

# uRiKA and Graph Analytics

uRiKA == universal RDF integration Knowledge Appliance



## Agenda





# uRiKA – YarcData Hardware and Software

### • Cray Hardware Engine

- Originally designed for deep analysis of large datasets
- Very large *scalable* shared memory
  - Architecture can support 512TB shared memory
  - Typical systems are 2 TB to 32 TB
- Multithreading
  - Unique highly multithreaded architecture
  - 128 hardware threads per processor
  - Extreme parallelism, hides memory latency

## • Multithreaded Graph Database

- Highly parallel in-memory RDF quad store
- High performance inference engine
- High performance parallel I/O

### Industry Standard Front End

- Based on Jena open source semantic DB
- WS02 application framework
- All standard SuSE Linux infrastructure and languages



### Hybrid Graph Appliance: Ease of Use AND Performance!



- Proven Cray infrastructure
- Cray XT5 3D Torus High Speed Interconnect



- Service Nodes:
  - AMD Opteron Processors, SuSe Linux
  - Open environment based on Jena
- Threadstorm Processors:
  - Multithreaded processors, MTK OS
  - Pre-programmed by YarcData



# **Emerging Web 3.0 Standards: RDF and SPARQL**

### Resource Description Framework (RDF)

- Designed to enable semantic web searching and integration of disparate data sources
- W3C standard formats
- Every datum represented as subject/predicate/object
  - Ideally with each of those expressed with a URI
- Standard ontologies in some domains
  - e.g., Open Biological and Biomedical Ontologies (OBO)
- Examples:



<ncbitax:NCBITaxon\_840261> <ncbitax:NCBITaxon\_195644> <ncbitax:NCBITaxon\_816681> rdf:type rdfs:subClassOf rdfs:label owl:Class <ncbitax:NCBITaxon\_185881> "Characiformes sp. BOLD:AAG5151"@en



## Semantic Database of RDF Triples

- RDF triples databases are inherently graphical
- Some researchers call semantic databases "semantic graph databases"



# Emerging Web 3.0 Standards: RDF and SPARQL

### SPARQL Protocol and RDF Query Language (SPARQL)

- Enables matching of graph patterns in the semantic DB
- Reminiscent of SQL

<pre># Lehigh University BenchMark (LUBM) Query 9 PREFIX rdf: <http: 02="" 1999="" 22-rdf-syntax-ns#="" www.w3.org=""></http:></pre>	PREFIX == shorthand for a URI
PREFIX ub: <http: 0401="" 2004="" td="" univ-benc<="" www.lehigh.edu="" ~zhp2=""><td>h.owl#&gt;</td></http:>	h.owl#>
SELECT ?X, ?Y, ?Z	variables to be returned from
WHERE	the query
{?X rdf:type ub:Student .	
?Y rdf:type ub:Faculty .	"find sets of (X, Y, Z) with a
?Z rdf:type ub:Course .	subject X of type Student, a
?X ub:advisor ?Y .	subject Y of type Faculty, and a
?Y ub:teacherOf ?Z .	subject Z of type Course,
?X ub:takesCourse ?Z}	where X is an advisee of Y, Y
	teaches course Z, and X takes



Company Confidential – Do Not Distribute

course Z"

### **uRiKA Software Stack**



### Industry-standard, Opensource Software Stack

• Linux, Java, Apache, WS02, Gadgets, Mashups...

### • Reusable Existing Skillsets

• OSGI, App Server, SOA, ESB, Web toolkit...

### No Lock-in

 All applications and artifacts built on uRiKA can be run on other platforms



## **uRiKA High Level Functional View**



### uRiKA complements existing Data Warehouse/Hadoop environment by offloading Graph Analytics





## Workflow





# Implementation





## Agenda





## **Getting Data into RDF Format**





# ... But Conversion is Not Enough

- Primary RDF goal is to be able to fuse data from different sources
- Ontologies (definitions of entities and relationships) must either be the same or be mappable onto each other
- Same: Use common ontologies
  - Generic: *e.g.*, RDF, OWL, Dublin Core
  - Social Networks: FOAF
  - Biology: *e.g.*, OBO, NCBI
- Map: Use ontology-mapping tools (e.g., Top Braid Composer)





# **Converting Data among Different RDF Formats**

- All files must be converted into .nt/.nq format before being Loaded into CQE
- In the fullness of time, the admin UI will convert from other formats as part of making the SDB
- Today, we use various conversion tools from the command line
  - rdf2rdf: Java tool, see <u>http://www.l3s.de/~minack/rdf2rdf/</u>
  - RIOT: runnable via bash scripts, see <u>http://incubator.apache.org/jena/documentation/io/riot.html</u>



# Extracting Data from an SQL Database

### • Use the D2R Server\*

- D2R can be used for an RDB in place; we extract the data once for ingestion into uRiKA
- See <u>D2R Server Data Extraction Configuration Guide</u> for details
- Two-step process
  - Generate a mapping from the target database
    - **e.g.,** generate-mapping -d jdbc:mysql://*server/database* -o map.n3
    - Many SQL DB types supported by D2R; MySQL, Oracle, and Postgres used so far with uRiKA
    - Edit the mapping file if (as typical) not all tables/fields are needed
  - Extract the data into an RDF file
    - e.g., dump-rdf -m map.n3
    - Move the result file into Lustre

\* http://www4.wiwiss.fu-berlin.de/bizer/d2rq/spec/#specification



### Importing Data from Data Stores e.g., Hadoop, Accumulo





## Agenda





# **Ingestion Data Flow**



- One or more .nt files (already resident on Lustre) can be combined into a single SDB
- A rules file may also be provided, which causes the ingestion step to run inferencing
- The output of the ingestion step is ~10 files in internal binary formats
- Those files are loaded directly into a uRiKA instance when it's initiated

20



.nt/.ng file

### **UI – Home**





**Company Confidential – Do Not Distribute** 

 $\frac{2}{2}$ 

## Ingestion UI – Import Data

Firefox			-	_		
( 302 Gadget Server	7:8080/carbon/dashboard/in	dex.isp			▼ Bina	
Most Visited 🚼 Google	C NYTimes F Facebook	MPR classical     RadioHe	artland 🏂 GMaps W Wikiped	ia O MPR News	59	Bookmarks
🔍 uRiK			Sigi	ned-in as admin   Sign-ou	ut   Help   Managem	nent Console
Home	Manage Data	Explore Data	Build a Query	Learn More	) 🖉 🕗	D 🕗
Tab id:0	Tab id:45	Tab id:46				
Import   <u>Build</u>	Load   Delete   St	atus				
Choose the type of	data you want to import. 🌶	To import from a relational	database (RDBMS) <u>click for ir</u>	structions.		=
Structured	Unstructo	ured Ru	es			
Structured I Choose your dai Loc Name Your Dai	Data File ta source: Local File ▼ al File ▲ Triple or Quad file ta File	es (.nt or .nq) are required.	Browse_			-
						×



**Company Confidential – Do Not Distribute** 

 $\frac{2}{2}$ 

## **Ingestion UI – Import Rules**

Firefox <b>T</b>			_	-		_				1 🕺
🚯 WSO2 Gadge	t Server	+	-	-	-					
کا 🔶	172.30.48.147:8080/carbor	n/dashboard/index.j	isp			ີ 🗸 🗸	C Bing			2
Most Visited	Google 🖲 NYTimes	f Facebook 💿	MPR classical 🧿 Rad	ioHeartla	nd 🏂 GMaps W Wil	kipedia 💿 MPR News			🗈 B	lookmarks
		ered By				Signed-in as admin	Sign-out   Help	o   Mana	gement Conso	ole
Home	Manage	e Data	Explore Data		Build a Query	Learn More		ø	o D o	
Tab id:0	Tab id:	45	Tab id:46							
Impor	L Build L Load LL	Delete   Ctatu								×
Character					hand (DDBMC) aliab	6				
Struct	ne type of data you war		Import from a relation	Pulos	abase (RDBMS) <u>click</u>	ror instructions.				E
Struct	ureu	Unstructured		Kules						
Info	roncing Buloc									
Choo	a the inferencing rules	you want to impo	t \Lambda For instructions		ting Rules click here					
Choo	Local File	you want to impo	re. 222 For instructions	soncrea	Browse					
	Name Your Rules									
	Set				import					
										-



**Company Confidential – Do Not Distribute** 

2 23

# **Ingestion UI – Build**

Fir	efox `	×	_			-	and the second		-				
🔂 V	VSO2 G	Gadget Se	erver	+		-	1000		10				
4	>	172	2.30.48.147:80	080/carbon/dashboard/index	.jsp			☆ ⊽ C ] <mark>ວ</mark> -	Bing				٩
P N	lost Vis	sited 🛂	Google 🖲	NYTimes 🛃 Facebook 🧿	MPR classical 🧿 RadioHeartl	and 🯂 GMa	ps W Wikipedia	MPR News				8	🖬 Book
Ę		u	<b>RiK</b> A	Powered By	TEX COMPANY		Signe	ed-in as admin   Sign-ou	t   Help	Mana	igeme	ent Co	onsole
F	lome			Manage Data	Explore Data	Build a Q	uery	Learn More		ø	0	Ð	<b>@</b>
	ab id	1:0		Tab id:45	Tab id:46								
													×
	Im	nort l	Build   I	oad   Delete   Stat	116								
	-		Dalla   I		<u>us</u>					_			
		Build	Knowled	gebase									
		Select th	he Data File(	(s) you want to build into a	Knowledgebase.								
		A 16		erest he built when a LC									
			viedgebases	cannot be built when a LC	AD is in progress.		Data	and Time					
		INC	ame				Date						
			dataset				02/02/2	2012 2:39AM	Â				
			nasa_datas	set			02/01/2	2012 9:46PM					
			mondial-na	la8443			01/05/2	2012 9:36PM	=				
			dbpedia_du	uisburg_essen			01/05/2	2012 9:36PM					
			dbpedia_ei	nstein_stuttgart			01/05/2	2012 9:36PM					
			dbpedia_lei	ipzig_berlin			01/05/2	2012 9:36PM					
			lubm0				01/05/2	2012 9:36PM	-				
			-inputDataN	0442			01/05/	2012 0.26DM					
		🔲 En	nable Inferer	ncing									
		Name Yo	our Knowled	gebase			Cancel Bui	ld					

 Can select which (Lustreresident) files to combine into a knowledge base

> 2 24

marks



## Ingestion UI – Load

Firefox VSO2 Ga	adget Serve	r +	Provide and	-		in the			<b>x</b> ⊽
( <b>{</b> )	172.30	.48.147:8080/carbon/dashboard/	index.jsp		∰ ⊽ C <mark></mark>	• Bing		م	
Most Visit	ited 🔧 Go	ogle 🖲 NYTimes 🛃 Faceboo	k 💿 MPR classical 💿 RadioHeartl	land 🏂 GMaps W Wikipedia	O MPR News			🖪 Bo	okmarks
	uR			Signe	ed-in as admin   Sign-ou	it   Help   Ma	nageme	nt Consol	e
Home		Manage Data	Explore Data	Build a Query	Learn More	6	9 🕗	D 🕺	
Tab id:	:0	Tab id:45	Tab id:46						
							1		٩
Imp	port   <u>B</u>	uild   Load   <u>Delete</u>   <u>9</u>	<u>Status</u>						
									E
	Load Kn	owledgebase	ad into your URIKA avetam						
		laebases cannot be loaded wh	en a BUU D is in progress						
	- Knowied	Name	en a boreb is in progress.	c.	ate and Time				
		E Charles I mondial I lubr	~0	02//	2/2012 2:00DM	<u>^</u>			
	0	±. ⊡Data Set	no	02/0	JZ/ 2012 3:00PM				
		🗄 🚞 Rule Set				=			
				02//	12/2012 4-24PM				
	٢			02/0	12/2012 4.24FM				
	$\odot$	<sup>⊞</sup> <mark>È</mark> NASA_POC		02/0	02/2012 1:53AM	~			
	Cancel	Load							

 Load starts the SDB instance (after stopping a current instance)

> 2 235



### **UI – Learn More**

inour outget server		+	Participation		_	
172.30.4	8.147:8080/carbon/d	lashboard/index.isi	p		<u>ר</u> קר פו נוסד Bin	
Most Visited 🚼 Good	gle 🖲 NYTimes 📕	Facebook O M	' 1PR classical 💿 RadioHea	rtland 🏂 GMaps W Wik	ipedia O MPR News	Bookmark
🔍 uRi	KA Powere		COMPANY		Signed-in as admin   Sign-out	Help   Management Console
Home	Manage I	Data	Explore Data	Build a Query	Learn More	ø 🛛 🗅 🕹
Tab id:0	Tab id:4	5	Tab id:46			
Quick Start Install The following quick st	ation and Config	<b>juration Guide</b> get you up and r	unning with uRiKA.			
Description	Q	uick Start Guid	e			
Front End (FE) envir assumptions, setup, installation instruction	and no	eploying uRiKA to ode environment	<u>o a service</u>			1
Quick Start User G	<b>uides</b> s are intended to j	ump start the use	er's experience. Further	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description	uides s are intended to j Quick Start Gui	ump start the use	er's experience. Further	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS	uides s are intended to j Quick Start Gui D2R User Guide	ump start the use ide More Ir D2R Doc D2R Maj	er's experience. Further <b>1fo</b> <u>cumentation</u> <u>oping Language</u>	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS Visual relationship finder	uides s are intended to j Quick Start Gui D2R User Guide RelFinder User G	ump start the use ide More Ir D2R Doc D2R May uide Visual D	ar's experience. Further fo cumentation pping Language ata Web	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS Visual relationship finder Learning SPARQL	uides s are intended to j Quick Start Gui D2R User Guide RelFinder User G	ump start the use ide More Ir D2R Dor D2R Maj uide Visual D SPAROL	ar's experience. Further fo cumentation oping Language ata Web Reference	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS Visual relationship finder Learning SPARQL Learning about WSO2	uides s are intended to j Quick Start Guid D2R User Guide RelFinder User G User and gadget permission management	ump start the use       de     More Ir       D2R Dor     D2R May       Uide     Visual D       SPAROL       WS02 G	ar's experience. Further fo cumentation pping Language ata Web Reference iadget Site	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS Visual relationship finder Learning SPARQL Learning about WSO2	uides s are intended to j Quick Start Guide D2R User Guide RelFinder User G User and gadget permission management	ide More Ir D2R Dor D2R May uide Visual D SPAROL WSO2 G	er's experience. Further fo cumentation oping Language ata Web Reference iadget Site	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.
Quick Start User G The quick start guide Description Extracting structured data from an RDBMS Visual relationship finder Learning SPARQL Learning about WSO2	uides s are intended to j Quick Start Guid D2R User Guide RelFinder User G User and gadget permission management	ump start the use ide More Ir D2R Dor D2R Maj uide Visual D SPAROL WSO2 G	ar's experience. Further fo cumentation oping Language ata Web Reference adget Site	details can be found by fo	ollowing the links to each tools' respect	ive documentation website.



**Company Confidential – Do Not Distribute** 

2 26

## Agenda





## **SPARQL** Query



#### Lehigh University Benchmark (LUBM) Query 9





# SPARQL roadmap

### • Current uRiKA release (0.9) covers SPARQL 1.0

- Released 2008
- Read-only database
- Basic pattern matching, filters, unions

### Next release will cover SPARQL 1.1

- Specification released 2011
- Property paths not included (spec under revision)
- Improved Database Administrator interface
- More input file formats supported

### • SPARQL 1.1 Update specification partially supported

In-memory INSERT DATA, DELETE DATA



# Inferencing

#### schema.ttl

```
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix : <http://example.org/vehicles/> .
```

:Vehicle a rdfs:Class . :Car rdfs:subClassOf :Vehicle . :SportsCar rdfs:subClassOf :Car .

#### data.ttl

```
@prefix ex: <http://example.org/vehicles/> .
@prefix : <http://myvehicledata.com/> .
```

```
:FordFiesta a ex:Car .
```

- :AudiA8 a ex:Car .
- :FerrariEnzo a ex:SportsCar .

#### query.rq

```
PREFIX ex: <http://example.org/vehicles/> .
PREFIX : <http://myvehicledata.com/> .
SELECT ?car
WHERE { ?car a ex:Car }
```

### • Without inferencing, query will return

• FordFiesta, AudiA8

### • With inferencing, query will return

• FordFiesta, AudiA8, FerrariEnzo

• We expect inferencing to be a uRiKA strength

Example from http://www.dotnetrdf.org/content.asp?pageID=Inference%20and%20Reasoning



## Agenda





# **Visualizing Results**

- RelFinder (built into the UI as Explore Data) does this with independently of the relationships in the data
- New visualization packages are being investigated for the next release
- Google Gadgets within the WSO2 framework are the primary interface for gadget development
  - These interfaces are not fully documented or exposed in 0.9









Ya

## Agenda





# Summary

- uRiKA is targeted at large-scale semantic-graph processing
- Some end-to-end ETL workflows need components beyond what uRiKA has been tested with today
- We expect to test and qualify third-party components
- Most customers have their own favorite ETL conponents- we will learn a lot by working with early customers
- ... and modify uRiKA based on those lessons



### Learn More

- SPARQL by Example tutorial, by Lee Feigenbaum, http://www.cambridgesemantics.com/2008/09/sparql-by-example/
- Search RDF data with SPARQL, by Philip McCarthy <u>http://www.ibm.com/developerworks/xml/library/j-sparql/</u>
- Semantic Web for the Working Ontologist, by Dean Allemang and James Hendler, ISBN 978-0123859655
- Learning SPARQL, by Bob DuCharme, O'Reilly, ISBN 978-1-449-30659-5
- RDF: <u>www.w3.org/RDF/</u>
- SPARQL: <u>www.w3.org/TR/rdf-sparql-query/</u>
- D2R: www4.wiwiss.fu-berlin.de/bizer/d2r-server
- WSO2: wso2.com/products/application-server
- Google Gadgets: www.google.com/webmasters/gadgets/



# Thank you!

James D. Maltby, Ph.D jmaltby@yarcdata.com



# **Backup Slides**





### **Data Sizes**

- Today, typically the input data should consume no more than 50-60% of the memory of the system, to leave space for temporary storage during queries
- Estimating memory consumption from the number of triples is inexact, but 100B:triple is a reasonable rule of thumb



# **New SPARQL 1.1 features**

### Aggregates and GROUP BY

Divide solution set into subgroups, can apply operator to subgroup

### Subqueries

Like RDBMS "view" => query against a previous result

### Negation

NOT EXISTS – does a pattern exist

## • SPARQL 1.1: New FILTER expressions

- EXISTS, NOT EXISTS
- IN, NOT IN
- COALESCE return RDF term value of the first term in the expression that evaluates without error
- IF evaluates the Effective Boolean Value of the first argument, returns 2nd or 3rd depending on whether EBV was true or fals



# More new features in SPARQL 1.1



• Query remote SPARQL endpoint

# • AVG(), MIN(), MAX(), COUNT()

Functions over returned values

### • MINUS – eliminate a specific match

### BIND

• Assign a value to a new variable in a group graph pattern



# **Dynamic in-memory Updates**

### • SPARQL 1.1 Update requests

- Subset of full specification
- Submitted over HTTP like a SPARQL query

### INSERT DATA QuadData

- Performed in-memory on compute nodes
- Forward static inferencing on inserted triples
- Copy of update stored on disk

### DELETE DATA QuadData

- Also in-memory, logged to disk
- Forward static inferencing on deleted triples
- No backward inferencing

