



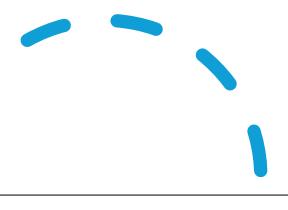
Just FTP It?

Chris Bontempi, University of Connecticut Health Center Network for Advanced NMR (NAN) - NSF Mid-scale RI-2 #1946970 CI Compass Virtual Workshop January 2025



Project Overview

- Nuclear Magnetic Resonance (NMR) Spectrometers
- Network for Advanced NMR (NAN)
 - NSF Mid-scale RI-2 Consortium: Network for Advanced NMR Award #1946970
 - Joint project between UCHC, UGA and UW-Madison
 - 19 spectrometers at 3 facilities, Including first publicly available 1.1GHz class spectrometers in the U.S.
 - Infrastructure/operation evaluated in 6-month workshop engagement with TrustedCl
- Objectives Democratize the use of NMR
 - · Archive NMR data
 - · Good data stewardship
 - FAIR principles
 - Promote NMR facilities and services
 - Provide knowledge bases for various levels of expertise
 - Management tools and a conduit to publishing results
 - Accommodate various fields of use
 - Biology and Biophysics
 - Chemistry
 - Material Science
 - Metabolomics
- PIs
 - Jeff Hoch University of Connecticut Health Center
 - Chad Rienstra University of Wisconsin, Madison (NMRFAM)
 - Kathrine Henzler-Wildman University of Wisconsin, Madison (NMRFAM)
 - Art Edison University of Georgia, Athens (CCRC)













Acquire

Process
data with NMRbox

Share

Participate
in challenges with NUScon

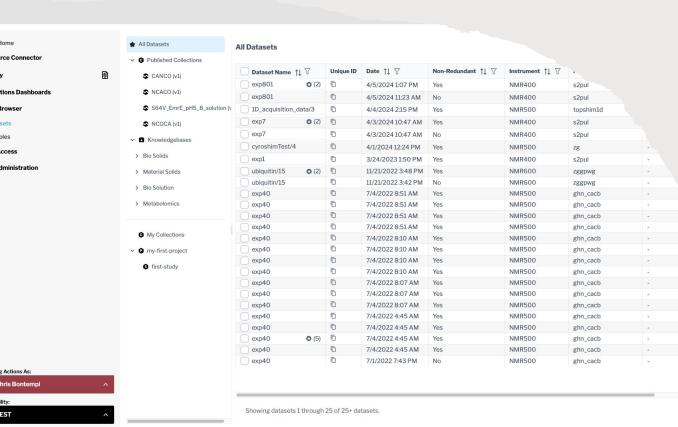
...provides seamless integrations between powerful NMR platforms

Early Considerations for NAN Data Ingest

- NAN is a supplement to facility data storage and backups, not a replacement for it
- We're more likely to get more data if NAN requires minimal extra user steps
 the stretch goal is no extra steps
- Should we keep it simple and just have the user FTP their data to the repository?
 - No automation provides
 - Fewer user steps
 - Consistency of content



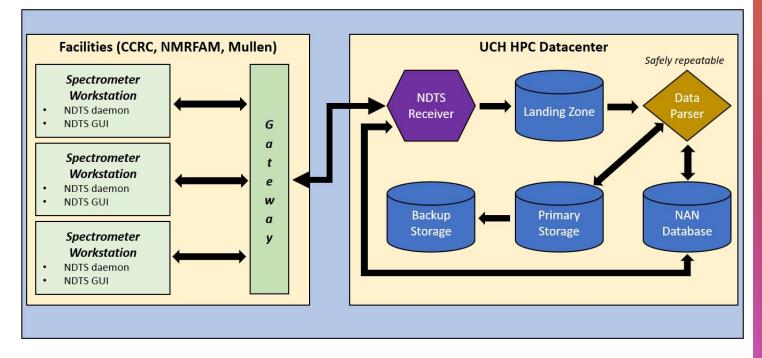
Data Lifecycle



- Pre-configuration settings
 - Specified by the facility manager
 - Specified by the user
- Data acquisition modes
 - Results of acquisitions harvested live from active spectrometer workstations in near real time
 - Historical results manually uploaded from active spectrometer workstations at any time
 - Coming soon archives of results manually uploaded from arbitrary locations
- Received and staged in landing zone
- Parsed into NAN repository
 - Identifier assigned
 - Contents interpreted
- Experimenter and PI of Experimenter can manage, modify, publish, share, and remove experiment data
- Original data is always preserved and archived on WORM devices
- 300K+ experiments harvested to date

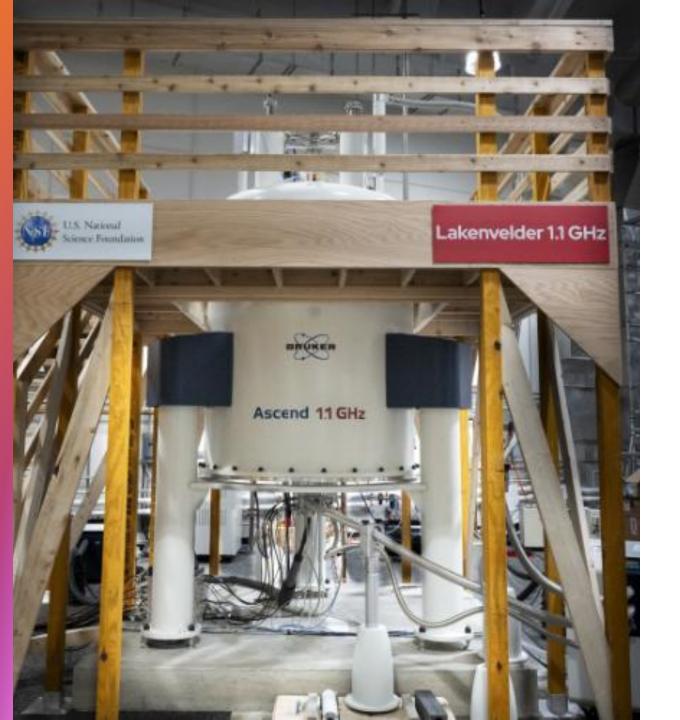
Data Acquisition to First Storage

- Detect/determine configuration and settings
- Detect acquisition start/finish on spectrometer workstation
 - Or trigger a manual upload
- Copy all relevant data to the site gateway
 - TCP socket-based copy
 - If the gateway is unavailable, copy to a temporary location on the spectrometer workstation
- Gateway copies all relevant data to the repository receiver
 - · Secure, reliable transmission
- Arrival of data triggers Parser operations
 - Moves data to Repository
 - Shared network storage device as Linux file system
 - Archiving and WORM takes place from here to Scality array
 - Write metadata to postgresql tables
 - Write visibility statistics to Elasticsearch
- Experiment is now available for download, sharing, publishing, and visibility/operational concerns



Challenges and Lessons Learned

- What is a discrete work unit?
- Unique identifier
- Workflow considerations
- Robust transmission
- Technical Minutiae
- Agile transformation
- Visibility



Identify Discrete Units

- No intrinsic, system-assigned identifier from the software that produces the results
- Of Interest:
 - Individual units single experiment matched set of parameters and results
 - Combined units matched parameters and results, but multiple different samples
 - E.g., metabolomics
 - Calibrating shifting parameters, combined results
 - Transients before and/or after the results of interest
 - User has wide latitude to design the acquisition process with pulse programs, making predicting the exact behavior somewhat challenging
- What are we identifying?
 - The acquisition what if it's multi-part? Does a single acquisition make sense on its own in that case?
 - The upload attempt what if it fails? Do we end up with multiple upload instances of the same acquisition?
 - What is the "real thing"? Are there different "real things" depending on the circumstances?

Considerations

Allow user-selected opt-out

Allow facility-manager-config ured opt-out

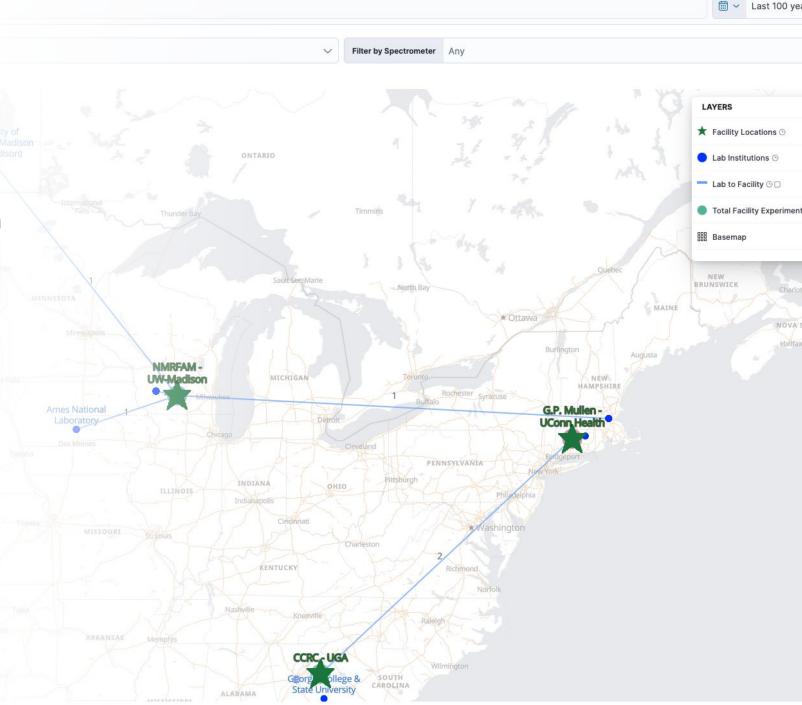
- Specify information that cannot be harvested
- Mapping operating system users to NAN users, samples, projects, etc.
- Which operating system user process is running the experiment?
- Required to detect experiment start/finish

Keeping multiple settings separate and organized (wrong owner problem) Thankfully only one user can run the spectrometer control subsystem at a time, which we can detect



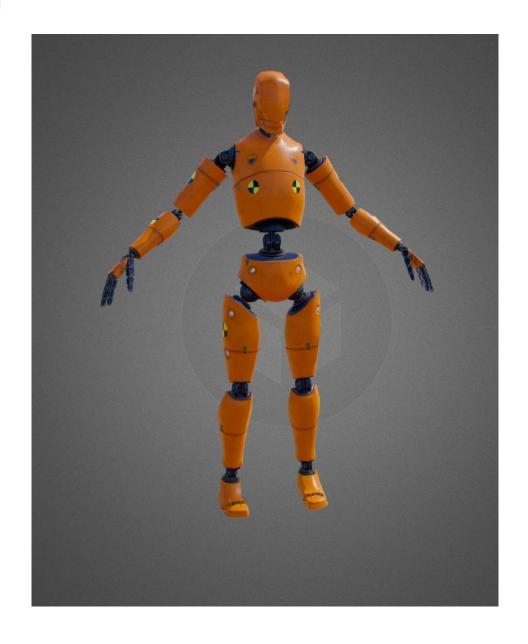
ve

- Ensure every bit gets to the repository
 - Read it before it is overwritten
 - Handle partial results identify broken parts
 - No bit flipping or network loss
- Secure transmission
 - No snooping
- Good citizen behavior
 - Don't attract/allow attacks, etc.
- What about extended downtime?
 - Detect issues that require manual intervention
- Retry until successful or called off
- Never INTERFERE WITH THE DEVICE, storage, the network, etc.
- Detect true duplicates in the landing zone/repository



Technical Minutiae

- Variations in workstation operating systems and supported standards (1990s to 2025)
- Receiver
- Gateway
- Workstation (NAN agent, everything else)
- · Any of the networks involved
- Repository file systems
 - latency
 - read-only
- Databases
 - Maintenance windows
 - Failures
- Program bugs
 - Unexpected conditions
 - Things that break file systems, databases, network protocols
 - E.g., had to tell rsync not to pre-interpret the directory name because of embedded spaces and special characters
- Parsing challenges
 - Missing stuff
 - Extra stuff
 - Unexpected character sets (e.g., latin-1)
 - Non-standard XML, JSON implementations



Recovery Challenges



- Resolved infrastructure failure, followed by a surge of backlogged results – spike in demand
- User expectations around visibility
 - How long until something unexpected is dealt with?
- Use manual transmit for recovery
 - Legitimate provenance?
 - Do parts match (parameter files match results)?
 - Has anything (important) been modified?
- Anticipate conditions that might make data unrecoverable
 - Facility backups

Evolving Transformation

Being able to Re-Parse some or all Experiments



New points of emphasis added



New usage perspectives

A different way of using NMR that we didn't anticipate

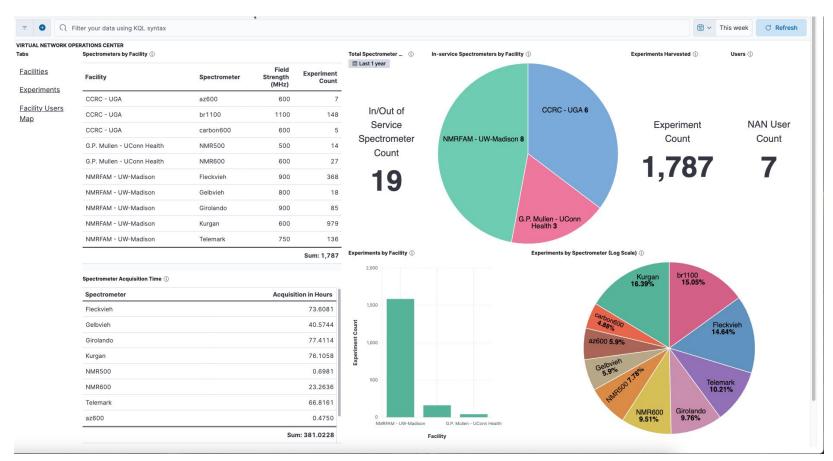


Refined calculations and formulas

Correcting mistakes



Newer workstation software provides new or better results and/or metadata



Visibility

- Instrument every operation/process
- Instrument your data
- Understand what you have
- Understand what you missed
- What is "useless"?
 - Parameters without results/results without parameters?
 - Know if what you have can't provide what most people want
 - But don't assume that everyone wants the same things
 - Look for unexpected patterns

Current and Future Challenges



- Unexpected experiment volume
 - Attention (but not blind adherence) to rational limits of usefulness – throttling?
- Complex experiment groupings
 - What can we safely demote?
- Scaling up to more facilities
 - Well-positioned
- Expanding to multiple repositories
 - Data separation (e.g., Europe and North America)
 - Localized visibility
 - Consolidated visibility
- Continue to root out single points of failure