

Migrating Data When Decommissioning PetaBytes of Storage

John Constable Informatics Support Group, ICT <u>jc18@sanger.ac.uk</u> @kript





Background

- 19PB of genomic data in 399 Resources on 76 resource servers over six Zones
- 41 servers need decommissioning this year, another 10 next year; aka ~10 PB across three types of hardware.
- Generating 10TB/data week, expecting to go up to 760TB if the scientists turn on all the PacBio/Nanopore sequencers they might buy for upcoming programs like 'Tree of Life'





https://docs.irods.org/4.2.6/system_overview/tips_and_tricks/#decom missioning-a-storage-resource

Advice is;

Determine which iRODS server will host the new device.
 Create a new iRODS resource that uses the new device.
 Add the new resource to the appropriate resource hierarchy (could be standalone).

4. Replicate data to the new resource.

- 5. Trim data from the to-be-retired resource.
 - 6. Remove the to-be-retired resource.
- 7. Safely disconnect the to-be-retired device.



You know how, in this over-heated, bullshit world, you just want someone to tell it straight? That's Ben Goldacre's job' CAITLIN MORAN



Ben Goldacre

Selected journalism and other stuff from the bestselling author of BAD SCIENCE and BAD PHARMA



4. Replicate data to the new resource.





Yak Shaving

Any apparently useless activity which, by allowing you to overcome intermediate difficulties, allows you to solve a larger problem.

"I was doing a bit of yak shaving this morning, and it looks like it might have paid off."

Definition credit <u>ghyston.com</u> Photo by Bryan Minear on Unsplash





Standing On The Shoulders Of Giants

This is mostly the work of my colleague Brett Hartley, with input from Terrell and the iRODS team





Solution One: iphymv within a single subtree

"Physically move a file in iRODS to another storage resource.

Note that if the source copy has a checksum value associated with it, a checksum will be computed for the replicated copy and compared with the source value for verification."

(from the man page)





Solution One: iphymv within a single subtree - REJECTED

<u>Issue 4010</u> - "repl to resource with existing replica does nothing"

"Nothing happens. Repl logic short-circuits resource plugins by detecting the good replica and determining that there is nothing to do."





Created by Gi studio from Noun Project



Solution Two: move resource out of hierarchy and then iphymv.

As a bonus, this would also stop new files being written to the resource!





Created by Qi studio from Noun Project



Solution Two: move resource out of hierarchy and then iphymv - REJECTED

In 4.1.x the resource location is stored as a string for each object, e.g.

ils -1 jc18_2G_20170710

jc18 0 root;replicate;seq-red;red4;irods-seq-i21-de 1744830464 2018-04-18.15:11
& jc18_2G_20170710

jc18 1 root;replicate;seq-green;green1;irods-seq-sr01-ddn-ra08-33-34-35 1744830464
2018-04-18.15:11 & jc18_2G_20170710

So every object would need an SQL UPDATE operation. We have hundreds of thousands of objects in each resource *and* it's a one-off, non-resumable operation.



ed by Gi studio Ioun Project



Solution Two: move resource out of hierarchy and then iphymv - REJECTED

Also, we were slightly spooked by

<u>#4402</u> - "renaming resource with substring affects all similarly named resources"





Created by Gi studio from Noun Project Greated by Gi st from Noun Proje



Solution Three: itrim everything off the resource, mark as down, then rebalance





Greated by Gi studio

Created by Gi a from Noun Proj



Solution Three: itrim everything off the resource, mark as down, then rebalance REJECTED

This leaves us with a period of time where each object only has 1 replica, which involves more risk than we were willing to accept.

Oh, and itrim *cowardly and unreasonably* refuses to trim below two objects, especially in a compound tree with two leaves below a replication resource





eated by Qi studio Greated by Qi studio Greater from

Created by Qi studio from Noun Project Greated by Gi studio from Noun Project







Created by Qi studio Cre

Created by Qi st from Noun Proje Created by Gi s from Noun Proje







Created by Gi studio Cre from Nous Project from Greated by Qi stu from Noun Project Greated by Gi studio from Noun Project



BUT!

Issue: 4212 - "iphymv doesn't move file in composite resource tree"

NOW we have Three Copies!

This could be something about our rulebase but...









Greated by Gi studio from Noun Project Greated by Gi studio from Noun Project Created by Gi stur from Noun Project



Member

...

Solution Four: iphymv out of the composite resource, then back in



trel commented on 10 Dec 2018

Whoa, that is exciting.

So, you targeted root;replicate;red;red3;irods-seq-i18-fg and you got a bonus replica in green9 with a new timestamp...





Greated by Gi studio from Nous Project Noun Project

lio Greated by Gi stu from Noun Projec Created by Gi s from Noun Proje



So we need a way to address the three replicas - Brett scripted a tool using the python API (including adding functionality as merge request <u>#162</u>!)



Created by Juan Pablo Bravo from Noun Project





Created by Gi studio

ated by Gi studio a Noun Project ted by Gi studio G Noun Project fr studio Great ject from

Greated by Gi studio from Noun Project



However, files are still being written to the resource, while we drain it.

Solution:

Set minimum_free_space_for_create_in_bytes (See <u>Using free_space check on</u> <u>unixfilesystem resources</u> in the manual) to be slightly larger than the filesystem backing the resource. This ensures that no files can be written to the resource, even once it is emptied.







If you don't already have one, find a resource outside of the composite resource which is large enough to hold the largest file in the retiring resource.

Fortunately, we can use the demoResc's on the IRES's, since even the largest files are *only* 600GB At the moment*, as long as we're careful about parallelisation...







So for each file all we need to do is;

iphymv -M -S \$retiringresourcehierarchy -R \$outsideresource \$file iphymv -M -S \$outsideresource -R root \$file irods-triple-replicas/triples.py \$file tee \$file >> movedfiles.log





Terrell came up with a one liner to do most all of this (adjusted for an attempt at readability)

#!/bin/bash

SIDECAR="demoResc"

HIER TO BE DRAINED="root;replicate;red;red3;irods-seq-i18-bc"

iquest "iphymv -M -S \"\${HIER_TO_BE_DRAINED}\" -R \"\${SIDECAR}\"
\"%s/%s\" && iphymv -M -S \"\${SIDECAR}\" -R "root" \"%s/%s\"; echo
%s/%s > trimmedfile; irods-triple-replicas/triples.py -f trimmedfile;
cat trimmedfile >> movedfiles; rm trimmedfile" "select COLL_NAME,
DATA_NAME, COLL_NAME, DATA_NAME, COLL_NAME, DATA_NAME where
DATA RESC HIER = '\${HIER TO BE DRAINED}'"









Disclaimers:

- We have tested this successfully on our development zone, but have yet to move production data.
- 2. No Yaks were harmed in the making of this talk





Thank you for staying awake listening!

Questions?

<u>Credits!</u>

Brett Hartley, ISG Helen Cousins, ISG for the Yak Photo's *in-situ* Terrell and the iRODS Team <u>Baffalo by Qi studio from the Noun Project</u> <u>Centaur by Eliricon from the Noun Project</u> <u>Superhero by Juan Pablo Bravo from the Noun Project</u> <u>Sidecar By DiabloTim, Oakland (from the Noun Project)</u> <u>Two Yaks Photo by DDP on Unsplash</u>

