



FONDAZIONE
BRUNO KESSLER



RDF_{pro}

an Extensible Tool for Building Stream-Oriented RDF Processing Pipelines

Riva del Garda, 19 October 2014

Marco Rospocher¹, Marco Amadori², Michele Mostarda², Francesco Corcoglioniti

⁽¹⁾ Data and Knowledge Management Unit, FBK-Irst, <http://dkm.fbk.eu/>

⁽²⁾ Web of Data Unit, FBK-Irst <http://wod.fbk.eu/>

<http://fracor.bitbucket.org/rdfpro>



The problem

perform simple RDF processing tasks

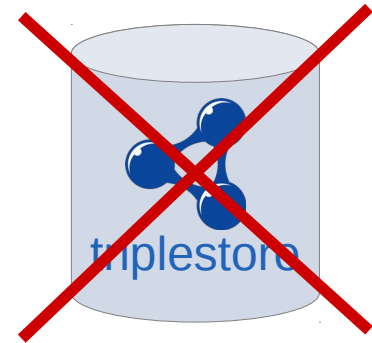
- filtering and transformation (quad-level)
- basic inference (RDFS)
- dataset merging → deduplication, owl:sameAs smushing
- simple statistics extraction (VOID+)
- ...

on large datasets

- LOD-sized: 100M+ triples
- quads, not just triples

on a single commodity machine

- no cluster / distributed computing
- no triplestore or other data index



The solution



RDF_{pro}

pro = processor (and not 'professional!')



- ~ Java command line tool ~
- ~ embeddable Java library ~
- ~ public domain code ~

<http://fracor.bitbucket.org/rdfpro/>

RDF_{pro} ingredients

① streaming

realized via the **RDF processor** abstraction



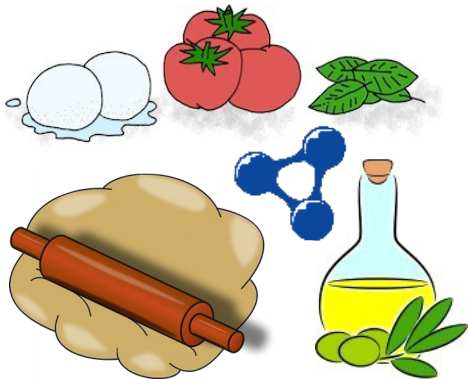
invocation syntax: `rdfpro @P args`

pro:

- natural model for many tasks
- $O(n)$ time complexity
→ fast, also due to sequential data access
- $O(1)$ space complexity (usually)
→ copes with arbitrarily large datasets

cons:

- restrictive model!



RDF_{pro} ingredients

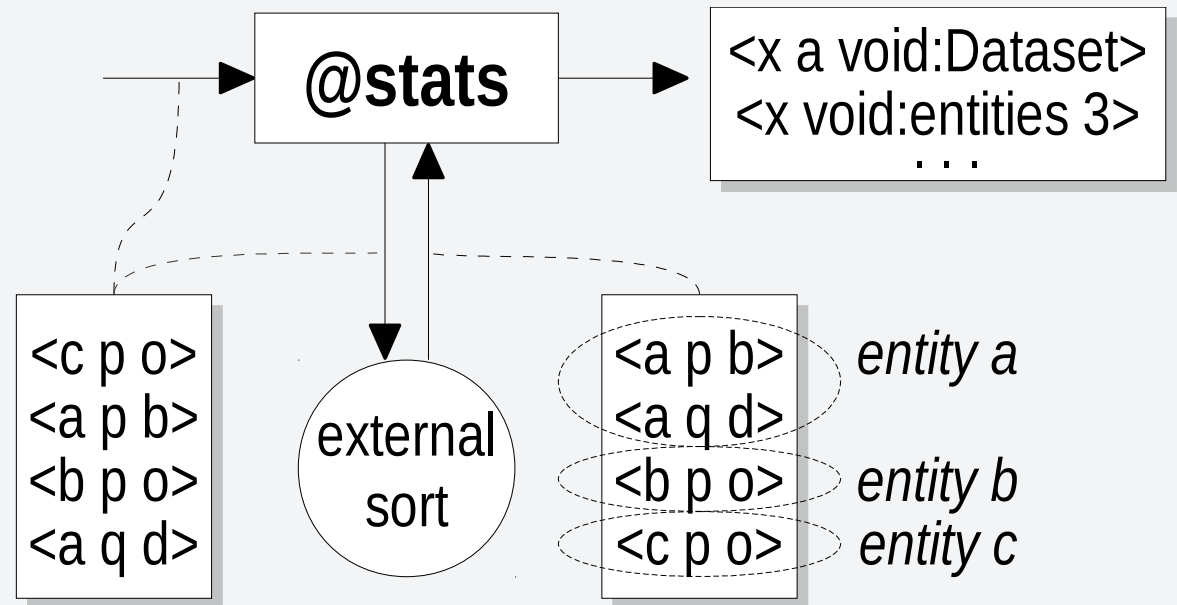
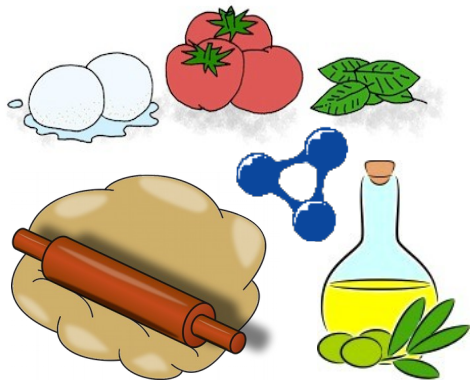
① streaming

② sorting

realized via **external sorting** (sort utility)

allows tasks not doable with pure streaming

- duplicate removal
- set operations (quad union, intersection, diff.)
- VOID statistics extraction
- ...

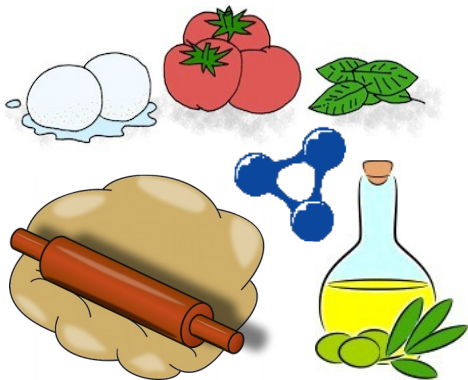


RDF_{pro} ingredients

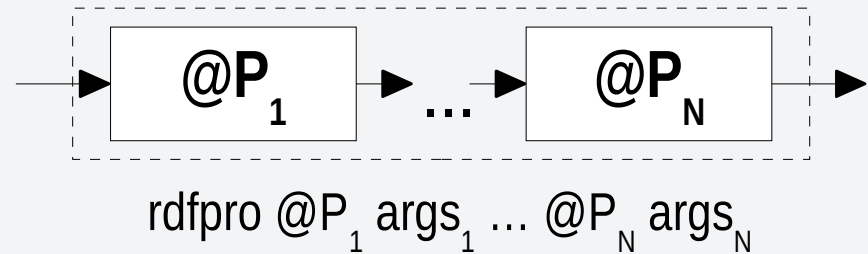
① streaming

② sorting

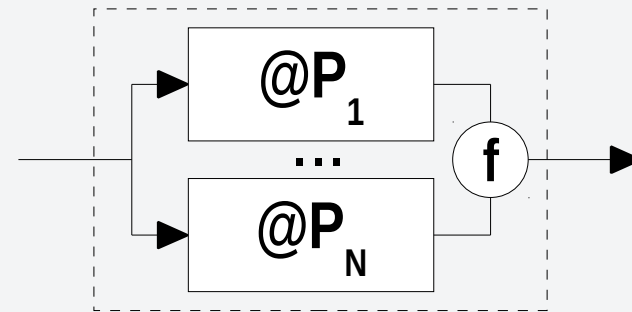
③ pipelining



① sequence composition



② parallel composition



pro:

- reduced I/O costs (less temporary files)
- reduced execution time (parallelism)

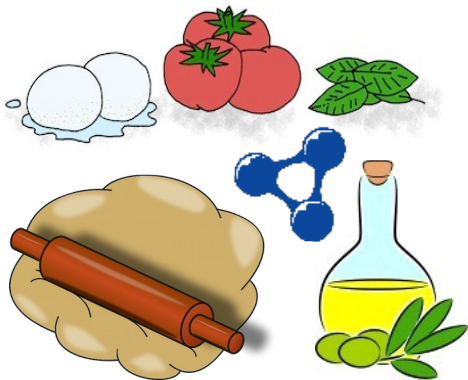
RDF_{pro} ingredients

① streaming

② sorting

③ pipelining

④ multi-threading



① inter-processor parallelism

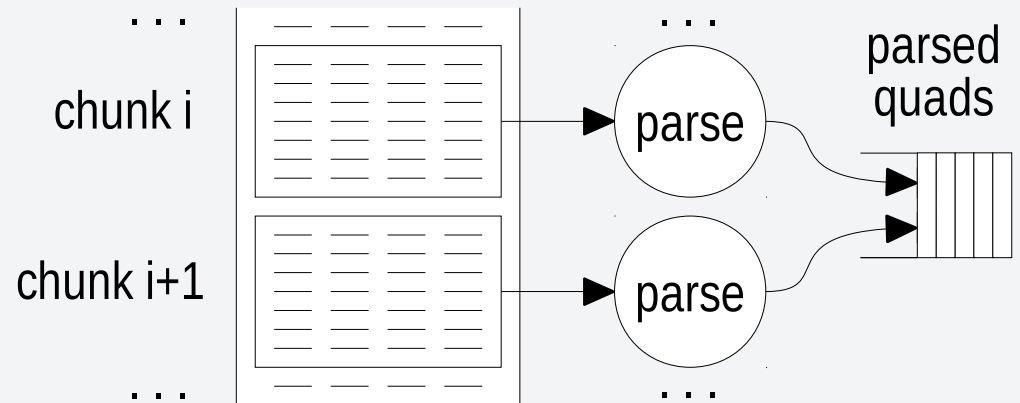
- multiple processors run in parallel

② intra-processor parallelism

- handleStatement() called concurrently

③ I/O parallelism

- multiple files read/written in parallel
- single files split in chunks processed in parallel (line-oriented RDF formats only)



Putting all together, you can ...

move data around

- @read / @write files
- @download from / @upload to SPARQL endpoints

transform data

- general purpose data @transform using Groovy
- @infer the RDFS closure
- @smush data, replacing URI aliases with canonical URIs
- extract @tbox and VOID @stats

compose these tasks freely

- also via set operations



A simple use case

integrate:

- **Freebase** (2014/07/10 dump, 2623 MQuads)
- **GeoNames** (2013/08/27 dump 125 MQuads)
- **DBpedia EN, ES, IT, NL** (subset of ver. 3.9, 271 MQuads)



performing:

- **filtering** (remove redundant quads & quads in unwanted languages)
- **smushing** (based on owl:sameAs links in DBpedia)
- **inference** (excluding `<X rdf:type rdfs:Resource>` stuff)
- **statistics extraction** (VOID with class & property partitions)

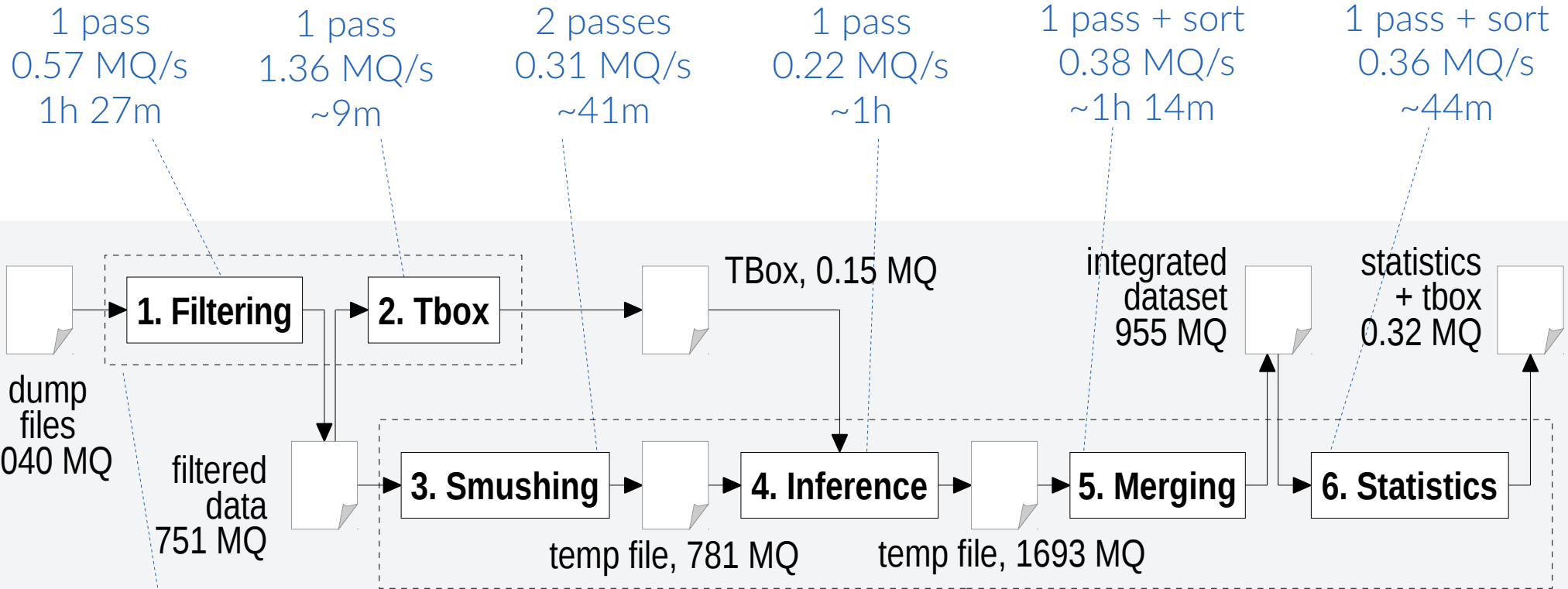
using:

- **a small workstation** (I7 860, 16 GB ram, 500 GB 7200 rpm hd)
- **RDF_{pro} + parallel sort + pigz + pbzip2**



A simple use case

tasks performed individually - 5h 16m total



1-2 aggregated:
1 pass, 0.56 MQ/s, 1h 29m

3-6 aggregated:
2 passes, 0.09 MQ/s, 2h 16m

aggregated tasks - 3h 46m total (-28%)

A simple use case

individual tasks

Task	Input size		Output size		Throughput		Time
	[MQuad]	[GB]	[MQuad]	[GB]	[MQuad/s]	[MB/s]	
1. Filtering	3019.89	29.31	750.78	9.68	0.57	5.70	1:27:46
2. TBox extraction	750.78	9.68	0.15	0.01	1.36	18.00	9:11
3. Smushing	750.78	9.68	780.86	10.33	0.31	4.04	40:53
4. Inference	781.01	10.34	1693.59	15.56	0.22	2.91	1:00:30
5. Deduplication	1693.59	15.56	954.91	7.77	0.38	3.61	1:13:33
6. Statistics	954.91	7.77	0.32	0.01	0.36	3.02	44:00
whole processing	3019.89	29.31	955.23	7.78	0.16	1.58	5:15:53

aggregated tasks

Task	Input size		Output size		Throughput		Time
	[MQuad]	[GB]	[MQuad]	[GB]	[Mquad/s]	[MB/s]	
1-2 aggregated	3019.89	29.31	750.92	9.69	0.56	5.60	1:29:23
3-6 aggregated	750.92	9.69	955.23	7.78	0.09	1.21	2:16:08
whole processing	3019.89	29.31	955.23	7.78	0.22	2.22	3:45:31

① download

<http://fracor.bitbucket.org/rdfpro> (or Google for it!)

RDF_{pro}

About ▾

Tool ▾

Libraries ▾

Maven Reports ▾

Links ▾

RDF_{pro}

An Extensible Tool for Building Stream-Oriented RDF Processing Pipelines

Download

About

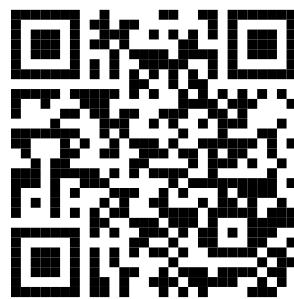
RDF_{pro} (RDF Processor) is a public domain, Java command line tool and library for **RDF processing**. RDF_{pro} offers a suite of stream-oriented, highly optimized **RDF processors** for common tasks that can be assembled in complex **pipelines** to efficiently process RDF data in one or more passes. RDF_{pro} originated from the need of a tool supporting typical **Linked Data integration tasks**, involving dataset sizes up to few billions triples.

Features

- RDF quad (triple + graph) filtering and replacement
- RDFS inference with selectable rules
- owl:sameAs [smushing](#)
- TBox and VOID statistics extraction
- RDF deduplication, intersection and difference
- data upload/download via SPARQL endpoints
- data read/write in multiple (compressed) formats (rdf, rj, jsonld, nt, nq, trix, trig, tqi, ttl, n3, brf)
- command line [tool](#) + [core](#), [tqi](#), [jsonld](#) libraries
- based on [Sesame](#).
- public domain software ([Creative Commons CC0](#))

News

- 2014-08-04 Version 0.2 has been released
- 2014-07-24 Version 0.1 has been released



LIBRARIES

[RDF_{pro}](#)
[TQL](#)
[JSONLD](#)
[Javadoc](#)

LINKS

[BitBucket project](#)
[Issue tracker](#)
[Contact authors](#)

RDF_{pro} is public domain software developed within:



[Back to top](#)



RDF_{pro} cookbook

① download

② install

check requirements:

- Java 1.7+ (Oracle, OpenJDK, whatever)
- gzip, bzip2, sort utilities available on PATH

extract the download tarball:

```
$ tar tf rdfpro-0.3.tar.gz
```

check that everything works:

```
$ cd rdfpro  
$ ./rdfpro -v  
RDF Processor Tool (RDFpro) 0.3  
Java 64 bit (Oracle Corporation) 1.7.0_67  
This is free software released into the public domain
```

suggestions:

- add rdfpro directory to PATH
- install and configure pigz and pbzip2 (see web site)

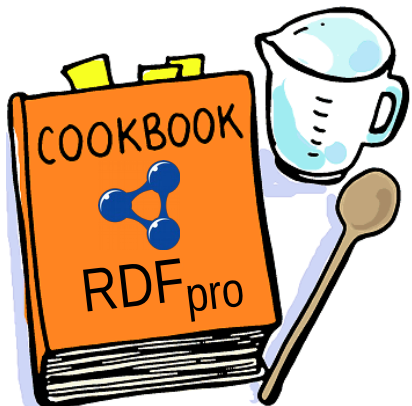


RDF_{pro} cookbook

- ① download
- ② install
- ③ try it out!

let's get and process some data from Dbpedia:

```
$ ./rdfpro \  
> @read http://dbpedia.org/resource/Riva_del_Garda \  
> http://it.dbpedia.org/resource/Riva_del_Garda \  
> @smush \  
> @infer http://downloads.dbpedia.org/3.9/dbpedia_3.9.owl.bz2 \  
> @transform "emitIf(t = rdf:type)" \  
> @unique \  
> @write riva_del_garda.ttl.gz
```



**That's all:
enjoy cooking triples with RDF_{pro} and...
happy eating !!**



for any question about the menu RDF_{pro}, contact
Francesco Corcoglioniti <corcoglio@fbk.eu>