

information technology

Research Data Management at RUG

Centre for Information Technology



RUG RDMS

The Research Data Management System at the University of Groningen (RUG RDMS): architecture, solution engines and challenges

A. Tsyganov, S. Stoica, M. Babai, V. Soancatl-Aguilar, J. Mc Farland, G. Strikwerda, M. Klein, V. Boxelaar, A. Pothaar, C. Marocico, J. van den Buijs

(rdms-support@rug.nl)



RUG RDMS



\checkmark Facts and figures

Key figures

Position in international rankings

Annual Reports University of Groningen

Institutional budget

Allocation of funds

Key figures

- Founded in 1614
- 11 faculties (1 in the Frisian capital of Leeuwarden)
- 6,250 fte staff (including University Medical Center Groningen, UMCG)
- More than 120 nationalities
- 140,000 alumni
- 1,000 million EUR budget

Education (2020)

- 34,000 students
- 8,250 international students (24%)
- > 120 English-taught Master's degree programmes
- > 45 Bachelor's degree programmes, of which >35 taught in English

Research (2020)

- 3,600 fte academic staff (39% international)
- 425 full professors
- 4,350 PhD candidates (51% international)
- 591 PhD theses
- 8,000 research publications (dissertations not included)
- 18 patent applications

Excellent prizes and grants (2015-2020)

Nobel Prize for Ben Feringa



 center for information technology

https://www.rug.nl/about-ug/profile/facts-and-figures/



RUG RDMS: Agenda

- Storage architecture
- Application design
- Web Interface
- iRODS rules
- Custom policies engine
- Metadata templates
- Auditing
- What's next?
- Q&A





RUG RDMS



Storage Architecture



5

- our own "Google Drive" for Research-data (keep data on-premises)
- object store instead of old-school filesystems
- self service: user can set ACL's on files without helpdesk
- self service: user can update/add metadata to files
- home-data, group-data, project-data
- files/access can be audited (rabbitmq/kibana/ELK-stack)
- files are replicated (2 copies, spread between 3 locations)





3 'Big' storage-servers (located in 3 different DC's)

Ubuntu/ZFS: fancy/hip filesystem (data-deduplication, compression, 'self'-healing, software RAID

root@gargantua0:~# zpool status | grep sd | wc -l 90

root@gargantua0:~# zpool list NAME SIZE ALLOC FREE EXPANDSZ FRAG CAP DEDUP HEALTH pool0 616T 97.1T 519T - 0% 15% 1.00x ONLINE





IRODS Zones:

- rugZone: big zone for current projects
- umcgZone: dedicated zone for UMCG (storage on datahandeling)
- sramZone: dedicated (test) zone for SRAM (external users test)
- testZone: dedicated test zone for user-test/workshops
- devZone: dedicated dev zone for developers/playground

More/extra Zones are VM's, so they can run on the same central Infrastructure.



RDMS backend should survive..

- broken disks (hot-swappable)
- broken volume (files are replicated)
- broken server(s)
- 1 downed data-centre
- small scale network outage
- maintenance...
- resc-vm's are small (move/migrated if needed)
- scalable workload (load-balancers, multi-site)
- High Performance Computing workloads





Authentication:

- PAM-based authentication stack
- login:email/university-password
- RUG-LDAP (rug-users)
- SRAM-LDAP (for external users)







Security best-practises:

- no root ssh login on systems
- SSL everything
- only 2 admins (admin-group) can 'sudo'
- ssh-login only from admin-network (BWP-lan)
- host-based firewall (iptables)
- network based central firewall (Palo Alto)
- stealth (intruder detection/file integrity)
- fail2ban on webdav
- todo/nice-to-have: fail2ban on iRODS
- monthly security-audit on systems by sec team
- external company pentesting (yearly)



university of

groningen



RUG RDMS



Application design



- Modularity
 - Different applications can make use of the system through an API
- Abstraction
 - Functionality is split hierarchically among the different layers: iRODS, mid-tier, front-end;
 - iRODS layer implements basic functionality and exposes it to higher layers through an API
- Containerization: Docker for services, VMs for iRODS/resource servers
- Robust: vertically and horizontally

























Metadata

Mutable metadata

Immutable metadata

- User generated metadata from Web interface
 User generated metadata from command line interface
- Automatic extracted metadata
- System metadata

System metadata is a number of metadata attributes in different formats that are used by the RDMS to control data workflows

- Mutable lists of metadata (id's)
- Object definition: project/dataset etc..





RUG RDMS



Web Interface



		5	4	C	Filt	lter d	lata											÷
0	Filter by name								rug/	home/a	a.tsyga	anov@rug	.nl/				0	+
	★ Starred	N	ame											Size	Creator		Modifi	ied
**	CBTest	CI CI	BTest											-	a.tsygano	v@rug.nl	Oct, 22	2 2019
	 test UMCG_PILOT1_SHARE 	🗖 te	est											ж	irods		Oct, 22	2 2019
	 Projects Team drives 	D U1	CG_PILOT	1_SHARE										-	a.tsygano	v@rug.nl	Jun, 02	2 2021
	• Shared	B st	ickers.j	iba										2.1 MB	a.tsygano	v@rug.nl	Oct, 29	9 2020
		🖹 te	est.txt											200 B	a.tsygano	v@rug.nl	Oct, 29	9 2020
の																		
t t t t t t t t t t t t t t t t t t t																		
t‡ € 8																		

	1	Main menu	Ð																			Ľ
	•		5	4	G	Filt	ilter data														₹	(
	8	Data browser							r	ug/ho	me/a.	tsyga	nov@ru	ug.nl/						0	+	
	**	Groups and Team drives		Name												Size	Creat	or		Modifie	ed	
		CBTest		CBTest												-	a.tsyg	anov@rug	.nl	Oct, 22	2019	
		Projects Mod_PILOT1_SHARE		test												-	irods			Oct, 22	2019	
		 Projects Paream drives 		UMCG_PILO	T1_SHARE	RE										-	a.tsyg	anov@rug	.nl	Jun, 02	2021	
		Shared		stickers.	jba											2.1 MB	a.tsyg	anov@rug	.nl	Oct, 29	2020	
		J. J	8	test.txt												200 B	a.tsyg	anov@rug	.nl	Oct, 29	2020	
	는 1 부부	Metadata Templates Settings																				
	.5	Backgroud																				
	Ç	processes/jobs																				
	\geq	Contact support																				
	θ	User profile																				
Jer .	0																	L	Logout			
0	https://re	search.web.rug.nl/rdmswebapp/#																				

ALC: N	My folder a.tsyganov@rug.nl	Ð							
•		5	Ļ	G	Filter data				÷
8	Filter by name				<pre>rug/home/a.tsyganov@rug.nl/</pre>			0	+
**	▼ ■My folder		Name			Size	Creator	Modifie	ed
	CBTest		CBTest			_	a.tsyganov@rug.nl	Oct, 22	201
	 Lest UMCG_PILOT1_SHARE 		test			-	irods	Oct, 22	201
	 Projects Team drives 		UMCG_PILC	DT1_SHARE		-	a.tsyganov@rug.nl	Jun, 02	202
	> Shared	6	stickers	jba		2.1 MB	a.tsyganov@rug.nl	Oct, 29	202
		B	test.txt			200 B	a.tsyganov@rug.nl	Oct, 29	202
症									
¢									
\geq									
θ									
0									

2

(1)



X

*	My folder a.tsyganov@rug.nl	Ð			~	Contra 100									
•	Filter by name	5	, Ly		G	Filter data								•	_
-	★ Starred		Name				rug/	nome/a.ts	yganoverug.ni/			Size	Creator	Modif	+ ied
**	 My folder CBTest 		CBTes	t								-	a.tsyganov@rug.nl	Oct, 2	2 2019
	 test UMCG_PILOT1_SHARE 		test									<u></u>	irods	Oct, 2	2 2019
	 Projects Param drives 		UMCG_	PILOTI	L_SHARE							-	a.tsyganov@rug.nl	Jun, (2 2021
	> C Shared	< E	stick	ers.jp	og		Show/Hide details			i *	0	2.1 MB	a.tsyganov@rug.nl	0ct, 2	9 2020
		Đ	test.	txt			Rename					200 B	a.tsyganov@rug.nl	0ct, 2	9 2020
							Make readonly (keep o	wnership)							
							Set inheritance optic	n							
							Add User				. 4 . 1				
							Add Project			Rigi		IICK			
							Add Group			 action 	ons	men	u		
							Add Metadata								
							Load metadata								
	Fxtr	act					Add metadata template								
站	moto		_				Extract Metadata								
φ	meta	uala	a												
\sim															
-															
θ															
0															
-															



	100	My folder a.tsyganov@rug.nl	(: stic	kers.jpg	file	×	Q [7]
	•	Territor I	r 4	C Filter data		4	Permissi	on C	OWN		
	8	Filter by name		rug/home/a.t	syganov@rug.nl/	• +	Inherite	nce f	false		
		* Starred	Nam	ie			$\langle \rangle$	A = B	0		
	**	My folder	CBTe	set							
		test					Name		Value		
		UMCG_PILOT1_SHARE	🗖 test	1			EXIF	ApertureValue	7983/3509		
		Projects		C DILOUI CUNDE			EXIF	BrightnessValue	2632/897		
		> 🔁 Team drives	L UMCC	PILOTI_SHARE			EXIF	ColorSpace	sRGB		
		Shared	stic	ckers.jpg			EXIF	ComponentsConfigurat	ion YCbCr		
			S				EXIF	DateTimeDigitized	2016:06:03 13:33:20		
			E test	t.txt			EXIF	DateTimeOriginal	2016:06:03 13:33:20		
							EXIF	ExifImageLength	3024		
							EXIF	ExifImageWidth	4032		
							EXIF	ExifVersion	0221		
							EXIF	ExposureBiasValue	0		
							EXIF	ExposureMode	Auto Exposure		
							EXIF	ExposureProgram	Program Normal		
							EXIF	ExposureTime	1/30		
							EXIF	FNumber	11/5		
							EXIF	Flash	Flash did not fire, com mode	pulsory flash 🖌	
							EXIF	FlashPixVersion	0100	1	
				Extracted	Ч		EXIF	FocalLength	83/20	1	
					u		EXIF	FocalLengthIn35mmFil:	.m 29	1	
				metadata	а		EXIF	ISOSpeedRatings	64	1	
	\$				-		EXIF	LensMake	Apple	/	
							EXIF	LensModel	iPhone 6s back camera 4	.15mm f/2.2 🖍	
	\sim						EXIF	LensSpecification	[83/20, 83/20, 11/5, 11	/5]	
	-						EXIF	MeteringMode	Pattern	1	
	•						EXIF	SceneCaptureType	Standard	1	
	0						EXIF	SceneType	Directly Photographed	1	
rer							EXIF	SensingMethod	One-chip color area	/	
1	0						EXIF	WhiteBalance	Auto	/ (Ċ
							GPS	GPSAltitudeRef	0	/	
							GPS	GPSI.atitudeRef	N	1	_

RUG RDMS



iRODS rules



RUG RDMS: iRODS rules

- Custom rules:
 - delayed rules and monitoring
 - copy data with permissions & metadata
 - automatic metadata extraction
 - custom policies engine
 - audit and permission control extra rules
- Rules to support tape library archiving





RUG RDMS



Custom policies engine















```
/testZone/home/Projects/project0_5n1 => one direction hash transformation =>
     => 2c7197f0a89e1c842180756537534a81a069be79e8ec6fa1473af21c
   Name:
     sysmdt rdms policy 2c7197f0a89e1c842180756537534a81a069be79e8ec6fa1473af21c
   Value:
     {"policy_name" : "project_user_participation_enddate",
      "policy creator" : "atsG",
      "input_parameters" :
                      "user name" : "atsG",
                      "end date" : "10/06/2021 11:52",
                      "date format" : "%d/%m/%Y %H:%M"
   Unit:
     POLICY PROJECT TORUN atsG
    university of
                    center for
                     information technology
    groningen
```

Policy metadata















ugrdmsdb=> \d RRREGISTERED_POLICIES_SYS

Table "public.rrregistered_policies_sys"

Column	I Туре	Collation	Nullable	Default
<pre>policy_id policy_verification_id</pre>	+ text text	+		
policy_status	character varying(255)	1	I I	
created	l timestamp without time zone	1		now()
changed	l timestamp without time zone	l	1	now()
policy_object	l text	l,	1	





- Code injection protection: if verification is turned on then only policies with the correct verification hash will be triggered.
- Protection from a manual change of the metadata for policies



Server cron job (as rods user)





Server cron job (as rods user)

icommands - get all metadata that are policies to run





Server cron job (as rods user)

icommands - get all metadata that are policies to run

get policy information:

- metadata id
- user that needs to run policy





Server cron job (as rods user)

icommands - get all metadata that are policies to run

get policy information:

- metadata id
- user that needs to run policy

set irods environment username and run asynchronously irules with metadata id as input

university of

groningen



center for information technology







groningen





from .policy_project_user_participation_enddate import run_policy_project_user_participation_enddate

```
C_AVAILABLE_POLICIES = {
    "project_user_participation_enddate" : run_policy_project_user_participation_enddate
}
# method to run policy that we've fetched from the metadata of the object
def run_policy(self):
    l function_name, l parameters = self.parce_policy() # parce json
    if self.policyIsValid():
        if self.c_namespace.C_AVAILABLE_POLICIES.has key(l function_name):
            self.c_namespace.C_AVAILABLE_POLICIES[l_function_name]( self, l_parameters) # execute code
```



RUG RDMS



Metadata Templates



RUG RDMS: Metadata Templates

Why our approach?

- Metadata is the necessary component to transform data into knowledge.
- A well-written metadata template is *crucial* to tag the data being stored
- Several *requests* for metadata template in different research domains:
 - social sciences, archeology, microbiology



RUG RDMS: Metadata Templates

{	
"Publication":	{
{	
"Title": "An overview of RDMS",	
"Authors": [
"RDMS Team Member 1",	
"RDMS Team Member 2"],	
"Date": "2021",	
"Publisher": "RUG",	
"Pages" : "4"	}
}	
}	





RUG RDMS: Metadata Templates

- There is no solution of one-size-fits-all
 - Each research domain has specific metadata-attributes
 - Within the same domain research questions may have different collection of tags
 - The technical skills of the users are very different



Provide the *User* with the possibility to *define* domain/research *specific* metadata templates



RUG RDMS: Metadata Templates - Technical solution





RUG RDMS: Metadata Templates - Technical solution

elf defined structur	e		
Elements	Description	Base type	
City	Enter description here	string	Self defined types
Camera	Enter description here	string	
Angle	Enter description here	Azimuth Angle	
Subject	Enter description here	string	•
Number of images	Enter description here	intger_between_10_and_20	





information technology

RUG RDMS: Metadata Templates - Technical solution

×

Self defined structur	e	Add Metadata Temp	late
		City*	Attributes are validated at input
Elements		Rotterdam	
Name	Description	Camera*	
City	Enter description here	Ladybug5+	
Company		Angle*	<u> </u>
Camera	Enter description nere	211	①
Angle	Enter description here	Ensure this value is less than or equal to 180.	
		Subject*	
Subject	Enter description here	Inspection of asfalt quality	
Number of images	Enter description here	Number of images*	
		10	9



RUG RDMS



Auditing



RUG RDMS: Auditing





RUG RDMS: Auditing



RUG RDMS: Auditing



RUG RDMS: Auditing - Index Lifecycle Management





RUG RDMS: Elastic Search - Index Lifecycle Management





RUG RDMS: Elastic Search - Index Lifecycle Management



RUG RDMS: Elastic Search - Index Lifecycle Management



RUG RDMS: What's next?

- Make code open source & improve deployment
- Integrate switch between different iRODS zones
- Finish customers use cases (currently there are 7 major project running from different faculties)
- Make iRODS NFS mount
- File level versioning
- Multi factor authentication & External users
- Fail2ban
- Windows ingest large files
- Add metadata via mount
- Next version of the metadata templates
- Auditing SLM snapshots





center for information technology



rdms-support@rug.nl

Q&A?



