



Robustness

1 **Contents**

2 **Robustness** **2**

3 Introduction 2

4 Requirements 2

5 Approach 2

6 Application data 2

7 Mitigating the effects of lack of disk space 10

8 Resource management 11

9 USB undervoltage 14

10 Risks 14

11 Design notes 15

12 BTRFS Overview 15

13 BTRFS robustness supporting features 16

14 **Robustness**

15 **Introduction**

16 This design identifies circumstances that, though undesired because of the risk
17 of loss of functionality, cannot be completely avoided and provides suggestions
18 for dealing with them in such a way that as little functionality as possible is
19 lost.

20 Note that improving D-Bus' robustness is a topic that will be covered in a later
21 stage in its own design document. About securing D-Bus services, please see
22 the security design.

23 **Requirements**

24 Minimize loss of data and loss of functionality due to data corruption in these
25 abnormal circumstances:

- 26 • Unexpected power loss
- 27 • Unexpected removal of storage devices
- 28 • Unexpected lack of disk space
- 29 • Physical damage to the media and other hardware errors

30 Minimize loss of functionality due to processes hogging these shared resources:

- 31 • CPU
- 32 • GPU
- 33 • I/O
- 34 • memory

- 35 • network queue
- 36 • D-Bus daemon

37 Approach

38 This section explains how to address the requirements in several specific cases,
39 taking into account different data sets and circumstances.

40 Application data

41 This section contains recommendations about how to robustly deal with data
42 generated by applications.

43 General guidelines

44 No software should assume that opening files will always succeed. Failure condi-
45 tions should be dealt with and the process will either continue running with as
46 little loss of functionality as possible, or will log a message and exit. Programs
47 should do the same when writing data (the filesystem may be full, or any other
48 mode of error might occur).

49 For example, if the browser application finds out at start up that the cookies
50 file is corrupted, it should move the old file away (or just delete it) and run as
51 usual other than past persistent cookies will have been lost. Or if there was
52 an error when writing a new persistent cookie to disk, the browser would keep
53 running with that cookie being transient (in memory only).

54 In order to reduce the effects of data corruption, regardless of the causes, it
55 would make sense to store different data sets in separate files. So that if cor-
56 ruption happens in, for example, the browser cookie store, it would not affect
57 unrelated functionality such as playlists.

58 For big data sets, Collabora recommends SQLite with either Write-Ahead Log-
59 ging (WAL¹) or [roll-back journal](http://www.sqlite.org/draft/lockingv3.html#rollback)². For smaller data sets, a robust method is
60 to write to a temporary file and rename it on top of the old one once finished.
61 This method is called “atomic overwrite-by-rename” and is mostly used when
62 editing a file in-place.

63 POSIX requires the atomicity of [overwrite-by-rename](http://pubs.opengroup.org/onlinepubs/009695399/functions/rename.html)³. Btrfs, Ext3 and Ext4
64 give atomic overwrite-by-rename guarantees, as well as atomic truncate guaran-
65 tees. The FAT filesystem guarantees neither.

66 SQLite

¹<http://www.sqlite.org/draft/wal.html>

²<http://www.sqlite.org/draft/lockingv3.html#rollback>

³<http://pubs.opengroup.org/onlinepubs/009695399/functions/rename.html>

67 For applications using SQLite for their storage, Collabora recommends using
68 either WAL or the rollback journal so that transactions are committed atomi-
69 cally. In addition, filesystem-specific tuning would be done by configuring the
70 SQLite system library for optimal performance.

71 WAL will be the best option in most cases, except when transactions will be
72 very big (involving more than 100 MB) and when writes are very seldom, then
73 the rollback journal would be preferred.

74 Collabora will run the [TCL test harness](#)⁴ for SQLite in LAVA, to detect any
75 issues in the specific configuration and software in the target platform. These in-
76 clude robustness tests that reproduce out-of-memory errors, input/output errors
77 and abnormal termination (crashes or power loss).

78 **Tracker**

79 Tracker stores data in SQLite files, so the robustness considerations that apply
80 to SQLite apply to Tracker as well. By default it uses WAL instead of the tradi-
81 tional rollback journal, which gives better performance for Tracker’s workload
82 with the same robustness guarantees.

83 **User settings**

84 For configuration settings in general, Collabora recommends using the [GSet-](#)
85 [tings](#)⁵ API from GLib with the [dconf](#)⁶ backend. When updating the database,
86 dconf will write the whole new contents to a new file, then atomically renaming
87 it on top of the old one.

88 For bigger pieces of data (individual settings whose data component exceeds 1
89 KB), Collabora recommends using plain files via a known-robust file-handling
90 library (such as [GKeyFile](#)⁷ from Glib, which is already a dependency) or SQLite.

91 **Media**

92 For media, the meta-data is stored in Tracker, with the actual data files in the
93 /home filesystem and in attached removable devices.

94 If the Tracker database that contains the meta-data has been corrupted, it
95 should be moved to the side (or deleted) and recreated again by indexing all
96 available media files. To minimize the chances of corruption, refer to [Tracker](#).

97 Software that reads the actual media files should assume media files may contain
98 invalid data and ignore them without further loss of functionality. Corrupted
99 media files should not be displayed in the UI.

⁴<http://www.sqlite.org/testing.html#tcl>

⁵<http://developer.gnome.org/gio/stable/GSettings.html>

⁶<https://live.gnome.org/dconf>

⁷<https://developer.gnome.org/glib/2.37/glib-Key-value-file-parser.html>

100 **Caches**

101 All software that uses a cache file should be ready to find that the cache is
102 unusable and cope with it without loss of functionality (temporary degradation
103 of performance is obviously expected in this case though the mechanism by
104 which the cache became corrupted will be treated by developers as a bug to
105 fix).

106 For example, if during start-up the Folks caches are found to be unreadable, lib-
107 folks would remove the corrupted cache files and recreate them, taking a longer
108 time to reply to queries. As the application using Folks would be executing
109 the queries asynchronously, the UI would keep being functional while the query
110 executes.

111 Examples of other components that use caches and that should cope with cache
112 corruption are the browser and the email client.

113 **Filesystems**

114 The reliability with which data is stored depends on both the storage medium
115 as well as the filesystem. In this section, we cover FAT32 and Btrfs. Ext4 is
116 mentioned, as it is a popular default filesystem on many Linux distributions –
117 however it doesn't suit the needs of the rollback system – either for system roll-
118 backs (See the System Update and Rollback Design) or for application rollbacks
119 (See the Applications Design).

120 The FAT32 filesystem is not robust under abnormal circumstances since it was
121 not made for devices which could be disconnected at any moment. In general, an
122 approach where writes to the device are tightly controlled and restricted to small
123 time-windows would help minimize the chances of corruption. See the *Media*
124 *and Indexing* design for a detailed explanation of the issues and suggestions.

125 The Ext4 filesystem is quite robust under power failure by default. It can be
126 made even more robust by [mounting it under data=journal](#)⁸ mode, but at a
127 large cost to performance.

128 Btrfs has been created on very robust principles, building upon the experience
129 of Ext4. Some brief technical details are provided at the end of this document
130 in [BTRFS overview](#).

131 **Filesystem options**

132 Filesystems usually have parameters that can be tuned to suit specific work-
133 loads. Some of them affect performance as well as robustness; either by trading
134 off between the two, or by taking advantage of specific hardware features avail-
135 able with the storage media.

⁸<http://kernel.org/doc/Documentation/filesystems/ext4.txt>

136 • **FAT32** is a simple filesystem that does not have many filesystem options
137 related to performance or robustness. Since we will not be creating any
138 FAT32 partitions ourselves, only mount-time options are interesting for
139 us. The recommended options are listed below:

140 – sync, flush
141 These filesystem options ensure that the kernel, as well as the filesys-
142 tem, flush data to the partition as soon as possible. This greatly
143 reduces the chances of data loss or filesystem corruption when USB
144 drives are yanked out by the user.

145 – ro (read only)
146 It is recommended that FAT32 partitions be mounted read-only to
147 avoid filesystem corruption, and other related problems as detailed in
148 the “*Media and Indexing Design*” in the section “*Indexing database*
149 *on removable device*”.

150 • **Btrfs** is relatively new, and so does not have many options relevant to
151 our needs of enhancing reliability on eMMC storage media. The available
152 options are listed below.

153 – *Mount-time options:*

154 * commit=number (default: 30)
155 Set the interval of periodic commit. This option is recent ([since](#)
156 [kernel 3.12](#))⁹.

157 * ssd
158 This option enables SSD-specific optimizations and disables some
159 optimisations specifically for rotating media. This option is en-
160 abled automatically on non-rotating storage.

161 * Recovery (default: off)
162 This option can be used to attempt recovery of a corrupted
163 filesystem (See [Repair and recovery](#)).

164 – *Filesystem creation options:*

165 * -s sector-size
166 This is the size of the filesystem blocks used for allocations. Ide-
167 ally, this should be the same size as the block size for the storage
168 medium.

169 * -M
170 This sets BTRFS to use “mixed block groups” - a mode that
171 stores data and metadata chunks together on disk for more effi-
172 cient space utilization for small filesystems – but incurs a perfor-
173 mance penalty on large ones. This option is not mature and will
174 be evaluated in the future.

⁹https://btrfs.wiki.kernel.org/index.php/Mount_options

175 The System Updates and Rollback Design describes the partition layout for
176 Apertis. Not all the partitions have the same requirements, so both the FAT32
177 and BTRFS filesystems are used. The partitions are configured as:

- 178 • **Factory Recovery** – This partition is never mounted read-write and must
179 be readable by the boot loader. Currently the boot loader for Apertis
180 – U-boot – does not support BTRFS. While patches exist to add that
181 functionality, they have not yet seen widespread testing. FAT32 will likely
182 be the filesystem chosen for the factory recovery image.
- 183 • **Minimal Boot partitions** – These partitions must also be readable by
184 the boot loader, and are currently FAT32. They are not normally mounted
185 at run-time, instead they are created, mounted, and populated by the
186 system update software once – and only ever accessed by the boot loader
187 afterwards. They will be mounted with the “sync” and “flush” flags.
- 188 • **System** - Since BTRFS provides an excellent snapshot mechanism to
189 assist system rollbacks (See [Cheap, fast, and atomic snapshots and roll-](#)
190 [back](#)), this partition will be populated with a BTRFS filesystem created
191 with the appropriate sector size for the storage device. It may be created
192 with mixed block groups to save storage space if that option does not lead
193 to instability. It will be mounted with the `ssd` option as well as `read-only`.
194 During a system update a single subvolume of the system subvolume will
195 be mounted read-write. The `repair mount` option will never be attempted
196 on the system partition, instead rollbacks or factory recovery will be used
197 to avoid potentially putting the system into an unknown state.
- 198 • **General Storage** – This partition shares similar requirements to the
199 system partition. It will be BTRFS, created with an appropriate sector
200 size and possibly mixed block groups. It will be mounted with the `ssd`
201 option. This is the only built-in non-volatile storage that will always be
202 mounted read-write. In the case of a damaged filesystem, repair may be
203 attempted on this partition.

204 Additionally, there are 2 partitions for raw status flag data that do not use
205 filesystems at all. See the System Updates and Rollback Design for more details.

206 Checksumming

207 Checksumming is used for detecting filesystem corruption due to any reason.
208 Different filesystems have different mechanisms for checksumming which give
209 us coverage for various different causes of filesystem corruption. Each mecha-
210 nism consumes I/O and CPU resources, and that must be weighed against the
211 advantages that it gives us.

212 It is important to note that checksumming does not protect us against corrup-
213 tion or help us in fixing the root cause of the corruption; it only allows us
214 to detect filesystem corruption when it happens. Hence, it is only useful as a

215 warning sign and recovering from data corruption is beyond the scope of this
216 feature.

- 217 • **FAT32** is a very old and simplistic filesystem, and it has no inbuilt facili-
218 ties for checksumming.
- 219 • **Btrfs** maintains a *checksum tree* for all the blocks that it allocates and
220 writes to. Hence, all file data and metadata is checksummed. This is the
221 default behaviour and the current checksum algorithm uses few resources.
222 This method of checksumming can detect all the ways in which corruption
223 can occur to data on the filesystem. See [Checksumming](#) for more detail.
- 224 • **Ext4** maintains checksums for journal data only, no checksumming of file
225 data takes place.

226 Alignment

227 The first piece of tuning that a filesystem on flash storage needs, is a proper
228 mapping of the filesystem blocks to the page size of the erase blocks on the flash.
229 This consists of two parts:

- 230 1. Ensuring that the filesystem and storage erase block sizes match using
231 filesystem creation options.
- 232 2. Aligning the block allocations in the filesystem with the storage blocks
233 by using the appropriate offsets while partitioning, or while creating the
234 filesystem.

235 If either of these is not satisfied, each filesystem block write will trigger two or
236 more flash block writes, and reduce the performance as well as reliability of the
237 MMC card.

238 The storage erase block size [can be read](#)¹⁰ from /proc/mtd or from U-Boot but
239 the flash storage can report something different than the real numbers. Some
240 sizes are available on the [Linaro wiki](#)¹¹. Linaro-image-tools is [now able](#)¹² to
241 generate images with a correct alignment.

242 Testing

243 Collabora will add tests to LAVA for testing how FAT32 and Btrfs behave on
244 the i.MX6 under stress, as well as for tuning the above mentioned parameters
245 for reliability and performance.

246 Root filesystem

247 The approach will be to mount as many parts of the root filesystem read-only
248 as possible such that the only writes to it would be during updates. This would

¹⁰http://processors.wiki.ti.com/index.php/Get_the_Flash_Erase_Block_Size

¹¹<https://wiki.linaro.org/WorkingGroups/KernelArchived/Projects/FlashCardSurvey>

¹²<https://bugs.launchpad.net/linaro-image-tools/bug/626907>

249 reduce the chances of catastrophic filesystem corruption in the event of power
250 failure and invalid system file modification by bugs in system or application
251 software. The only partition that is to be mounted writable is the user partition
252 that will be mounted in /home. All the other writable parts of the / filesystem
253 will be backed by tmpfs, located in RAM. We will avoid the lack of space
254 problem by only storing small files in tmpfs or files which don't take space (lock
255 files, socket files). Bigger files such as programs, libraries, configuration files will
256 remain on disk and available read-only.

257 See the *System Updates and Rollback* design for detailed information about the
258 robustness of the update process.

259 **Other filesystems**

260 The system should be able to function even if mounting one or more of the
261 non-essential file systems fails. Even if the system is able to keep running, it
262 would do so with reduced functionality, so some recovery action would need
263 to be taken in order to regain the lost functionality. The system should try
264 to recover automatically as far as possible. In the case of unrecoverable system
265 failure, the user can be instructed at system boot to request technical assistance
266 at a service shop.

267 **Main storage**

268 In case of power loss, the flash media can become corrupted due to how writes
269 are performed. Apertis will be notified via a GPIO signal 100 milliseconds before
270 power is completely lost, in order to give the flash controller time to commit to
271 non-volatile media what is in its cache.

272 Given the short time available and the general slowness of flash devices when
273 writing, we recommend that the signal is handled in the kernel, because
274 userspace will not have enough time to react (depending on the load and the
275 scheduler, it could take from 10 ms to 100 ms for the signal to start being
276 processed by a userspace process). A device driver should be written that,
277 when the GPIO signal is received:

- 278 1. stops flushing dirty pages to the drive,
- 279 2. tells the flash controller to flush its caches to permanent storage, and
- 280 3. starts the shutdown sequence.

281 The device driver will start handling the signal 10-100 μ s after the GPIO is
282 activated. In spite of this, if the device has big caches and is slow to write,
283 corruption of arbitrary data blocks can still happen.

284 In general, drive health data should be monitored so that the user can be notified
285 about disk failures which require a garage visit for hardware replacement.

286 As no more dirty pages will be flushed to the storage device when the GPIO
287 signal is received, the data in the page cache will be lost. To reduce the amount
288 of data that could be lost, eMMC reliable writes can be used, and the page cache
289 configuration can be tuned. But it has to be noted that use of reliable writes
290 and reducing the amount of in-flight data is a trade-off against performance
291 that can be quantified only on the final hardware configuration through direct
292 experimentation.

293 **Removable devices**

294 External devices that can be removed at any moment are not reliable for writ-
295 ing of critical data. In addition to the problem of corruption of files being
296 written, wear leveling by the controller might corrupt unrelated blocks which
297 might even contain the directory table or the file allocation table, rendering the
298 whole partition unusable.

299 The quality of external storage devices such as flash drives varies greatly, in some
300 cases the device will unexpectedly stop responding to commands, or data will
301 be lost. Applications that write to removable drives must be robust enough to
302 be able to continue in the face of such errors with minimal loss of functionality.

303 As mentioned in [Filesystems](#), the safest way to use removable drives is by re-
304 stricting the processes that can write to the drive, and minimizing the time-
305 window for the writes. For that to be practical, there should be a system ser-
306 vice that is the only one allowed to write to removable devices and that would
307 accept requests from applications, remount the device read-write, write the new
308 contents, then remount read-only again.

309 Since, for interoperability reasons, the filesystem used in removable devices is
310 FAT32, in addition to the issues mentioned in this section, the robustness con-
311 siderations that were explained earlier in [Filesystems](#) also apply.

312 **Mitigating the effects of lack of disk space**

313 In order to reduce the chances that the system will find itself in a situation
314 where lack of disk space is problematic, it is recommended that available disk
315 space is monitored and applications notified so they can react and modify their
316 behavior accordingly. Applications may chose to delete unused files, delete or
317 reduce cache files or purge old data from their databases.

318 The recommended mechanism for monitoring available disk space is for a dae-
319 mon running in the user session to call *statvfs* (2) periodically on each mount
320 point and notify applications with a D-Bus signal. [Example code](#)¹³ can be
321 found in the GNOME project, which uses a similar approach (polling every 60
322 seconds).

¹³<http://git.gnome.org/browse/gnome-settings-daemon/tree/plugins/housekeeping/gsd-disk-space.c#n693>

323 Additionally, so error messages can be stored also in low-space conditions, it
324 is recommended that *journal* is configured to leave an amount of free space
325 smaller than the reserved blocks of the filesystem that backs the log files. This
326 way, applications will still be able to log messages after applications have con-
327 sumed all the space available to them.

328 In case applications cannot be trusted to properly delete non-essential files, a
329 possibility would be for them to state in their manifest where such files will be
330 stored, so the system can delete them when needed.

331 In order to make sure that malfunctioning applications cannot cause disruption
332 by filling filesystems, it would be required that each application writes to a
333 separate filesystem.

334 It may be worth noting that temporary directories should be emptied on reboot.

335 **Resource management**

336 The robustness goal of resource management is to prevent one or more applica-
337 tions from disrupting basic functionality due to excessive resource consumption.
338 The basic mechanism for this is to allocate resources in such a way that applica-
339 tions cannot starve services in the base system. This is to be achieved firstly by
340 changing the resource allocation policy to give higher priority to services, and
341 secondly by limiting the maximum amount of resources that an application can
342 consume at a time.

343 Resource limits are capable of helping ensure a process does not render the
344 whole system unresponsive. However, some design decisions also play an im-
345 portant role here. If the user has no way to kill the process that became too
346 slow or unresponsive, the user experience will suffer. The same goes for the
347 case in which an application gets stuck into a failing scenario, such as a web
348 browser automatically loading pages that were open when the browser closed
349 unexpectedly. For these reasons care must be exercised while designing the user
350 interactions for both the system chrome and applications to be sure such cases
351 are addressed.

352 If, despite throttling, some processes still impact the overall user experience
353 negatively because of excessive resource usage, there is the option of identifying
354 those processes and terminating them. Collabora recommends against this be-
355 cause it is very difficult to automatically distinguish between processes that use
356 large amounts of resources due to malfunction or maliciousness and processes
357 that use excessive resources for legitimate purposes. Killing the wrong process
358 may free up resources but is likely to be perceived by the user as a severe defect
359 in the overall user experience.

360 As a general recommendation, for optimal responsiveness, applications should
361 not block the UI thread when calling anything that is not assured to return
362 almost immediately, which includes all local or remote I/O operations. When

363 the potential duration of an operation is a considerable portion of the commonly-
364 considered maximum acceptable response time (100 ms), it should be done
365 asynchronously. GLib contains [asynchronous APIs](#)¹⁴ for I/O in its [file](#)¹⁵ and
366 [streaming](#)¹⁶ classes.

367 CPU

368 To make sure that important processes have available CPU cycles even when mal-
369 functioning or malicious applications monopolise the CPU, it is recommended to
370 set task scheduler priorities according to the importance of processes. Systemd
371 can do this for services by setting the [CPUSchedulingPriority](#)¹⁷ property in the
372 service unit file of the process. When the process described by the service unit
373 file starts new processes, they stay in the same cgroups and they keep the same
374 CPUSchedulingPriority.

375 At present (Q1 2014), systemd manages the user session on target images but
376 not on the SDK. With the user session managed by systemd, the priorities of ap-
377 plications are no longer set by the application launcher using [sched_setscheduler](#)
378 [\(2\)](#)¹⁸.

379 If there are processes that need real-time capabilities, or that should have very
380 low CPU access, the CPUSchedulingPolicy property can be used to change to
381 the rr (real-time) or idle scheduling policies. Real-time access for a process
382 should be carefully considered and tested because it can have a negative impact
383 on the process and even the entire system.

384 For identifying processes that use an excessive amount of CPU, the [cpuacct](#)¹⁹
385 cgroups controller can be used.

386 Though it is not recommended to automatically terminate local applications
387 with excessive CPU usage, it makes sense for web pages. Web pages are not
388 screened before they execute on the system, hence it is important to ensure that
389 their ability to disrupt system functionality is minimised. For this, WebKit can
390 detect when a block of JavaScript code has been executing for too long, pause
391 it, and give the embedding application the possibility of canceling the execution
392 of this block of code. Collabora has added API to WebKit-Clutter for this.

393 I/O

394 Similar to CPU usage, Collabora recommends giving priority to important pro-
395 cesses when there is contention for I/O bandwidth. Collabora recommends that
396 important services have a value for the property IOSchedulingPriority lower
397 than 4 (the default). If, for any reason, some applications need priorities other

¹⁴<http://developer.gnome.org/gio/stable/async.html>

¹⁵http://developer.gnome.org/gio/stable/file_ops.html

¹⁶<http://developer.gnome.org/gio/stable/streaming.html>

¹⁷<http://0pointer.de/public/systemd-man/systemd.exec.html>

¹⁸http://www.kernel.org/doc/man-pages/online/pages/man2/sched_setscheduler.2.html

¹⁹<http://www.kernel.org/doc/Documentation/cgroups/cpuacct.txt>

398 than the default, the application launcher can use the `ioprio_set`²⁰ (2) syscall
399 to change their priority. When the process described by the service unit file
400 starts new processes, they stay in the same cgroups and they keep the same
401 `IOSchedulingPriority`.

402 **Memory**

403 Collabora recommends putting a single limit on the amount of memory that the
404 whole application set can allocate so a fair reserve is left for the base software.
405 This limit should be just big enough so that the Apertis instance never reaches
406 the “out of memory” (OOM²¹) condition at the system level. For example, if
407 the total of memory available for processes is 1GB, there is no swap, and we
408 know that the services in the base system should need a maximum of 300MB,
409 then all applications should belong to a cgroup that is limited to 700MB of
410 memory.

411 In specific cases, it may make sense to put a different limit on a specific appli-
412 cation, but it can easily be counterproductive and cause a waste of memory.

413 Something else worth doing is to make sure that the OOM killer²² selects ap-
414 plications for killing instead of system services. For this, the systemd prop-
415 erty `OOMScoreAdjust` can be used to reduce the chances that a service will be
416 killed. For applications, it is recommended that the application launcher sets
417 its `/proc/<pid>/oom_score_adj` (see [here](#)²³) to be higher than 0. The ideal
418 value may vary depending upon the importance of each application.

419 With the example setup mentioned before, the OOM killer will terminate the
420 bulkiest application when one of these conditions are met:

- 421 • The total memory taken by applications all together is going to increase
422 over 700MB.
- 423 • The total memory taken by all processes (services plus applications) is
424 going to increase over 1GB.

425 To make better use of the available memory, it’s recommended that applications
426 listen to the cgroup notification `memory.usage_in_bytes`²⁴ and when it gets
427 close to the limit for applications, start reducing the size of any caches they
428 hold in main memory. It may be good to do this inside the SDK and provide
429 applications with a glib signal that they can listen for.

430 **Network queue**

431 Processes would be classified into cgroup classes such as:

²⁰http://www.kernel.org/doc/man-pages/online/pages/man2/ioprio_set.2.html

²¹http://en.wikipedia.org/wiki/Out_of_memory

²²<http://lwn.net/Articles/317814/>

²³<http://www.kernel.org/doc/Documentation/filesystems/proc.txt>

²⁴<http://www.kernel.org/doc/Documentation/cgroups/memory.txt>

- 432 • Interactive (VoIP, internet radio)
- 433 • Semi-interactive (web pages, maps)
- 434 • Asynchronous (mail, app notifications, etc)
- 435 • Bulk (downloads, system updates)

436 Cgroup controllers are only used for classification of outgoing packets. [NET-](#)
437 [PRIO_CGROUP](#)²⁵ and [NET_CLS_CGROUP](#)²⁶ would be used for setting the
438 priority, and for classifying processes into cgroups. By thus tagging packets
439 with the cgroup of applications and services, *tc*²⁷ can be used to set limits to
440 the rate at which processes send packets (<http://lartc.org/howto/>).

441 Bandwidth rate-limiting would be required to ensure interactive streams do not
442 get starved by lower priority streams.

443 There is little we can do about latency for applications like VoIP, since even when
444 the bandwidth is sufficient, the bottlenecks are the hardware buffers, queues,
445 and scheduling on various devices outside the control of our system. This is an
446 open problem in networking, and a large part of it is related to [Bufferbloat](#)²⁸.

447 Note that there's no robustness issue that can be prevented by limiting the rate
448 at which processes receive incoming packets.

449 GPU

450 As explained in the WebGL design, the [GL_EXT_robustness](#)²⁹ extension pro-
451 vides a mechanism by which the watchdog in the GL implementation can reset
452 the GPU, invalidating all GL contexts and thus stopping all GPU activity.

453 Unfortunately, this only prevents denial of service (DoS) conditions caused by
454 WebGL, because processes must opt-in to use this extension. Thus, applications
455 may intentionally or unintentionally ignore the extension and continue monopo-
456 lising the GPU. Within the web browser, scripts that use WebGL and take over
457 the GPU will be interrupted and terminated by the browser.

458 If it runs its own GL implementation, then it could monitor GPU resource
459 usage and reset those contexts that seem to be disrupting the rest of the sys-
460 tem. It could notify processes via the [GL_EXT_robustness](#) extension and even
461 terminate them if they ignore the context reset notifications.

462 Accounting

²⁵<http://lwn.net/Articles/474695/>

²⁶[http://docs.fedoraproject.org/en-US/Fedora/16/html/Resource_Management_Guide/
sec-net_cls.html](http://docs.fedoraproject.org/en-US/Fedora/16/html/Resource_Management_Guide/sec-net_cls.html)

²⁷<http://lartc.org/manpages/tc.txt>

²⁸<http://en.wikipedia.org/wiki/Bufferbloat>

²⁹http://www.khronos.org/registry/gles/extensions/EXT/EXT_robustness.txt

463 Besides setting limits on resources, cgroups also allows to retrieve resource us-
464 age metrics. As examples, for CPU usage the *cpucct* cgroup controller con-
465 tains the *usage*, *stat* and *usage_percpu* reports; the *memory* controller provides
466 usage data in its *stat* report; the *blkio* controller has *throttle.io_serviced* and
467 *throttle.io_service_bytes*.

468 **USB undervoltage**

469 In the case that the system momentarily isn't able to power connected USB de-
470 vices such as MP3 players or smartphones due to voltage drops, the system will
471 power off and on again these devices, so that the connection gets reestablished
472 and the user experience gets affected as little as possible.

473 **Risks**

- 474 • FAT32 is fundamentally unreliable, specially on removable devices.
- 475 • Robustness of flash media varies greatly and the user may not be able
476 to distinguish failures caused by the hardware from failures due to the
477 software.
- 478 • Excessively-low resource limits for applications can lead to resource waste;
479 excessively-high may be less effective in avoiding DoS. There may not exist
480 a good middle point.
- 481 • Heuristics used to determine when to kill a process with excessive resource
482 usage are not perfect and can cause major failure from the user point of
483 view.
- 484 • If Vivante does not implement `GL_EXT_robustness` properly, web pages
485 could DoS the whole system.
- 486 • Bugs in the OpenGL implementation can lead to instability, data loss and
487 privacy breaches that can be triggered from web pages.
- 488 • If the flash media loses power while a block is open for writing, it is possible
489 that several random blocks elsewhere in the same drive will be corrupted.
490 This can affect other filesystems, even if they are mounted read-only.

491 **Design notes**

492 The following items have been identified for future investigation and design work
493 later in the project and are thus not addressed in this design:

- 494 • Vulnerability to DoS attacks in D-Bus and proposed solutions.
- 495 • Optimization of the SQLite configuration parameters for the specific
496 filesystems in use in Apertis.

497 No updates as of March 2014.

498 **BTRFS Overview**

499 The most powerful feature of Btrfs³⁰ is the fact that all information (data +
500 metadata) is stored in the same basic data structures, and all modification of
501 these data structures is performed in a copy-on-write (CoW) fashion.

502 Since all information on disk is stored using the same type of data structure, this
503 allows metadata and data to share features such as checksumming and striping.

504 This combined with the fact that Btrfs uses CoW while modifying all informa-
505 tion, means that in theory, the filesystem is always consistent if the storage
506 device supports “Force Unit Access³¹” correctly. However, in practice, filesys-
507 tem bugs, a lack of maturity in the code, and other (unforeseen) problems may
508 prevent this.

509 **BTRFS robustness supporting features**

510 **Cheap, fast, and atomic snapshots and rollback**

511 All snapshots in Btrfs are CoW copies of the subvolume being snapshotted with
512 an incremented reference count for the blocks. As a result, creating snapshots
513 is very fast, and they take up a negligible amount of space. Just like every
514 other operation, the snapshot is created atomically by the use of transactions
515 and sequenced flushes. Further, all snapshots are actually just subvolumes, and
516 hence can be mounted on their own.

517 Unlike LVM2, which creates snapshots in the form of block devices that can be
518 mounted, Btrfs creates snapshots in the form of subvolumes, which are repre-
519 sented as subdirectories.

520 Even though snapshots are displayed in a subdirectory they are not “owned” by
521 that subvolume. Snapshots and subvolumes are identical in Btrfs, and are first-
522 class citizens with respect to other subvolumes. This means that the default
523 subvolume can be set at any time. The change will be made the next time the
524 filesystem is mounted. All the subvolumes, except the top-level subvolume, can
525 also be deleted; irrespective of their relationships with each other.

526 **Repair and recovery**

527 If, for any reason, the root node or the superblock gets corrupted and the filesys-
528 tem cannot be mounted, mounting in recovery mode will make btrfs check the
529 superblock (or alternate superblocks if the superblock is also corrupted) for
530 alternate roots from previous transactions. This is possible because all modifi-
531 cations to the Btrfs trees are done in a CoW manner and existing roots are not
532 deleted. The filesystem stores the last four roots as a backup for the recovery
533 option.

³⁰<https://btrfs.wiki.kernel.org/>

³¹http://en.wikipedia.org/wiki/SCSI_Write_Commands#Write_.2810.29

534 **Checksumming**

535 The header of every chunk of space in Btrfs has space for 32-bytes of check-
536 sums of the chunk itself. In addition, there is a checksum tree which maintains
537 checksums for each block of data. Since the data as well as the metadata blocks
538 are referenced in the checksum tree, all information in the filesystem is check-
539 summed.

540 Currently, Btrfs uses the CRC-32 checksum algorithm, but there are plans to
541 upgrade that, and add the option to set the checksum algorithm when the
542 filesystem is created.