

## Towards Dataset Dynamics: Change Frequency of Linked Open Data Sources

J. Umbrich, M. Hausenblas, A. Hogan, A. Polleres, S. Decker





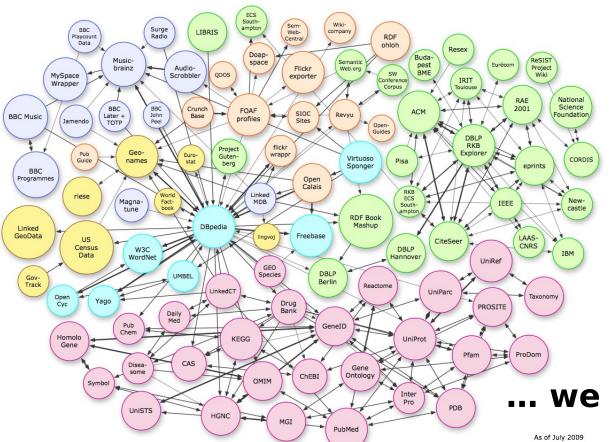
#### **Motivation (Situation)**



**Digital Enterprise Research Institute** 

www.deri.i

е



In late 2009: we have

over **100** open datasets

providing over 13.7

billion RDF triples,

interlinked by some **142** 

million links ...

we can expect/have
many changes

(new resources, added links, updates)

courtesy of Richard Cyganiak and Anja Jentzsch





#### **Motivation (Dataset Dynamics)**

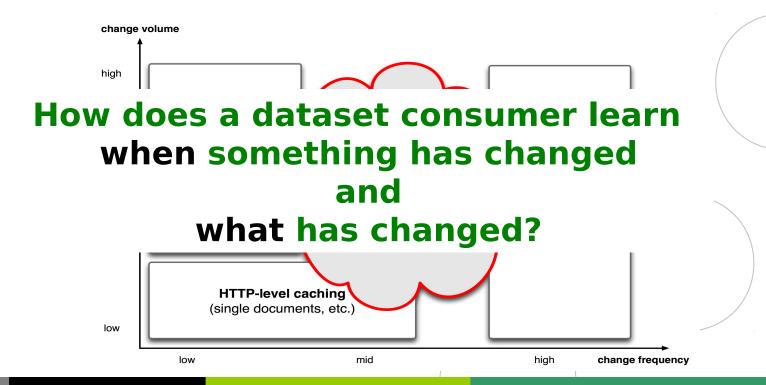


**Digital Enterprise Research Institute** 

www.deri.

-

- Outcomes/Findings/Lessons learned from
  - □ LDC09: **D**ataset **C**hange **M**anager/**W**atch-Dog demo
  - Dataset Dynamics Group







### **Motivation** (Dynamics in Web of

#### **Documents**)



**Digital Enterprise Research Institute** 

www.deri.

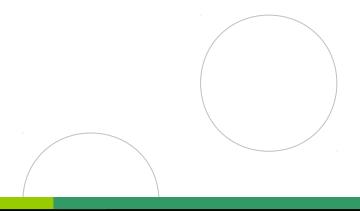
e

- Web of Documents
  - □ Change frequency of Web documents can be modeled as a Poisson Process [Cho at al. VLDB 2000]
  - □ Improvement of Estimators for uncertainty [Cho at al. ACM Journal 2001]
  - □ What is new on the Web? (search engine perspective) [Ntoulas et al., WWW2004]
- Applications
  - □ Web crawling and caching
  - ☐ Maintaining link integrity
  - □ replication and synchronisation
    - Servicing of continuous queries





## How dynamic is the Web of Data?





# Detection Level: A Matter of Granularity Digital Enterprise Research Institute

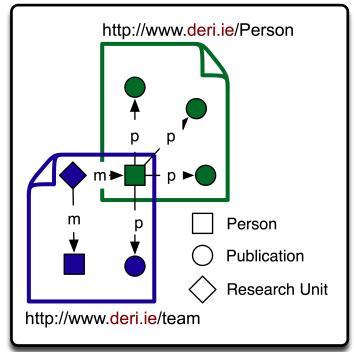


www.deri.i

e

- Document-centric
  - ☐ HTTP Get body content

- Entity-centric
  - Entity-per-document change
  - □ (Global) Entity-change
    - Reuse of entities among sources (square)
    - wrt. to specific documentse.g. per namespace, pay-level-doma...







#### **Change Detection Mechanism**



**Digital Enterprise Research Institute** 

www.deri.

е

- Content monitoring fetching the entire content
- HTTP caching (response header as of RFC 2616)
- Notification active notification

	Content	HTTP	Notificatio n
availability	+	+/-	+/-
reliability	+	+/-	unknown
costs	high	low	unknown
scalability	high	high	unknown
documents	yes	yes	yes





# Preliminary Results & Lessons learned





#### **Experimental Setup**



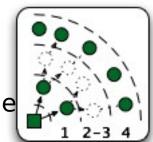
**Digital Enterprise Research Institute** 

www.deri.i

e

8.18%)

- Dataset Web data (passive monitoring)
  - □ weekly snapshots over 24 weeks (start Nov. 08)
  - 4 hop neighborhood from Tim Berners-Lee FOAF file
  - □ 550K RDF/XML docs, 3.3M unique entities



7.12%)

(16.75%)

(67.95%)

- Change detection
  - skolemise blank nodes within a d
  - □ pairwise comparison of statements by scanning sorted list

Only ETag

Both

None

Only Last-Modified

- Detection level
  - □ Document-centric
  - ☐ Entity-change (pay-level-domain)



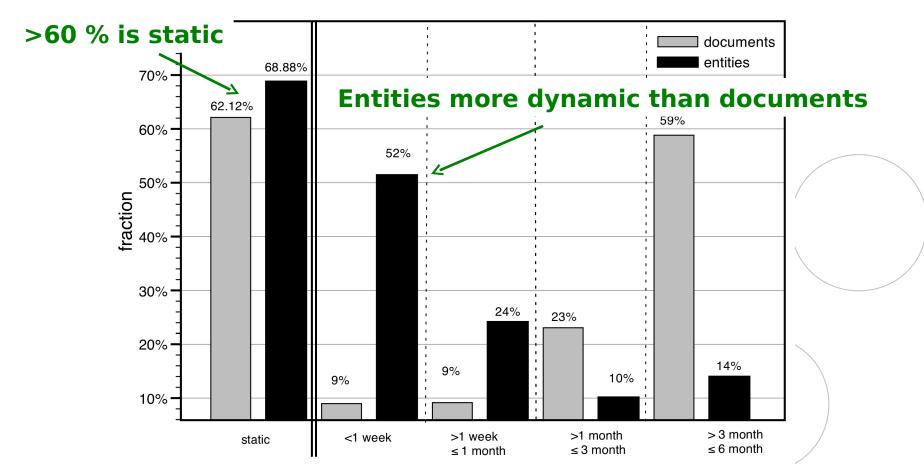
#### **Average Change Frequency**



**Digital Enterprise Research Institute** 

www.deri.

€



average change frequency





#### **Entity Changes per Document**



**Digital Enterprise Research Institute** 

	U	A	D	(UA   UD   AD)	Total
U	76.88 %	9.46 %	7.08 %	3.87 %	97.29 %
A	9.46 %	0.19 %	2.29 %	3.87 %	15.81 %
D	7.08 %	2.29 %	0.23 %	3.87 %	13.5 %

#### Legend:

 $\mathbf{U} \neq \mathsf{Update}$ 

 $\mathbf{A} = Add$ 

**D** = Delete

of an entity





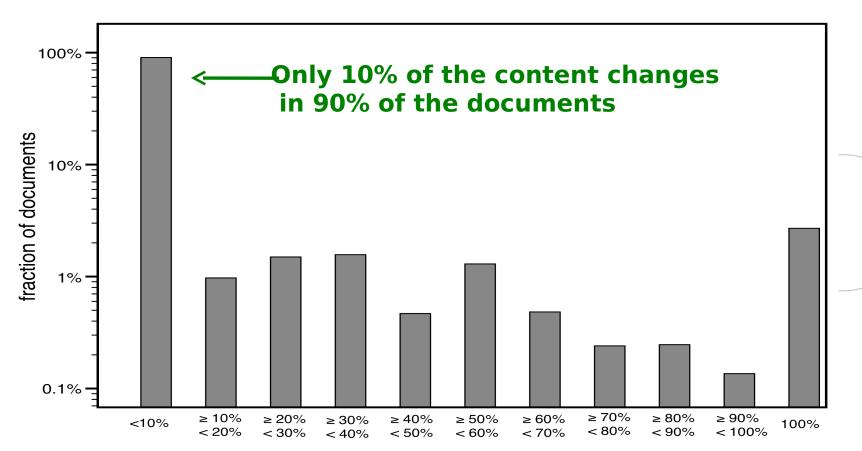
# Fraction of Changes per Document



**Digital Enterprise Research Institute** 

www.deri.i

6



average entity fraction which changes per document





- Web of Data has similar characteristics compared to the Web of Document concerning change frequencies
  - □ Change frequency of Web documents can be modeled as a Poisson Process [Cho at al. VLDB 2000]
- Document-centric change detection not helpful
  - □ different focus (entity-centric)
- Still needs more experimental verification





**Digital Enterprise Research Institute** 

www.deri.i

6

- Large scale experiment over a long period
- Study of what and how much changes
- Dataset Dynamics Group: http://groups.google.com/group/dataset-dynamics

Meet-up during WWW2010 on 29 April: http://esw.w3.org/Camps:LODCampW3CTrack#breakout

**Questions?** 



