t) as positive and $\frac{n}{2}$ as negative
- Can be modeled with other functions such as entropy

EAGLE Algorithm Step 4: Active Learning

BOBBIT AKSW

- Compute $\delta((s, t))$ for all (s, t) returned by a LS
- Pick k most informative
- Require labeling from user
- Update list of positive and negative examples



May 25, 2017 28 / 68

EAGLE Algorithm Step 5: Remove least fit elements

Mohamed Ahme



• Fitness = F-Measure w.r.t known data

(InfAI & FORTH



TUTORIAL:

LD

Accura

May 25, 2017 29 / 68

EAGLE Algorithm Step 5: Remove least fit elements

Mohamed



• Fitness = F-Measure w.r.t known data

(InfAI & FORTH



TUTORIAL:

LD

ACCUR.

May 25, 2017 29 / 68

EAGLE Algorithm Step 5: Remove least fit elements



• Fitness = F-Measure w.r.t known data



• If termination conditions not met, goto Step 2

(INFAI & FORTH

• Else terminate and pick fittest LS

Mohamed Ahmed Sherif et al.

EAGLE Algorithm

BOBBIT AKSW

- Measure degree of monogamy of links [NIK+12]
- Only works for 1-1 relations, e.g., owl:sameAs

$$\mathcal{P}(M) = rac{|\{s| \exists t : (s,t) \in M\}|}{\sum\limits_{s} |\{t : (s,t) \in M\}|},$$
 $\mathcal{R}(M) = rac{|\{t| \exists s : (s,t) \in M\}|}{\sum\limits_{t} |\{s : (s,t) \in M\}|},$
 $\mathcal{F}^{\beta}(M) = (1+\beta^2) rac{\mathcal{P}_d(M)\mathcal{R}_d(M)}{\beta^2 \mathcal{P}_d(M) + \mathcal{R}_d(M)}$





Experiments and Results

EAGLE

Experimental Setup:

- Compared batch learning and genetic programming
- Used 3 different data sets
 - Dailymed-Drugbank (LATC)
 - OBpedia-LinkedMDB (LATC)
 - OBLP-ACM
- Compared different sizes of population (20,100)
- Compared random annotation with active learning
- Mutation and crossover rates = 0.6
- Maximal number of iterations = 50

(INFAI & FORTH
EAGLE Experiments and Results (Dailymed-Drugbank)



(INFAI & Mohamei

EAGLE Experiments and Results (DBpedia-LinkedMDB)



Mohamed Ahmed Sherif et al. (InfAI & FORTH

D TUTORIAL: ACCURACY

MAY 25, 2017 33 / 68

EAGLE Experiments and Results (DBLP-ACM)



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

AY 25, 2017 34 / 68



- EAGLE Experiments and Results
 - Larger population leads to
 - Better results, yet
 - Longer runtimes
 - For most datasets, population size of 100 seems sufficient for most linked data sets
 - EAGLE is more time-efficient than state of the art
 - 337s for ACM-DBLP (n=100) vs.

(INFAI & FORTH)

- 1553s for Marlin (ADTree)
- 2196s for Marlin (SVM)
- 4320s for Febrl (SVM)
- Active learning clearly outperforms random annotation

Table of Contents





Mohamed Ahmed Sherif et al. (InfAI & FORTH) LD Tutorial: Accuracy

May 25, 2017 36 / 68

COALA Correlation-Aware Active Learning of Link Specifications

(INFAI & FORTH)

Mohamed Ahmed Sherif





LD TUTORIAL: ACCURAC

May 25, 2017 37 / 68





- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH

LD TUTORIAL:

ACCURA



May 25, 2017 38 / 68





- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH)



LD TUTORIAL: ACCURAC

May 25, 2017 38 / 68





2017

- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH)



LD TUTORIAL: ACCURAC





- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH

LD TUTORIAL:

ACCURA



May 25, 2017 39 / 68





- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH)



LD TUTORIAL: ACCURAC

May 25, 2017 39 / 68





May 25, 2017

- Supervised (mostly active, e.g., RAVEN, EAGLE, SILK)
- Unsupervised (e.g., KnoFuss, EUCLID, EAGLE)

(INFAI & FORTH)

Mohamed



LD TUTORIAL: ACCURAC

COALA Learning Complex Specifications



Insight

Mohamed Ahmed

• Choice of right example is key for learning

(INFAI & FORTH)

LD TUTORIAL: ACCURAC

• So far, only use of informativeness

May 25, 2017 40 / 68

COALA Learning Complex Specifications



Insight

- Choice of right example is key for learning
- So far, only use of informativeness

Question

• Can we do better by using more information?

(INFAI & FORTH

Accu

May 25, 2017 40 / 68



COALA

Insight

- Choice of right example is key for learning
- So far, only use of informativeness

Question

• Can we do better by using more information?

(INFAI & FORTH

• Higher F-measure

Learning Complex Specifications

• Often slower

Coala Approach

Basic Idea



• Use similarity of link candidates when selecting most informative examples

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: ACCURACY

May 25, 2017 41 / 68

COALA Approach

Basic Idea



• Use similarity of link candidates when selecting most informative examples



May 25, 2017 41 / 68

COALA Approach Basic Idea



• Use similarity of link candidates when selecting most informative examples





- COALA Similarity of Candidates
 - Link candidate x = (s, t) can be regarded as vector $(\sigma_1(x), \ldots, \sigma_n(x)) \in [0, 1]^n$.
 - Similarity of link candidates x and y:

$$sim(x, y) = \frac{1}{1 + \sqrt{\sum_{i=1}^{n} (\sigma_i(x) - \sigma_i(y))^2}}.$$
(1)

• Allows exploiting both intra- and inter-class similarity

(INFAI & FORTH

- Rationale: Use intra-class similarity
- Approach

COALA

Graph Clustering

- Cluster elements of S^+ and S^- independently
- Choose one element per cluster as representative
- Present oracle with most informative representatives

(INFAI & FORTH

LD TUTORIAL: ACCURA



COALA BorderFlow



• $G = (V, E, \omega)$ with $V = S^+$ or $V = S^-$

(INFAI & FORTH

TUTORIAL:

ACCURA

- $\omega(x,y) = sim(x,y)$
- Keep best ec edges for each $x \in V$







May 25, 2017

- Seed-based algorithm
- Goal: Maximize borderflow ratio $bf(X) = \frac{\Omega(b(X),X)}{\Omega(b(X),n(X))}$

(InfAI & FORTH

LD TUTORIAL:

ACCURA





- Seed-based algorithm
- Goal: Maximize borderflow ratio $bf(X) = \frac{\Omega(b(X), X)}{\Omega(b(X), n(X))}$



May 25, 2017 45 / 68



🚓 HOBBIT (AKSW)

• Seed-based algorithm

Mohamed Ahmed Sherif et al.

• Goal: Maximize borderflow ratio $bf(X) = \frac{\Omega(b(X),X)}{\Omega(b(X),n(X))}$



LD

http://sourceforge.net/projects/cugar-framework/

(INFAI & FORTH)

May 25, 2017 45 / 68





- Seed-based algorithm
- Goal: Maximize borderflow ratio $bf(X) = \frac{\Omega(b(X),X)}{\Omega(b(X),n(X))}$

FORTH



TUTORIAL

May 25, 2017 46 / 68

COALA BorderFlow

HOBBIT AKSW

- Seed-based algorithm
- Goal: Maximize borderflow ratio $bf(X) = \frac{\Omega(b(X),X)}{\Omega(b(X),n(X))}$



http://sourceforge.net/projects/cugar-framework/ Mohamed Ahmed Sherif et al. (INFAI & FORTH) LD TUTORIAL: ACCURACY

May 25, 2017 46 / 68



• Rationale: Use both inter- and intra-class similarity

(INFAI & FORTH

TUTORIAL:

Accura

• Approach

Spreading Activation

COALA

- $\mathcal{M}_0: m_{ij} = sim(x_i, x_j)$ with $(x_i, x_j) \in (S^+ \cup S^-)^2$
- \mathcal{A}_0 : $a_i = ifm(x_i)$

May 25, 2017 47 / 68



- COALA Spreading Activation
 - Rationale: Use both inter- and intra-class similarity
 - Approach
 - $\mathcal{M}_0: m_{ij} = sim(x_i, x_j)$ with $(x_i, x_j) \in (S^+ \cup S^-)^2$
 - \mathcal{A}_0 : $a_i = ifm(x_i)$
 - $\mathcal{A}_t = \mathcal{A}_{t-1} + \mathcal{M}_{t-1}\mathcal{A}_{t-1}$ (spread activation)
 - $\mathcal{A}_t = \mathcal{A}_t / \max(\mathcal{A}_t)$ (normalize) $\mathcal{M}_t = \mathcal{M}_{t-1}^{(\mathbf{f})}$ (weight decay)

(INFAI & FORTH

TUTORIAL:

ACCUR



- COALA Spreading Activation
 - Rationale: Use both inter- and intra-class similarity
 - Approach
 - $\mathcal{M}_0: m_{ij} = sim(x_i, x_j)$ with $(x_i, x_j) \in (S^+ \cup S^-)^2$
 - \mathcal{A}_0 : $a_i = ifm(x_i)$
 - $\mathcal{A}_t = \mathcal{A}_{t-1} + \mathcal{M}_{t-1}\mathcal{A}_{t-1}$ (spread activation)
 - $\mathcal{A}_t = \mathcal{A}_t / \max(\mathcal{A}_t)$ (normalize) $\mathcal{M}_t = \mathcal{M}_{t-1}^{(\mathbf{f})}$ (weight decay)







- Used EAGLE as active learning approach
 - Mutation and crossover rate = 0.6
 - Selection rate = 0.7
 - Not deterministic \Rightarrow Ran each experiment 5 times
 - 5 queries to oracle per iteration
 - 10 iterations overall
 - 2 populations sizes: 20 and 100
 - 50 generations between iterations
- Two real-world and three synthetic datasets

(INFAI & FORTH

• Single thread of a server (JDK1.7, Ubuntu 10.0.4, AMD Opteron 2GHz, 2GB/Experiment)

COALA Evaluation Parameters for WD



- Ran experiments on DBLP-ACM
- Population = 20
- $r \in \{2, 4, 8, 16, 32\}$



May 25, 2017 49 / 68

COALA Evaluation Parameters for CL_____



- Ran experiments on DBLP-ACM
- Population = 20
- $ec \in \{1, 2, 3, 4, 5\}$



May 25, 2017 50 / 68





- Population = 100, final values
- Better results, yet unclear when to use WD or CL

(INFAI & FORTH

DataSet	EAGLE	WD	CL
Abt	0.19±0.04	0.25±0.04	0.23±0.04
DBLP	$0.91{\pm}0.03$	$0.96{\pm}0.01$	$0.96{\pm}0.02$
Person1	$0.86{\pm}0.02$	$0.89{\pm}0.01$	$0.81{\pm}0.18$
Person2	$0.74{\pm}0.03$	$0.71{\pm}0.08$	0.77±0.03
Restaurant	$\textbf{0.89}{\pm 0.0}$	$0.86{\pm}0.02$	$\textbf{0.89}{\pm 0.0}$

LD TUTORIAL: ACCURACY

Table of Contents

Mohamed Ahmed Sherif

(INFAI & FORTH





LD TUTORIAL: ACCURAC

May 25, 2017 52 / 68



- WOMBAT Motivation
- - Need for automatic LD for evolving datasets
 - Mostly positive examples on the Web of Data
 - Negative examples rarely to be found
 - Missing links cannot be regarded as negative examples (Open World Assumption)



Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL: ACCURACY

May 25, 2017 53 / 68

WOMBAT

Mohamed Ahmed Sherif et

(INFAI & FORTH

Idea





LD TUTORIAL: ACCURAC

May 25, 2017 54 / 68


TORIAL: ACCURA

(INFAI & FORTH



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

) Tutorial: Accuracy

May 25, 2017 56 / 68

Simple Operator

HOBBIT AKSW



Complete Operator

HOBBIT AKSW



- uses a more sophisticated expansion strategy
- allows learning arbitrarily nested LS
- is an upward complete refinement operator

(INFAI & FORTH

• is improved using pruning

 $\psi(L) = \begin{cases} \{A_{i_1} \setminus A_{j_1} \sqcap \cdots \sqcap A_{i_m} \setminus A_{j_m} | A_{i_k}, A_{j_k} \in \mathbf{A} \\ \text{for all } 1 \leq k \leq m \} \\ \{L \sqcup A_i \setminus A_j | A_i \in \mathbf{A}, A_j \in \mathbf{A} \} \\ \{L_1 \} \cup \{L \sqcup A_i \setminus A_j | A_i \in \mathbf{A}, A_j \in \mathbf{A} \} \\ \{L_1 \sqcap \cdots \sqcap L_{i-1} \sqcap L' \sqcap L_{i+1} \sqcap \cdots \sqcap L_n | L' \in \psi(L_i) \} \\ \cup \{L \sqcup A_i \setminus A_j | A_i \in \mathbf{A}, A_j \in \mathbf{A} \} \\ \{L_1 \sqcup \cdots \sqcup L_{i-1} \sqcup L' \sqcup L_{i+1} \sqcup \cdots \sqcup L_n | L' \in \psi(L_i) \} \\ \cup \{L \sqcup A_i \setminus A_j | A_i \in \mathbf{A}, A_j \in \mathbf{A} \} \end{cases}$

if $L = \bot$ if $L \in A$ if $L = L_1 \setminus L_2$ if $L = L_1 \Box \cdots \Box L_n (n \ge 2)$ if $L = L_1 \sqcup \cdots \sqcup L_n (n \ge 2)$



May 25, 2017 58 / 68

Pruning Procedure

AKSW

ψ is an upward refinement operator

The set of links generated by a child node is a superset of or equal to the set of links generated by its parent

- r_{max} is bounded by the most general constructable LS
- p_{max} is bounded as false positives cannot disappear during generalisation

•
$$F_{max} = \frac{2p_{max}r_{max}}{p_{max}+r_{max}}$$

• Prune all nodes in the search tree with $F_{max} < F_{best}$



Pruning Procedure

BOBBIT AKSW

ψ is an upward refinement operator

The set of links generated by a child node is a superset of or equal to the set of links generated by its parent

- r_{max} is bounded by the most general constructable LS
- p_{max} is bounded as false positives cannot disappear during generalisation

•
$$F_{max} = \frac{2p_{max}r_{max}}{p_{max}+r_{max}}$$

• Prune all nodes in the search tree with $F_{max} < F_{best}$



Evaluation Setup

AKSW

- 8 benchmark datasets (5 real-world, 3 synthetic)
- 2.80 GHz PC running OpenJDK 64-Bit Server 1.7.0_75 on Ubuntu 14.04.2 LTS
- 7 GB RAM
- Wombat
 - Similarity measures: jaccard, trigrams, cosine and qgrams
 - Termination: F = 1 or max number of refinement tree depth of 10

LD TUTORIAL: ACCURAC

• Properties coverage threshold = 0.6

(INFAI & FORTH)



BOBBIT AKSW

WOMBAT Evaluation: 10-Fold Cross Validation F-Measure

Dataset	Wombat Simple	WOMBAT Complete	Euclid Linear	EUCLID Conjunction	EUCLID Disjunction	EAGLE
Person 1	1.00	1.00	0.64	0.97	1.00	0.99 ± 0.004
Person 2	1.00	0.99	0.22	0.78	0.96	0.94 ± 0.032
Restaurants	0.98	0.97	0.97	0.97	0.97	0.97 ± 0.024
DBLP-ACM	0.97	0.98	0.98	0.98	0.98	0.98 ± 0.007
Abt-Buy	0.60	0.61	0.06	0.06	0.52	0.65 ± 0.025
Amazon-GP	0.70	0.67	0.59	0.71	0.73	0.71 ± 0.033
DBP-LMDB	0.99	1.00	0.99	0.99	0.99	0.99 ± 0.004
DBLP-GS	0.94	0.94	0.90	0.91	0.91	$\textbf{0.93} \pm \textbf{0.006}$
Average	0.90	0.90	0.67	0.80	0.88	$\textbf{0.90} \pm \textbf{0.017}$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

May 25, 2017 61 / 68

BOBBIT AKSW

WOMBAT <u>Evaluation:</u> Amazon – Google Products Dataset



WOMBAT Evaluation: Pruning Procedure

HOBBIT AKSW

Dataset	10%	20%	30%	40%	50%	60%	70%	80%
Person 1	1.57	2.13	1.85	2.13	2.13	2.13	2.13	2.13
Person 2	1.29	1.29	1.57	1.57	1.57	1.57	1.57	1.57
Restaurant	1.17	1.45	1.17	1.45	1.45	1.45	1.45	1.45
DBLP-ACM	6.23	5.58	6.79	6.85	6.85	6.85	6.79	6.79
Abt-Buy	3.38	3.00	3.00	3.39	3.39	3.39	1.79	3.39
Amazon-GP	1.14	1.38	1.33	1.37	1.38	1.45	1.54	1.59
DBP-LMDB	1.00	1.86	2.86	1.86	1.86	2.33	2.36	2.36
DBLP-GS	1.79	1.93	2.01	2.36	2.45	1.66	2.44	2.26

 $\begin{array}{l} \mbox{Pruning factor:} & \mbox{number of searched nodes (search tree size + pruned nodes)} \\ \hline \mbox{Max. size of the search tree (2000 nodes in this set of experiments)} \end{array}$

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

May 25, 2017 63 / 68

WOMBAT Evaluation: Training with only 2%

Mohamed Ahmed Sherif et al.

(INFAI & FORTH)



Dataset	Pessimistic	Re-weighted	Simple	Complete
Persons 1	1.00	1.00	1.00	1.00
Persons 2	0.97	1.00	0.80	0.84
Restaurants	0.95	0.94	0.98	0.88
DBLP-ACM	0.93	0.95	0.94	0.94
Amazon-GP	0.39	0.43	0.53	0.45
Abt-Buy	0.36	0.37	0.37	0.36
Average	0.77	0.78	0.77	0.74

LD TUTORIAL: ACCURACY



HOBBIT AKSW

- Conclusion
 - \bullet Presented $\mathrm{WOMBAT},$ the first approach to learn LS from positive examples
 - $\bullet~\mathrm{WOMBAT}$ is based on generalisation over the space of LS
 - Presented 2 operators to achieve this goal
 - Evaluated WOMBAT against SOTA
 - $\bullet~{\rm WOMBAT}$ outperforms SOTA by 11% on average
- Future work
 - Parallelize WOMBAT
 - Try more aggressive pruning techniques for better scalability
 - Apply active learning strategies

(INFAI & FORTH)

• Unsupervised WOMBAT



Table of Contents





Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: ACCURACY

 ${\rm May}\ 25,\ 2017 \qquad 66\ /\ 68$

Summary and Conclusion

🚓 HOBBIT 🗛 📕

May 25, 2017



• Large number of challenges to learning accurate specifications

- Reduce labeling effort
 ⇒ Active learning
- Learn complex specifications
 Genetic programming
- Is the second se
- Challenges include
 - Determinism
 - Oeep learning
 - Self-checking
 - 4 . . .

CHALLENGE YOURSELF AND HAVE FUN!

Tutorial: Accurac

A cknowledgment

(INFAI & FORTH

HOBBIT AKSW



This work was supported by grants from the EU H2020 Framework Programme provided for the project HOBBIT (GA no. 688227).

Link Discovery Tutorial Part IV: Benchmarking

Mohamed Ahmed Sherif⁽¹⁾, Irini Fundulaki⁽²⁾, Axel-Cyrille Ngonga Ngomo⁽¹⁾

(1) Institute for Applied Informatics, Germany (2) FORTH, Greece



May 28th, 2017 Portoroz, Slovenia.

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 24, 2017 1 / 36

The Questions(s)

HOBBIT AKSW

Instance matching research has led to the development of various systems.

- What are the problems that I wish to solve?
- What are the relevant key performance indicators?
- What is the behavior of the existing engines w.r.t. the key performance indicators?



Which are the tool(s) that I should use for my data and for my use case?

D Tutorial:Intro

May 24, 2017 2 / 36

Importance of Benchmarking

BOBBIT AKSW

Benchmarks exist

- To allow adequate measurements of systems
- To provide evaluation of engines for real (or close to real) use cases

Provide help

- Designers and Developers to assess the performance of their tools
- Users to compare the different available tools and evaluate suitability for their needs
- Researchers to compare their work to others

Leads to improvements:

- Vendors can improve their technology
- Researchers can address new challenges
- Current benchmark design can be improved to cover new necessities and application domains

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D Tutorial:Intro

May 24, 2017 3 / 36



The Answer Benchmark your engines!

Instance Matching Benchmark comprises of

(INFAI & FORTH

- **Datasets:** The raw material of the benchmarks. These are the source and the target dataset that will be matched together to find the links
- Test Cases: Address heterogeneities (structural, value, semantic) of the datasets to be matched
- Gold Standard (Ground Truth / Reference Alignment): The "correct answer sheet" used to judge the completeness and soundness of the instance matching algorithms
- Metrics: The performance metric(s) that determine the systems behaviour and performance

Benchmark Datasets Characteristics

BOBBIT AKSW

Nature

- Real Datasets: Widely used datasets from a domain of interest
 - Realistic conditions for heterogeneity problems
 - Realistic distributions
 - Error prone, hard to create Reference Alignment
- **Synthetic Datasets:** Produced with a data generator (that hopefully produces data with interesting characteristics
 - Fully controlled test conditions
 - Accurate, Easy to create Reference Alignments
 - Unrealistic distributions
 - Systematic heterogeneity problems
- Schema
 - Datasets to be matched have the same or different schemas
- Domain
 - Datasets come from the same or different domains

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 5 / 36

Benchmark Test Cases
Variations

HOBBIT AKSW

• Value

- Name style abbreviations, Typographical errors, change format (date/gender/number), synonym change, language change (multilinguality)
- Structural
 - Change property depth, Delete/add property, split property values, transformation of object/data to data/object type property
- Semantics
 - class deletion/modification, invert property assertions, change class/property hierarchy, assert class disjointness
- Combinations of Variations

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 6 / 36





- The "correct answer sheet" used to judge the completeness and soundness of the instance matching algorithms
- Characteristics
 - Existence of errors / missing alignments

(INFAI & FORTH

• Representation: owl:sameAs and skos:exactMatch

LD TUTORIAL:INTRO

May 24, 2017 7 / 36

Metrics

BOBBIT AKSW

- Precision $P = \frac{tp}{(tp+fn)}$ Recall $R = \frac{tp}{(tp+fp)}$ F-measure $F = 2 \times P \times \frac{R}{(P+R)}$



Instance Matching Benchmarks Desirable Attributes



Systematic Procedure	matching tasks should be reproducible and the execution must
	be comparable
Availability	benchmark should be available
Quality	precise evaluation rules and high quality ontologies must be pro- vided
Equity	evaluation process should not privilege any system
Dissemination	benchmark should be used to evaluate instance matching systems
Volume	dataset size
Gold Standard	gold standard should exist and be as accurate as possible

(INFAI & FORTH

LD TUTORIAL:INTRO

May 24, 2017

What about Benchmarks?

(INFAI & FORTH)

HOBBIT AKSW

Instance matching techniques have, until recently, been benchmarked in an ad-hoc way.

There is no standard way of benchmarking the performance of the systems, when it comes to Linked Data.

LD TUTORIAL:INTRO

May 24, 2017 10 / 36





- IM benchmarks have been mainly driven forward by the **Ontology Alignment Evaluation** Initiative (OAEI)
 - organizes annual campaign for ontology matching since 2005
 - hosts independent benchmarks

(INFAI & FORTH)

- In 2009, OAEI introduced the Instance Matching (IM) Track
 - focuses on the evaluation of different instance matching techniques and tools for Linked Data

LD TUTORIAL:INTRO

May 24, 2017 11 / 36

Instance Matching Benchmarks

(INFAI & FORTH



- Bechmark Generators
- Synthetic Benchmarks
- Real Benchmarks

LD TUTORIAL:INTRO

May 24, 2017 12 / 36

SWING [FMN+11]

S HOBBIT AKSW



Lance [SDF+15b]

(INFAI & FORTH)

HOBBIT AKSW

Flexible, generic and domain-independent benchmark generator which takes into consideration RDFS and OWL constructs in order to evaluate instance matching systems.

LD TUTORIAL:INTRO

May 24, 2017 14 / 36

Lance [SDF+15b]

HOBBIT AKSW

Lance provides support for:

- Semantics-aware transformations
 - Complex class definitions (union, intersection)
 - Complex property definitions (functional properties, inverse functional properties)
 - Disjointness (properties)
- Standard value and structure based transformations
- Weighted gold standard based on tensor factorization
- Varying degrees of difficulty and fine-grained evaluation metrics

Available at http://github.com/jsaveta/Lance

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D Tutorial:Intro

May 24, 2017 15 / 36

Lance Architecture

HOBBIT AKSW



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 16 / 36

Synthetic Benchmarks

HOBBIT AKSW





TUTORIAL:INTRO

May 24, 2017 17 / 36

Synthetic IM Benchmarks $O_{Verview (1)}$



	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015
Systematic Procedure	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Availability	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	-	-	\checkmark	\checkmark
Quality	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Equity	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Dissemination	6	3	6	1	3	4	4	5	5
Volume	0.2K	1.4K	0.86K	4K	0.375K	1.5K	0.43K	2.650K	10K
Gold Standard	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 24, 2017 18 / 36

Synthetic IM Benchmarks Overview (2)

🔊 HOBBIT AKSW

	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015
Value Variations	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Structural Variations	\checkmark	\checkmark	\checkmark	\checkmark	-	-	-	+	+
Logical Variations	\checkmark	\checkmark	-	\checkmark	-	\checkmark	-	-	-
Multilinguality	-	-	-	-	-	-	\checkmark	\checkmark	\checkmark

	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015
Blind Evaluations	-	-	-	-	-	-	\checkmark	\checkmark	\checkmark
1-n Mappings	-	-	\checkmark	-	-	-	\checkmark	\checkmark	-

(INFAI & FORTH Mohamed

LD TUTORIAL:INTRO

May 24, 2017 19

Synthetic IM Benchmarks _{Overview (3)}

BOBBIT AKSW

	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015	Lance 2015
Systematic Procedure	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Availability	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	-	-	\checkmark	\checkmark	\checkmark
Quality	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Equity	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Dissemination	6	3	6	1	3	4	4	5	5	2
Volume	0.2K	1.4K	0.86K	4K	0.375K	1.5K	0.43K	2.650K	10K	> 1 <i>M</i>
Gold Standard	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 20 / 36

Synthetic IM Benchmarks _{Overview (4)}



	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015	Lance 2015
Value Variations	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Structural Variations	\checkmark	\checkmark	\checkmark	\checkmark	-	-	-	+	+	+
Logical Variations	\checkmark	\checkmark	-	\checkmark	-	\checkmark	-	-	-	+
Multilinguality	-	-	-	-	-	-	\checkmark	\checkmark	\checkmark	\checkmark

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 24, 2017 21 / 36
Synthetic IM Benchmarks _{Overview (5)}



	IIMB 2009	IIMB 2010	PR 2010	IIMB 2011	Sandbox 2012	IIMB 2012	RDFT 2013	ID-REC 2014	Author Task 2015	Lance 2015
Blind Evaluations	-	-	-	-	-	-	\checkmark	\checkmark	\checkmark	\checkmark
1-n Mappings	-	-	\checkmark	-	-	-	\checkmark	\checkmark	-	-

Mohamed Ahmed Sherif et al. (InfAI & FORTH

LD TUTORIAL:INTRO

May 24, 2017 22 / 36

Real Benchmarks

HOBBIT AKSW

ARS

(OAEI 2009)

Data Interlinking

(OAEI 2010)

Data Interlinking – New York Times

(OAEI 2011)

May 24, 2017 23 / 36

Real IM Benchmarks Overview (1)



	ARS	DI 2010	DI 2011
Systematic Procedure	\checkmark	\checkmark	\checkmark
Availability	\checkmark	\checkmark	-
Quality	\checkmark	\checkmark	\checkmark
Equity	\checkmark	\checkmark	\checkmark
Dissemination	5	2	3
Volume	100K	6K	NA
Gold Standard	\checkmark	\checkmark	+

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 24 / 36

Real IM Benchmarks Overview (2)



	ARS	DI 2010	DI 2011
Value Variations	\checkmark	\checkmark	\checkmark
Structural Variations	\checkmark	\checkmark	-
Logical Variations	-	-	-
Multilinguality	-	-	-
Blind Evaluations	-	-	-

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 25 / 36



S HOBBIT AKSW



HOBBIT AKSW



Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 27 / 36

HOBBIT AKSW



HOBBIT AKSW

Logical Variations



LD TUTORIAL:INTRO

May 24, 2017 29 / 36

HOBBIT AKSW





Mohamed Ahmed

(INFAI & FORTH)

HOBBIT AKSW

Scalability & Expressiveness

LD TUTORIAL:INTRO

May 24, 2017 32 / 36



- Only one benchmark that tackles both, combination of variations and scalability issues
- $\bullet\,$ Not enough IM benchmark using the full expressiveness of RDF/OWL language

D TUTORIAL:INTRO

May 24, 2017 33 / 36

HOBBIT AKSW

- Systems can handle the value variations, the structural variation, and the simple logical variations separately.
- More work needed for complex variations (combination of value, structural, and logical)
- More work needed for structural variations
- Enhancement of systems to cope with the clustering of the mappings (1-n mappings)

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D Tutorial:Intro

May 24, 2017 34 / 36



- Many instance matching benchmarks have been proposed
- Each of them answering to some of the needs of instance matching systems.
- It is essential to start creating benchmarks that will "show the way to the future"
- Extend the limits of existing systems.

(INFAI & FORTH)

LD TUTORIAL:INTRO

May 24, 2017 35 / 36

A cknowledgment

(INFAI & FORTH

HOBBIT AKSW

36



This work was supported by grants from the EU H2020 Framework Programme provided for the project HOBBIT (GA no. 688227).

LD TUTORIAL:INTRO

Link Discovery Tutorial Part V: Hands-On

Mohamed Ahmed Sherif⁽¹⁾, Irini Fundulaki⁽²⁾, Axel-Cyrille Ngonga Ngomo⁽¹⁾

(1) Institute for Applied Informatics, Germany (2) FORTH, Greece



May 28th, 2017 Portoroz, Slovenia.

LD TUTORIAL: HANDS-ON

Iohamed Ahmed Sherif et al. (InfAI & FORTH

fay 24, 2017 1 / 19

Table of Contents

HOBBIT AKSW

- 1 Test Dataset
- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration
- 4 Task III: Use Limes GUI
- **5** Open Challenges

(INFAI & FORTH) LD TUTORIAL: HANDS-O

May 24, 2017 2 / 19

Table of Contents

HOBBIT AKSW

1 Test Dataset

- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration

(INFAI & FORTH)

- 4 Task III: Use LIMES GUI
- **6** Open Challenges

LD TUTORIAL: HANDS-ON

May 24, 2017 3 / 19

🚓 HOBBIT 🗛 🖌

Semantic Web Dog Food Corpus

Data exposed: Metadata (papers, presentations, people) for several semantic web related conferences and workshops, including the most recent ISWC, ESWC and WWW events.

- Dumps: http://data.semanticweb.org/dumps
- Endpoint: http://data.semanticweb.org/sparql
- DataHub: https://datahub.io/dataset/semantic-web-dog-food

Table of Contents

HOBBIT AKSW

1) Test Datase

- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration

(INFAI & FORTH)

- 4 Task III: Use LIMES GUI
- **6** Open Challenges

LD TUTORIAL: HANDS-ON

May 24, 2017 5 / 19

Task I: Deduplication Execute the given Configuration File



- Decentralized nature of LOD
- Data contain duplicates
- How to efficiently detect similar resources?

(INFAI & FORTH

LD TUTORIAL: HANDS-ON

May 24, 2017 6 / 19

Task I: Deduplication Execute the given Configuration File



- Decentralized nature of LOD
- Data contain duplicates
- How to efficiently detect similar resources?

(INFAI & FORTH

Example

Find duplicate authors in Semantic Web Dog Food Corpus?

LD TUTORIAL: HANDS-ON

May 24, 2017 6 / 19

LIMES Configuration File $_{PREFIXes}$



<PREFIX>
 <NAMESPACE>http://www.w3.org/2000/01/rdf-schema#</NAMESPACE>
 <LABEL>rdfs</LABEL>
<//PREFIX>

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 7 / 19

LIMES Configuration File SOURCE and TARGET datasets



<SOURCE>

<ID>SDF1</ID>
<ENDPOINT>semanticDogFood.nt</ENDPOINT>
<VAR>?x</VAR>
<PAGESIZE>-1</PAGESIZE>
<RESTRICTION>?x a foaf:Person</RESTRICTION>
<PROPERTY>rdfs:label</PROPERTY>
<TYPE>NT</TYPE>
</SOURCE>

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 8 / 19

LIMES Configuration File



<METRIC>Levenshtein(x.rdfs:label , y.rdfs:label)</METRIC>

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 9 / 19

LIMES Configuration File ACCEPTANCE and REVIEW



<ACCEPTANCE>
 <THRESHOLD>0.9</THRESHOLD>
 <FILE>similarAuthor.nt</FILE>
 <RELATION>ov:similarTo</RELATION>
</ACCEPTANCE>
<REVIEW>
 <THRESHOLD>0.5</THRESHOLD>
 <FILE>similarAuthor_review.nt</FILE>
 <RELATION>owl:sameAs</RELATION>

</REVIEW>

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 10 / 19

LIMES Configuration File

(INFAI & FORTH)



<OUTPUT>TTL</OUTPUT>

Mohamed

LD TUTORIAL: HANDS-ON

May 24, 2017 11 / 19

Task I Run Limes



Run Limes

Mohamed Ahmed Sheri

java -jar limes-core-1.0.0.jar task1.xml

(INFAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 12 / 19

Task I Using Machine Learning

Mohamed Ahmed Sherif et al



<MLALGORITHM> <NAME>wombat simple</NAME> <TYPE>unsupervised</TYPE> </MLALGORITHM>

(INFAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 13 / 19

Table of Contents

HOBBIT AKSW

1 Test Datase

- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration

(INFAI & FORTH)

- 4 Task III: Use LIMES GUI
- **6** Open Challenges

LD TUTORIAL: HANDS-ON

May 24, 2017 14 / 19



Find publications

Find similar publications

Task II

- with similar keyword,
- but do not link any publication to itself

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 15 / 19



Task II Find similar publications

Find publications

- with similar keyword,
- but do not link any publication to itself

(INFAI & FORTH)

1. Find publications

- ?x a swrc:InProceedings
- ?y a swrc:InProceedings

LD TUTORIAL: HANDS-ON

May 24, 2017 15 / 19



Task II Find similar publications

Find publications

- with similar keyword,
- but do not link any publication to itself

1. Find publications

- ?x a swrc:InProceedings
- ?y a swrc:InProceedings

2. with similar keyword, but do not link any publication to itself

MINUS(jaccard(x.swrc:listKeyword , y.swrc:listKeyword)|0.3, ExactMatch(x.swrc:listKeyword , y.swrc:listKeyword)|1.0)

Mohamed Ahmed Sherif et al. (InfAI & FORTH)

D TUTORIAL: HANDS-ON

May 24, 2017 15 / 19

Table of Contents

HOBBIT AKSW

1 Test Datase

- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration

(INFAI & FORTH)

LD TUTORIAL: HANDS-ON

4 Task III: Use Limes GUI

(5) Open Challenges

May 24, 2017 16 /

Task III Use Limes GUI



Run Limes GUI

Mohamed Ahmed Sherin

java -jar limes-core-1.0.0.jar -g

(INFAI & FORTH)

LD TUTORIAL: HANDS-ON

May 24, 2017 17 / 19

Table of Contents

HOBBIT AKSW

1 Test Datase

- 2 Task I: Execute given LIMES Configuration
- 3 Task II: Create your first LIMES Configuration

(INFAI & FORTH

- 4 Task III: Use LIMES GUI
- **5** Open Challenges

LD Tutorial: Hands-On

May 24, 2017 18 / 19
A cknowledgment

(INFAI & FORTH

HOBBIT AKSW

19



This work was supported by grants from the EU H2020 Framework Programme provided for the project HOBBIT (GA no. 688227).

TUTORIAL: HANDS-ON