

第十八讲：文件系统实例

第 2 节：EXT4 文件系统

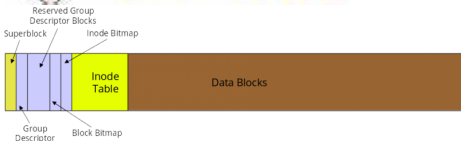
向勇、陈渝

清华大学计算机系

xyong,yuchen@tsinghua.edu.cn

2020 年 4 月 19 日

历史



Ext 2 3 4

Ext4 File System

把

Linux 文件系统历史

尽管 EXT[1-4] 文件系统家族是为 Linux 编写的，但它的根源是 Minix 操作系统和 Minix 文件系统，它们早于 Linux 大约五年，于 1987 年首次发布。如果我们回顾一下 EXT 文件系统家族历史和技术从其 Minix 根源的演变，了解 EXT4 文件系统要容易得多。

Reference:

OSTEP txtbook: Crash Consistency: FSCheck and Journaling; Ext4 slides from Mingming Cao; An introduction to Linux's EXT4 filesystem by David Both

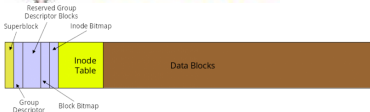


- Linus Torvalds 刚开始编写 Linux 内核时，他需要一个文件系统，但是不想编写一个文件系统。
- 拿来主义：直接采用 Andrew S. Tanenbaum 编写 Minix 文件系统



MINIX FS

- 引导扇区：在其上安装了硬盘的第一个扇区。引导块包括一个很小的引导记录和一个分区表。
- 超级块：每个分区中的第一个块是一个超级块，其中包含元数据，该元数据定义了其他文件系统结构，并将它们定位在分配给该分区的物理磁盘上。



- Linus Torvalds 在编写原始的 Linux 内核时，他需要一个文件系统，但是不想编写一个文件系统。
- 拿来主义：直接采用 Andrew S. Tanenbaum 编写 Minix 文件系统



MINIX FS

- 索引节点区位图：表明哪些索引节点已使用和哪些索引节点空闲。
- 索引节点区：它在磁盘上有自己的空间。每个 inode 包含有关一个文件的信息，包括数据块的位置，即属于该文件的区域。
- 数据区位图：来跟踪使用和免费的数据区。
- 数据区：保存实际存储数据。

Ext 2 3 4

- 最初的 EXT 文件系统（扩展）由 RémyCard 编写，并于 1992 年随 Linux 发行，以克服 Minix 文件系统的某些大小限制。

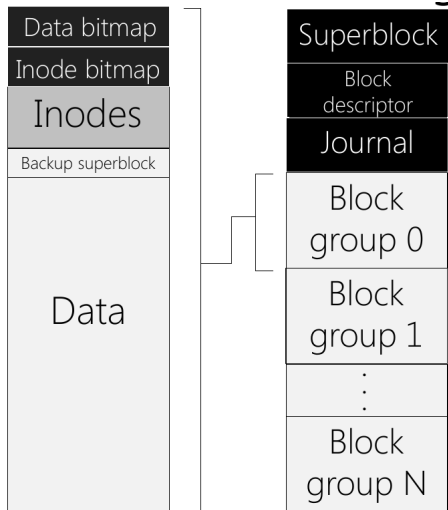
EXT FS

- 主要的结构更改是基于 Unix 文件系统（UFS）（也称为 Berkeley 快速文件系统（FFS））的文件系统的元数据。很少有关于 EXT 文件系统的公开信息可以验证，显然是因为它存在重大问题，并很快被 EXT2 文件系统所取代。

Ext 2 3 4

- EXT3 的唯一目标是克服花费大量时间恢复异常文件系统的时间。EXT3 文件系统增加了 journal 机制，它预先记录了将对文件系统执行的更改。其余磁盘结构与 EXT2 中的相同。

EXT3 FS



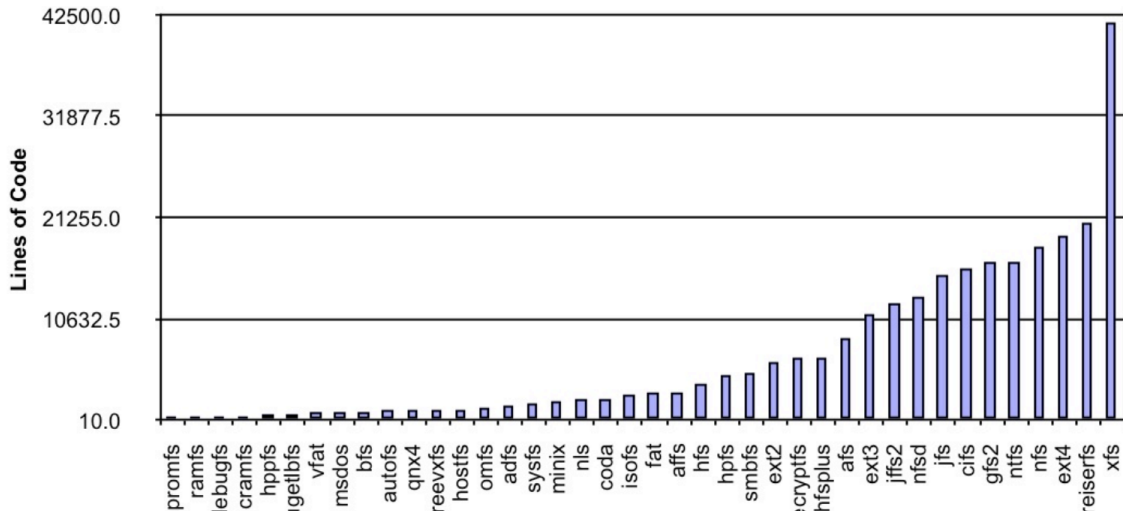


- EXT4 文件系统主要进一步改善了性能，可靠性和大容量支持

EXT4 FS

- 为了提高可靠性，添加了元数据和日志校验和。
- 为了满足各种关键任务要求，文件系统时间戳得到了改进，增加了纳秒精度时间间隔。
- 为提高容量和访问大容量文件的性能，把数据分配从固定块更改为扩展区

LOC for all File Systems in Linux 2.6.27



细节 – 支持大容量存储

EXT4 FS: 支持大容量存储

Ext4
File System

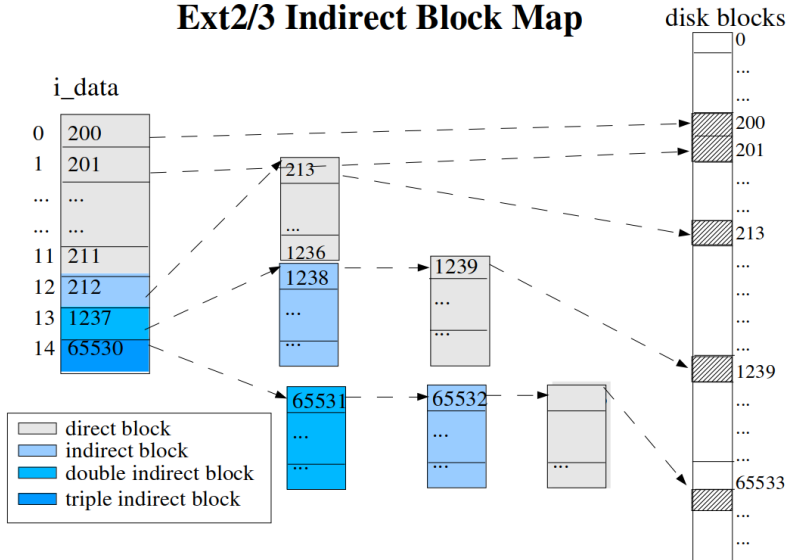
- EXT4 的重要特征：
支持大容量存储和
快速恢复异常状态



- 一部蓝光 8K 3D 电影，512GB~1TB
- 1EB 文件系统大小，16TB 文件大小，子目录个数无限制

细节 – 支持大容量存储

Ext2/3 Indirect Block Map



细节 – 支持大容量存储

extent: 一段连续的储存块

logical	length	physical
0	1000	200

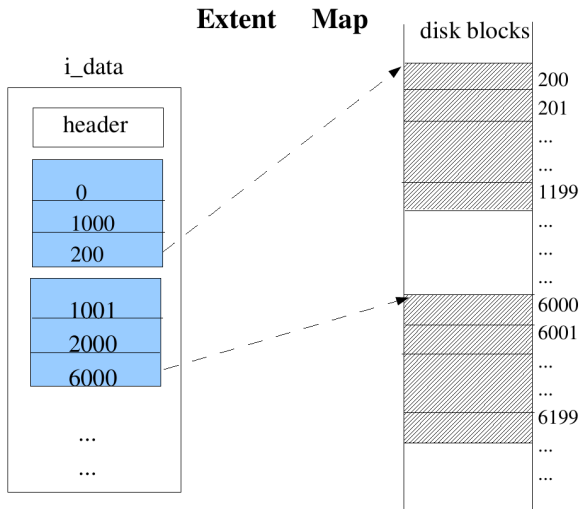
```
struct ext4_extent {  
    __le32 ee_block;    /* first logical block extent covers */  
    __le16 ee_len;      /* number of blocks covered by extent */  
    __le16 ee_start_hi; /* high 16 bits of physical block */  
    __le32 ee_start;    /* low 32 bits of physical block */  
};
```

细节 – 支持大容量存储

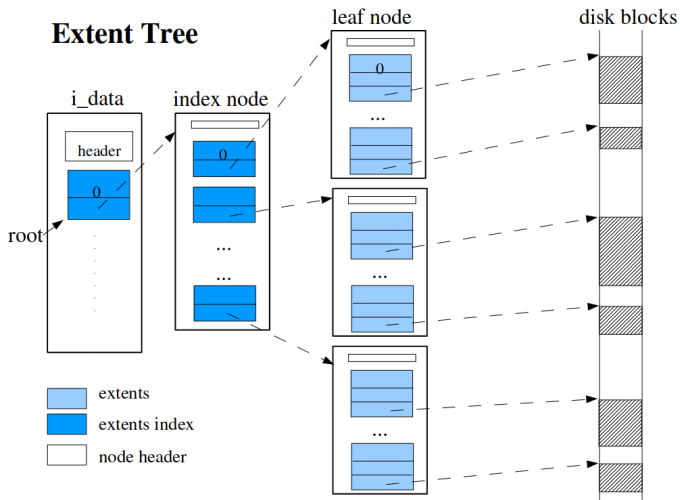
extent: 一段连续的储存块

logical	length	physical
0	1000	200

```
struct ext4_extent {
    __le32 ee_block; /* first logical block extent covers */
    __le16 ee_len; /* number of blocks covered by extent */
    __le16 ee_start_hi; /* high 16 bits of physical block */
    __le32 ee_start; /* low 32 bits of physical block */
};
```

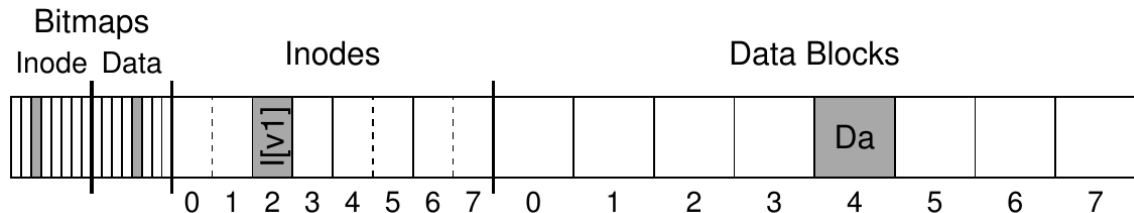


细节 – 支持大容量存储



细节 – 支持恢复异常

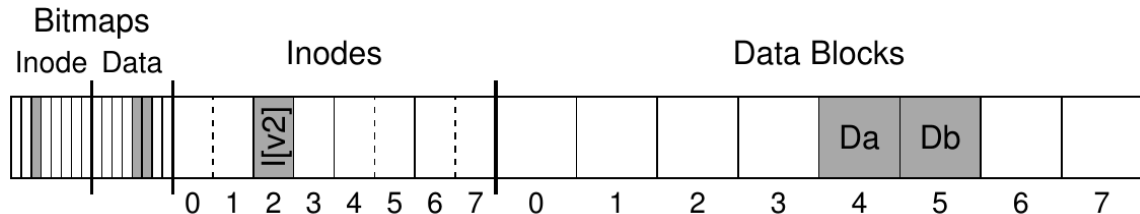
EXT4: 恢复异常文件系统 for crash-consistency problem



permissions : read-write
size : 1
pointer : 4
pointer : null
pointer : null
pointer : null

细节 – 支持恢复异常

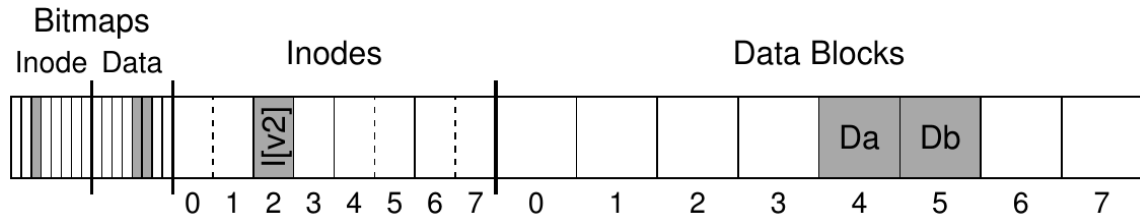
EXT4: 恢复异常文件系统 for crash-consistency problem



```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```


细节 – 支持恢复异常

EXT4: 恢复异常文件系统 for crash-consistency problem



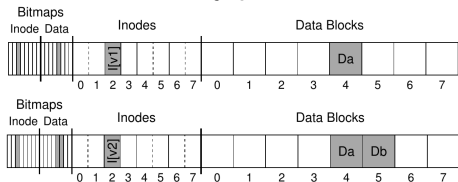
permissions : read-write
size : 2
pointer : 4
pointer : 5
pointer : null
pointer : null

3 个写操作

- inode ($I[v2]$)
- bitmap ($B[v2]$)
- data block (Db)

细节 – 支持恢复异常 – Crash 场景

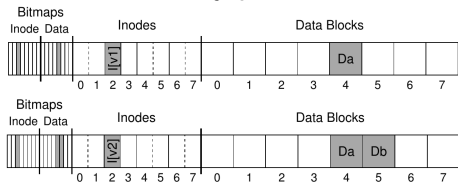
EXT4: 恢复异常文件系统 for crash-consistency problem



```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

细节 – 支持恢复异常 – Crash 场景

EXT4: 恢复异常文件系统 for crash-consistency problem



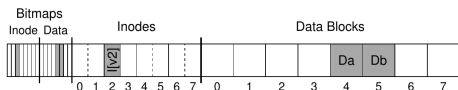
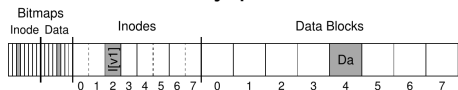
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

Crash 场景

- 只有 data block (Db) 写入磁盘
- 只有 inode (I[v2]) 写入磁盘
- 只有 bitmap (B[v2]) 写入磁盘

细节 – 支持恢复异常 – Crash 场景

EXT4: 恢复异常文件系统 for crash-consistency problem



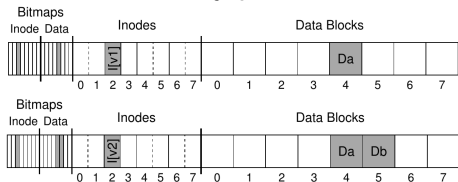
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

Crash 场景

- 只有 data block (Db) 写入磁盘
- 只有 inode (I[v2]) 写入磁盘
- 只有 bitmap (B[v2]) 写入磁盘
- inode(I[v2]) 和 bitmap (B[v2]) 写入磁盘
- inode(I[v2]) 和 data block(Db) 写入磁盘
- bitmap(B[v2]) 和 data block(Db) 写入磁盘

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



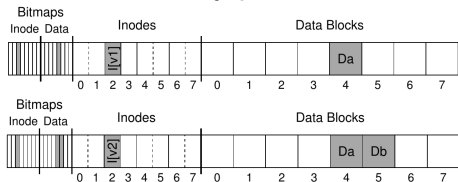
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1：File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



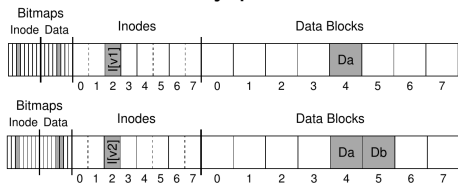
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1 : File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



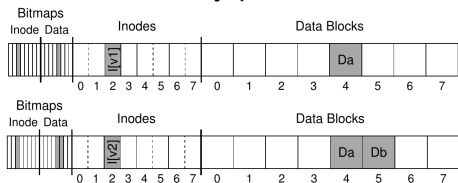
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1 : File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks
- Inode state/links: 有效的类型字段/链接计数

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



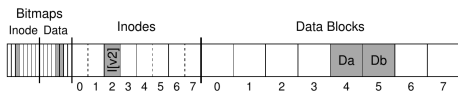
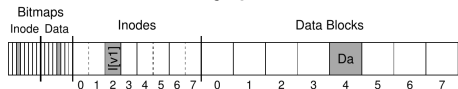
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1 : File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks
- Inode state/links: 有效的类型字段/链接计数
- Duplicates: 两个 inode 指向同一 data

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



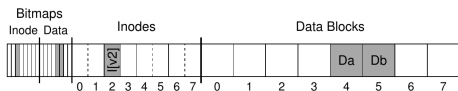
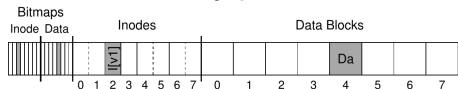
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1 : File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks
- Inode state/links: 有效的类型字段/链接计数
- Duplicates: 两个 inode 指向同一 data
- Bad blocks: 指针显然指向超出其有效范围

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



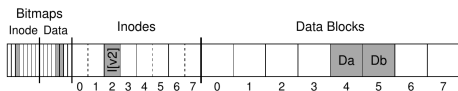
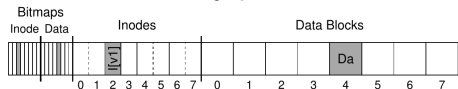
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 1：File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks
- Inode state/links: 有效的类型字段/链接计数
- Duplicates: 两个 inode 指向同一 data
- Bad blocks: 指针显然指向超出其有效范围
- Directory: “.” 和 “..” 是头两项

细节 – 支持恢复异常 – FSCK

EXT4: 恢复异常文件系统 for crash-consistency problem



```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

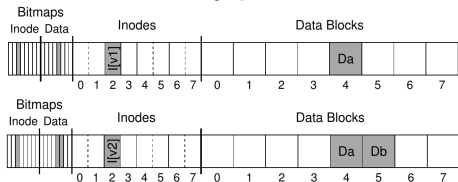
解决方法 1：File System Checker (fsck)

- Superblock: 文件系统大小大于已分配的块数
- Free blocks: bitmap v.s. Free blocks
- Inode state/links: 有效的类型字段/链接计数
- Duplicates: 两个 inode 指向同一 data
- Bad blocks: 指针显然指向超出其有效范围
- Directory: “.” 和 “..” 是头两项

太慢

细节 – 支持恢复异常 – Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



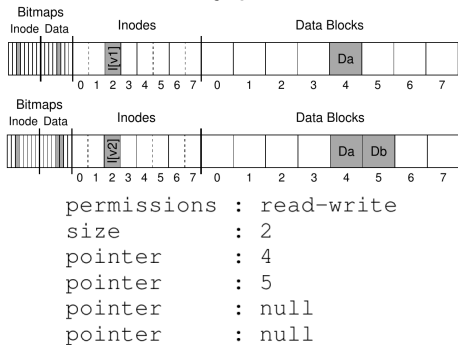
```
permissions : read-write
size        : 2
pointer     : 4
pointer     : 5
pointer     : null
pointer     : null
```

解决方法 2：日志 Journaling (Write-Ahead Logging)

- 从数据库管理系统的世界中窃取一个想法
- 最早（1987）在 Cedar 文件系统中出现
- 出现在 EXT3/4, JFS, XFS, NTFS 等

细节 – 支持恢复异常 – Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



解决方法 2：日志 Journaling (Write-Ahead Logging)

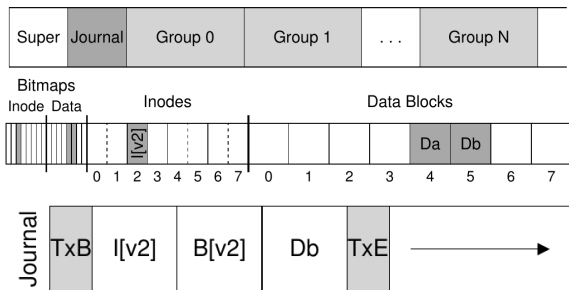
- 从数据库管理系统的世界中窃取一个想法
- 最早（1987）在 Cedar 文件系统中出现
- 出现在 EXT3/4, JFS, XFS, NTFS 等

基本思路

- 更新磁盘时，在覆盖相关结构之前，先写下一点日志（在磁盘上某个设定好的其他位置），以描述要执行的操作。

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem

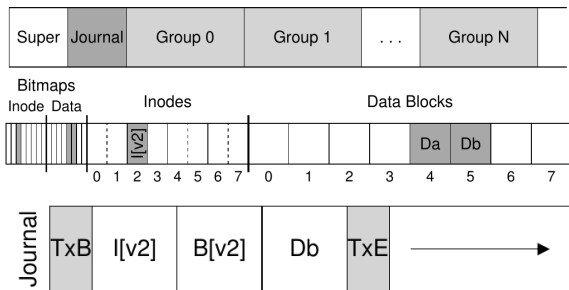


方法 2：Data Journaling

- TxB: transaction 开始
- Tx E: transaction 结束
- logical logging: 中间 3 块数据

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



方法 2：Data Journaling

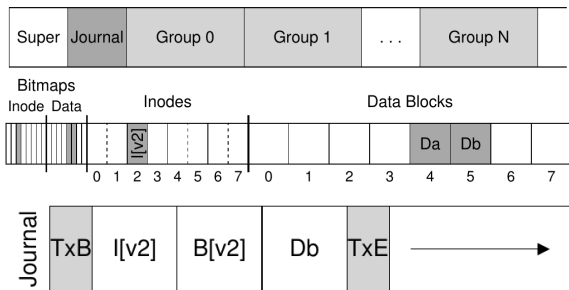
- TxB: transaction 开始
- TxE: transaction 结束
- logical logging: 中间 3 块数据

上述 transaction 写到磁盘上后, 更新磁盘, 覆盖相关结构 (checkpoint)

- I[v2]
- B[v2]
- Db

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem

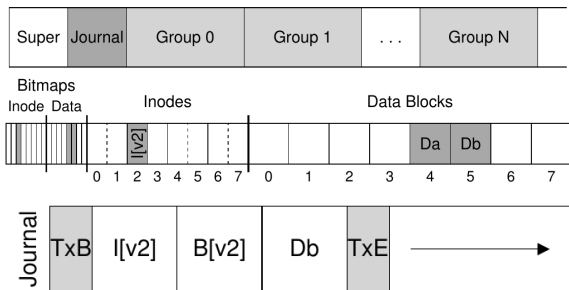


方法 2：Data Journaling

- TxB: transaction 开始
- Tx E: transaction 结束
- logical logging: 中间 3 块数据

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



方法 2：Data Journaling

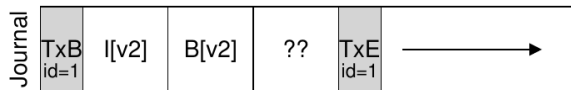
- TxB: transaction 开始
- TxE: transaction 结束
- logical logging: 中间 3 块数据

上述 transaction 写到磁盘上后, 更新磁盘, 覆盖相关结构 (checkpoint)

- I[v2]
- B[v2]
- Db

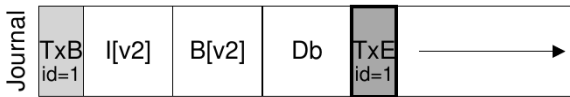
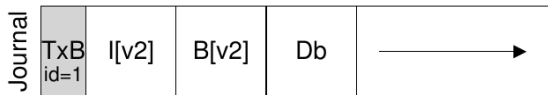
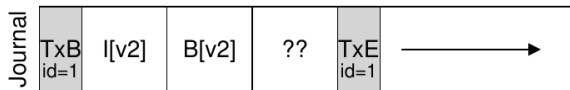
细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for
crash-consistency problem



细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for
crash-consistency problem

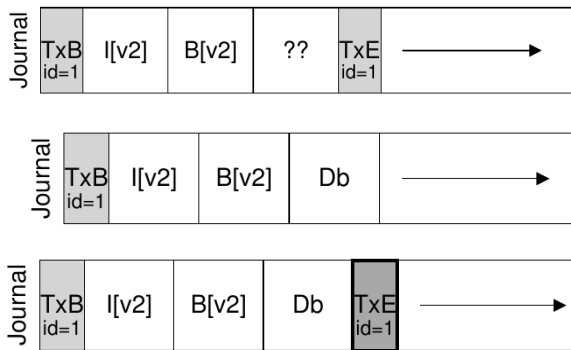


方法 2：Data Journaling

- Journal write
- Journal commit
- Checkpoint

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



方法 2：Data Journaling

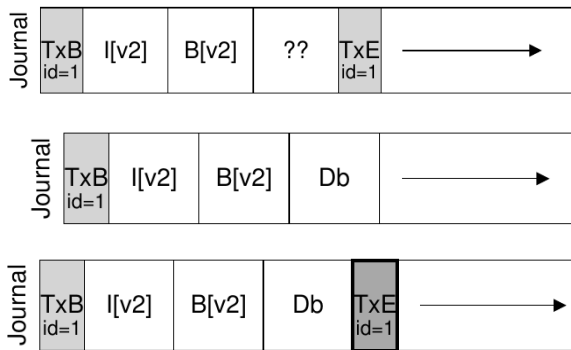
- Journal write
- Journal commit
- Checkpoint

恢复 (Recovery): 在此更新序列期间的任何时间都可能发生崩溃。

- 如果崩溃发生在将事务安全地写入日志之前
- 如果崩溃是在事务提交到日志之后但在检查点完成之前发生

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



方法 2：Data Journaling

- Journal write
- Journal commit
- Checkpoint

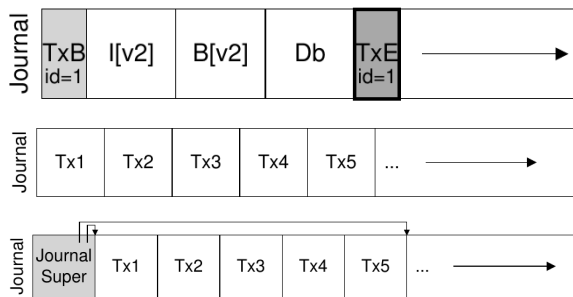
恢复 (Recovery): 在此更新序列期间的任何时间都可能发生崩溃。

- 如果崩溃发生在将事务安全地写入日志之前
- 如果崩溃是在事务提交到日志之后但在检查点完成之前发生

太多写，慢！

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem

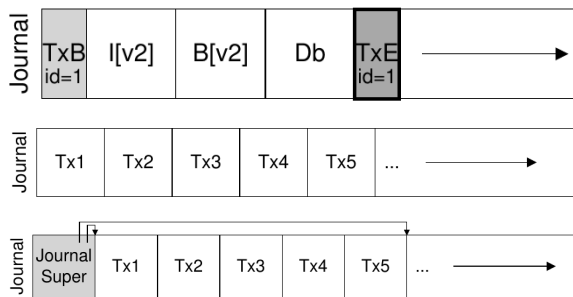


方法 2：Data Journaling: 提高速度

- 批处理日志更新
- 使日志有限：循环日志
- 日志超级块 journal superblock

细节 – 支持恢复异常 – Data Journaling

EXT4: 恢复异常文件系统 for crash-consistency problem



方法 2 : Data Journaling: 提高速度

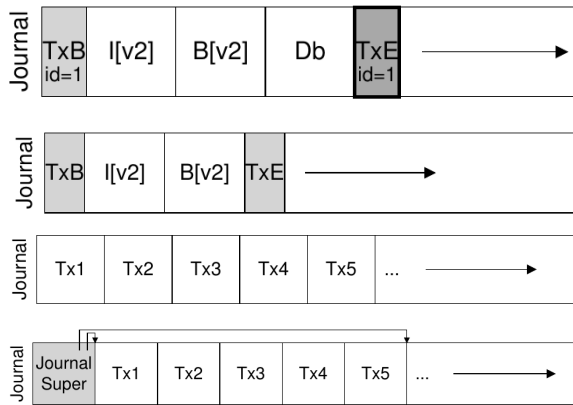
- 批处理日志更新
- 使日志有限：循环日志
- 日志超级块 journal superblock

新的更新过程

- Journal write
- Journal commit
- Checkpoint
- Free: 一段时间后，通过更新日记帐超级块将交易记录标记为空闲

细节 – 支持恢复异常 – Metadata Journaling

EXT4: 恢复异常文件系统 for
crash-consistency problem

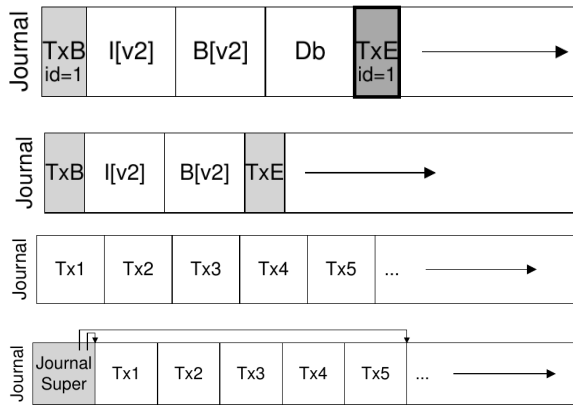


方法 2：Metadata Journaling: 进一步提高速度

- 我们什么时候应该将数据块 Db 写入磁盘？

细节 – 支持恢复异常 – Metadata Journaling

EXT4: 恢复异常文件系统 for
crash-consistency problem

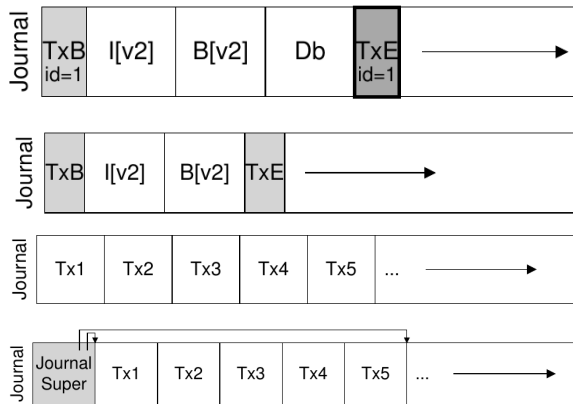


方法 2：Metadata Journaling: 进一步提高速度

- 我们什么时候应该将数据块 Db 写入磁盘？
- 事实证明，数据写入的顺序对于仅元数据的日记记录确实很重要。
- 如果在事务（包含 I [v2] 和 B [v2]）完成后将 Db 写入磁盘，这样有问题吗？

细节 – 支持恢复异常 – Metadata Journaling

EXT4: 恢复异常文件系统 for
crash-consistency problem



方法 2：Metadata Journaling: 进一步提高速度
新的更新过程

- Data write
- Journal metadata write
- Journal commit
- Checkpoint metadata
- Free

通过强制首先写入数据，文件系统可以保证指针永远不会指向垃圾数据。

细节 – 支持恢复异常 – Metadata Journaling

Data Journaling 时间线

TxB	Journal Contents		TxE	File System	
	(metadata)	(data)		Metadata	Data
issue complete	issue complete	issue complete			
			issue complete		
				issue complete	issue complete

Metadata Journaling 时间线

TxB	Journal Contents		TxE	File System	
	(metadata)			Metadata	Data
issue complete	issue complete				issue complete
			issue complete		
				issue complete	
					issue complete