Google Git

kernel / pub / scm / fs / xfs / xfs-documentation / master / . / design /
XFS_Filesystem_Structure / **allocation_groups.asciidoc**

```
blob: c0ba16a82d489fedc10bfafd90c26409938b7730 [file] [log] [blame]
```

```
  1   [[Allocation_Groups]]
  2   = Allocation Groups
  3
  4   As mentioned earlier, XFS filesystems are divided into a number of equally
  5   sized chunks called Allocation Groups. Each AG can almost be thought of as an
  6   individual filesystem that maintains its own space usage. Each AG can be up to
  7   one terabyte in size (512 bytes × 2^31^), regardless of the underlying device's
  8   sector size.
  9
 10   Each AG has the following characteristics:
 11
 12           * A super block describing overall filesystem info
 13           * Free space management
 14           * Inode allocation and tracking
 15           * Reverse block-mapping index (optional)
 16           * Data block reference count index (optional)
 17
 18   Having multiple AGs allows XFS to handle most operations in parallel without
 19   degrading performance as the number of concurrent accesses increases.
 20
 21   The only global information maintained by the first AG (primary) is free space
 22   across the filesystem and total inode counts. If the
 23   +XFS_SB_VERSION2_LAZYSBCOUNTBIT+ flag is set in the superblock, these are only
 24   updated on-disk when the filesystem is cleanly unmounted (umount or shutdown).
 25
 26   Immediately after a +mkfs.xfs+, the primary AG has the following disk layout;
 27   the subsequent AGs do not have any inodes allocated:
 28
 29   .Allocation group layout
 30   image::images/6.png[]
 31
 32   Each of these structures are expanded upon in the following sections.
 33
 34   [[Superblocks]]
 35   == Superblocks
 36
 37   Each AG starts with a superblock. The first one, in AG 0, is the primary
 38   superblock which stores aggregate AG information. Secondary superblocks are
 39   only used by xfs_repair when the primary superblock has been corrupted.  A
 40   superblock is one sector in length.
```

```
41
42      The superblock is defined by the following structure. The description of each
43      field follows.
44
45      [source, c]
46      ----
47      struct xfs_sb
48      {
49              __uint32_t              sb_magicnum;
50              __uint32_t              sb_blocksize;
51              xfs_rfsblock_t          sb_dblocks;
52              xfs_rfsblock_t          sb_rblocks;
53              xfs_rtblock_t           sb_rextents;
54              uuid_t                  sb_uuid;
55              xfs_fsblock_t           sb_logstart;
56              xfs_ino_t               sb_rootino;
57              xfs_ino_t               sb_rbmino;
58              xfs_ino_t               sb_rsumino;
59              xfs_agblock_t           sb_rextsize;
60              xfs_agblock_t           sb_agblocks;
61              xfs_agnumber_t          sb_agcount;
62              xfs_extlen_t            sb_rbmblocks;
63              xfs_extlen_t            sb_logblocks;
64              __uint16_t              sb_versionnum;
65              __uint16_t              sb_sectsize;
66              __uint16_t              sb_inodesize;
67              __uint16_t              sb_inopblock;
68              char                    sb_fname[12];
69              __uint8_t               sb_blocklog;
70              __uint8_t               sb_sectlog;
71              __uint8_t               sb_inodelog;
72              __uint8_t               sb_inopblog;
73              __uint8_t               sb_agblklog;
74              __uint8_t               sb_rextslog;
75              __uint8_t               sb_inprogress;
76              __uint8_t               sb_imax_pct;
77              __uint64_t              sb_icount;
78              __uint64_t              sb_ifree;
79              __uint64_t              sb_fdblocks;
80              __uint64_t              sb_frextents;
81              xfs_ino_t               sb_uquotino;
82              xfs_ino_t               sb_gquotino;
83              __uint16_t              sb_qflags;
84              __uint8_t               sb_flags;
85              __uint8_t               sb_shared_vn;
86              xfs_extlen_t            sb_inoalignmt;
87              __uint32_t              sb_unit;
88              __uint32_t              sb_width;
89              __uint8_t               sb_dirblklog;
```

```
 90            __uint8_t                sb_logsectlog;
 91            __uint16_t               sb_logsectsize;
 92            __uint32_t               sb_logsunit;
 93            __uint32_t               sb_features2;
 94            __uint32_t               sb_bad_features2;
 95
 96            /* version 5 superblock fields start here */
 97            __uint32_t               sb_features_compat;
 98            __uint32_t               sb_features_ro_compat;
 99            __uint32_t               sb_features_incompat;
100            __uint32_t               sb_features_log_incompat;
101
102            __uint32_t               sb_crc;
103        xfs_extlen_t             sb_spino_align;
104
105        xfs_ino_t                sb_pquotino;
106        xfs_lsn_t                sb_lsn;
107        uuid_t                   sb_meta_uuid;
108        xfs_ino_t                sb_rrmapino;
109    };
110    ----
```

111    *sb_magicnum*::
112    Identifies the filesystem. Its value is +XFS_SB_MAGIC+ ``XFSB'' (0x58465342).
113
114    *sb_blocksize*::
115    The size of a basic unit of space allocation in bytes. Typically, this is 4096
116    (4KB) but can range from 512 to 65536 bytes.
117
118    *sb_dblocks*::
119    Total number of blocks available for data and metadata on the filesystem.
120
121    *sb_rblocks*::
122    Number blocks in the real-time disk device. Refer to
123    xref:Real-time_Devices[real-time sub-volumes] for more information.
124
125    *sb_rextents*::
126    Number of extents on the real-time device.
127
128    *sb_uuid*::
129    UUID (Universally Unique ID) for the filesystem. Filesystems can be mounted by
130    the UUID instead of device name.
131
132    *sb_logstart*::
133    First block number for the journaling log if the log is internal (ie. not on a
134    separate disk device). For an external log device, this will be zero (the log
135    will also start on the first block on the log device).  The identity of the log
136    devices is not recorded in the filesystem, but the UUIDs of the filesystem and
137    the log device are compared to prevent corruption.
138

```
139    *sb_rootino*::
140    Root inode number for the filesystem.  Normally, the root inode is at the
141    start of the first possible inode chunk in AG 0.  This is 128 when using a 4KB
142    block size.
143
144    *sb_rbmino*::
145    Bitmap inode for real-time extents.
146
147    *sb_rsumino*::
148    Summary inode for real-time bitmap.
149
150    *sb_rextsize*::
151    Realtime extent size in blocks.
152
153    *sb_agblocks*::
154    Size of each AG in blocks. For the actual size of the last AG, refer to the
155    xref:AG_Free_Space_Management[free space] +agf_length+ value.
156
157    *sb_agcount*::
158    Number of AGs in the filesystem.
159
160    *sb_rbmblocks*::
161    Number of real-time bitmap blocks.
162
163    *sb_logblocks*::
164    Number of blocks for the journaling log.
165
166    *sb_versionnum*::
167    Filesystem version number. This is a bitmask specifying the features enabled
168    when creating the filesystem. Any disk checking tools or drivers that do not
169    recognize any set bits must not operate upon the filesystem. Most of the flags
170    indicate features introduced over time. If the value of the lower nibble is >=
171    4, the higher bits indicate feature flags as follows:
172
173    .Version 4 Superblock version flags
174    [options="header"]
175    |=====
176    | Flag                          | Description
177    | +XFS_SB_VERSION_ATTRBIT+      |
178    Set if any inode have extended attributes.  If this bit is set; the
179    +XFS_SB_VERSION2_ATTR2BIT+ is not set; and the +attr2+ mount flag is not
180    specified, the +di_forkoff+ inode field will not be dynamically adjusted.
181    See the section about xref:Extended_Attribute_Versions[extended attribute
182    versions] for more information.
183
184    | +XFS_SB_VERSION_NLINKBIT+     | Set if any inodes use 32-bit di_nlink values.
185    | +XFS_SB_VERSION_QUOTABIT+     |
186    Quotas are enabled on the filesystem. This
187    also brings in the various quota fields in the superblock.
```

```
188
189  | +XFS_SB_VERSION_ALIGNBIT+      | Set if sb_inoalignmt is used.
190  | +XFS_SB_VERSION_DALIGNBIT+     | Set if sb_unit and sb_width are used.
191  | +XFS_SB_VERSION_SHAREDBIT+     | Set if sb_shared_vn is used.
192  | +XFS_SB_VERSION_LOGV2BIT+      | Version 2 journaling logs are used.
193  | +XFS_SB_VERSION_SECTORBIT+     | Set if sb_sectsize is not 512.
194  | +XFS_SB_VERSION_EXTFLGBIT+     | Unwritten extents are used. This is always set.
195  | +XFS_SB_VERSION_DIRV2BIT+      |
196  Version 2 directories are used. This is always set.
197
198  | +XFS_SB_VERSION_MOREBITSBIT+   |
199  Set if the sb_features2 field in the superblock contains more flags.
200  |=====
201
202  If the lower nibble of this value is 5, then this is a v5 filesystem; the
203  +XFS_SB_VERSION2_CRCBIT+ feature must be set in +sb_features2+.
204
205  *sb_sectsize*::
206  Specifies the underlying disk sector size in bytes.  Typically this is 512 or
207  4096 bytes. This determines the minimum I/O alignment, especially for direct I/O.
208
209  *sb_inodesize*::
210  Size of the inode in bytes. The default is 256 (2 inodes per standard sector)
211  but can be made as large as 2048 bytes when creating the filesystem.  On a v5
212  filesystem, the default and minimum inode size are both 512 bytes.
213
214  *sb_inopblock*::
215  Number of inodes per block. This is equivalent to +sb_blocksize / sb_inodesize+.
216
217  *sb_fname[12]*::
218  Name for the filesystem. This value can be used in the mount command.
219
220  *sb_blocklog*::
221  log~2~ value of +sb_blocksize+. In other terms, +sb_blocksize = 2^sb_blocklog^+.
222
223  *sb_sectlog*::
224  log~2~ value of +sb_sectsize+.
225
226  *sb_inodelog*::
227  log~2~ value of +sb_inodesize+.
228
229  *sb_inopblog*::
230  log~2~ value of +sb_inopblock+.
231
232  *sb_agblklog*::
233  log~2~ value of +sb_agblocks+ (rounded up). This value is used to generate inode
234  numbers and absolute block numbers defined in extent maps.
235
236  *sb_rextslog*::
```

```
237    log~2~ value of +sb_rextents+.

238

239    *sb_inprogress*::
240    Flag specifying that the filesystem is being created.

241

242    *sb_imax_pct*::
243    Maximum percentage of filesystem space that can be used for inodes. The default
244    value is 5%.

245

246    *sb_icount*::
247    Global count for number inodes allocated on the filesystem. This is only
248    maintained in the first superblock.

249

250    *sb_ifree*::
251    Global count of free inodes on the filesystem. This is only maintained in the
252    first superblock.

253

254    *sb_fdblocks*::
255    Global count of free data blocks on the filesystem. This is only maintained in
256    the first superblock.

257

258    *sb_frextents*::
259    Global count of free real-time extents on the filesystem. This is only
260    maintained in the first superblock.

261

262    *sb_uquotino*::
263    Inode for user quotas. This and the following two quota fields only apply if
264    +XFS_SB_VERSION_QUOTABIT+ flag is set in +sb_versionnum+. Refer to
265    xref:Quota_Inodes[quota inodes] for more information.

266

267    *sb_gquotino*::
268    Inode for group or project quotas. Group and project quotas cannot be used at
269    the same time on v4 filesystems.  On a v5 filesystem, this inode always stores
270    group quota information.

271

272    *sb_qflags*::
273    Quota flags. It can be a combination of the following flags:

274

275    .Superblock quota flags
276    [options="header"]
277    |=====
278    | Flag                        | Description
279    | +XFS_UQUOTA_ACCT+           | User quota accounting is enabled.
280    | +XFS_UQUOTA_ENFD+           | User quotas are enforced.
281    | +XFS_UQUOTA_CHKD+           | User quotas have been checked.
282    | +XFS_PQUOTA_ACCT+           | Project quota accounting is enabled.
283    | +XFS_OQUOTA_ENFD+           | Other (group/project) quotas are enforced.
284    | +XFS_OQUOTA_CHKD+           | Other (group/project) quotas have been checked.
285    | +XFS_GQUOTA_ACCT+           | Group quota accounting is enabled.
```

```
286    | +XFS_GQUOTA_ENFD+                | Group quotas are enforced.
287    | +XFS_GQUOTA_CHKD+                | Group quotas have been checked.
288    | +XFS_PQUOTA_ENFD+                | Project quotas are enforced.
289    | +XFS_PQUOTA_CHKD+                | Project quotas have been checked.
290    |=====
291
292    *sb_flags*::
293    Miscellaneous flags.
294
295    .Superblock flags
296    [options="header"]
297    |=====
298    | Flag                         | Description
299    | +XFS_SBF_READONLY+           | Only read-only mounts allowed.
300    |=====
301
302    *sb_shared_vn*::
303    Reserved and must be zero (``vn'' stands for version number).
304
305    *sb_inoalignmt*::
306    Inode chunk alignment in fsblocks.  Prior to v5, the default value provided for
307    inode chunks to have an 8KiB alignment.  Starting with v5, the default value
308    scales with the multiple of the inode size over 256 bytes.  Concretely, this
309    means an alignment of 16KiB for 512-byte inodes, 32KiB for 1024-byte inodes,
310    etc.  If sparse inodes are enabled, the +ir_startino+ field of each inode
311    B+tree record must be aligned to this block granularity, even if the inode
312    given by +ir_startino+ itself is sparse.
313
314    *sb_unit*::
315    Underlying stripe or raid unit in blocks.
316
317    *sb_width*::
318    Underlying stripe or raid width in blocks.
319
320    *sb_dirblklog*::
321    log~2~ multiplier that determines the granularity of directory block allocations
322    in fsblocks.
323
324    *sb_logsectlog*::
325    log~2~ value of the log subvolume's sector size. This is only used if the
326    journaling log is on a separate disk device (i.e. not internal).
327
328    *sb_logsectsize*::
329    The log's sector size in bytes if the filesystem uses an external log device.
330
331    *sb_logsunit*::
332    The log device's stripe or raid unit size. This only applies to version 2 logs
333    +XFS_SB_VERSION_LOGV2BIT+ is set in +sb_versionnum+.
334
```

```
335   *sb_features2*::
336   Additional version flags if +XFS_SB_VERSION_MOREBITSBIT+ is set in
337   +sb_versionnum+. The currently defined additional features include:
338
339   .Extended Version 4 Superblock flags
340   [options="header"]
341   |=====
342   | Flag                              | Description
343   | +XFS_SB_VERSION2_LAZYSBCOUNTBIT+ |
344   Lazy global counters. Making a filesystem with this bit set can improve
345   performance. The global free space and inode counts are only updated in the
346   primary superblock when the filesystem is cleanly unmounted.
347
348   | +XFS_SB_VERSION2_ATTR2BIT+      |
349   Extended attributes version 2. Making a filesystem with this optimises the
350   inode layout of extended attributes.  If this bit is set and the +noattr2+
351   mount flag is not specified, the +di_forkoff+ inode field will be dynamically
352   adjusted.  See the section about xref:Extended_Attribute_Versions[extended
353   attribute versions] for more information.
354
355   | +XFS_SB_VERSION2_PARENTBIT+    |
356   Parent pointers. All inodes must have an extended attribute that points back to
357   its parent inode. The primary purpose for this information is in backup systems.
358
359   | +XFS_SB_VERSION2_PROJID32BIT+ |
360   32-bit Project ID.  Inodes can be associated with a project ID number, which
361   can be used to enforce disk space usage quotas for a particular group of
362   directories.  This flag indicates that project IDs can be 32 bits in size.
363
364   | +XFS_SB_VERSION2_CRCBIT+       |
365   Metadata checksumming.  All metadata blocks have an extended header containing
366   the block checksum, a copy of the metadata UUID, the log sequence number of the
367   last update to prevent stale replays, and a back pointer to the owner of the
368   block.  This feature must be and can only be set if the lowest nibble of
369   +sb_versionnum+ is set to 5.
370
371   | +XFS_SB_VERSION2_FTYPE+        |
372   Directory file type.  Each directory entry records the type of the inode to
373   which the entry points.  This speeds up directory iteration by removing the
374   need to load every inode into memory.
375   |=====
376
377   *sb_bad_features2*::
378   This field mirrors +sb_features2+, due to past 64-bit alignment errors.
379
380   *sb_features_compat*::
381   Read-write compatible feature flags.  The kernel can still read and write this
382   FS even if it doesn't understand the flag.  Currently, there are no valid
383   flags.
```

```
384
385    *sb_features_ro_compat*::
386    Read-only compatible feature flags.  The kernel can still read this FS even if
387    it doesn't understand the flag.
388
389    .Extended Version 5 Superblock Read-Only compatibility flags
390    [options="header"]
391    |=====
392    | Flag                            | Description
393    | +XFS_SB_FEAT_RO_COMPAT_FINOBT+ |
394    Free inode B+tree.  Each allocation group contains a B+tree to track inode chunks
395    containing free inodes.  This is a performance optimization to reduce the time
396    required to allocate inodes.
397
398    | +XFS_SB_FEAT_RO_COMPAT_RMAPBT+ |
399    Reverse mapping B+tree.  Each allocation group contains a B+tree containing
400    records mapping AG blocks to their owners.  See the section about
401    xref:Reconstruction[reconstruction] for more details.
402
403    | +XFS_SB_FEAT_RO_COMPAT_REFLINK+ |
404    Reference count B+tree.  Each allocation group contains a B+tree to track the
405    reference counts of AG blocks.  This enables files to share data blocks safely.
406    See the section about xref:Reflink_Deduplication[reflink and deduplication] for
407    more details.
408
409    | +XFS_SB_FEAT_RO_COMPAT_INOBTCNT+ |
410    Inode B+tree block counters.  Each allocation group's inode (AGI) header
411    tracks the number of blocks in each of the inode B+trees.  This allows us
412    to have a slightly higher level of redundancy over the shape of the inode
413    btrees, and decreases the amount of time to compute the metadata B+tree
414    preallocations at mount time.
415
416    |=====
417
418    *sb_features_incompat*::
419    Read-write incompatible feature flags.  The kernel cannot read or write this
420    FS if it doesn't understand the flag.
421
422    .Extended Version 5 Superblock Read-Write incompatibility flags
423    [options="header"]
424    |=====
425    | Flag                            | Description
426    | +XFS_SB_FEAT_INCOMPAT_FTYPE+ |
427    Directory file type.  Each directory entry tracks the type of the inode to
428    which the entry points.  This is a performance optimization to remove the need
429    to load every inode into memory to iterate a directory.
430
431    | +XFS_SB_FEAT_INCOMPAT_SPINODES+ |
432    Sparse inodes.  This feature relaxes the requirement to allocate inodes in
```

```
433   chunks of 64.  When the free space is heavily fragmented, there might exist
434   plenty of free space but not enough contiguous free space to allocate a new
435   inode chunk.  With this feature, the user can continue to create files until
436   all free space is exhausted.
437
438   Unused space in the inode B+tree records are used to track which parts of the
439   inode chunk are not inodes.
440
441   See the chapter on xref:Sparse_Inodes[Sparse Inodes] for more information.
442
443   | +XFS_SB_FEAT_INCOMPAT_META_UUID+ |
444   Metadata UUID.  The UUID stamped into each metadata block must match the value
445   in +sb_meta_uuid+.  This enables the administrator to change +sb_uuid+ at will
446   without having to rewrite the entire filesystem.
447
448   | +XFS_SB_FEAT_INCOMPAT_BIGTIME+ |
449   Large timestamps.  Inode timestamps and quota expiration timers are extended to
450   support times through the year 2486.  See the section on
451   xref:Timestamps[timestamps] for more information.
452
453   | +XFS_SB_FEAT_INCOMPAT_NEEDSREPAIR+ |
454   The filesystem is not in operable condition, and must be run through
455   xfs_repair before it can be mounted.
456
457   | +XFS_SB_FEAT_INCOMPAT_NREXT64+ |
458   Large file fork extent counts.  This greatly expands the maximum number of
459   space mappings allowed in data and extended attribute file forks.
460
461   |=====
462
463   *sb_features_log_incompat*::
464   Read-write incompatible feature flags for the log.  The kernel cannot recover
465   the FS log if it doesn't understand the flag.
466
467   .Extended Version 5 Superblock Log incompatibility flags
468   [options="header"]
469   |=====
470   | Flag                               | Description
471   | +XFS_SB_FEAT_INCOMPAT_LOG_XATTRS+      |
472   Extended attribute updates have been committed to the ondisk log.
473
474   |=====
475
476   *sb_crc*::
477   Superblock checksum.
478
479   *sb_spino_align*::
480   Sparse inode alignment, in fsblocks.  Each chunk of inodes referenced by a
481   sparse inode B+tree record must be aligned to this block granularity.
```

```
482
483    *sb_pquotino*::
484    Project quota inode.
485
486    *sb_lsn*::
487    Log sequence number of the last superblock update.
488
489    *sb_meta_uuid*::
490    If the +XFS_SB_FEAT_INCOMPAT_META_UUID+ feature is set, then the UUID field in
491    all metadata blocks must match this UUID.  If not, the block header UUID field
492    must match +sb_uuid+.
493
494    *sb_rrmapino*::
495    If the +XFS_SB_FEAT_RO_COMPAT_RMAPBT+ feature is set and a real-time
496    device is present (+sb_rblocks+ > 0), this field points to an inode
497    that contains the root to the
498    xref:Real_time_Reverse_Mapping_Btree[Real-Time Reverse Mapping B+tree].
499    This field is zero otherwise.
500
501    === xfs_db Superblock Example
502
503    A filesystem is made on a single disk with the following command:
504
505    ----
506    # mkfs.xfs -i attr=2 -n size=16384 -f /dev/sda7
507    meta-data=/dev/sda7                isize=256    agcount=16, agsize=3923122 blks
508             =                         sectsz=512   attr=2
509    data      =                        bsize=4096   blocks=62769952, imaxpct=25
510             =                         sunit=0      swidth=0 blks, unwritten=1
511    naming    =version 2               bsize=16384
512    log       =internal log            bsize=4096   blocks=30649, version=1
513             =                         sectsz=512   sunit=0 blks
514    realtime =none                     extsz=65536  blocks=0, rtextents=0
515    ----
516
517    And in xfs_db, inspecting the superblock:
518
519    ----
520    xfs_db> sb
521    xfs_db> p
522    magicnum = 0x58465342
523    blocksize = 4096
524    dblocks = 62769952
525    rblocks = 0
526    rextents = 0
527    uuid = 32b24036-6931-45b4-b68c-cd5e7d9a1ca5
528    logstart = 33554436
529    rootino = 128
530    rbmino = 129
```

```
531   rsumino = 130
532   rextsize = 16
533   agblocks = 3923122
534   agcount = 16
535   rbmblocks = 0
536   logblocks = 30649
537   versionnum = 0xb084
538   sectsize = 512
539   inodesize = 256
540   inopblock = 16
541   fname = "\000\000\000\000\000\000\000\000\000\000\000\000"
542   blocklog = 12
543   sectlog = 9
544   inodelog = 8
545   inopblog = 4
546   agblklog = 22
547   rextslog = 0
548   inprogress = 0
549   imax_pct = 25
550   icount = 64
551   ifree = 61
552   fdblocks = 62739235
553   frextents = 0
554   uquotino = 0
555   gquotino = 0
556   qflags = 0
557   flags = 0
558   shared_vn = 0
559   inoalignmt = 2
560   unit = 0
561   width = 0
562   dirblklog = 2
563   logsectlog = 0
564   logsectsize = 0
565   logsunit = 0
566   features2 = 8
567   ----
568
569
570   [[AG_Free_Space_Management]]
571   == AG Free Space Management
572
573   The XFS filesystem tracks free space in an allocation group using two B+trees.
574   One B+tree tracks space by block number, the second by the size of the free
575   space block. This scheme allows XFS to find quickly free space near a given
576   block or of a given size.
577
578   All block numbers, indexes, and counts are AG relative.
579
```

```
580   [[AG_Free_Space_Block]]
581   === AG Free Space Block
582
583   The second sector in an AG contains the information about the two free space
584   B+trees and associated free space information for the AG. The ``AG Free Space
585   Block'' also knows as the +AGF+, uses the following structure:
586
587   [source, c]
588   ----
589   struct xfs_agf {
590       __be32                  agf_magicnum;
591       __be32                  agf_versionnum;
592       __be32                  agf_seqno;
593       __be32                  agf_length;
594       __be32                  agf_roots[XFS_BTNUM_AGF];
595       __be32                  agf_levels[XFS_BTNUM_AGF];
596       __be32                  agf_flfirst;
597       __be32                  agf_fllast;
598       __be32                  agf_flcount;
599       __be32                  agf_freeblks;
600       __be32                  agf_longest;
601       __be32                  agf_btreeblks;
602
603       /* version 5 filesystem fields start here */
604       uuid_t                  agf_uuid;
605       __be32                  agf_rmap_blocks;
606       __be32                  agf_refcount_blocks;
607       __be32                  agf_refcount_root;
608       __be32                  agf_refcount_level;
609       __be64                  agf_spare64[14];
610
611       /* unlogged fields, written during buffer writeback. */
612       __be64                  agf_lsn;
613       __be32                  agf_crc;
614       __be32                  agf_spare2;
615   };
616   ----
617
618   The rest of the bytes in the sector are zeroed. +XFS_BTNUM_AGF+ is set to 3:
619   index 0 for the free space B+tree indexed by block number; index 1 for the free
620   space B+tree indexed by extent size; and index 2 for the reverse-mapping
621   B+tree.
622
623   *agf_magicnum*::
624   Specifies the magic number for the AGF sector: ``XAGF'' (0x58414746).
625
626   *agf_versionnum*::
627   Set to +XFS_AGF_VERSION+ which is currently 1.
628
```

629    *agf_seqno*::
630    Specifies the AG number for the sector.
631
632    *agf_length*::
633    Specifies the size of the AG in filesystem blocks. For all AGs except the last,
634    this must be equal to the superblock's +sb_agblocks+ value. For the last AG,
635    this could be less than the +sb_agblocks+ value. It is this value that should
636    be used to determine the size of the AG.
637
638    *agf_roots*::
639    Specifies the block number for the root of the two free space B+trees and the
640    reverse-mapping B+tree, if enabled.
641
642    *agf_levels*::
643    Specifies the level or depth of the two free space B+trees and the
644    reverse-mapping B+tree, if enabled.  For a fresh AG, this value will be one,
645    and the ``roots'' will point to a single leaf of level 0.
646
647    *agf_flfirst*::
648    Specifies the index of the first ``free list'' block. Free lists are covered in
649    more detail later on.
650
651    *agf_fllast*::
652    Specifies the index of the last ``free list'' block.
653
654    *agf_flcount*::
655    Specifies the number of blocks in the ``free list''.
656
657    *agf_freeblks*::
658    Specifies the current number of free blocks in the AG.
659
660    *agf_longest*::
661    Specifies the number of blocks of longest contiguous free space in the AG.
662
663    *agf_btreeblks*::
664    Specifies the number of blocks used for the free space B+trees. This is only
665    used if the +XFS_SB_VERSION2_LAZYSBCOUNTBIT+ bit is set in +sb_features2+.
666
667    *agf_uuid*::
668    The UUID of this block, which must match either +sb_uuid+ or +sb_meta_uuid+
669    depending on which features are set.
670
671    *agf_rmap_blocks*::
672    The size of the reverse mapping B+tree in this allocation group, in blocks.
673
674    *agf_refcount_blocks*::
675    The size of the reference count B+tree in this allocation group, in blocks.
676
677    *agf_refcount_root*::

678     Block number for the root of the reference count B+tree, if enabled.

679

680     *agf_refcount_level*::

681     Depth of the reference count B+tree, if enabled.

682

683     *agf_spare64*::

684     Empty space in the logged part of the AGF sector, for use for future features.

685

686     *agf_lsn*::

687     Log sequence number of the last AGF write.

688

689     *agf_crc*::

690     Checksum of the AGF sector.

691

692     *agf_spare2*::

693     Empty space in the unlogged part of the AGF sector.

694

695     [[AG_Free_Space_Btrees]]

696     === AG Free Space B+trees

697

698     The two Free Space B+trees store a sorted array of block offset and block

699     counts in the leaves of the B+tree. The first B+tree is sorted by the offset,

700     the second by the count or size.

701

702     Leaf nodes contain a sorted array of offset/count pairs which are also used for

703     node keys:

704

705     [source, c]

706     ----

707     struct xfs_alloc_rec {

708         __be32                      ar_startblock;

709         __be32                      ar_blockcount;

710     };

711     ----

712

713     *ar_startblock*::

714     AG block number of the start of the free space.

715

716     *ar_blockcount*::

717     Length of the free space.

718

719     Node pointers are an AG relative block pointer:

720

721     [source, c]

722     ----

723     typedef __be32 xfs_alloc_ptr_t;

724     ----

725

726     * As the free space tracking is AG relative, all the block numbers are only

727   32-bits.
728   * The +bb_magic+ value depends on the B+tree: ``ABTB'' (0x41425442) for the block
729   offset B+tree, ``ABTC'' (0x41425443) for the block count B+tree.  On a v5
730   filesystem, these are ``AB3B'' (0x41423342) and ``AB3C'' (0x41423343),
731   respectively.
732   * The +xfs_btree_sblock_t+ header is used for intermediate B+tree node as well
733   as the leaves.
734   * For a typical 4KB filesystem block size, the offset for the +xfs_alloc_ptr_t+
735   array would be +0xab0+ (2736 decimal).
736   * There are a series of macros in +xfs_btree.h+ for deriving the offsets,
737   counts, maximums, etc for the B+trees used in XFS.
738
739   The following diagram shows a single level B+tree which consists of one leaf:
740
741   .Freespace B+tree with one leaf.
742   image::images/15a.png[]
743
744   With the intermediate nodes, the associated leaf pointers are stored in a
745   separate array about two thirds into the block. The following diagram
746   illustrates a 2-level B+tree for a free space B+tree:
747
748   .Multi-level freespace B+tree.
749   image::images/15b.png[]
750
751   [[AG_Free_List]]
752   === AG Free List
753
754   The AG Free List is located in the 4^th^ sector of each AG and is known as the
755   AGFL. It is an array of AG relative block pointers for reserved space for
756   growing the free space B+trees. This space cannot be used for general user data
757   including inodes, data, directories and extended attributes.
758
759   With a freshly made filesystem, 4 blocks are reserved immediately after the free
760   space B+tree root blocks (blocks 4 to 7). As they are used up as the free space
761   fragments, additional blocks will be reserved from the AG and added to the free
762   list array.  This size may increase as features are added.
763
764   As the free list array is located within a single sector, a typical device will
765   have space for 128 elements in the array (512 bytes per sector, 4 bytes per AG
766   relative block pointer). The actual size can be determined by using the
767   +XFS_AGFL_SIZE+ macro.
768
769   Active elements in the array are specified by the
770   xref:AG_Free_Space_Block[AGF's] +agf_flfirst+, +agf_fllast+ and +agf_flcount+
771   values. The array is managed as a circular list.
772
773   On a v5 filesystem, the following header precedes the free list entries:
774
775   [source, c]

```
776    ----
777    struct xfs_agfl {
778        __be32              agfl_magicnum;
779        __be32              agfl_seqno;
780        uuid_t              agfl_uuid;
781        __be64              agfl_lsn;
782        __be32              agfl_crc;
783    };
784    ----
785
786    *agfl_magicnum*::
787    Specifies the magic number for the AGFL sector: "XAFL" (0x5841464c).
788
789    *agfl_seqno*::
790    Specifies the AG number for the sector.
791
792    *agfl_uuid*::
793    The UUID of this block, which must match either +sb_uuid+ or +sb_meta_uuid+
794    depending on which features are set.
795
796    *agfl_lsn*::
797    Log sequence number of the last AGFL write.
798
799    *agfl_crc*::
800    Checksum of the AGFL sector.
801
802    On a v4 filesystem there is no header; the array of free block numbers begins
803    at the beginning of the sector.
804
805    .AG Free List layout
806    image::images/16.png[]
807
808    The presence of these reserved blocks guarantees that the free space B+trees
809    can be updated if any blocks are freed by extent changes in a full AG.
810
811    ==== xfs_db AGF Example
812
813    These examples are derived from an AG that has been deliberately fragmented.
814    The AGF:
815
816    ----
817    xfs_db> agf 0
818    xfs_db> p
819    magicnum = 0x58414746
820    versionnum = 1
821    seqno = 0
822    length = 3923122
823    bnoroot = 7
824    cntroot = 83343
```

```
825   bnolevel = 2
826   cntlevel = 2
827   flfirst = 22
828   fllast = 27
829   flcount = 6
830   freeblks = 3654234
831   longest = 3384327
832   btreeblks = 0
833   ----
834
835   In the AGFL, the active elements are from 22 to 27 inclusive which are obtained
836   from the +flfirst+ and +fllast+ values from the +agf+ in the previous example:
837
838   ----
839   xfs_db> agfl 0
840   xfs_db> p
841   bno[0-127] = 0:4 1:5 2:6 3:7 4:83342 5:83343 6:83344 7:83345 8:83346 9:83347
842               10:4 11:5 12:80205 13:80780 14:81496 15:81766 16:83346 17:4 18:5
843               19:80205 20:82449 21:81496 22:81766 23:82455 24:80780 25:5
844               26:80205 27:83344
845   ----
846
847   The root block of the free space B+tree sorted by block offset is found in the
848   AGF's +bnoroot+ value:
849
850   ----
851   xfs_db> fsblock 7
852   xfs_db> type bnobt
853   xfs_db> p
854   magic = 0x41425442
855   level = 1
856   numrecs = 4
857   leftsib = null
858   rightsib = null
859   keys[1-4] = [startblock,blockcount]
860             1:[12,16] 2:[184586,3] 3:[225579,1] 4:[511629,1]
861   ptrs[1-4] = 1:2 2:83347 3:6 4:4
862   ----
863
864   Blocks 2, 83347, 6 and 4 contain the leaves for the free space B+tree by
865   starting block. Block 2 would contain offsets 12 up to but not including 184586
866   while block 4 would have all offsets from 511629 to the end of the AG.
867
868   The root block of the free space B+tree sorted by block count is found in the
869   AGF's +cntroot+ value:
870
871   ----
872   xfs_db> fsblock 83343
873   xfs_db> type cntbt
```

```
874   xfs_db> p
875   magic = 0x41425443
876   level = 1
877   numrecs = 4
878   leftsib = null
879   rightsib = null
880   keys[1-4] = [blockcount,startblock]
881           1:[1,81496] 2:[1,511729] 3:[3,191875] 4:[6,184595]
882   ptrs[1-4] = 1:3 2:83345 3:83342 4:83346
883   ----
884
885   The leaf in block 3, in this example, would only contain single block counts.
886   The offsets are sorted in ascending order if the block count is the same.
887
888   Inspecting the leaf in block 83346, we can see the largest block at the end:
889
890   ----
891   xfs_db> fsblock 83346
892   xfs_db> type cntbt
893   xfs_db> p
894   magic = 0x41425443
895   level = 0
896   numrecs = 344
897   leftsib = 83342
898   rightsib = null
899   recs[1-344] = [startblock,blockcount]
900           1:[184595,6] 2:[187573,6] 3:[187776,6]
901           ...
902           342:[513712,755] 343:[230317,258229] 344:[538795,3384327]
903   ----
904
905   The longest block count (3384327) must be the same as the AGF's +longest+ value.
906
907   [[AG_Inode_Management]]
908   == AG Inode Management
909
910   [[Inode_Numbers]]
911   === Inode Numbers
912
913   Inode numbers in XFS come in two forms: AG relative and absolute.
914
915   AG relative inode numbers always fit within 32 bits. The number of bits actually
916   used is determined by the sum of the xref:Superblocks[superblock's] +sb_inoplog+
917   and +sb_agblklog+ values. Relative inode numbers are found within the AG's inode
918   structures.
919
920   Absolute inode numbers include the AG number in the high bits, above the bits
921   used for the AG relative inode number. Absolute inode numbers are found in
922   xref:Directories[directory] entries and the superblock.
```

```
923
924    .Inode number formats
925    image::images/18.png[]
926
927    [[Inode_Information]]
928    === Inode Information
929
930    Each AG manages its own inodes. The third sector in the AG contains information
931    about the AG's inodes and is known as the AGI.
932
933    The AGI uses the following structure:
934
935    [source, c]
936    ----
937    struct xfs_agi {
938        __be32              agi_magicnum;
939        __be32              agi_versionnum;
940        __be32              agi_seqno
941        __be32              agi_length;
942        __be32              agi_count;
943        __be32              agi_root;
944        __be32              agi_level;
945        __be32              agi_freecount;
946        __be32              agi_newino;
947        __be32              agi_dirino;
948        __be32              agi_unlinked[64];
949
950        /*
951         * v5 filesystem fields start here; this marks the end of logging region 1
952         * and start of logging region 2.
953         */
954        uuid_t              agi_uuid;
955        __be32              agi_crc;
956        __be32              agi_pad32;
957        __be64              agi_lsn;
958
959        __be32              agi_free_root;
960        __be32              agi_free_level;
961
962        __be32              agi_iblocks;
963        __be32              agi_fblocks;
964
965    }
966    ----
967    *agi_magicnum*::
968    Specifies the magic number for the AGI sector: ``XAGI'' (0x58414749).
969
970    *agi_versionnum*::
971    Set to +XFS_AGI_VERSION+ which is currently 1.
```

```
 972
 973   *agi_seqno*::
 974   Specifies the AG number for the sector.
 975
 976   *agi_length*::
 977   Specifies the size of the AG in filesystem blocks.
 978
 979   *agi_count*::
 980   Specifies the number of inodes allocated for the AG.
 981
 982   *agi_root*::
 983   Specifies the block number in the AG containing the root of the inode B+tree.
 984
 985   *agi_level*::
 986   Specifies the number of levels in the inode B+tree.
 987
 988   *agi_freecount*::
 989   Specifies the number of free inodes in the AG.
 990
 991   *agi_newino*::
 992   Specifies AG-relative inode number of the most recently allocated chunk.
 993
 994   *agi_dirino*::
 995   Deprecated and not used, this is always set to NULL (-1).
 996
 997   *agi_unlinked[64]*::
 998   Hash table of unlinked (deleted) inodes that are still being referenced. Refer
 999   to xref:Unlinked_Pointer[unlinked list pointers] for more information.
1000
1001   *agi_uuid*::
1002   The UUID of this block, which must match either +sb_uuid+ or +sb_meta_uuid+
1003   depending on which features are set.
1004
1005   *agi_crc*::
1006   Checksum of the AGI sector.
1007
1008   *agi_pad32*::
1009   Padding field, otherwise unused.
1010
1011   *agi_lsn*::
1012   Log sequence number of the last write to this block.
1013
1014   *agi_free_root*::
1015   Specifies the block number in the AG containing the root of the free inode
1016   B+tree.
1017
1018   *agi_free_level*::
1019   Specifies the number of levels in the free inode B+tree.
1020
```

```
1021   *agi_iblocks*::
1022   The number of blocks in the inode B+tree, including the root.
1023   This field is zero if the +XFS_SB_FEAT_RO_COMPAT_INOBTCNT+ feature is not
1024   enabled.
1025
1026   *agi_fblocks*::
1027   The number of blocks in the free inode B+tree, including the root.
1028   This field is zero if the +XFS_SB_FEAT_RO_COMPAT_INOBTCNT+ feature is not
1029   enabled.
1030
1031   [[Inode_Btrees]]
1032   == Inode B+trees
1033
1034   Inodes are traditionally allocated in chunks of 64, and a B+tree is used to
1035   track these chunks of inodes as they are allocated and freed. The block
1036   containing root of the B+tree is defined by the AGI's +agi_root+ value.  If the
1037   +XFS_SB_FEAT_RO_COMPAT_FINOBT+ feature is enabled, a second B+tree is used to
1038   track the chunks containing free inodes; this is an optimization to speed up
1039   inode allocation.
1040
1041   The B+tree header for the nodes and leaves use the +xfs_btree_sblock+ structure
1042   which is the same as the header used in the xref:AG_Free_Space_Btrees[AGF
1043   B+trees].
1044
1045   The magic number of the inode B+tree is ``IABT'' (0x49414254).  On a v5
1046   filesystem, the magic number is ``IAB3'' (0x49414233).
1047
1048   The magic number of the free inode B+tree is ``FIBT'' (0x46494254).  On a v5
1049   filesystem, the magic number is ``FIB3'' (0x46494254).
1050
1051   Leaves contain an array of the following structure:
1052
1053   [source,c]
1054   ----
1055   struct xfs_inobt_rec {
1056        __be32                      ir_startino;
1057        __be32                      ir_freecount;
1058        __be64                      ir_free;
1059   };
1060   ----
1061
1062   *ir_startino*::
1063   The lowest-numbered inode in this chunk.
1064
1065   *ir_freecount*::
1066   Number of free inodes in this chunk.
1067
1068   *ir_free*::
1069   A 64 element bitmap showing which inodes in this chunk are free.
```

```
1070
1071    Nodes contain key/pointer pairs using the following types:
1072
1073    [source,c]
1074    ----
1075    struct xfs_inobt_key {
1076          __be32                          ir_startino;
1077    };
1078    typedef __be32 xfs_inobt_ptr_t;
1079    ----
1080
1081    The following diagram illustrates a single level inode B+tree:
1082
1083    .Single Level inode B+tree
1084    image::images/20a.png[]
1085
1086
1087    And a 2-level inode B+tree:
1088
1089    .Multi-Level inode B+tree
1090    image::images/20b.png[]
1091
1092
1093    === xfs_db AGI Example
1094
1095    This is an AGI of a freshly populated filesystem:
1096
1097    ----
1098    xfs_db> agi 0
1099    xfs_db> p
1100    magicnum = 0x58414749
1101    versionnum = 1
1102    seqno = 0
1103    length = 825457
1104    count = 5440
1105    root = 3
1106    level = 1
1107    freecount = 9
1108    newino = 5792
1109    dirino = null
1110    unlinked[0-63] =
1111    uuid = 3dfa1e5c-5a5f-4ca2-829a-000e453600fe
1112    lsn = 0x1000032c2
1113    crc = 0x14cb7e5c (correct)
1114    free_root = 4
1115    free_level = 1
1116    ----
1117
1118    From this example, we see that the inode B+tree is rooted at AG block 3 and
```

```
1119    that the free inode B+tree is rooted at AG block 4.  Let's look at the
1120    inode B+tree:
1121
1122    ----
1123    xfs_db> addr root
1124    xfs_db> p
1125    magic = 0x49414233
1126    level = 0
1127    numrecs = 85
1128    leftsib = null
1129    rightsib = null
1130    bno = 24
1131    lsn = 0x1000032c2
1132    uuid = 3dfa1e5c-5a5f-4ca2-829a-000e453600fe
1133    owner = 0
1134    crc = 0x768f9592 (correct)
1135    recs[1-85] = [startino,freecount,free]
1136            1:[96,0,0] 2:[160,0,0] 3:[224,0,0] 4:[288,0,0]
1137            5:[352,0,0] 6:[416,0,0] 7:[480,0,0] 8:[544,0,0]
1138            9:[608,0,0] 10:[672,0,0] 11:[736,0,0] 12:[800,0,0]
1139            ...
1140            85:[5792,9,0xff80000000000000]
1141    ----
1142
1143    Most of the inode chunks on this filesystem are totally full, since the +free+
1144    value is zero.  This means that we ought to expect inode 160 to be linked
1145    somewhere in the directory structure.  However, notice that 0xff80000000000000
1146    in record 85 -- this means that we would expect inode 5847 to be free.  Moving
1147    on to the free inode B+tree, we see that this is indeed the case:
1148
1149    ----
1150    xfs_db> addr free_root
1151    xfs_db> p
1152    magic = 0x46494233
1153    level = 0
1154    numrecs = 1
1155    leftsib = null
1156    rightsib = null
1157    bno = 32
1158    lsn = 0x1000032c2
1159    uuid = 3dfa1e5c-5a5f-4ca2-829a-000e453600fe
1160    owner = 0
1161    crc = 0x338af88a (correct)
1162    recs[1] = [startino,freecount,free] 1:[5792,9,0xff80000000000000]
1163    ----
1164
1165    Observe also that the AGI's +agi_newino+ points to this chunk, which has never
1166    been fully allocated.
1167
```

```
1168   [[Sparse_Inodes]]
1169   == Sparse Inodes
1170
1171   As mentioned in the previous section, XFS allocates inodes in chunks of 64.  If
1172   there are no free extents large enough to hold a full chunk of 64 inodes, the
1173   inode allocation fails and XFS claims to have run out of space.  On a
1174   filesystem with highly fragmented free space, this can lead to out of space
1175   errors long before the filesystem runs out of free blocks.
1176
1177   The sparse inode feature tracks inode chunks in the inode B+tree as if they
1178   were full chunks but uses some previously unused bits in the freecount field to
1179   track which parts of the inode chunk are not allocated for use as inodes.  This
1180   allows XFS to allocate inodes one block at a time if absolutely necessary.
1181
1182   The inode and free inode B+trees operate in the same manner as they do without
1183   the sparse inode feature; the B+tree header for the nodes and leaves use the
1184   +xfs_btree_sblock+ structure which is the same as the header used in the
1185   xref:AG_Free_Space_Btrees[AGF B+trees].
1186
1187   It is theoretically possible for a sparse inode B+tree record to reference
1188   multiple non-contiguous inode chunks.
1189
1190   Leaves contain an array of the following structure:
1191
1192   [source,c]
1193   ----
1194   struct xfs_inobt_rec {
1195        __be32                    ir_startino;
1196        __be16                    ir_holemask;
1197        __u8                      ir_count;
1198        __u8                      ir_freecount;
1199        __be64                    ir_free;
1200   };
1201   ----
1202
1203   *ir_startino*::
1204   The lowest-numbered inode in this chunk, rounded down to the nearest multiple
1205   of 64, even if the start of this chunk is sparse.
1206
1207   *ir_holemask*::
1208   A 16 element bitmap showing which parts of the chunk are not allocated to
1209   inodes.  Each bit represents four inodes; if a bit is marked here, the
1210   corresponding bits in ir_free must also be marked.
1211
1212   *ir_count*::
1213   Number of inodes allocated to this chunk.
1214
1215   *ir_freecount*::
1216   Number of free inodes in this chunk.
```

```
1217
1218   *ir_free*::
1219   A 64 element bitmap showing which inodes in this chunk are not available for
1220   allocation.
1221
1222   === xfs_db Sparse Inode AGI Example
1223
1224   This example derives from an AG that has been deliberately fragmented.  The
1225   inode B+tree:
1226
1227   ----
1228   xfs_db> agi 0
1229   xfs_db> p
1230   magicnum = 0x58414749
1231   versionnum = 1
1232   seqno = 0
1233   length = 6400
1234   count = 10432
1235   root = 2381
1236   level = 2
1237   freecount = 0
1238   newino = 14912
1239   dirino = null
1240   unlinked[0-63] =
1241   uuid = b9b4623b-f678-4d48-8ce7-ce08950e3cd6
1242   lsn = 0x600000ac4
1243   crc = 0xef550dbc (correct)
1244   free_root = 4
1245   free_level = 1
1246   ----
1247
1248   This AGI was formatted on a v5 filesystem; notice the extra v5 fields.  So far
1249   everything else looks much the same as always.
1250
1251   ----
1252   xfs_db> addr root
1253   magic = 0x49414233
1254   level = 1
1255   numrecs = 2
1256   leftsib = null
1257   rightsib = null
1258   bno = 19048
1259   lsn = 0x50000192b
1260   uuid = b9b4623b-f678-4d48-8ce7-ce08950e3cd6
1261   owner = 0
1262   crc = 0xd98cd2ca (correct)
1263   keys[1-2] = [startino] 1:[128] 2:[35136]
1264   ptrs[1-2] = 1:3 2:2380
1265   xfs_db> addr ptrs[1]
```

```
1266  xfs_db> p
1267  magic = 0x49414233
1268  level = 0
1269  numrecs = 159
1270  leftsib = null
1271  rightsib = 2380
1272  bno = 24
1273  lsn = 0x600000ac4
1274  uuid = b9b4623b-f678-4d48-8ce7-ce08950e3cd6
1275  owner = 0
1276  crc = 0x836768a6 (correct)
1277  recs[1-159] = [startino,holemask,count,freecount,free]
1278          1:[128,0,64,0,0]
1279          2:[14912,0xff,32,0,0xffffffff]
1280          3:[15040,0,64,0,0]
1281          4:[15168,0xff00,32,0,0xffffffff00000000]
1282          5:[15296,0,64,0,0]
1283          6:[15424,0xff,32,0,0xffffffff]
1284          7:[15552,0,64,0,0]
1285          8:[15680,0xff00,32,0,0xffffffff00000000]
1286          9:[15808,0,64,0,0]
1287          10:[15936,0xff,32,0,0xffffffff]
1288  ----
```

Here we see the difference in the inode B+tree records.  For example, in record 2, we see that the holemask has a value of 0xff.  This means that the first sixteen inodes in this chunk record do not actually map to inode blocks; the first inode in this chunk is actually inode 14944:

```
xfs_db> inode 14912
Metadata corruption detected at block 0x3a40/0x2000
...
Metadata CRC error detected for ino 14912
xfs_db> p core.magic
core.magic = 0
xfs_db> inode 14944
xfs_db> p core.magic
core.magic = 0x494e
----
```

The chunk record also indicates that this chunk has 32 inodes, and that the missing inodes are also ``free''.

[[Real-time_Devices]]
== Real-time Devices

The performance of the standard XFS allocator varies depending on the internal state of the various metadata indices enabled on the filesystem.  For

```
1315   applications which need to minimize the jitter of allocation latency, XFS
1316   supports the notion of a ``real-time device''.  This is a special device
1317   separate from the regular filesystem where extent allocations are tracked with
1318   a bitmap and free space is indexed with a two-dimensional array.  If an inode
1319   is flagged with +XFS_DIFLAG_REALTIME+, its data will live on the real time
1320   device.  The metadata for real time devices is discussed in the section about
1321   xref:Real-time_Inodes[real time inodes].
1322
1323   By placing the real time device (and the journal) on separate high-performance
1324   storage devices, it is possible to reduce most of the unpredictability in I/O
1325   response times that come from metadata operations.
1326
1327   None of the XFS per-AG B+trees are involved with real time files.  It is not
1328   possible for real time files to share data blocks.
```

Powered by Gitiles | Privacy | Terms                                      txt    json