

A Framework for Telescope Data Quality Management

N. M. Radziwill and R. F. DuPlain¹

*National Radio Astronomy Observatory, P.O. Box 2, Green Bank, WV,
24944*

Abstract. Whether an astronomer performs his or her own observations, or uses previously observed data extracted from an archive, the question "Can I meet my scientific objectives using this data?" must be answered. To do this, four categories of quality should be evaluated: observing system quality, observability, raw data quality and derived data quality. Devices and individual software modules must be operating properly to complete the observation when it is scheduled. The raw data must meet global standards for completeness and accuracy, and specific standards set by the astronomer. Additionally, the datasets must be accessible and sufficiently complete for the science to be performed. A framework for managing quality based on data quality rules is presented, derived from an 18-month prototyping exercise at the Robert C. Byrd Green Bank Telescope (GBT) in Green Bank, WV. A primary goal was to retain the astronomer's control over his or her evaluations. These quality management principles could be extended to any telescope or data production facility for which an interpretation of data is critical.

1. Introduction

Data quality is a pivotal issue for telescopes, influencing various practices including service observing, automated data processing, proposal planning, time allocation, and archive-based research. When data is taken on behalf of an observer, it must meet specified quality objectives to be a valuable use of the observer's allocated time. Producing automated "science quality" data products depends on accurate characterization and assessment of data quality objectives. A telescope's capacity to produce data with desired characteristics impacts the content of a proposal, the time requested, and the time ultimately allocated.

A researcher using an archive may have a different scientific intent than the original observer, and must be able to assess the quality of a dataset according to his or her own objectives. As a result, a dataset may be high quality to one researcher, but low quality to another. According to Silva & Peron (2004), "the usefulness and quality of any given science data product is in the eye of the beholder, ie. it strongly depends on the technical needs and science objectives of the end-user astronomer or data analysis team."

¹University of Cincinnati, 2600 Clifton Avenue, Cincinnati, OH, 45221

Despite their fundamental role, quality issues are typically addressed in a fragmented manner; control system trending and data quality assessment, for example, are treated independently. Quality management takes the needed systemic approach to reduce variation, improve productivity, and streamline processes, applying quality principles simultaneously to telescope operations and observational datasets. Data quality is relative to the researcher's intent, and impacts many parts of the "end to end" process from proposal preparation through offline processing. From the perspective of operations, how do we gather and use this information in the most effective and value-adding way?

2. Goals

The definition of quality as "fitness for use" (Juran & Gryna 1988) was employed because it relates to operational goals: ensuring that the dataset is fit for use by the observer and that the telescope was fit for use at the scheduled time. Furthermore, the dataset should potentially be fit for use from the perspective of a future researcher using an archive. Operationally, quality assessments should feed back into the planning process to help satisfy additional goals, including increasing the operational availability of the telescope (recovering from faults quickly by identifying and repairing root causes of problems efficiently). By addressing root cause problems instead of symptoms, support resources are used more effectively. Lending objectivity to downtime metrics also helps prioritize fixes. The cycle time for implementing the most value-adding updates to the instrument's hardware and software can be shortened, and the time to deliver new scientific capabilities can be minimized. The researcher must be able to characterize and assess data quality according to many different objectives, possibly trying on different strategies to assess the same dataset. Automated assessments must not take control away from the scientist, who may not be able to succinctly define or describe complex quality characteristics.

3. Quality Categorization

In the case of a perfectly operating instrument (ie. no unnecessary latencies), the ultimate measure of its scientific productivity will be a) throughput of observations and b) the quality of the output data. If a dataset can meet the objectives of many observers (in addition to the original observer) its value as a data product is enhanced. To evaluate whether a dataset will satisfy the scientist's intent, quality measures in four areas must be evaluated. A similar taxonomy outlined by Hanuschik et al. (2002) did not include the observability category.

1. **Observing Systems Quality** All equipment needed for the observation must be installed, fully operable, and communications paths must be operational between the hardware and the software. This is a system dependent criterion, and is resolved by a) maximizing operational availability while b) minimizing the time required to identify and solve root causes of faults.
2. **Observability** Even if the system is equipped to execute an observation at a particular time, this is immaterial if the source is not above the horizon, if pointing is inadequate (as in the case of high winds), etc. This is a

time dependent criterion, and is resolved by a) holding an appropriately-sized inventory of executable observations and b) scheduling policies and algorithms that execute these observations at times that they will produce data at the desired quality level.

3. **Raw Data Quality** This information is needed to control software behavior, e.g. not launching a pipeline process if the data is insufficient. Availability, completeness, and precision are all characteristics of the raw data. Evaluations are device or observing-mode specific, and stewardship of the quality requirements can remain with the telescope scientific staff.
4. **Derived Data Quality** These criteria all involve scientific evaluation of the data, for example, rms of spectra, signal to noise, and dynamic range. Quality checks at this level are crafted based on the observer's intent.

4. A Data Quality Management Framework

The discipline of quality management involves establishing the structures, policies and guidelines that enable an organization to meet its quality goals. The following structures for a telescope quality management program were identified:

1. **Establish Data Quality Policies.** For example, a policy describing what data should be irreversibly blanked, what should be masked or reversibly blanked, what should be reversibly flagged, and at which stages of the raw data production process each type of editing should be applied.
2. **Understand Dependencies in the Data Production Process.** To solve the right operational problems at the most appropriate times, a framework for root cause analysis should be in place. For this GBT, this involved constructing a two-dimensional model of component dependencies versus dependencies in classifications of problems encountered.
3. **Implement a Data Store for Quality Information.** Many software applications require quality checks, especially of the raw data, to determine appropriate behaviors. In many telescope systems, quality checks are duplicated in the software. This redundancy is eliminated when quality information is centrally managed. A quality database should be independent of any databases used during the observing process to minimize fault potential.
4. **Provide Continuous Monitoring of Critical Infrastructure and Control System Monitor Points.** By detecting instrument or observing mode failures before they occur, and adjusting schedules or scheduling algorithms appropriately, a loss of telescope availability can be prevented. Several ESO telescopes already implement this practice successfully.
5. **Apply Rulesets for Evaluating the Quality of Data During Production.** The end result of this step is that a dataset exists for which a set of global and observation-specific assessments are true. Pipeline heuristics fall in this category; they seek to automatically translate observer intent into data quality rules (Loshin 2001).
6. **Allow Researcher to Apply Additional Rulesets Based on Intent/Apply Algorithms to Assess Derived Data Quality.** A library of additional rules (based on a science data model) should be available to

enable viewing the data through different filters of quality, enabling the subjective assessments that will always be required from experts.

7. **Close the Loop Between Data & the Production Process.** In addition to revealing faults in the production process, quality knowledge will continually increase throughout the lifetime of an instrument and its observing modes. The goal of an autonomic system whereby knowledge of how to improve quality is automatically learned and applied to the production process represents a leap forward in continuous quality improvement.

5. Conclusions and Future Research

Quality control involves applying data quality attributes to datasets to determine their viability; quality management measures, controls, and automates the continuous improvement process to generate higher quality science products more readily. Data quality is a systemwide issue, not limited to the science products or raw datasets. The control system must be functioning properly and the raw data must meet quality objectives to ensure accurate software behavior. Systemic management of data quality can yield **operational efficiencies**; in software, this is the ability to code less (because quality checks are performed centrally) and fix bugs faster (because the root causes of errors are more readily determined). Data quality assessment has a subjective component, which can be addressed by designing multiple tiers of rulesets. Some of these will be applied by the system (e.g. by the pipeline) and some manually by the researcher. Raw data diagnostics should be implemented first because the quality of operations can often be inferred by indicators within the raw data. Ongoing work is examining a) the benefits realized from feedback between raw data quality diagnostics and its production by the control system, and b) the utility of stratifying rules (global, observation dependent, and user-defined). Quantitatively understanding the interactions between quality information in the control system, pipeline, and archive is the goal. By shifting the focus from quality control to quality management, the next wave of innovations for optimizing the scientific productivity of a telescope could be identified.

References

- Hanuschik, R., Smoker, J., Kaufer, A., Palsa, R., Kiesgen, M., 2002. Quality Control of VLT FLAMES/GIRAFFE Data. Proc. SPIE 5493, Optimizing Scientific Return for Astronomy through Information Technologies 2002, 564-573, ed. P. J. Quinn & A. Bridger (Bellingham, WA: SPIE)
- Juran, J. M. & Gryna, F. M., 1988. Juran's Quality Control Handbook, 4th Ed. (New York: McGraw-Hill, Inc.)
- Loshin, D., 2001. Enterprise Knowledge Management: The Data Quality Approach. (San Diego: Academic Press).
- Silva, D. & Peron, M. VLT Science Products Produced by Pipelines. The Messenger, **118**, Dec. 2004.