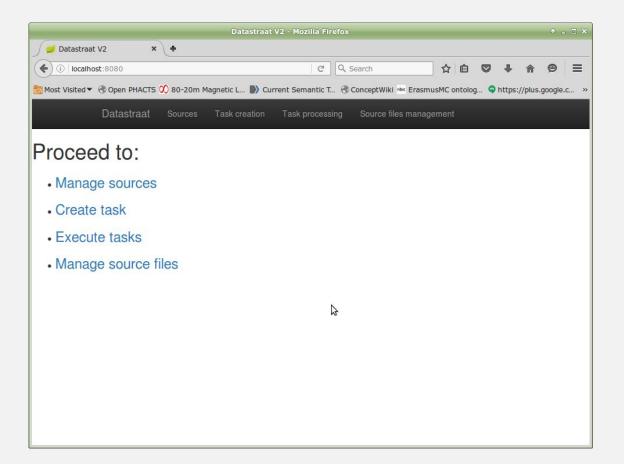


DPS V2

Our starting point for DPS

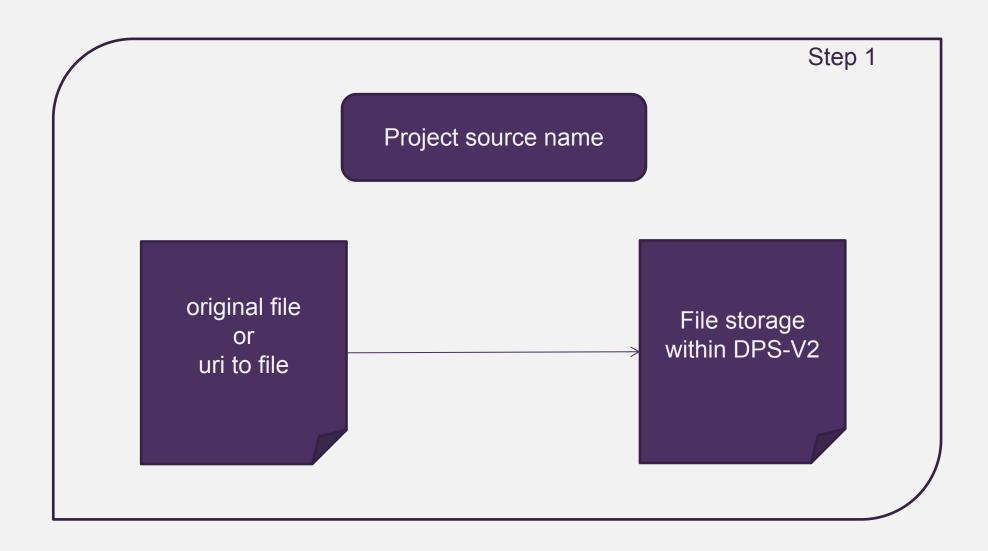


- **E** ssh -f dps@178.63.49.197 -L 8082:127.0.0.1:8082 -N
- E browser uri;
 - ⊢ localhost:8082



Step 1; loading data into DPS





Uploading / downloading my data

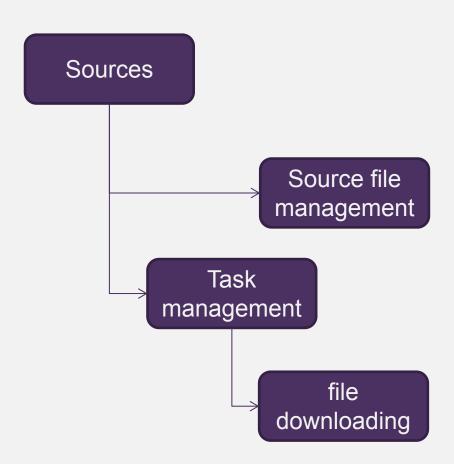


E Sources

- → give your project a reference name
 - Its a handle we use during other steps in subsequent processes
- E Decide how you upload your data file
 - source file management, select input file, enter destination file name, upload

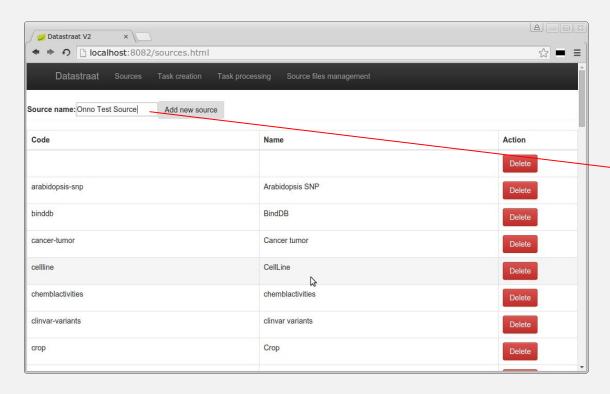
OR

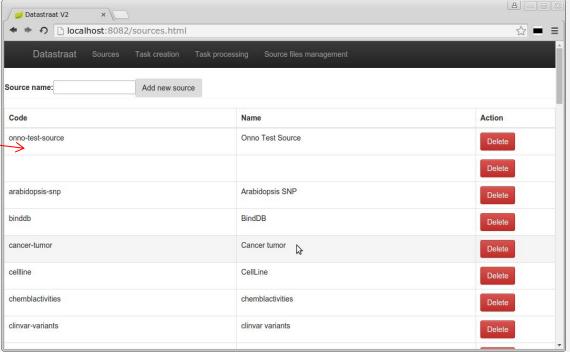
task management, file downloading through url definition



Add new source (project name)

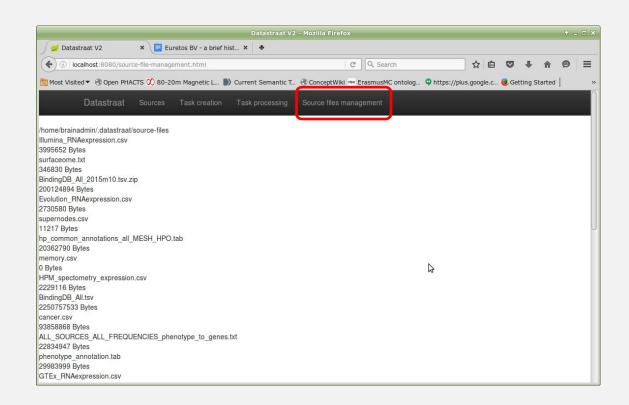


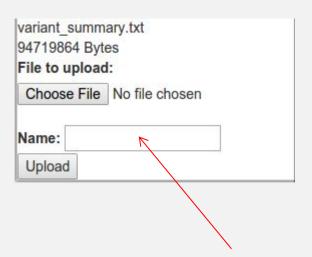




Source file management, uploading your file







Define the filename as used to store within DPS-V2





atastraat V2	× ⟨ ♠ Robomongo host:8082/task-creation	n.html					۲۵
Datastraat	Sources Task creation	Task processing	Source files	management			
s Drafts							
< Get back							
Source:	8						
Entrezgene							
Processing step	configuration:						
file-downloading							
File Destionation I	Path?:						
/disk-ba	ackup-2/inputdata/EntrezGene						
File URL(2):							
ftp://ftp.	ncbi.nih.gov/gene/DATA/mim2	2gene_medgen					
Save Save	as draft						
				Please	e note: or	nly single	
						upported	



Excercise step 1

Excercise step 1 (use your own name)



E Add a new source "Onno Test Source"

□ Input file "onno.txt"

#subjectName, predicateName, objectName, pmid "smoking", "causes", "cancer", "24950063" smoking, relates, cancer, 12907484 "excercise", "prevents", "cancer", 23665342

E Import it to the DPS server



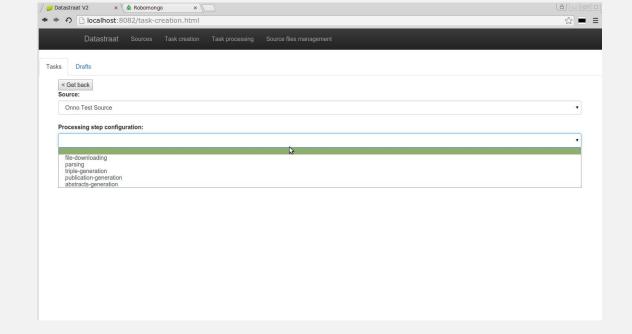
Task management

Task management



- **E** file downloading
- **E** parsing
 - → step 2
- **E** triple-generation
 - → step 3
- **E** publication-generation
- **E** abstracts-generation
 - step 5

Training coverage



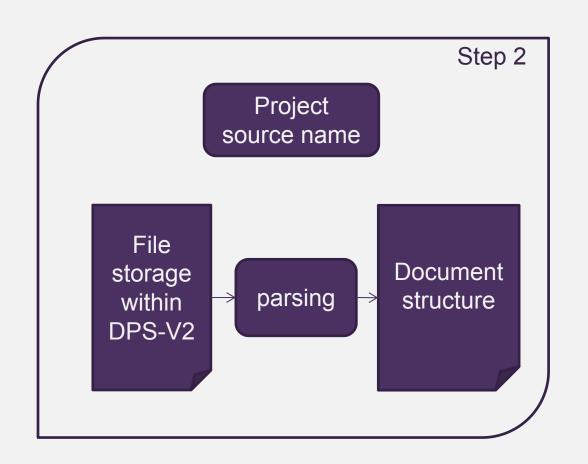


Parsing logic

Step 2; parsing data within DPS



- E The data file we stored on the server requires parsing to normalize into document format
- E Its important to understand that all values are parsed as strings
- E Parsing only transforms input to structured document format
- E There is only one cleanup function, which is limited to removing header and / or tail of string line.

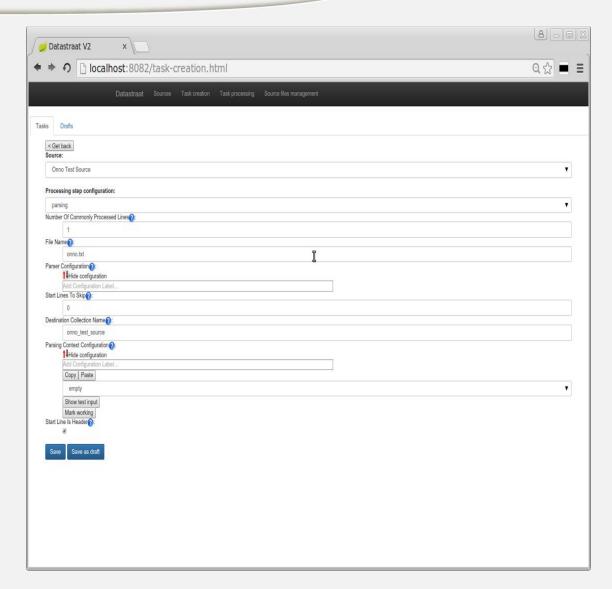


Creating the parser (step 2-1)



- **E** select source (project)
- E processing step; parsing
- E select file name
- E confirm first row is header true/false
- parsing context configuration; used for workflows, set to empty

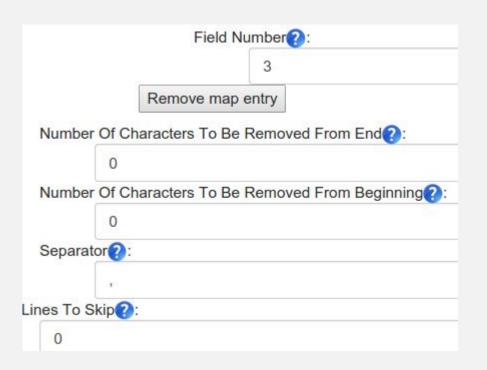
Validation function does not work for parsing

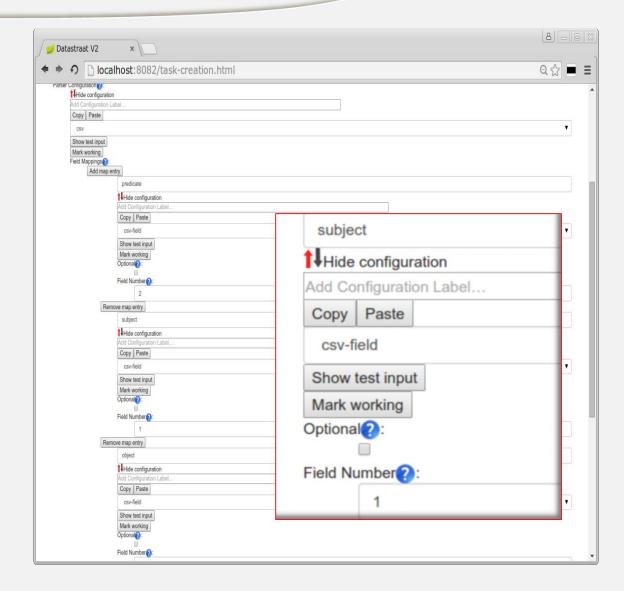


Creating the parser (step 2-2)



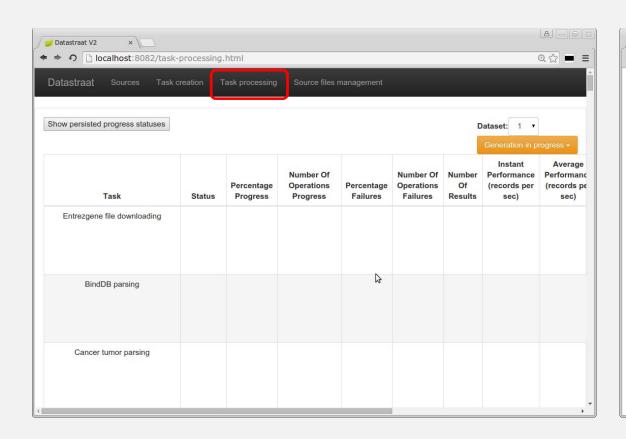
- **E** Field Mapping definition
 - → add map entry for each field

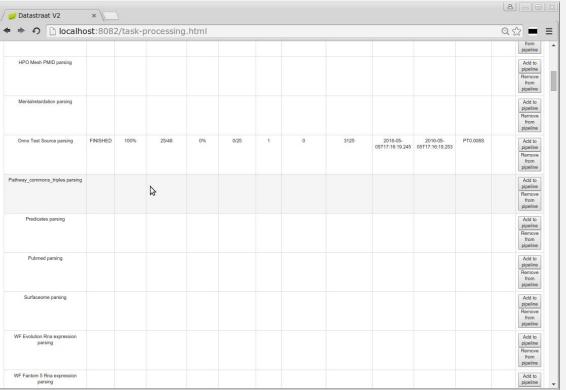




Execute the parser (step 2-3)

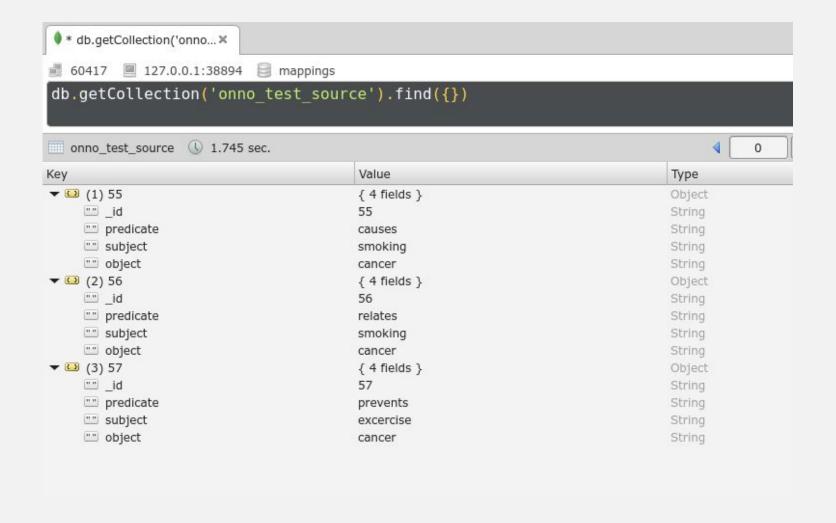






Resulting collection within documentstore







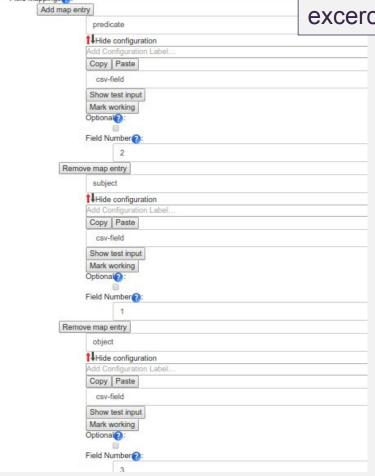
Summary parsing

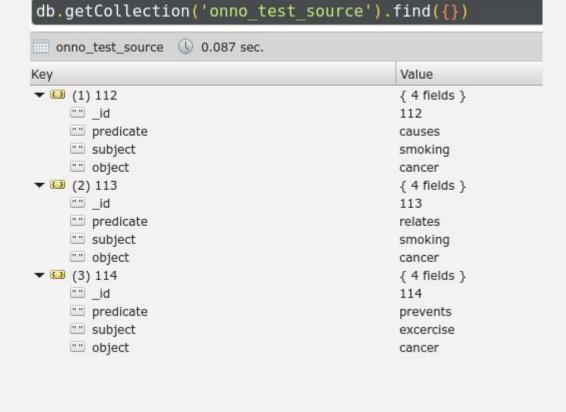
CSV



#header stuff smoking,causes,cancer smoking,relates,cancer excercise,prevents,cancer

This is document after we execute the task





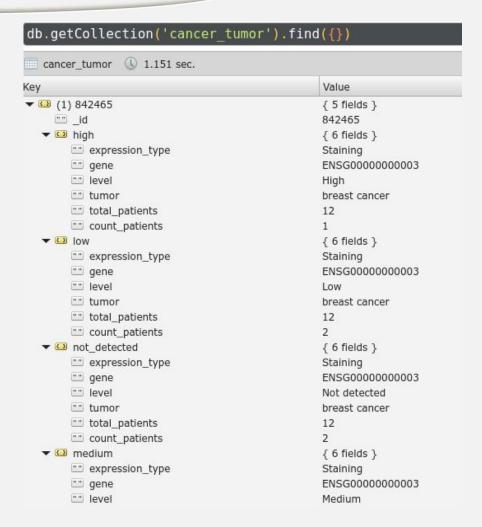
csv-multi-line



"Gene","Gene name","Tumor","Level","Count patients","Total patients","Expression type"
"ENSG0000000003","TSPAN6","breast cancer","High","1","12","Staining"
"ENSG00000000003","TSPAN6","breast cancer","Medium","7","12","Staining"
"ENSG00000000003","TSPAN6","breast cancer","Low","2","12","Staining"
"ENSG00000000003","TSPAN6","breast cancer","Not detected","2","12","Staining"
"ENSG00000000003","TSPAN6","carcinoid","High","0","4","Staining"
"ENSG00000000003","TSPAN6","carcinoid","Medium","1","4","Staining"
"ENSG00000000003","TSPAN6","carcinoid","Low","1","4","Staining"
"ENSG000000000003","TSPAN6","carcinoid","Not detected","2","4","Staining"



example parser logic; cancer tumor

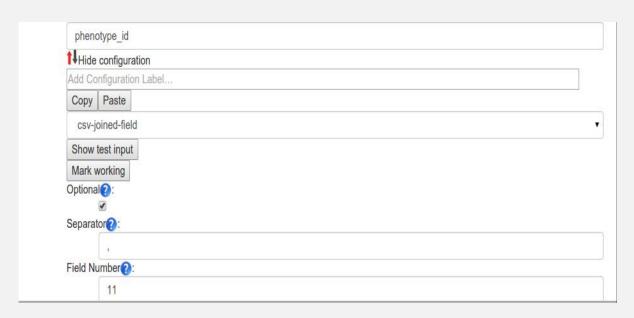


This is document after we execute the task.

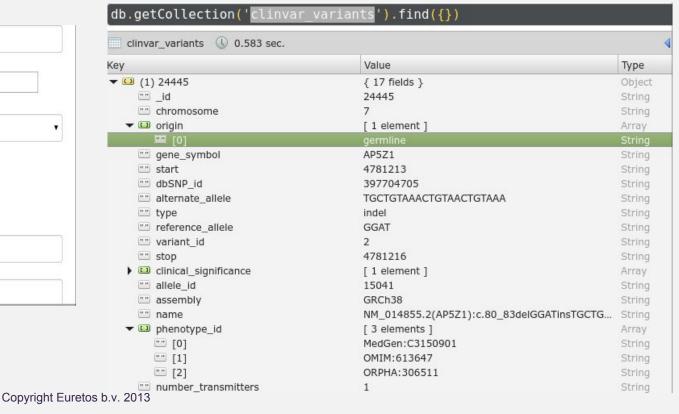
csv-joined-field



15041 NM_014855.2(AP5Z1):c.80_83delGGATinsTGCTGTAAACTGTAACTGTAAA (p.Arg27_Ala362de indel linsLeuLeuTer) AP5Z1 Pathogenic 397704705 9907 RCV00000012 MedGen:C3150901,OMIM:613647,ORPHA:306511 germline GRCh37 7 4820844 4820 no assertion criteria provided NM_014855.2:c.80_83delGGATinsTGCTGTAAACTGTAA 847 CTGTAAA NP_055670.1:p.Arg27_Ala362delinsLeuLeuTer 1 Jun 29, 2010 OMIM Allelic Variant:613653.0001 2 GGAT TGCTGTAAACTGTAACTGTAAA 1 NC 000007.13



example parser logic; clinvar_variants





Excercise step 2

Excercise step 2 (use your own project)



- **E** Create your own parsing
- **E** Execute as task
- E Check resulting document using Robomongo

Configuration robomongo

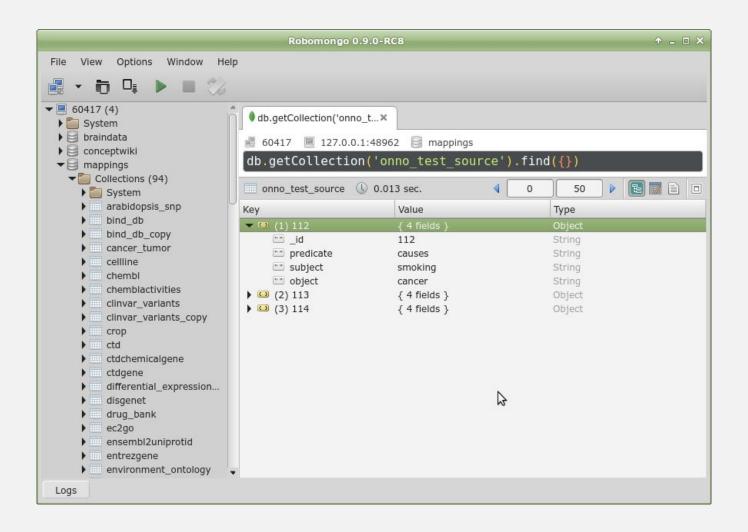






Robomongo view

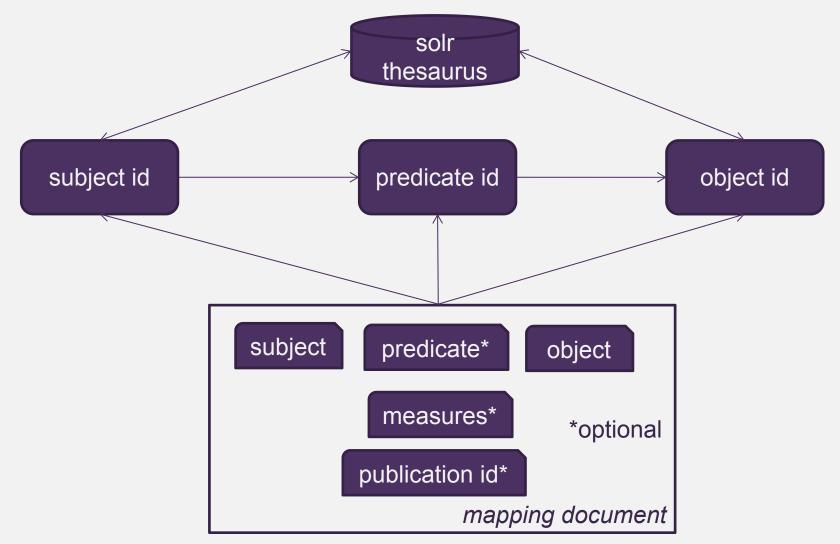






Creating triple generator logic



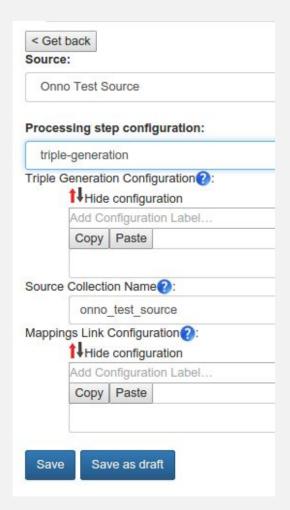


Copyright Euretos b.v. 2013

Main configuration elements

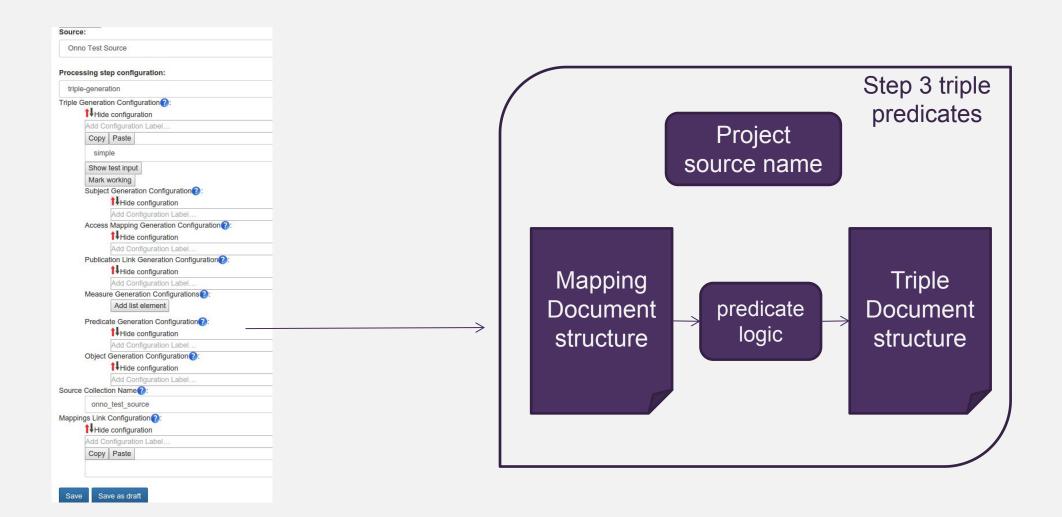


- Confirm the source collection name
- E Mapping link configuration:



Predicate

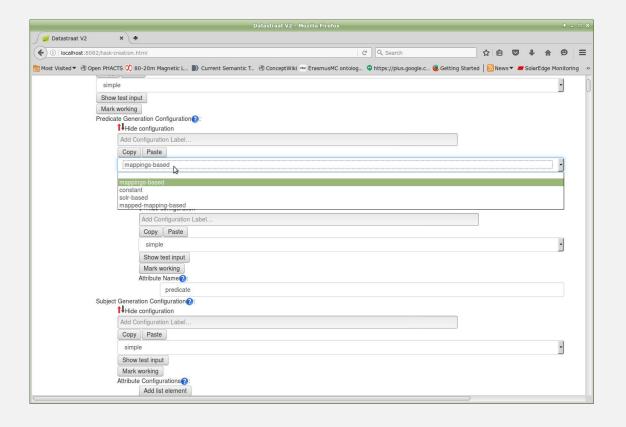




Predicate assignment

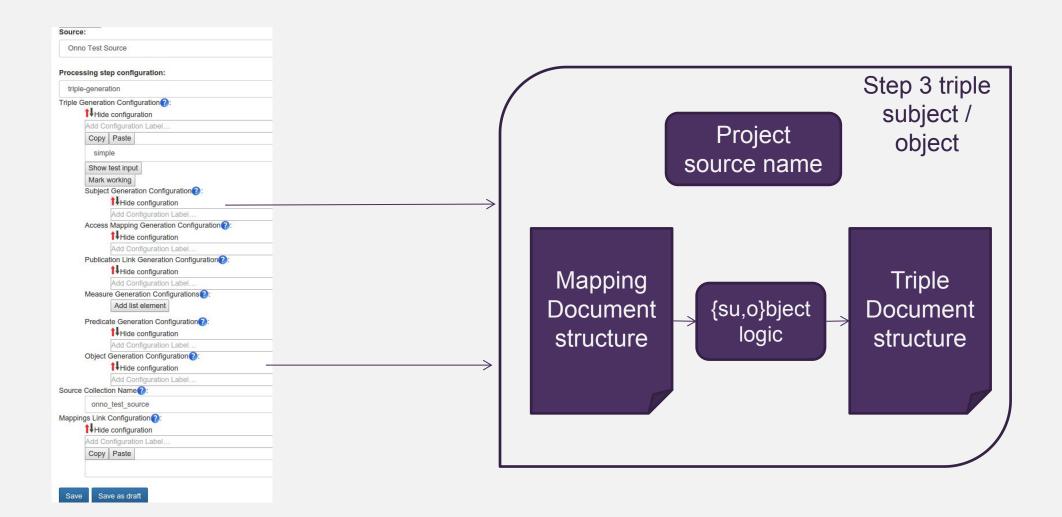


- E Assignment / creation of predicate can be based on;
 - mapping-based
 - retrieve predicate name from mapping collection
 - → constant
 - statically assigned by user into predicate name field
 - solr-based
 solr-based
 - solr request / response determines which predicate will be assigned
 - mapped-mapping-based
 - take a value from mapping collection and depending on value returns predicate name



Subject / Object

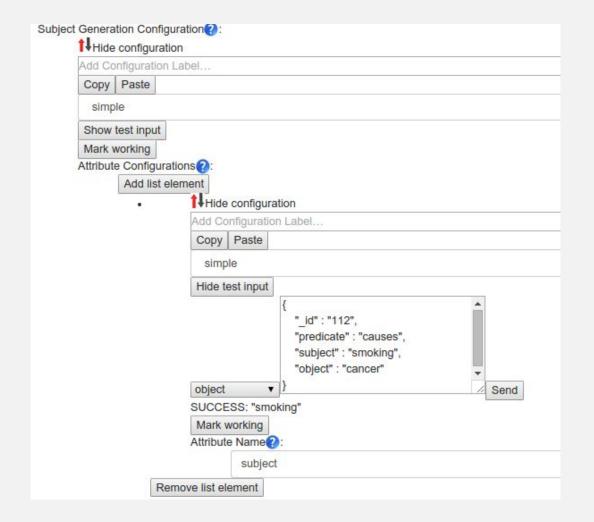




Subject / object -> retrieve values

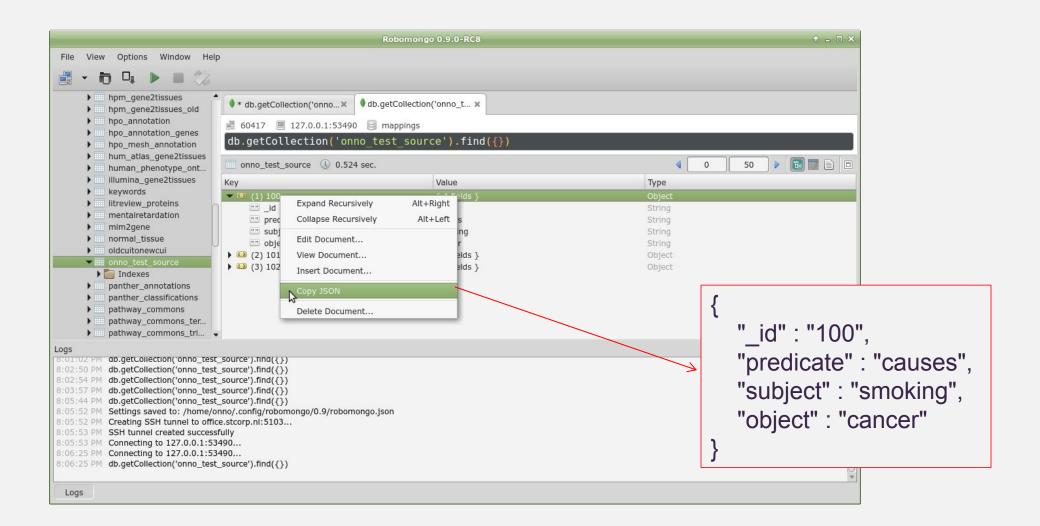


- E use simple to keep it simple
- E attribute name is subject or object depending case
- you can test mapping field through copy / paste json object (see next slide)



Robomongo (copy json)



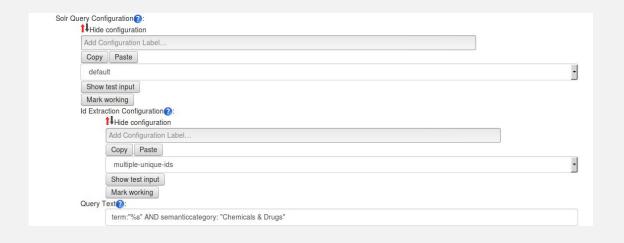


Subject / Object -> Solr query



E Example solr queries:

- → term:"%s"
- term:"%s" AND semantictype:"28" AND taxonomies:"arabidopsis thaliana"
- (term:"%s" AND semantictype:"116")
- (term:%s" AND (source:"umls" OR source:"uniprot"))
- E Good idea to use your API access to collect key, value pairs?

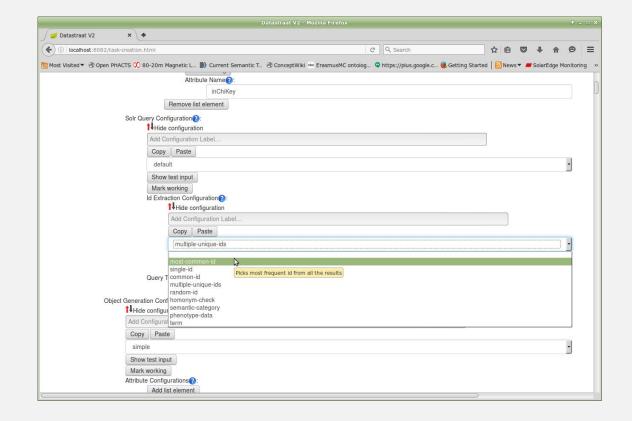


Note: keep in mind that semantic type value is numerical. This is different than api solr semantic type value

Subject / Object -> Solr response

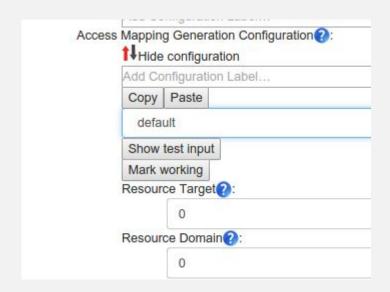


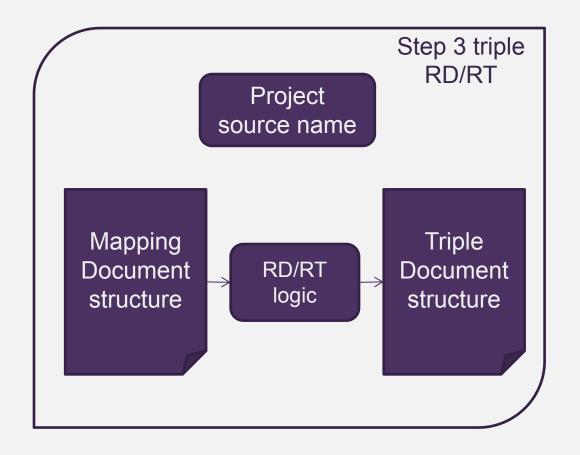
- what to do if (multiple) concept(s) are returned?
- E most-common-id
 - in case multiple ids are found, use the id with most documents
 - when tie, we pick one
- single-id
 - we expect single id to use.
 - if multiple, skip resolution
- multiple-unique-ids
 - in case multiple ids are found, use them all
- random-id
 - in case multiple ids are found, pick one
- homonym-check
 - in case multiple ids are found, pick the most likely one
 - when tie, skip resolution
- Special cases (dont use);
- semantic-category
 - when match, return semantic category
- E phenotype-data
 - when match, return phenotype data
- E term
 - when match, return preferred term



Access Mapping

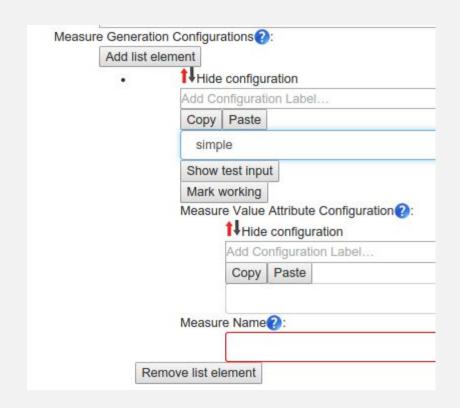


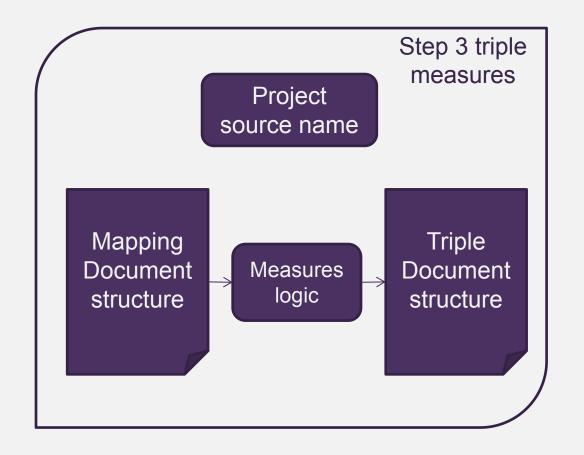




Measures

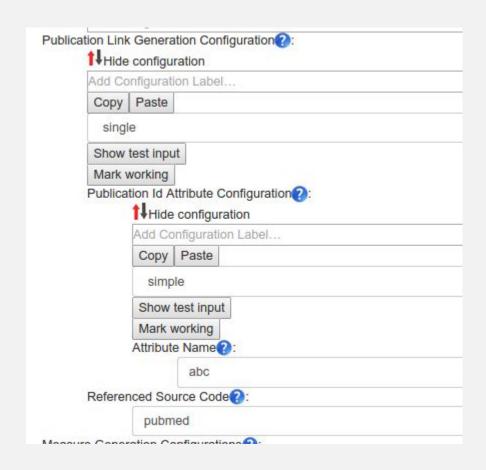


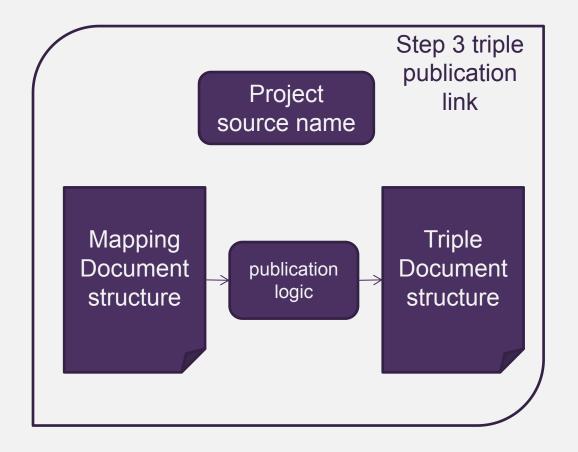




Publication links









Publication generation