



Research Data Migrations and Challenges

Will Schmied, Storage Architect

Opinions expressed are solely my own and do not express the views or opinions of my employer.

St. Jude 101

Cost to family: \$0



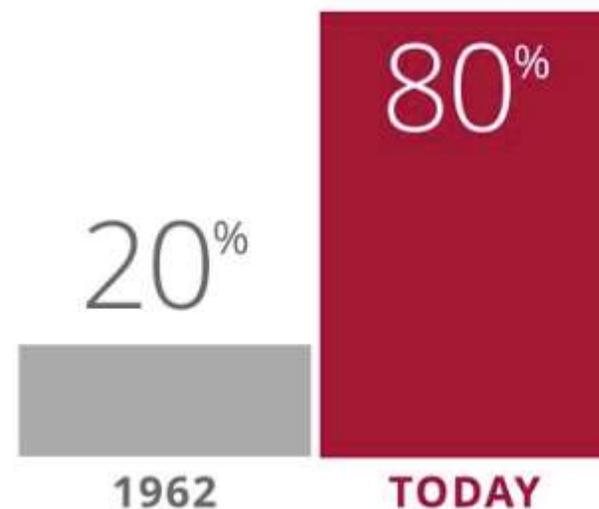
Families never receive a bill from St. Jude for treatment, travel, housing or food – because all a family should worry about is helping their child live.

Saving kids worldwide



St. Jude has treated children from all 50 states and **from around the world.**

Our goal: 100% survival



Treatments invented at St. Jude have helped push the overall **childhood cancer survival rate from 20% to more than 80%** since it opened more than 50 years ago.

A Brief History of Research Storage at St. Jude

- St. Jude has used GPFS / Spectrum Scale since 2011.
 - Before that, various Linux “roll your own” based solutions.
 - This first GPFS system was ~ 1.5 PB usable capacity initially and later grew to a whopping ~2.0 PB.
- Currently on our third generation of Spectrum Scale clusters.
 - The fourth generation is in the installation and configuration phase currently.
- Hardware currently single vendor (DDN).
 - Supporting systems (for LSF, Bright) from Dell/EMC.
- Software licensing directly from IBM.
- I've been managing Spectrum Scale since May 2013.

We Started Out Small...

- In the beginning, there was SoNAS.
 - St. Jude was doing real HPC compute work, but the options in 2011 were few.
 - At one point in time, we had FOUR different SoNAS systems installed on campus.
- By early 2013, we knew we needed native HPC, not HPC over NFS.
 - We were not sold on the ESS at that time, so we built our own TWO clusters on top of DCS 3700 storage.
 - These would serve very specific computational workflows for the *Pediatric Cancer Genome Project*.

...But Quickly Grew (and Grew More!)

- By early 2015, the limitations of our DCS 3700 filesystems were very apparent.
 - Metadata on its own pool needed to be addressed.
 - We also wanted to get out of managing server/storage bare-metal provisioning.
 - These first two GRIDScaler systems were to be an HPC specific cluster and a NAS specific cluster.
 - These were to replace the “research” clusters, first and second generation.
 - Two of the SoNAS were still on site.

The Research Storage Environment Today

- We currently have SEVEN DDN GRIDScaler clusters running in production.
 - SFA 7700, 12K, 14K, 7990, 18K.
 - All are running SS v4.2.3.
 - Largest cluster is ~22.5 PB (37 servers, 297 clients).
 - Total is ~48 PB, 90 servers, 383 native clients, ~4,500 protocols clients.
- Fourth generation cluster is currently being installed and configured, “Jude”.
 - It will replace and consolidate two largest third generation clusters.
 - Planned to be ~38 PB usable by end of CY 2020.
- Remaining smaller SS v4 clusters have residual “capital lifetime”.
 - These will be upgraded in place to SS v5* for mainstream support.
 - More forklift upgrades in the future to get a clean SS v5 filesystem.

A Tale of Way (more than two) Migrations

- We find ourselves migrating often:
 - To meet research compute needs as data sets grow exponentially.
 - Due to undesirable conditions or restrictions within current system.
- Some reasons for previous forklift migrations:
 - *Before* GPFS → SoNAS (Gen1): different foundational technologies
 - SoNAS (Gen1) → Gen2 / Gen3: restricted / closed source environment
 - V3.5 (Gen2) → v4.1 (Gen3): no desire to update problematic source cluster
 - V4.2 (Gen3) → v5.0.4 (Gen4): ensure full v5 functionality to prevent future constraints
- Ironically, our goal since “generation 2” has been to upgrade the filesystem in place and refresh the backing storage via NSD disk migration.
 - **We have yet to accomplish this goal.**

This is What Migrations Usually Feel Like...



Migration \neq Fun

- There are many "migration tools" available.
 - Not all of them are actually *useful* (or safe).
- Some thoughts on the past methods:
 - *Before* GPFS \rightarrow SoNAS: distribution provided rsync. Migration of <1.0 PB.
 - SoNAS (first evacuation) \rightarrow DCS 3700: AFM IW. A brand-new technology! Experienced data loss and corruption. Migration of ~1.5 PB.
 - DCS 3700 (first evacuation) \rightarrow GRIDScaler: AFM SW. Better, not great. Left behind corruption on target. Migration of ~1.5 PB.
 - SoNAS (second evacuation) \rightarrow GRIDScaler: Many parallel rsync jobs (custom rsync binary). Better yet, not spectacular. Slow, but predictable. ~6.0 PB.
 - DCS 3700 (second evacuation) \rightarrow GRIDScaler: Many parallel rsync jobs (as above). Migration of ~1.0 PB.

What's Next for Migration?

- Gen3 → Gen4: Atempo Digital Archive (e.g. Miria) with supplemental (custom binary) rsync in special cases.
 - Good: Easy to manage like rsync, but much faster to transfer data.
 - Good: Choose between keeping source system ACLs or writing new ACLs on destination via inheritance (we have use cases for both options).
 - Good: No additional extended attributes added (as with AFM).
 - Good: No worries about managing gateway node queues and memory (AFM).
 - Mixed: Compute node cutover still must occur at some point. Not an actual problem since they need to be reimaged from SS v4 to SS v5 anyhow.
 - Mixed: Early migration data sets mounted on source cluster over NFS. We will have the opportunity to remote mount the v4 system on v5 later though.

Questions?

