

Design of the Data Management System for the protoDUNE Experiment

M. Potekhin^a, B. Viren^a, S. Fuess^b, O. Gutsche^b, R. Illingworth^b,
M. Mengel^b, A. Norman^b

^a*Brookhaven National Laboratory, Upton NY*

^b*Fermi National Accelerator Laboratory, Batavia IL*

Abstract

The protoDUNE experimental program is designed to test and validate the technologies and final designs that will be applied to the construction of the DUNE detector at the Sanford Underground Research Facility (SURF). The protoDUNE detectors will be run in a dedicated beam line at the CERN SPS accelerator complex.

Keywords: DAQ, Data Management, DUNE

1. ProtoDune Program and Detectors

2 The protoDUNE program will help validate various DUNE technology as-
3 pects before proceeding with the construction of the principal DUNE detec-
4 tors at SURF. It is designed for measurements with a test beam provided
5 by a dedicated target and beamline system at the CERN SPS accelerator
6 complex. It also has the potential to be an important platform for realistic
7 LArTPC detector characterization (e.g. PID, shower response, etc.) uti-
8 lizing controlled conditions of a test-beam experimental setup. The name
9 protoDUNE currently applies to two full-scale LArTPC prototypes based on
10 two different technologies. The full-scale designation is used to describe the
11 fact that the prototypes contain important (and large) structural and read-
12 out elements built according to the specifications (including the size) of the
13 eventual full detector.

14 The single-phase (SP) LArTPC functions without amplification in the
15 medium (liquid Argon) and is in essence a very large ionization chamber
16 equipped with a large number of readout electrodes (wires), each with its

17 own electronics chain. In this design, the front-end electronics is situated
18 within the cryostat in order to minimize noise (the so-called cold electronics
19 design). In the dual-phase (DP) TPC ionization electrons are extracted
20 from the liquid into the gaseous phase of Argon, and drift in Argon gas
21 towards a specially designed 2D structure on top of the detector where they
22 multiply according to principles of proportional chamber operation. The two
23 designs are complementary in the sense they explore different approaches to
24 optimization of the Liquid Argon detector characteristics.

25 In December of 2015 the dual-phase prototype was given the official design-
26 nation as a CERN experiment NP02, and the single-phase was designated as
27 NP04. Both are to be deployed at CERN in 2017 and scheduled to take data
28 in 2018. The prototypes will be placed in a specially constructed large-scale
29 extension of the existing experimental hall located in the CERN North Area.
30 Each prototype will be provided a dedicated optical fiber network connection
31 to the CERN central storage facilities located in the West Area campus of
32 CERN. The nominal bandwidth of these dedicated network connections will
33 be 20 Gbps for each experiment. The motivations for this specific choice of
34 nominal bandwidth will be presented in the following sections.

35 *1.1. protoDUNE Data Characteristics*

36 In order to provide the necessary precision for reconstruction of the ion-
37 ization patterns in the LArTPC, both single-phase and dual-phase designs
38 share the same fundamental characteristics:

- 39 • High spatial granularity of readout (e.g. the electrode pattern), and
40 the resulting high channel count

- 41 • High digitization frequency (which is essential to ensure a precise po-
42 sition measurement along the drift direction)

43 Another common factor in both designs is the relatively slow drift ve-
44 locity of electrons in Liquid Argon, which is of the order of millimeters per
45 microsecond, depending on the drift volume voltage and other parameters.
46 This leads to a substantial readout window (of the order of milliseconds)
47 required to collect all of the ionization in the Liquid Argon volume due the
48 event of interest. Even though the readout times are substantially different
49 in the two designs, the net effect is similar. The high digitization frequency
50 in every channel (as explained above) leads to a considerable amount of data

51 per event. Each event is comparable in size to a high-resolution digital pho-
52 tograph.

53 As will be shown in the following sections, it is foreseen that the total
54 amount of data to be produced by the protoDUNE detectors will be of the
55 order of a couple of petabytes (including commissioning runs with cosmic
56 rays). Instantaneous and average data rates in the data transmission chain
57 are expected to be substantial. For these reasons, capturing data streams
58 generated by the protoDUNE DAQ systems, buffering of the data, perform-
59 ing fast QA analysis, and transporting the data to sites external to CERN
60 for processing (e.g. FNAL, BNL, etc.) requires significant resources and
61 adequate planning.

62 *1.2. Prioritization*

63 All of the many elements in the chain of data acquisition, storage, distri-
64 bution and processing are critically important to derive physics results from
65 the data. At the same time, certain components of the data chain need to
66 be prioritized over others in order to perform the measurements during a
67 potentially limited time period.

68 The priority components are the DAQ and the Raw Data Management
69 System, which includes capturing the data coming out of the DAQ, trans-
70 porting the data to persistent mass storage and prompt Quality Assurance
71 which is required to ensure corrective action can be taken if the detector or
72 certain system problems are identified in the QA process. The latter can be
73 thought of as sophisticated monitoring done in near-time, which implies a
74 “few minutes scale of processing.

75 **2. DAQ Interface to Data Handling System**

76 The plan is to have adequate buffering capability in the DAQ for both
77 NP02 and NP04. In this case, adequate indicates in part conforming to a
78 CERN requirement that the experiment must be able to keep taking data for
79 a least 3 days at nominal rate, even if there is an occasional problem with the
80 data link between the experiment site and CERN storage facilities, an issue
81 with central storage or any type of similar outage. At the time of writing,
82 there are differences in two respective approaches:

- 83 • Buffer depth in NP02 is larger, in order to make possible some process-
84 ing right in the data room of the experiment. A number of middleware

85 options are being explored for this storage solution, in particular the
86 BeeGFS file system.

- 87 • In NP04 the emphasis is made on a more lightweight and fault tolerant
88 setup which satisfies the general throughput requirement. No extensive
89 processing is foreseen on the experiment site. Among the technical
90 options for the buffer farm file access is xrootd.

91 It is this outer layer of the data acquisition system in either experiment that
92 will need to be interfaced with protoDUNE's raw data management complex.

93 **3. Requirements**

94 The following is a summary of basic requirements for the protoDUNE data
95 management system:

- 96 • Transfer raw data files from both online disk buffer farms of the DP and
97 SP prototype detectors (NP02 and NP04 respectively) to CERN EOS
98 disk and from there to CERN tape (CASTOR), FNAL tape (Enstore)
99 and other end-points.
- 100 • Ensure that the throughput is adequate and there are no bottlenecks
101 for the Data Acquisition System given the expected data rates over the
102 nominal two month running (see table)
- 103 • Record metadata about file status and outcome of file operations
- 104 • Operate at CERN and FNAL with support for initial setup and ongoing
105 operations
- 106 • Provide monitoring of overall system health, alerts on error and support
107 debugging of problems.
- 108 • Provide triggers to perform file operations (copy, delete) based on con-
109 figurable rules
- 110 • Support express lane process at CERN and other institutions.

111 Table 1 contains performance measures driven by the extreme of each DP
112 and SP detector and which the file handling system must accommodate on
113 the assumption that they apply to both detectors.

114 Table 2 The file handling system assumes the following limits will not be
115 exceeded

Performance Benchmark	Single Phase	Dual Phase
Raw data volume	2.5 PB	
Raw file volume	2 M	
Data Rate	20 Gbps	
Latency to reach EOS	10 min	
Latency to reach express lane processing	10 min	

Table 1: Expected protoDUNE performance characteristics.

Performance Benchmark	Single Phase	Dual Phase
Pending files in FTS dropbox	1 PB	
Simultaneous active files in transfer	50,000	
File registrations rate	3600/hr	
File registrations	200,000/day	

Table 2: Performance requirements for the data handling system used in the protoDUNE-experiment.

116 4. System Topology

117 The baseline design for the online data handling system is shown schemat-
118 ically in Fig. 1. This design leverages the technology of the Fermilab File
119 Transfer Service (F-FTS) running in two positions within the CERN com-
120 puting environment. The primary F-FTS systems are homed within or strad-
121 dling the border of each of the detectors DAQ domains (one for single phase
122 and one for dual phase readout detectors) and is configured to run on a system
123 which is able to access the DAQs buffer disks either through a POSIX filesys-
124 tem or a protocol layer (e.g. XROOTD, gridFTP or others). The primary
125 FTS system is configured on this system to look at one or more input direc-
126 tories or storage locations, commonly referred to as a dropboxes. The FTS
127 daemon performs periodic scans of the configured dropboxes. The system will
128 operate asynchronously and independently of the protoDUNEDAQ systems.
129 When the protoDUNEDAQ has produced a file that it wishes to have passed
130 off to the storage system, it will move (perform a filesystem/storage system
131 atomic move operation or the equivalent) the file to the dropbox location.
132 When a new file is located within one of the dropboxes, the system will ini-
133 tiate the file registration and transfer operations. In the protoDUNE model,
134 the primary F-FTS will initiate a copy operation from the online disk buffer
135 farms into the EOS storage system. The system will use the 3rd party copy

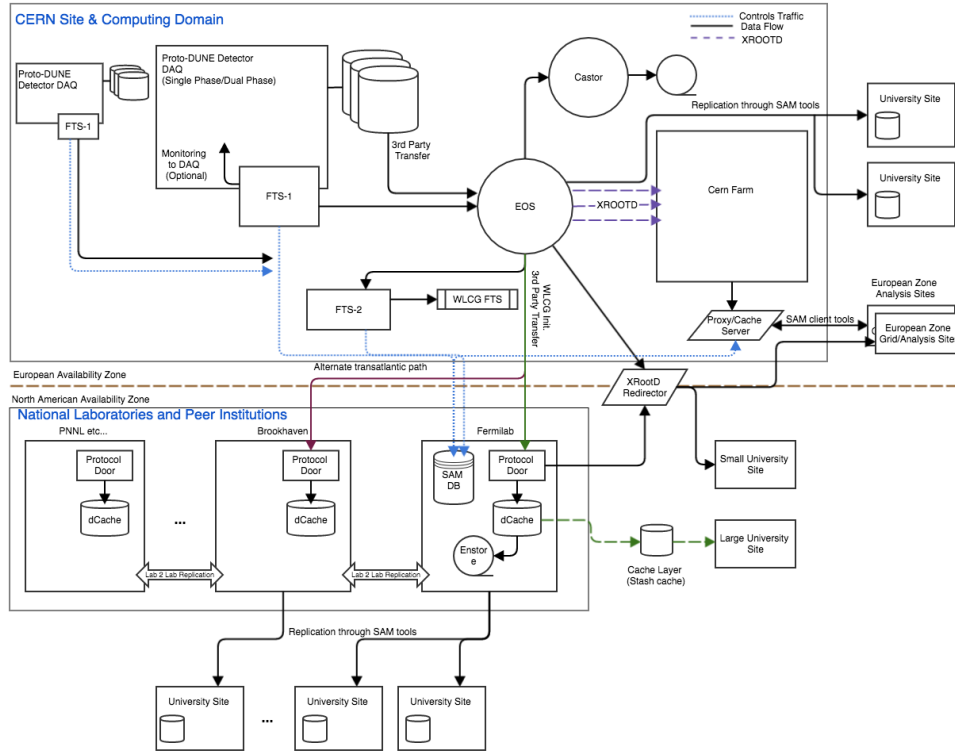


Figure 1: protoDUNE data handling system topology

136 support that is provided by the XROOTD protocol to allow for optimized
 137 transfers into EOS. Upon completion of the initial copy into EOS, the F-
 138 FTS will initiate a chained copy (a copy that is dependent on the initial copy
 139 into EOS) of the data from EOS to the Castor archive system. The F-FTS
 140 system will register each files metadata records (containing both basic and
 141 physics metadata, defined below) in the SAM data handling system. Upon
 142 completion of the replication of the data to Castor, the F-FTS will enter
 143 a monitoring/polling state of its operations to determine when the file has
 144 been successfully written to the archival media.

145 The primary F-FTS will handle the cleanup of its input dropbox. The
 146 F-FTS performs configurable periodic cleanup passes through its current file
 147 sets. The baseline cleanup logic is shown in Fig. 2. This cleanup process
 148 ensures that all files are successfully transferred and archived before they are

149 deleted and provides an age criterion on files so that files can be retained
150 within the DAQ environment for a period of time. This allows operations
151 like online/nearline processing or log files to be examined within the DAQ
152 environment after the actual transfers to archival storage have completed.

153 The secondary F-FTS service runs outside of the detector DAQ domains,
154 but within the CERN computing sphere for maximum efficiency. This F-FTS
155 system is configured to monitor, as its input dropboxes, the EOS locations
156 that the primary DAQ F-FTS systems are using as output endpoints. The
157 system operates in the same manner as the primary F-FTS, scanning for
158 new files and then initiating copy requests of the data to a set of one or more
159 destination endpoints. This F-FTS system will initiate the copy request
160 between the European availability zone, and the North American availabil-
161 ity zone. At least one of the endpoint destinations for this service will be
162 the FNAL-based dCache/Enstore system, where a second archival copy of
163 the data will be recorded. Additional copies across the availability zones
164 can be configured based on available resources and available transatlantic
165 bandwidth (e.g. direct transfer from CERN to BNLS dCache or to PNNLs
166 storage facilities). To perform the transfers, the F-FTS will interface with
167 the WLCG-FTS (as a supported transfer protocol) to schedule the actual
168 transfer between the storage elements. Underlying the F-FTS systems will
169 be a fully featured data management layer (SAM) which will provide the
170 metadata and replica catalogs. All data being handled by the transfer sys-
171 tems will have corresponding records in the data handling catalog so that
172 the content, locations and provenance of the data can be fully tracked. The
173 primary catalog systems will reside at FNAL with proxy/cache layers in the
174 European and North American zones to ensure high speed connections be-
175 tween the servers and the [offline analysis] clients that will query the catalogs
176 from these zones. Similar scalable proxy and cache layers can be instantiated
177 for other identified zones which may require high speed or optimized access
178 to the catalog systems.

179 Replication between collaborating peer institutions in the same geographic
180 zone (i.e. national labs and large universities with significant computing/storage
181 resources in the North American zone) will be provided through standard
182 dataset replication tools that are a part of the SAM data handling suite.
183 Data access at smaller university sites or opportunistic computing resources
184 (e.g. Open Science Grid (OSG) affiliated sites) is provided and optimized
185 through support for streaming protocols such as XROOTD with redirector
186 services and through cache layers such as the stash cache system employed

187 by OSG.

188 **5. Technical Requirements and Specifications**

189 The data management and transfer systems requirements are enumerated
190 in table XX according to the category to which they belong.

191 *5.1. Data Throughput*

192 The goal of the data handling system is to be performant at a level that
193 allows for full exploitation of the underlying hardware and networking infras-
194 tructure upon which it is running. In the case of the proto-Dune experiment
195 the FTS system will be capable of operating at a sustained data processing
196 and transfer rate that is greater than 80% of the maximum theoretical band-
197 width available between the primary DAQ storage systems and the EOS file
198 storage system. In the current design this network bandwidth consists of two
199 20 Gb/s ethernet links. The RAID arrays or distributed file system archi-
200 tectures to which the data will be written to and read from are assumed to
201 have similar or lower available bandwidth. The actual observed read/write
202 bandwidth will be determined by the actual choice of hardware that is de-
203 ployed to the detectors. The primary DAQ FTS system will be required to
204 operate at a sustained effective bandwidth of the lesser of: two times 16 Gb/s
205 (80% of theoretic network bandwidth) or 80% of the measured read access
206 bandwidth from the storage arrays [during simultaneous write operations,
207 if the DAQ computing models requires this mixed access mode] under the
208 standard proto-DUNE operations.

209 *5.2. Transaction Throughput*

210 The goal of the data handling system is to be performant at a level that
211 allows for full exploitation of the underlying data catalog and database in-
212 frastructure while at the same time meeting or exceeding the rate at which
213 files or other data objects are generated by the proto-DUNE detectors. In the
214 case of the proto-DUNE experiment that FTS systems that run on the border
215 of the DAQ domain will be able to support a minimum transaction rate of
216 3600 file registration per hour (1 Hz) and an aggregated rate of 100,000 file
217 registration operations per day. This level of performance has been demon-
218 strated in other online and offline systems where the data handling system
219 has been deployed.

220 *5.3. Transaction Tracking*

221 The goal of the online components of the data handling system are to
222 be able to perform simultaneous end-to-end tracking of all files that are un-
223 dergoing active transport through the system prior to being written/verified
224 in the archival storage systems, without impacting or interrupting data tak-
225 ing. To achieve this goal the online data handling systems for proto-DUNE
226 will be able to support a total number of in flight files (per FTS instance)
227 equal to the 100,000 times the maximum supported duration (in days) of
228 an outage that can be sustained by the DAQ system under outages in the
229 network or storage systems downstream of the DAQ domain. In the case of
230 proto-DUNE the FTS systems will support at least 700,000 in flight files at
231 any given time, where at least 10,000 of these files are undergoing active reg-
232 istration/transport through the system at any given time (meaning at least
233 10k of the files are being actively managed while the remainder are waiting
234 in the designated dropbox location.

235 *5.4. Transfer Latency*

236 The file transfer and management systems are designed to operate in an
237 asynchronous manner to other components of the proto-DUNE DAQ. Due to
238 the polling models employed in this asynchronous model, many operations
239 do not have fixed temporal relations to other events in the DAQ system,
240 but rather will occur/complete within a well defined time window or with a
241 certain latency. In the case of the proto-DUNE experiment, the data man-
242 agement systems will be capable of operating within the following parameters

243 In the first stage of the file transfers, the latency between the completion
244 of the DAQ writing a file and the hand off/detection of the file by the FTS
245 system will be controlled by a configurable delay parameter (in minutes)
246 within the FTS which controls the intervals at which the FTS looks for new
247 files in its dropbox locations. Under normal (steady state) operation, the
248 average latency will be 1-2 polling intervals. The actual latency between
249 when the DAQ finishes writing a given file and when the FTS picks up and
250 begins to actively manage that file can be delayed by the number of other new
251 or pending file currently in the system. In this case the order of handling of
252 files will be handled via an internal queuing algorithm that efficiently handles
253 the files but does not provide any user specified prioritization.

254 In the second stage of file transfers, the latency between when a file
255 is generated by the DAQ and when it is transferred to the EOS system
256 and then written/verified to archival media (via Castor) is dependent on

257 first stage latencies (new file detection and file registration) and then the
258 details of the archival storage system and its characteristics. In particular
259 the DAQ to EOS to Castor path will constitute a set of chained dependencies
260 with an independent polling interval for each stage. Under normal operating
261 conditions that latency for the file to be transferred to EOS and be available
262 through the data handling system is expected to be 1-2 of these polling
263 cycles, contingent on write congestion in the EOS storage system. For full
264 registration of the files in Castor and onto type, the FTS/SAM systems will
265 support latencies in this recording of 3 hours to 30 days and will support a
266 configurable timeout parameter to indicate failures in this transition.

267 The latency between when files are written by the DAQ and when the files
268 are available on storage in the North American zone, will be determined by
269 the second stage latency of the files appearing on the EOS system from the
270 DAQ and then a secondary latency will be incurred based on the polling for
271 files that require transatlantic endpoints. Similar to the other stages this is
272 configurable through a polling interval. Once queued for handling the actual
273 file transport will be off loaded to the WLCG FTS system which will schedule
274 the files for transmission and will properly balance and throttle the site to
275 site traffic and will properly conform to the CERN computing environment.
276 Latencies at this stage are well understood based on the experience that the
277 LHC experiments have with WLCG FTS.

278 **6. Storage System Interface Specifications**

279 The module design of the data handling system provide considerable flex-
280 ibility in its ability to adapt to different requirements for both data input
281 and output storage systems. For the proto-Dune experiment these major in-
282 terface points are considered to be at the interface between the experiments
283 core DAQ components and at the interfaces between the mass storage and
284 archival storage systems at both CERN and FNAL. The interface choices
285 below detail these interfaces.

286 *6.1. DAQ to Data Handling System Interface*

287 The interface between the core DAQ system and the data handling system
288 is made at the DAQs disk buffer farm. When the final stage event builder
289 have completed the assembly of a file, the file will be handed off to the data
290 handling system by the DAQ by placing the file into a designated dropbox
291 area on the target file system. For POSIX compliant file systems, this can be

292 accomplished through an atomic move operation (i.e. Unix mv equivalent)
293 which minimizes the IO overhead associated with the hand off. No signals
294 need be emitted from the DAQ, nor will any other form of message be required
295 in the direction of the DAQ to the data handling system. The F-FTS will
296 detect new files through files becoming visible in configured dropbox locations
297 of the storage system (i.e. a new file appearing in a directory)

298 The F-FTS systems running as part of the data handling system will use
299 either a POSIX (or near POSIX) compliant file system view of the disk buffer
300 farm, or an API that provides access to standard meta information regarding
301 the storage systems dropbox area (directory) and contents (i.e. filenames, file
302 size, create/modify times etc) Additionally the FTS will require a file read
303 API for the purpose of metadata extraction from the file and for checksum
304 computations (if not provided through some other means). The system will
305 require a delete API and privilege, if the automated cleanup options are
306 desired/enabled. If the DAQ storage elements support third party transfer
307 protocols, the FTS can use these to reduced overhead in performing the
308 actual file transfer operations.

309 The organization of the FTS dropbox area(s) will be configured to sup-
310 port the volume of data that would be generated during a sustained outage
311 (of non-DAQ systems) of no less than 3 days plus the time required to reg-
312 ister/process the volume of data that would be generated during the outage
313 (i.e. the recovery time required to clear a backlog).

314 *6.2. Data Handling to EOS Interface*

315 The interface between the data handling system and the EOS storage
316 system will be made through the APIs provided by the standard protocols
317 already supported by both the EOS system and the SAM data handling
318 system. In particular the xrootd protocol will be the primary API used for
319 the interaction between the system, with secondary support for the gridftp,
320 webdav and srm protocol APIs. The EOS system will be declared, and
321 configured in the SAM system as a standard storage element and will have
322 files stored on it registered and mapped within SAM the replica catalog to
323 the well defined access URIs.

324 *6.3. Data Handling to CASTOR Interface*

325 The interface between the data handling system and the CASTOR archival
326 storage system for file ingest will be made through the APIs provided by the
327 standard protocols already supported by both the CASTOR system and the

328 SAM data handling system. In particular these system both already sup-
329 port the xrootd protocol, gridftp and srm. The CASTOR API to query the
330 tape archive interfaces (to determine if a file has been successfully archived
331 to tape and which tape it is located on) will be integrated into the SAM
332 data handling system in a manner similar to other tape systems that SAM
333 is already aware of (i.e. Enstore) The CASTOR system will be declared,
334 and configured in the SAM system as a standard tape storage system and
335 will have files stored on it registered and mapped within SAM the replica
336 catalog to the well defined access URIs along with additional tape location
337 information to permit optimized retrieval. Custom modules/utilities will be
338 developed as needed as part of the SAM data handling suite to provide ad-
339 ditional support for the CASTOR system in performing certain operations
340 (i.e. bulk queries or monitoring operations). These modules/utilities will be
341 distributed as part of the core SAM distribution.

342 *6.4. Data Handling to dCache/Enstore Interface*

343 The interface between the SAM data handling system and the dCache/Enstore
344 archival storage system is already fully defined and supported. The interface
345 fully supports standard access protocols including xrootd, gridftp, webdav
346 and srm as well as the proprietary dcap protocol and a specialized NFS
347 implementation that provides limited basic read access for hosts with local
348 mounts of the dCache system. The SAM system has modules that opti-
349 mize data access based on Enstore tape location information. This interface
350 is in wide scale production use across many experiments including NOvA,
351 MicroBooNE, Dune 35t, Minos and others.

352 *6.5. Data Handling to SAN/NAS Interface*

353 Generalized SAN and NAS systems when acting as data sources (input)
354 will be integrated to the data handling system through their POSIX style
355 interface if available. When enabled as data sinks (output) for the data
356 handling system, a front end data server(s) running standard access protocols
357 in the form of xrootd or gridftp will be utilized. These will then be mapped
358 into SAM as standard storage locations.

359 *6.6. Generalized Protocol Support Interfaces*

360 The data handling system provides a file delivery and transfer layer which
361 acts to provide protocol abstraction to the end users (i.e. provides a con-
362 sistent command interface) and provides a protocol bridge when performing

363 transfer operations between dissimilar storage systems (i.e. cross protocol
364 transfers). The Intensity Frontier Data Handling tool (IFHD) for proto-
365 DUNE will support protocol integration for:

366 XROOTD Full support will be provided for xrootd in the Fermi-SAM
367 data catalog, F-FTS and IFDH layers. Currently these tools
368 support read/write access methods. Full support for additional
369 xrootd features (listings, permissions modification, etc...) is
370 currently under development.

371 WLCG FTS SAM, F-FTS and IFDH will support WLCG-FTS as a 3rd
372 party or proxy transfer mechanism. In this mode outbound
373 file transfers originated from the CERN domain can offloaded
374 to WLCG-FTS so that the data traffic across the CERN site
375 can be properly scheduled and balanced. Support for this mode
376 of operation will be integrated into IFDH or directly into the
377 sam_cp layer of the F-FTS.

378 GridFTP This protocol is fully supported for most of the components
379 involved, but is not the preferred or documented interface for
380 any of the CERN storage components.

381 SRM This protocol is fully supported for most of the component sys-
382 tems involved, but is not the preferred or documented interface
383 specification for the CERN storage components.

384 HTTP Fermi-SAM uses HTTP for internal communication between
385 Fermi-FTS and the Fermi-SAM instances, and for client pro-
386 grams wanting file location and metadata information. It is
387 also a supported file transfer interface into DCache.

388 CVMFS The data handling suite includes support for distributing its
389 client tools and homing other parts of the suite on a CVMFS
390 read-only filesystem. This support allows the experiments to
391 provide widespread, efficient code distribution over a HTTP
392 based protocol while providing a file-system layer to the end
393 user applications. This is a standard distribution method used
394 to to distribute the DUNE and LARSOF T software, as well as
395 Fermi-SAM client utilities.

396 7. Data Replication

397 The data handling systems provide for multiple replicas of files to exist
398 across storage systems and locations. It will be necessary to replicate data
399 files or data sets between sites or storage locations. The data handling sys-
400 tems for protoDUNE will handle this replication in different ways depending
401 on geographic proximity and bandwidth. In particular in the protoDUN-
402 Earchitecture at least two Primary Zones exists in the form of the European
403 zone and the North American zone. The European zone encompasses CERN
404 and other institutions near to CERN geographically and in network band-
405 width proximity, while the North American zone would encompass FNAL,
406 BNL, the other institutions in close network proximity to the major labora-
407 tories and universities in the DUNE collaboration.

408 Replication within the primary zones will be handled in the following
409 manners by the data handling systems:

- 410 • European Zone – Replication and data transfers within the CERN do-
411 main will be handled by F-FTS instances configured with site specific
412 endpoints rules (i.e. the EOS and CASTOR rule sets) which will pro-
413 vide for the full data sets to available in both EOS and CASTOR.
414 Transfer and replication operations outside the CERN domain to insti-
415 tutions with registered storage elements will be performed through the
416 SAM suites replication tools (`sam_clone_dataset`). Transfer or replica-
417 tions to smaller institutions or analysis sites without registered stor-
418 age elements can be performed through the SAM/ifdh client tools
419 (`ifdh_fetch`). The SAM client tools are which are available through
420 CVMFS.
- 421 • North American Zone – Initial replication of the data to FNAL, BNL
422 and other collaborating labs or large university sites will be automated
423 through F-FTS instances at CERN and in North America which will
424 be used to optimize the data transfer paths to each of the institutions,
425 subject to bandwidth constraints of the host institutions. Transfers
426 elsewhere in the North American side would be performed using the
427 same SAM tool suites and strategy for handling registered and unreg-
428 istered storage, as specified above for the European zone.
- 429 • Tertiary Institutions – Other institutions wishing to host a set of datasets
430 or partial data sets, will be able to register their storage elements with

431 the SAM data catalog and will then be able to use the replication tools
432 to clone the relevant data to their site.

433 8. File Registration

434 Registering files or other forms of data with the SAM catalog require a set
435 of meta information about the files. This information is used to allow detailed
436 searches to be performed to select specific data for analysis. This metadata
437 is broken into two general types Base Metadata and Physics Metadata which
438 can be provided at the time of file registration, or modified later. These
439 metadata are composed of:

- 440 • Base Metadata – For each file the the Fermi-SAM database requires a
441 unique filename, the file size, and a file type string. It also allows one
442 or more checksums - each consisting of an algorithm type name and a
443 value - and a description of the application used to create the file. The
444 file can also have zero or more parents (which much be files already
445 existing in the catalog) which creates a tree of relationships between
446 files.
- 447 • Physics Metadata – There are a number of predefined physics metadata
448 fields such as the detector configuration, the run number, the data tier,
449 the event count, the start time and end time, and the data stream. It is
450 also possible to define arbitrary key names for other values with either
451 integer, float (64 bit), or string types.

452 The experiment must provide a method for determining or generating the
453 metadata (i.e. an external program that can be run on a file, or a python
454 module that extracts information from a filename, etc) and that method will
455 be invoked by the data handling system when it encounters appropriate files.
456 The output of the method used must conform to a suitable format supported
457 by the F-FTS and SAM systems. The systems support JSON formatted files
458 and python dictionaries for direct upload to the data catalogs.

459 9. Data and Replica Catalogs

460 The data management system for protoDUNE will leverage the DUNE
461 experiments SAM data catalog. The underlying resources of this catalog
462 and its software are capable of, and sized to support, a file inventory in

463 excess of 150 million files (demonstrated by the D experiments catalog which
464 uses the same technology). The data catalog provides a command line toolset
465 for both client and administrative functions, as well as a web based API for
466 interacting with the catalog. The web tools provide additional guidance for
467 assisting users in finding and classifying data through the indexes.

468 The data management system will also leverage the same DUNE instance
469 of the SAM catalog to maintain its replica catalog. This instance will provide
470 a unique namespace for the DUNE/protoDUNEexperiments and will contain
471 the appropriate storage identifiers and paths to enumerate the full locations
472 of DUNE/protoDUNEdata on all supported sites. The SAM replica catalog
473 itself is protocol neutral, but a mapping layer within the data handling system
474 performs the translation of locations into actual access URLs of the default,
475 preferred or requested protocol schema for a given storage system (i.e. the
476 system maps the location onto the method for how you retrieve the file).

477 **10. Data Transfer Technology**

478 The data management system for protoDUNE will support, through its
479 use the SAM tool suite, data transfers to and from CERN, FNAL and other
480 sites storage elements using a combination of protocols supported by the spe-
481 cific sites storage resources. In particular the system will natively support
482 the xrootd, gsiftp (gridftp), webdav(http), and srm protocols for access to
483 EOS, Castor and dCache/Enstore (as appropriate). The system can sup-
484 port multi-channel transfers and 3rd party transfers using these protocols to
485 achieve high bandwidth or low overhead transfers where needed. The sys-
486 tem also supports operations specific to local access methods on POSIX file
487 systems.

488 The file transport layer uses a modular design with an abstraction layer
489 so that new protocols or access methods can be added to the system without
490 changes to the other interface layers. The system can also select or prioritize
491 the replica source location and access protocol based on destination charac-
492 teristics (i.e. it can prefer a local cache copy of a file to a replica that is at
493 a transatlantic location) The system will also support offloading of transfers
494 to 3rd party transport systems (e.g. WLCG FTS) to properly balance the
495 resources of different WAN connections and site resources.

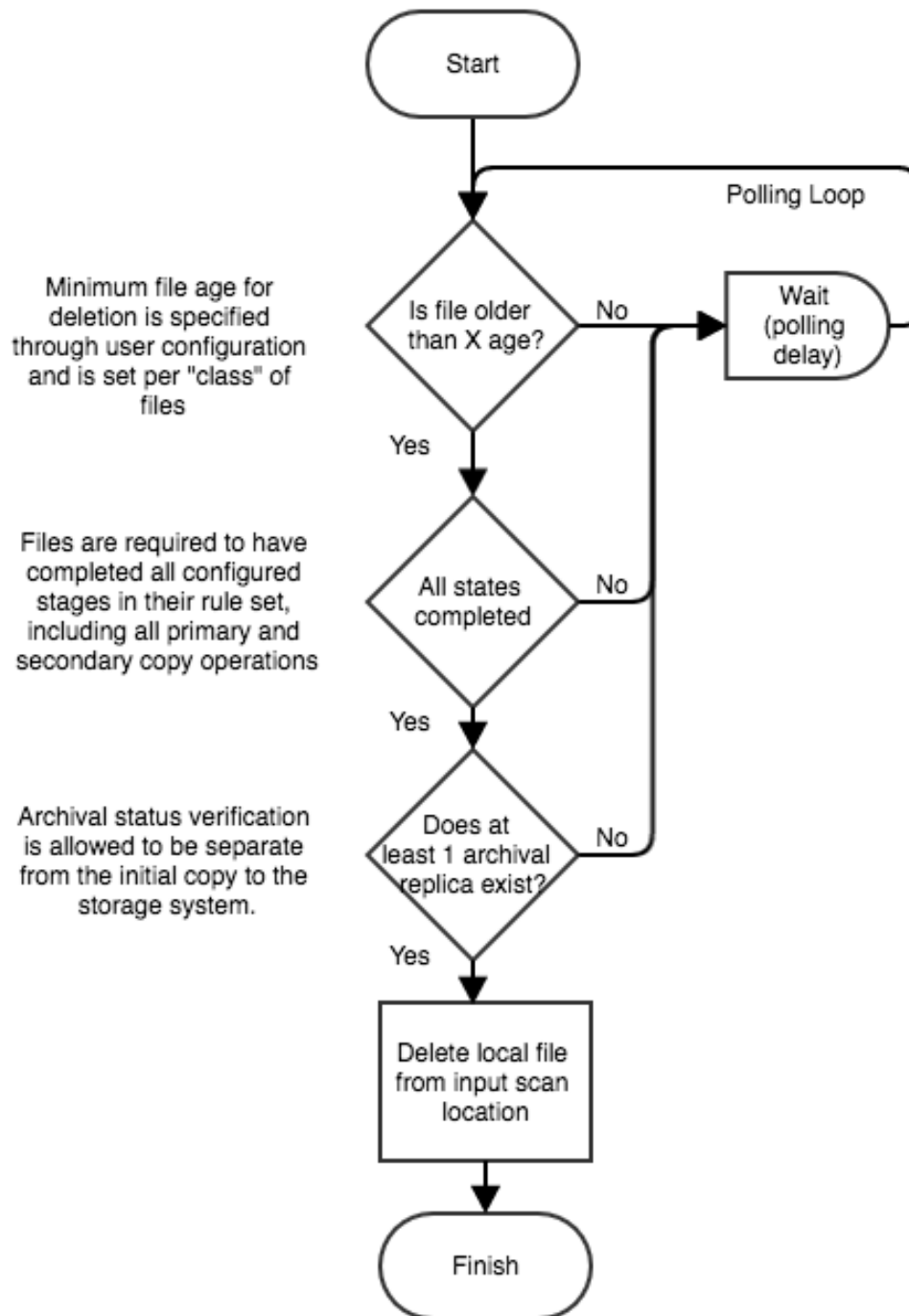


Figure 2: protoDUNEF-FTS File deletion logic and cleanup Flowchart. The F-FTS performs a configurable multi-stage validation procedure to determine if files are eligible for deletion/cleanup from the input area. In particular the system can verify the archival status of files (ensuring they have been written to tape and that tape has been closed/unloaded properly) prior to performing any file deletion.