

Christian Hummert  
Dirk Pawlaszczyk *Editors*

# Mobile Forensics – The File Format Handbook

Common File Formats and File Systems  
Used in Mobile Devices

OPEN ACCESS

 Springer

# Mobile Forensics – The File Format Handbook

Christian Hummert • Dirk Pawlaszczyk  
Editors

# Mobile Forensics – The File Format Handbook

Common File Formats and File Systems  
Used in Mobile Devices



Springer

*Editors*

Christian Hummert  
Agentur für Innovation in der Cybersicherheit  
Halle (Saale), Germany

Dirk Pawlaszczyk  
Fachgruppe Informatik  
Hochschule Mittweida  
Mittweida, Germany



ISBN 978-3-030-98466-3

ISBN 978-3-030-98467-0 (eBook)

<https://doi.org/10.1007/978-3-030-98467-0>

© The Editor(s) (if applicable) and The Author(s) 2022. This book is an open access publication.

**Open Access** This book is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this book are included in the book's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the book's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG  
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Preface

One of the biggest challenges in digital forensics is to gain a deep understanding of file systems and file formats. This knowledge is needed to recover files from corrupt file systems or reveal artefacts of former states of a computer system. For most file systems, it is easy to find resources that describe file systems at a high level. But for the detailed knowledge on the Hex-code level that is needed for digital forensics there are only few sources.

Brian Carrier did a great job in his book: “File System Forensic Analysis” [10]. This book was definitely a model for this work. When I started digital forensics, I really devoured this book. Reading Carrier I understood in detail how files are stored on computers and how to recover deleted files. “File System Forensic Analysis” gives deep insights into FAT, NTFS, EXT2, EXT3 and UFS. However, there are more file systems to discover: In my further work, I learned about EXT4, HFS+ and APFS. Concentrating on mobile forensics, there are even more file systems to explain.

In addition, there is more than file systems. Especially in mobile forensics, there are new file formats to encounter, which have a broader and more universal scope. File Formats like the SQLite database format are used in nearly every mobile device by millions of different Apps.

In January 2018, I started to write the proposal for the EU project FORMOBILE<sup>1</sup>. I aimed to provide an end-to-end mobile forensic investigation chain. I succeeded to build an outstanding consortium with 19 partners from 11 countries. Together we created a work plan that delivers novel tools to support mobile forensics, builds a new standard for mobile forensics and offers novel training for the forensic experts in this area. Happily, the project was funded and started in May 2020.

At the beginning of the proposal creation phase, we agreed to write a File Format Handbook that summarizes knowledge about various file formats and file systems common in mobile devices. I am more than happy to provide this file format handbook as a deliverable to the European Commission and a broader audience of forensic experts.

---

<sup>1</sup> <https://formobile-project.eu>

However, this book is not only a toolbox for experienced investigators that have learned about digital investigations from real cases and using analysis tools. The other target audience is students. Moreover, the book is aimed at people who are new to digital forensics and are more interested in the general theory of file recovery and file systems. It has to be admitted that this work is not a tutorial on how to use a specific tool but has a broader idea.

## Roadmap

This book is organized into two distinct parts (Fig. 1). Part I describes several different file systems that are commonly used in mobile devices. APFS is the file system that is used in all modern Apple devices. This includes the iPhones, iPads but also the Apple Computers like the MacBook Series. At the same time, Ext4 is very common in Android devices. Ext4 is the successor of the Ext2 and Ext3 file systems that were commonly used on Linux-based computers. The Flash-Friendly File System (F2FS) is a Linux system designed explicitly for NAND Flash memory. This type of memory is common in removable storage devices and mobile devices. Samsung Electronics developed the system in 2012. The QNX6 filesystem is present in Smartphones delivered by Blackberry (e.g. Devices that are using Blackberry 10) and modern vehicle infotainment systems that use QNX as their operating system.

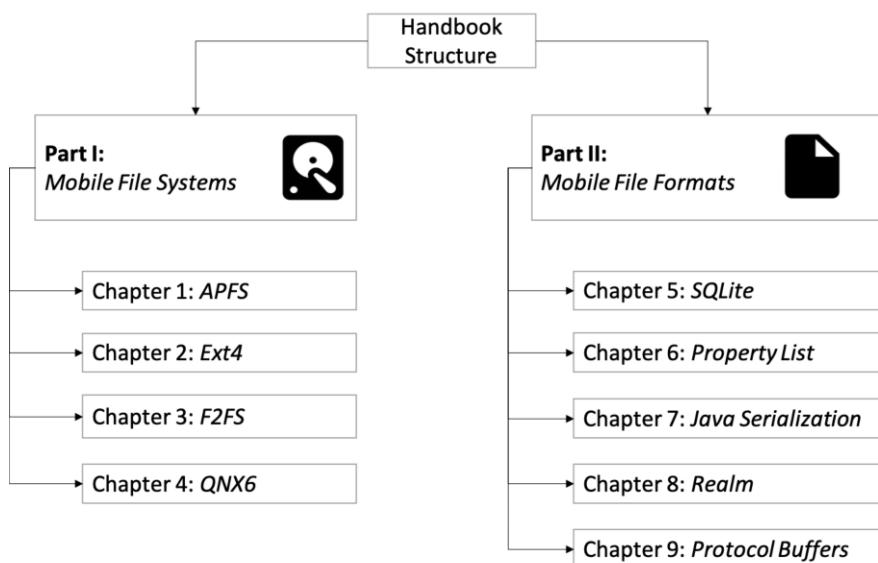


Fig. 1: Structure of the book

Part II describes five different file formats that are commonly used on mobile devices. SQLite is nearly omnipresent in mobile devices. The overwhelming majority of all mobile applications (Apps) store their data in such databases. Another important file format in the mobile world are Property Lists. They are especially frequent on Apple devices. Java Serialization is a popular technique for storing object states in the Java programming language. In the field of mobile forensics, we come across such artefacts. App developers very often resort to this technique to make their application state persistent. The Realm database format has emerged over recent years as a possible successor to the now ageing SQLite format and has a growing use on mobile devices. Protocol Buffers provide a format for taking compiled data and serializing it by turning it into bytes represented in decimal values. This technique is also often used in mobile devices.

## Scope of the Book

After the Roadmap shows the forthcoming chapters' names, it is time to summarize what is not included in this book. The book summarizes four file systems and five commonly used file formats. Next to the fundamental description of the formats, there are hints about the forensic value of possible artefacts and tools that can decode the files or file systems.

This book is not a step-by-step guide to investigating mobile devices, reconstructing file systems, or decoding file formats. In addition, the book does not describe what files a specific OS or application creates. So you want to gather information about which files to examine if a specific app or OS is given.

The book is appropriate for forensic experts who need knowledge about a specific file system or file format or students who want to become forensic experts. This book requires some knowledge about computers, mobile devices, file systems and file formats, so this is not an absolute beginners guide.

## Conventions Used in This Book

We cite numerous books, articles, and websites throughout the book. These citations appear in the text using square brackets [999]. All references can be found at the end of this book. The reader should further pay special attention to annotations that are marked with the following symbols from the text:

### ! Attention

Certain things should not be done to avoid mistakes. At one point or another, the reader will be warned of possible mistakes during a forensic investigation.

## > Important

This hint field is used whenever we think something is significant.

---

## Tips

If you find such a box in the text, we would like to give the reader an important hint. Taking these tips to heart can save much time in practice.

---

## Acknowledgements

The project has received funding from the European Union’s Horizon 2020 Research and Innovation Program under Grant Agreement No. 832800.

We would wish to thank everyone who has participated in any way and made this book possible. In particular, we would like to thank Georgina Humphries, without whose help the project would certainly not have gotten this far for their valuable comments and corrections, as well as Phil Cobley and Chris Currier, who have proofread the text.

Mittweida, Germany  
December 2021

*Christian Hummert  
Dirk Pawlaszczyk*

# Contents

## Part I Mobile File System Formats

<b>1 APFS.....</b>	<b>3</b>
Rune Nordvik	
1.1 Introduction .....	3
1.2 APFS - File system category .....	4
1.2.1 Finding the APFS container .....	4
1.2.2 Object header .....	5
1.2.3 Superblocks .....	8
1.2.4 Checkpoint mapping .....	11
1.2.5 Volumes .....	15
1.3 APFS - Metadata Category .....	26
1.4 APFS - File Name category .....	32
1.5 APFS - Content Category .....	34
1.6 APFS - Application Category .....	38
1.7 Comparing our results with a commercial tool .....	38
<b>2 Ext4.....</b>	<b>41</b>
Rune Nordvik	
2.1 Introduction .....	41
2.2 Ext4 - File system category .....	42
2.3 Superblock.....	43
2.3.1 Temporary data about the File system .....	43
2.3.2 Supported features .....	45
2.3.3 The group descriptor .....	51
2.4 Ext4 - Metadata Category .....	55
2.4.1 The inode .....	55
2.4.2 User privileges and type of file .....	56
2.4.3 Temporary metadata describing inodes.....	57
2.4.4 Temporary metadata manipulations.....	58
2.4.5 Links count .....	59

2.5	Ext4 - File Name category .....	65
2.6	Ext4 - Content Category .....	66
2.6.1	Recovery of files .....	66
2.6.2	Generic metadata time carving .....	67
2.6.3	Additional file content .....	67
2.7	Ext4 - Application Category .....	68
<b>3</b>	<b>The Flash-Friendly File System (F2FS) .....</b>	<b>69</b>
	Chris Currier	
3.1	Introduction .....	69
3.1.1	NAND (Not And) Flash Memory .....	69
3.1.2	Flash Translation Layer (FTL) .....	71
3.2	Flash Filesystems .....	71
3.2.1	The Log-Structured File System (LSFS) or (LFS) .....	72
3.2.2	Flash-Friendly File System (F2FS): Enter F2FS .....	72
3.2.3	Wandering Tree Problem .....	73
3.3	On-Disk Layout of F2FS .....	73
3.3.1	Creation of F2FS partitions with Mkfs.f2fs .....	75
3.3.2	F2FS on Disk .....	76
3.4	File Structure of F2FS .....	81
3.4.1	Node Structure .....	81
3.4.2	File Creation and Management .....	83
3.4.3	Fsck.f2fs Identifying Files .....	85
3.4.4	Metadata .....	86
3.4.5	Multi-Head Logging .....	87
3.4.6	Cleaning .....	88
3.5	Forensic Analysis .....	91
3.5.1	F2FS Sample Dataset .....	91
3.5.2	F2FS and Windows .....	92
3.5.3	Data-Extraction with XRY .....	93
3.5.4	Superblock Examination .....	94
3.5.5	Examine NAT, SIT & SSA with Linux .....	95
3.5.6	Carving for artefacts with XAMN .....	101
3.5.7	Node Allocation Table (NAT) Comparisons .....	105
3.6	F2FS - Application fields .....	108
3.7	Conclusion .....	108
<b>4</b>	<b>QNX6 .....</b>	<b>109</b>
	Conrad Meyer	
4.1	Introduction .....	109
4.2	QNX6 Filesystem Structure .....	110
4.2.1	Superblock .....	111
4.2.2	Bitmap .....	114
4.2.3	Inode .....	115
4.2.4	Directories .....	117
4.2.5	Long Filenames Inode .....	119

4.3	Example: Construction of a file . . . . .	119
4.4	Deleted Files . . . . .	122
4.5	Forensic Tools supporting QNX6 filesystems . . . . .	125

## Part II Mobile File Formats

<b>5</b>	<b>SQLite . . . . .</b>	129
	Dirk Pawlasczyk	
5.1	Introduction . . . . .	129
5.2	The SQLite File Structure . . . . .	130
5.2.1	The Database Header . . . . .	132
5.2.2	Storage Classes, Serial Types and Varint-Encoding . . . . .	135
5.2.3	Decoding The SQLite_Master Table . . . . .	136
5.2.4	Page Structure . . . . .	138
5.2.5	Recovering Data Records . . . . .	141
5.3	Accessing The Freelist . . . . .	144
5.4	More Artefacts . . . . .	146
5.4.1	Temporary File Types . . . . .	146
5.4.2	Rollback Journals . . . . .	148
5.4.3	Write-Ahead Logs . . . . .	151
5.5	Conclusions . . . . .	154
<b>6</b>	<b>Property Lists . . . . .</b>	157
	Christian Hummert and Georgina Louise Humphries	
6.1	Introduction . . . . .	157
6.2	Binary plist Structure . . . . .	158
6.3	Example . . . . .	161
6.4	Forensic Tools Supporting plists . . . . .	163
6.5	Conclusions . . . . .	165
<b>7</b>	<b>Java Serialization . . . . .</b>	167
	Dirk Pawlasczyk	
7.1	Introduction . . . . .	167
7.2	Object Serialization in Java . . . . .	168
7.2.1	Serialization Techniques in Java . . . . .	168
7.2.2	Serialization by Example . . . . .	169
7.3	Java Object Serialization Protocol Revealed . . . . .	172
7.4	Pitfalls and Security Issues . . . . .	177
7.4.1	Hands on Serialized Objects . . . . .	178
7.4.2	Beware of Gadget Chains . . . . .	178
7.5	Conclusions . . . . .	180

<b>8</b>	<b>Realm</b>	181
Phil Cobley and Ginger Geneste		
8.1	Organisation of this Chapter	181
8.2	Introduction	182
8.3	SQLite, It is Not!	183
8.3.1	Relational Databases	183
8.3.2	SQLite as a Relational Database	185
8.3.3	SQLite Schema	186
8.3.4	Temporary SQLite Files	186
8.3.5	SQLite File Format	188
8.4	How Realm Works	189
8.4.1	Realm Database Fundamentals	189
8.4.2	Common Concepts and Terminology	189
8.5	File Storage and Structures	196
8.5.1	Realm Files and Folders	196
8.5.2	The Realm File	196
8.5.3	Creating Realm Test Instance	198
8.5.4	The Realm Database File Structure	204
8.5.5	Realm File Header	205
8.5.6	Realm File Arrays	210
8.5.7	Realm Array Header	211
8.5.8	Checksum	212
8.5.9	Flags	212
8.5.10	Size	215
8.5.11	Realm Array Payload	216
8.5.12	Size Calculation Example	217
8.5.13	Array Example - Header	218
8.5.14	Array Example - Flags	218
8.5.15	Array Example - Size	219
8.6	Conclusion	220
<b>9</b>	<b>Protocol Buffers</b>	223
Chris Currier		
9.1	Introduction	223
9.1.1	What is a Protocol Buffer?	223
9.1.2	Why are Protocol Buffers Used?	225
9.2	Using Protocol Buffers	226
9.2.1	The Schema Definition	231
9.2.2	Compiling Your Protocol Buffer	238
9.2.3	Creation of a Protobufs with Python	242
9.2.4	Reversing Proto Buffer Messages	245
9.3	Practical Analysis of different Proto Buffers	248
9.3.1	Mobile Device Artifact Examples	249
9.3.2	Yet another example: Apply Property List (PLIST) Files	256
9.3.3	Suggested Examination Process of a File	257

<b>Contents</b>	<b>xiii</b>
9.3.4    Tools .....	259
9.4    Conclusion .....	260
<b>References .....</b>	<b>261</b>
<b>Index .....</b>	<b>267</b>

# **Part I**

## **Mobile File System Formats**

File system analysis examines one volume of a disk respective disk image. The data contained in the volume is interpreted as a file system. Typically an interpreted file system offers a listing of files organized in directories (at least this is true for the four described file systems in this book). One aim of file system analysis is usually recovering deleted files. In contrast to files that have been restored by a carving process and lack all formerly available metadata, files that have been undeleted via the file system may contain at least some metadata.

In this part of the book, the general design of four common file systems in mobile devices is described, and different analysis techniques are presented. This part abstractly approaches the topic and is not limited to how a specific tool analyzes a file system.

File systems are used by the operating system and provide mechanisms to store data in a hierarchy of files and directories. File systems organize metadata and user data such that the operating system can use this. Some operating systems rely on one specific file system, whereas others can handle many different file systems. The described file systems in this part are typically used in mobile devices.

# Chapter 1

## APFS



Rune Nordvik

**Abstract** The Apple File System (APFS) has been the standard FS for Apple devices since 2017. At that time, no digital forensic tools supported it, leaving tool dependent digital or mobile forensic investigators without the ability to investigate this file system properly. The APFS was first enabled on iOS, the operating system of iPhone, and later that same year on MacOS. APFS replaced the HFS+ FS. This chapter will document the important metadata structures of APFS, which is based on state of the art research, and we are focusing on the investigative meaning of the structures.

### 1.1 Introduction

Apple developed the APFS , and the main architect was Dominic Giampaolo [33]. Hansen and Toolan [33] started reverse engineering of the APFS for investigation purposes using pre-releases of the APFS as early as 2016, which were included for educational and development purposes in MacOS v 10.12 (Sierra) in September 2016. In March 2017, APFS was deployed on iPhone and iPad [33]. Recently, Apple has released developer documentation for APFS [4], and we use this accurate documentation. In order to develop offset tables that can be used when interpreting hex dumps, we scrutinize the C-structures from the APFS Developer documentation [4] and the type of the fields found in these structures.

We found research on how the APFS uses encryption on the user data partition of an iPhone and that users can not disable the encryption [17]. In this chapter, we have decided to use an image of an iPhone using iOS v 13.3 to see if there could be other partitions that are not encrypted.

---

Rune Nordvik

The Norwegian Police University College (Politihøgskolen), Slemdalsveien 5, 0369 Oslo, Norway,  
e-mail: [rune.nordvik@phs.no](mailto:rune.nordvik@phs.no)

## > Important

Other partitions can contain information that are relevant for the investigation. There could be logs that describe activated features, or information about when a device was rebooted, etc.

## 1.2 APFS - File system category

A new feature for APFS is how it structures volumes, and it does not use a typical partition table to divide the storage into partitions, each with its own FS volume. Instead, it uses storage or a partition to set up a container. This container will contain both container metadata, metadata for snapshots and volumes, and data blocks. This is illustrated in Fig. 1.1.

From an investigative perspective, this means it is not enough to document the partition systems like the Master Boot Record (MBR) and the Globally Unique Identifier (GUID) Partition Table (GPT). Now it is also a need to scrutinise the APFS container.

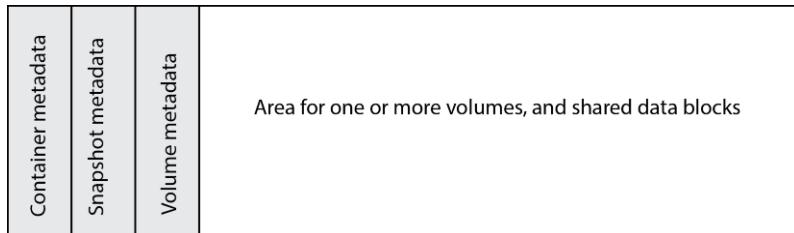


Fig. 1.1: APFS Container.

## ! Attention

Files located on one volume may share the identical data blocks with files located on another volume within the same APFS container.

### 1.2.1 Finding the APFS container

In order to find the APFS container, it is necessary to parse the GUID partition table (GPT) and a partition with the type 7C3457EF-0000-11AA-AA11-00306543ECAC (APFS\_GPT\_PARTITION\_UUID) is an APFS container. How to read the GPT is

described by Nikkel, 2009 [51]. The easiest way is to skip directly to the partition table starting on sector 2 of the disk. Each partition described in this table is 128 bytes (0x80), and it starts with the partition type (its first 16 bytes), followed by a globally unique identifier for this specific partition (also 16 bytes). From relative offset 32 (0x20), we find the start sector of the partition, and it is from that location we find the container.

We used an image from an iPhone 7 running iOS v 13.3, which should have one of the first implementations of the APFS. The first thing we noticed is that the default sector size is 4096 bytes, not the usual 512 bytes. In Fig. 1.2 we show with green background the GUID type of this partition as it is shown in a hex dump. We need to read this data in a special way in order to compare it to the APFS partition type GUID (7C3457EF-0000-11AA-AA11-00306543ECAC). The first four bytes need to be read as little-endian (LE), and therefore ef57347c is read from right to left (backwards) as 7c3457ef, which matches the first part of the APFS partition type. Then we continue with the next two bytes, and they are only zeros, meaning it does not change just because we read it backwards. The next two bytes are aa11, and must be read as 11aa. The next two bytes are not a multi-byte, and the endianness does not matter. They need to be read as single bytes, meaning aa11, then the final six bytes is also not a multi-byte, and must be read as they are, meaning 00306543ECAC. Now we have identified this partition as an APFS container.

00000000:	ef57	347c	0000	aa11	aa11	0030	6543	ecac	W4 .....0eC..
00000010:	a487	6696	dac6	3543	8401	4f98	27f0	0790	..f...5C..0.'...
00000020:	0800	0000	0000	0000	47d6	dc01	0000	0000	.....G.....
00000030:	0000	0000	0000	0000	4300	6f00	6e00	7400	.....C.o.n.t.
00000040:	6100	6900	6e00	6500	7200	0000	0000	0000	a.i.n.e.r.....
00000050:	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000060:	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000070:	0000	0000	0000	0000	0000	0000	0000	0000	.....

Fig. 1.2: APFS Container within GUID partition table.

The field with black background is the eight bytes describing the start sector of this APFS container, and since this is a multi-byte field, it must be read as LE, meaning 0x8 (8 in decimal). That is the sector we need to show in order to find the APFS Container. Highlighted with the yellow background, we can see the Unicode string "Container", which is the name of this partition.

### 1.2.2 Object header

Every object in APFS has an object header, shown in Fig. 1.3, which consists of the fields as described in Table 1.1.

00000000:	15fd	6ff8	2da1	de43	0100	0000	0000	0000	...o.-..C.....
00000010:	0e81	5800	0000	0000	0100	0080	0000	0000	...X.....

Fig. 1.3: Object Header.

Table 1.1: obj\_phys\_t

Offset	Size	Name	Description
0x0	0x8	o_cksum	Fletcher 64 bit checksum
0x8	0x8	o_oid	The object id
0x10	0x8	o_xid	The transaction id
0x18	0x4	o_type	The object type
0x1C	0x4	o_subtype	The object subtype

The object type field low 16 bits indicate a specific object type, while the high 16 bits are used for object type flags:

### Object type, some examples

- OBJECT\_TYPE\_NX\_SUPERBLOCK, 0x00000001
- OBJECT\_TYPE\_BTREE, 0x00000002
- OBJECT\_TYPE\_OMAP, 0x0000000b
- OBJECT\_TYPE\_FS, 0x0000000d
- OBJECT\_TYPE\_FSTREE 0x0000000e
- OBJECT\_TYPE\_INVALID, 0x00000000

### Object type masks

- OBJECT\_TYPE\_MASK, 0x0000ffff
- OBJECT\_TYPE\_FLAGS\_MASK, 0xffff0000

### Object type flags

- OBJ\_VIRTUAL, 0x00000000
- OBJ\_EPHEMERAL, 0x80000000
- OBJ\_PHYSICAL, 0x40000000
- OBJ\_NOHEADER, 0x20000000
- OBJ\_ENCRYPTED, 0x10000000
- OBJ\_NONPERSISTENT, 0x08000000

In Fig. 1.3 we can see the object header, which is 32 bytes in size. The first eight bytes 0x43dea12df86ffd15 (LE) is the *Fletcher checksum* (highlighted in dark blue), the *object id* is 0x1 (highlighted in light green), and the *transaction id* is 0x58810e (highlighted in yellow). The object id and transaction id combined specify a specific state in time.

$$\text{ObjectType} = \text{o\_type} \& \text{ OBJECT\_TYPE\_MASK}$$

$$\text{ObjectTypeFlags} = o\_type \& \text{OBJECT\_TYPE\_FLAGS\_MASK}$$

The object type is 80000001 (highlighted in red), and after computing the object type and flags, we found that the object type value is 0x00000001, and the object type flags is 0x80000000 (OBJ\_EPHEMERAL). We can also see the subtype has the value 0x0000 (highlighted in purple), and here we can use the same approach computing the subtype and flags. However, when the o\_subtype value is 0x0 it means no subtype. Based on our header interpretation, we can see that this superblock is an ephemeral object.

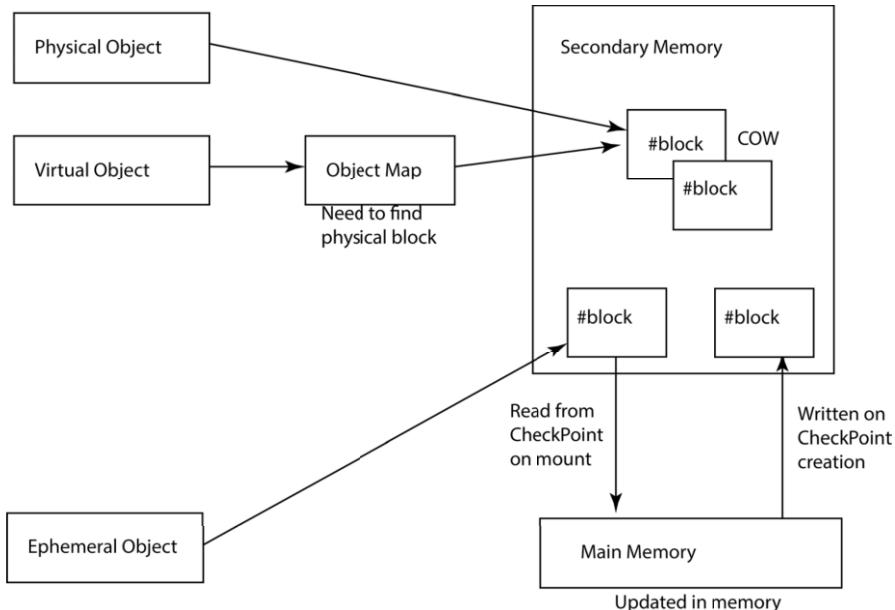


Fig. 1.4: APFS use of the ephemeral, physical, and virtual objects.

## Ephemeral Objects

Objects that should be in memory (ephemeral) and are changed in memory when needed, and will be written to disk as a part of a checkpoint.

## Physical Objects

Objects that are stored in a known block address, and change needs to be written to another location because of the Copy-On-Write (COW) feature. The *object id* is the

same as the block address, and therefore any change means saving to a new block address, this also means a new object id.

## Virtual Objects

Virtual objects stored at a block address that can be found by looking up in a *object map* (often a B-Tree). However, the object id is the same after updating a virtual object. When we look up a virtual object, we use its object id and the transaction id to specify the object at a specific time. This means that when a virtual object changes, it will be written to another physical block (COW), but the virtual object id is still the same in the object header o\_oid field.

### 1.2.3 Superblocks

APFS uses different kinds of superblocks, and the first we find is the Container Superblock (CSB), a nx\_superblock\_t structure. Addresses that point to locations on the disk are in 64bits, meaning eight bytes. These addresses point to ephemeral, physical or virtual objects. The latter objects need an object map in order to find the physical address where the object is located. The difference between these object types is illustrated in Fig. 1.4. In Fig. 1.6 we show a simplified overview of how we plan to go through the File System Category for the APFS.

00000000:	15fd	6ff8	2da1	de43	0100	0000	0000	0000	...	o	..C	.....
00000010:	0e81	5800	0000	0000	0100	0080	0000	0000	...	X	.....	
00000020:	4e58	5342	0010	0000	40d6	dc01	0000	0000	NXSB	...	@	.....
00000030:	0000	0000	0000	0000	0000	0000	0000	0000	...	..	.....	
00000040:	0200	0000	0000	0000	1a4b	a4a7	6d9f	4803	...	K	m.H	.....
00000050:	9fa1	27a6	d9c6	56c1	e044	7b00	0000	0000	...	'	V.D{	.....
00000060:	0f81	5800	0000	0000	1801	0000	186c	0000	...	X	..l.	.....
00000070:	0100	0000	0000	0000	1901	0000	0000	0000	...	..	.....	
00000080:	5d00	0000	db51	0000	0000	0000	0000	0000	l.	....Q	.....	
00000090:	0000	0000	0000	0000	0004	0000	0000	0000	...	..	.....	
000000a0:	2563	d501	0000	0000	0104	0000	0000	0000	%C	..	.....	
000000b0:	0000	0000	6400	0000	0204	0000	0000	0000	...	d	.....	
000000c0:	93d0	5600	0000	0000	6116	0400	0000	0000	...	V	..a	.....
000000d0:	427e	7a00	0000	0000	b57e	7a00	0000	0000	B~z	..	~z	.....

Fig. 1.5: APFS first superblock in block 0 of the partition.

In Fig. 1.5 we find the magic key NXSB (0x4e585342), which identifies this as a superblock. The next field (4 bytes) with blue background is the block size used, here 0x1000 (4096) bytes<sup>1</sup>. With black background, we have an 8-byte field describing

<sup>1</sup> The minimum block size is 4096 (default), and the maximum is 65536

how many blocks this container contains, here 0x1dcd640 (31249984). We can multiply this with the block size if we want the size in bytes. We want the size in GiB<sup>2</sup>, and use this computation.

$$GiB = \frac{31249984 * 4096}{1024^3} = 119.2$$

In light green background we can see the GUID, which uniquely identifies this container. It has the value A7A44B1A-9F6D-0348-9FA1-27A6D9C656C1.

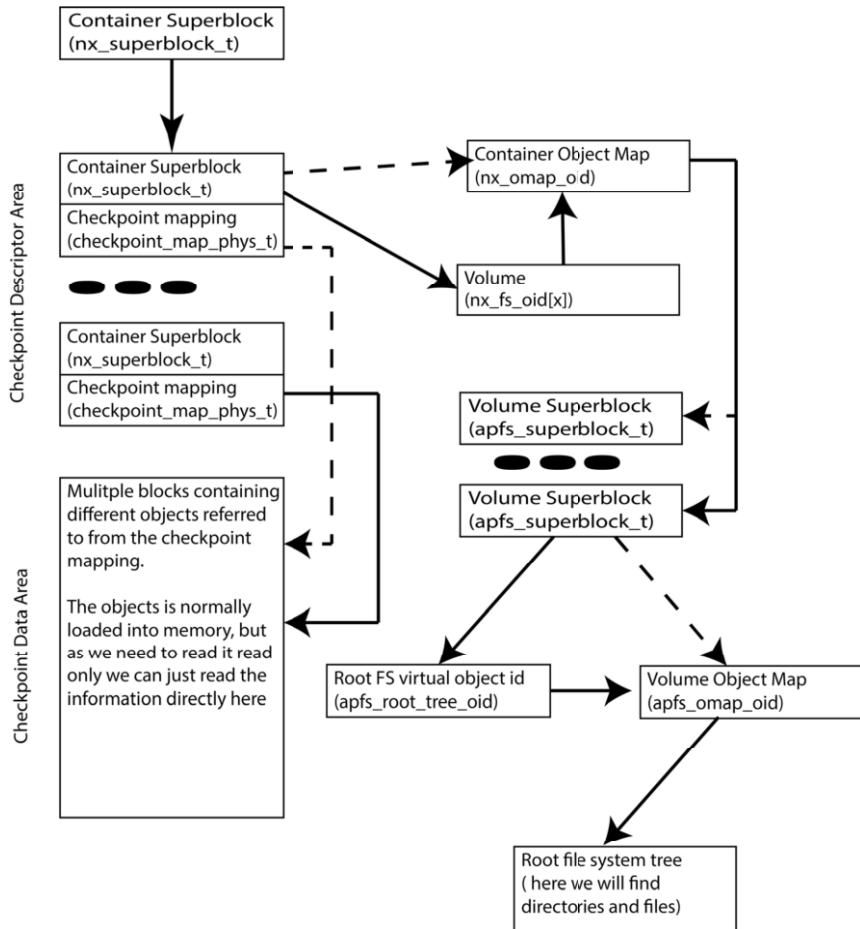


Fig. 1.6: Overview of how to manually read an APFS file system.

<sup>2</sup> The smallest supported container size is 1 MiB

Table 1.2: nx\_superblock

Offset	Size	Name	Description
0x20	0x4	magic	Container Magic
0x24	0x4	nx_block_size	Block Size in bytes
0x28	0x8	nx_block_count	Block count
0x30	0x8	nx_features	Features
0x38	0x8	nx_read_only_compatible_features	Read only compatible features
0x40	0x8	nx_incompatible_features	Incompatible features
0x48	0x10	nx_uuid	Container UUID
0x58	0x8	nx_next_oid	Next Object ID (OID)
0x60	0x8	nx_next_xid	Next Transaction ID (XID)
0x68	0x4	nx_xp_desc_blocks	Blocks used by the Checkpoint Descriptor Area
0x6C	0x4	nx_xp_data_blocks	Blocks used by Checkpoint Data Area
0x70	0x8	nx_xp_desc_base	Base address of Checkpoint Descriptor Area or Physical OID
0x78	0x8	nx_xp_data_base	Base address of Checkpoint Data Area or Physical OID
0x80	0x4	nx_xp_desc_next	Next index for Checkpoint Descriptor Area
0x84	0x4	nx_xp_data_next	Next index for Checkpoint Data Area
0x88	0x4	nx_xp_desc_index	Index for first item in Checkpoint Descriptor Area
0x8C	0x4	nx_xp_desc_len	Number of blocks used in Checkpoint Descriptor Area
0x90	0x4	nx_xp_data_index	Index for first item in Checkpoint Data Area
0x94	0x4	nx_xp_data_len	Number of blocks used in Checkpoint Data Area
0x98	0x8	nx_spaceman_oid	Space Manager Object ID (OID)
0xA0	0x8	nx_omap_oid	Container Object Map Object ID (OID)
0xA8	0x8	nx_reaper_oid	Reaper Object ID (OID)
0xB0	0x4	nx_test_type	Reserved for testing
0xB4	0x4	nx_max_file_systems	Maximum number of volumes in this container
0xB8	0x8	nx_fs_oid[0]	Start of array of OIDs for volumes in this container

We need to find the latest checkpoint superblock, an eight-byte address starting from offset 0x70 (the `nx_xp_desc_base` field), here in yellow background. Here we can see that this points to block 1, relative to the start of this container. At that location, we will either find a new superblock object, or a B-tree map in case this checkpoint superblock is not contiguous. This may or may not be the latest checkpoint.

The documentation from Apple [4] describes that we need to parse through all the blocks in the *Checkpoint Descriptor Area*, and find the block with the highest transaction id (XID) with the same object id (OID). If this block includes the magic key and the *Fletcher checksum* can be verified, then this block is the latest checkpoint.

In our example, the `nx_xp_desc_blocks` field highlighted in light brown at offset 0x68 has the value 0x118 (280). This means the checkpoint descriptor area consists of 280 blocks.

When going through each of these 280 blocks, we noticed that the magic key was only present in every second block, and this is correct since each additional superblock had the value 0x2 in the field `nx_xp_desc_len` at offset 0x8C. This field describes the number of blocks this checkpoint used in the Checkpoint Descriptor Area. We found that the superblock with the highest value in XID was, in our example, located at block 19 relative to the start of the container.

Fig. 1.7 is similar to the superblock we found in the first sector (sector 0) of the container. However, it has the highest transaction number of all the checkpoint

00000000:	0a2d e3f6 ac07 6b45	0100 0000 0000 0000	.....kE.....
00000010:	4b81 5800 0000 0000	0100 0080 0000 0000	K.X.....
00000020:	4e58 5342 0010 0000	40d6 dc01 0000 0000	NXSB....@.....
00000030:	0000 0000 0000 0000	0000 0000 0000 0000	.....
00000040:	0200 0000 0000 0000	1a4b a4a7 6d9f 4803	.....K.m.H.
00000050:	9fa1 27a6 d9c6 56c1	2745 7b00 0000 0000	...'.V.'E{.....
00000060:	4c81 5800 0000 0000	1801 0000 186c 0000	L.X.....l.....
00000070:	0100 0000 0000 0000	1901 0000 0000 0000	.....
00000080:	d700 0000 4553 0000	d500 0000 0200 0000	...ES.....
00000090:	4053 0000 0500 0000	0004 0000 0000 0000	@S.....
000000a0:	9868 d501 0000 0000	0104 0000 0000 0000	..h.....
000000b0:	0000 0000 6400 0000	0204 0000 0000 0000	...d.....
000000c0:	93d0 5600 0000 0000	6116 0400 0000 0000	..V.....a.....
000000d0:	427e 7a00 0000 0000	b57e 7a00 0000 0000	B~z.....~z.....

Fig. 1.7: Latest checkpoint Superblock in the Checkpoint Descriptor Area.

superblocks, and is, therefore, the current (latest) checkpoint. Two other requirements must be fulfilled:

- The NXSB magic must be found.
- The checksum must verify, or else there is something wrong with the checkpoint superblock.

#### ➤ Important

Finding the latest checkpoint is important since this is the last state for this file system. However, the other previous checkpoints may be interesting in order to recover files that are deleted in the latest checkpoint.

#### 1.2.4 Checkpoint mapping

We need to scrutinise the second block in this checkpoint descriptor. This block starts in the block after the superblock, as seen in Fig. 1.8.

There could be multiple checkpoint mapping blocks in a checkpoint, but in our example case it was just one superblock and one checkpoint mapping block for each checkpoint.

The 0x20 (32) bytes highlighted in yellow is the usual object header, but here the field is called `cpm_o`, and `cpm` is an abbreviation for checkpoint mapping. If we interpret the `o_type`, we can see that the object type is 0x0000000C (`OBJECT_TYPE_CHECKPOINT_MAP`), and this object type flag is 0x40000000 (`OBJ_PHYSICAL`). The last mapping block is always marked as the last, using the `cpm_flags` field, highlighted in blue. In this case, it is 0x01, since we only have one such mapping block. In our mapping, we have 0x5 records in an array (highlighted in

00000000:	e70b	0047	c524	a1f8	d800	0000	0000	0000	0000	.....G.	\$.....
00000010:	4c81	5800	0000	0000	0c00	0040	0000	0000	0000	L.X.....	@.....
00000020:	0100	0000	0500	0000	0500	0080	0000	0000	0000	.....	.....
00000030:	0010	0000	0000	0000	0000	0000	0000	0000	0000	.....	.....
00000040:	0004	0000	0000	0000	5e54	0000	0000	0000	0000	.....	^T.....
00000050:	0200	0080	0900	0000	0010	0000	0000	0000	0000	.....	.....
00000060:	0000	0000	0000	0000	0304	0000	0000	0000	0000	.....	.....
00000070:	5f54	0000	0000	0000	0200	0080	0900	0000	0000	_T.....	.....
00000080:	0010	0000	0000	0000	0000	0000	0000	0000	0000	.....	.....
00000090:	0504	0000	0000	0000	6054	0000	0000	0000	0000	.....	T.....
000000a0:	1100	0080	0000	0000	0010	0000	0000	0000	0000	.....	.....
000000b0:	0000	0000	0000	0000	0104	0000	0000	0000	0000	.....	.....
000000c0:	6154	0000	0000	0000	1200	0080	0000	0000	0000	aT.....	.....
000000d0:	0010	0000	0000	0000	0000	0000	0000	0000	0000	.....	.....
000000e0:	f346	0400	0000	0000	6254	0000	0000	0000	0000	.F.....	bT.....

Fig. 1.8: Checkpoint mapping block in the Checkpoint Descriptor Area.

Table 1.3: checkpoint\_map\_phys\_t

Offset	Size	Name	Description
0x0	0x20	cpm_o	Object header
0x20	0x4	cpm_flags	Checkpoint flags
0x24	0x4	cpm_count	Records in this mapping
0x28	var	cpm_map[cpm_count]	Array of Checkpoint mappings

green). The records are 0x28 bytes long, and we have only highlighted the important fields of the first record. Then we highlight the next records with either grey or white background. This shows we have a total of five records.

Table 1.4: checkpoint\_mapping\_t

Offset	Size	Name	Description
0x0	0x4	cpm_type	Low 16 bits for object type, and high 16 bits for object type flags
0x4	0x4	cpm_subtype	The object's subtype
0x8	0x4	cpm_size	The object size in bytes
0xC	0x4	cpm_pad	Not in use, padding
0x10	0x8	cpm_fs_oid	Virtual FS OID that this object is associated with
0x18	0x8	cpm_oid	Ephemeral object id
0x20	0x8	cpm_paddr	The address in the checkpoint data area where this object is stored

The first record cpm\_type highlighted in black colour is after mapping 0x00000005 (OBJECT\_TYPE\_SPACEMAN), and the object type flag is 0x80000000(OBJ\_Ephemeral). The size of the object is described by the cpm\_size (highlighted in orange) is 0x1000 (4096), or the same as the size of a block. The virtual object id of the fs volume that this object is associated with is defined in the field

cpm\_fs\_oid, and has the value 0x0 (highlighted in red). In the fields cpm\_oid, we will find the ephemeral object id, 0x400 (highlighted in purple). Finally, we find the field cpm\_paddr, which contains the value 0x545e (21598), the address to the checkpoint data area where this object is stored (highlighted in dark blue).

Table 1.5: Actual checkpoint mapping

Record	Type	Subtype	Ephemeral OID	Phys Address
1	SPACEMAN	Not used	0x400	0x545e (21598)
2	BTREE	SPACEMAN_FREE_QUEUE	0x403	0x545f (21599)
3	BTREE	SPACEMAN_FREE_QUEUE	0x405	0x5460 (21600)
4	NX_REAPER	Not used	0x401	0x5461 (21601)
5	NX_REAP_LIST	Not used	0x0	0x446f3 (280307)

All the records describe an object with the size 4096, and they all use the FS virtual object 0x0. The other values for these five records are listed in Table 1.5. We can follow the address of the first record to get statistics about the container and its internal pool bitmap. The bitmap describes which blocks are allocated (used) or unallocated (free). Table 1.6 can be used to interpret the space manager.

Table 1.6: spaceman\_phys\_t

Offset	Size	Name	Description
0x0	0x20	sm_o	Object header
0x20	0x4	sm_block_size	Block size
0x24	0x4	sm_blocks_per_chunk	Blocks per chunk
0x28	0x4	sm_chunks_per_cib	Chunks per cib
0x2C	0x4	sm_cibs_per_cab	Cibs per cab
0x30	0x60	spacdev[2]	Special structure
0x90	0x4	sm_flags	Flags
0x94	0x4	sm_ip_bm_tx_multiplier	Bitmap multiplier
0x98	0x8	sm_ip_block_count	Block count
0xA0	0x4	sm_ip_bm_size_in_blocks	Bitmap size in Blocks
0xA4	0x4	sm_ip_bm_block_count	Bitmap block count
0xA8	0x8	sm_ip_bm_base	Address to Bitmap base
0xB0	0x8	sm_ip_base	Address to ip base
0xB8	0x8	sm_fs_reserve_block_count	FS reserved block count
0xC0	0x8	sm_fs_reserve_alloc_count	FS reserved allocation count
0xC8	0x78	spacemanfreequeue[3]	Free queues
0x140	0x2	sm_ip_bm_free_head	bitmap free head
0x142	0x2	sm_ip_bm_free_tail	bitmap free tail
0x144	0x4	sm_ip_bm_xid_offset	Transaction id offset
0x148	0x4	sm_ip_bitmap_offset	bitmap offset
0x14C	0x4	sm_ip_bm_free_next_offset	Next bitmap free offset
0x150	0x4	sm_version	Spacemanager version

00000000:	fa45	8765	b8d5	cff5	0004	0000	0000	0000	0000	E.e.....
00000010:	4c81	5800	0000	0000	0500	0080	0000	0000	0000	L.X.....
00000020:	0010	0000	0080	0000	7e00	0000	fb01	0000	0000	.....~.....
00000030:	40d6	dc01	0000	0000	ba03	0000	0000	0000	0000	@.....
00000040:	0800	0000	0000	0000	4024	4901	0000	0000	0000	.....@I.....
00000050:	080a	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000060:	0000	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000070:	0000	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000080:	480a	0000	0000	0000	0000	0000	0000	0000	0000	H.....
00000090:	0100	0000	1000	0000	460b	0000	0000	0000	0000	.....F.....
000000a0:	0100	0000	1000	0000	316d	0000	0000	0000	0000	.....1m.....
000000b0:	416d	0000	0000	0000	f509	0000	0000	0000	0000	Am.....
000000c0:	8500	0000	0000	0000	1000	0000	0000	0000	0000	.....
000000d0:	0304	0000	0000	0000	4a81	5800	0000	0000	0000	.....J.X.....
000000e0:	0300	0000	0000	0000	0000	0000	0000	0000	0000	.....
000000f0:	2d00	0000	0000	0000	0504	0000	0000	0000	0000	.....
00000100:	4a81	5800	0000	0000	0002	0000	0000	0000	0000	J.X.....
00000110:	0000	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000120:	0000	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000130:	0000	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000140:	0c00	0a00	d809	0000	e009	0000	e809	0000	0000	.....
00000150:	0100	0000	d809	0000	c700	0000	0000	0000	0000	.....

Fig. 1.9: SPACEMAN block, for finding the internal pool bitmap.

We use Table 1.6 to interpret the Fig. 1.9. The easiest way to identify the block where the current bitmap starts is to add the fields sm\_ip\_bm\_base (0x6D31), which is pointed to by the first spaceman device in the field sm\_addr\_offset at 0x50 (0xa08, relative to the start of the block), and the sm\_ip\_bitmap\_offset (0x9E0) and then subtract with sm\_ip\_bm\_size\_in\_blocks (0x1). This gives the block 0x7710. The fields mentioned above depend on the number of spaceman devices, each occupying 0x30 bytes (in our case, there were two spaceman devices).

00000000:	00fe	1ffe	e1ff	0100	0000	e001	001e	0000	0000	.....
00000010:	1ee0	e101	1e00	e001	e001	fe1f	feff	ffff	ffff	.....
00000020:	ffff	ffff	ffff	ff3f	0000	0000	0000	0000	0000	.....?
00000030:	0000	00c0	ffff	ffff	ffff	ffff	1f00	0000	e0ff	.....
00000040:	ffff	.....								
00000050:	ffff	.....								
00000060:	ffff	.....								
00000070:	ffff	.....								
00000080:	ffff	.....								
00000090:	ffff	.....								
000000a0:	ffff	ffff	0700	0000	0000	0000	0000	0000	0000	.....
000000b0:	0000	ffff	.....							
000000c0:	ffff	.....								
000000d0:	ffff	ffff	ffff	ffff	ff07	0000	0000	fcff	fcff	.....

Fig. 1.10: Internal pool bitmap, if bit is 1 then the block is allocated, or if it is 0 then the block is free.

In Fig. 1.10 we see the start of the bitmap. Every byte must be converted to binary, and each bit represents the allocation status of one block (1=allocated, 0=unallocated). The first byte is 0x0, meaning the block 0-7 is defined as not allocated. The following byte is 0xFE (binary: 1111 1110). We start reading from the least significant bit. Moreover, since this is the second byte, the first bit represents block 8 (not allocated) but blocks 9-15 is allocated.

### ! Attention

The SPACEMAN is poorly documented in the APFS developer documentation, which means there could be more accurate methods to identify the bitmap.

## 1.2.5 Volumes

The maximum number of possible volumes are defined in field `nx_max_file_systems` found in the checkpoint superblock, and in this case, the number is 0x64 (100)<sup>3</sup>. However, not all of them are in use. From offset 0xB8 we find an array of fs volume virtual object ids, of which the first has the value 0x402. Only a few of them have a value different from zero. Each of these are virtual object ids that eventually will lead us to a file system tree. In order to identify the location of the volume boot record (VBR) for a file system, we need to map the virtual object id with the one found in the Container Object Map, where we can find the block, which is the address to the VBR. The Container Object Map object id can be found in the checkpoint superblock at byte offset 0xA0 (`nx_omap_oid`) and is a 64-bit value, here 0x1D56898 (30763160). However, this block consists of a B-Tree that needs to be parsed. Before we do the actual mapping, we need to learn the structure of the B-Tree.

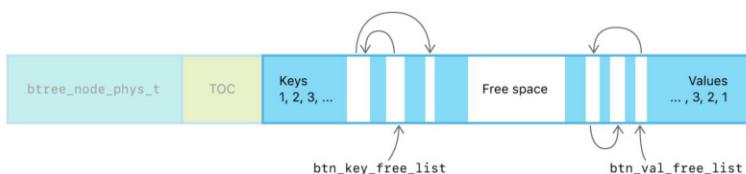


Fig. 1.11: The structure of a B-Tree node block (source: Apple File System Reference)

We can see from Fig. 1.11 that the first part of the block contains the `btree_node_phys_t` structure, which we also have described in Table 1.7, and both include the object header (0x20 in size) and the node header (0x18 in size). The data area is everything after the headers. This data area contains the table of content (TOC), keys, free space and the values.

<sup>3</sup> The maximum number of volumes is defined as 100

00000000:	518d	eb7c	eea0	3022	9868	d501	0000	0000	0.. ..0".h.....
00000010:	c382	5800	0000	0000	0200	0040	0b00	0000	..X.....@.....
00000020:	0700	0000	0500	0000	0000	c001	6000	200d	.....
00000030:	3000	1000	4000	1000	2000	3000	0000	1000	0..@...0.....
00000040:	5000	6000	4000	5000	1000	2000	1000	2000	P..@.P.....
00000050:	0000	0000	0000	0000	0000	0000	0000	0000	.....
...									
000001f0:	0000	0000	0000	0000	6116	0400	0000	0000	.....a.....
00000200:	1281	5800	0000	0000	b57e	7a00	0000	0000	..X.....~z.....
00000210:	0981	5800	0000	0000	0204	0000	0000	0000	..X.....
00000220:	0e81	5800	0000	0000	ffff	1000	0000	0000	..X.....
00000230:	c282	5800	0000	0000	427e	7a00	0000	0000	..X.....B~z.....
00000240:	2481	5800	0000	0000	93d0	5600	0000	0000	\$..X.....V.....
00000250:	c382	5800	0000	0000	0000	0000	0000	0000	..X.....
00000260:	0000	0000	0000	0000	0000	0000	0000	0000	.....
...									
00000f60:	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000f70:	0000	0000	0000	0000	0000	0000	0010	0000	.....
00000f80:	6068	d501	0000	0000	0000	0000	0010	0000	`h.....
00000f90:	7b72	d501	0000	0000	ffff	1000	0010	0000	{r.....
00000fa0:	9866	d501	0000	0000	0000	0000	0010	0000	.f.....
00000fb0:	f262	d501	0000	0000	0000	0000	0010	0000	.b.....
00000fc0:	4165	d501	0000	0000	0000	0000	0010	0000	Ae.....
00000fd0:	f178	d501	0000	0000	1200	0000	0010	0000	.x.....
00000fe0:	1000	0000	1000	0000	1000	0000	1000	0000	.....
00000ff0:	0500	0000	0000	0000	0100	0000	0000	0000	.....

Fig. 1.12: Container FS Volume Object Map.

Table 1.7: btree\_node\_phys\_t

Offset	Size	Name	Description
0x0	0x20	btn_o	Object header of the B-Tree block
0x20	0x2	btn_flags	Flags for this B-Tree
0x22	0x2	btn_level	Number of child levels under this node
0x24	0x4	btn_nkeys	Number of keys stored in this node
0x28	0x4	btn_table_space	Location(16 bits offset, 16 bits length) of Table of Content(TOC)
0x2C	0x4	btn_free_space	Location for shared free space for keys and values
0x30	0x4	btn_key_free_list	A linked list that tracks free key space
0x34	0x4	btn_val_free_list	A linked list that tracks free value space
0x38	0x8	btn_data[var]	The nodes storage area (toc, keys, free space, and values)

From the B-tree object header in Figure 1.12 we can see that the o\_type describe a physical (0x40000000) B-tree (0x2). The o\_subtype is a object map (0xb). Then we

interpret the B-tree node header:

- btn\_flags: 0x7,  
BTNODE\_ROOT, BTNODE\_LEAF, BTNODE\_FIXED\_KV\_SIZE

Table 1.8: B-Tree Node Flags

Define Name	Define Value	Description
BTNODE_ROOT	0x0001	The root node
BTNODE_LEAF	0x0002	A leaf node
BTNODE_FIXED_KV_SIZE	0x0004	Only use the offset for keys-value pairs
BTNODE_HASHED	0x0008	Contains child hashes
BTNODE_NOHEADER	0x00010	Object header consist of zeros
BTNODE_CHECK_KOFF_INVAL	0x8000	Will never appear on disk

- btn\_level: 0x0, There is no level under this one.
- btn\_nkeys: 0x05, there are 5 records.
- btn\_table\_space: 0x00 offset, 0x1c0, meaning TOC starts after the node header at 0x38, and is 0x1c0 in length. This also means the key are starts at 0x1f8, directly after the TOC.
- The shared free space starts at 0x60 in the key area, meaning  $0x1f8 + 0x60 = 0x258$ , and it is 0xd20 in length, meaning it end at  $0x258+0xd20=0xf78$ , where it meet the last part of the value area (the top of value area, where the free space ended).

From the TOC, we can find the key-value pairs. If BTNODE\_FIXED\_KV\_SIZE flag is set, only offsets to keys and values are used. If not, both offset and length are used. All offsets for keys are relative to the start of the key area, and all offsets for values are relative from the end of the value area (the bottom of the value area). At the end of the Root node block, we have the btree\_info\_t structure, which can be interpreted using Table 1.9. Only Root nodes should have this additional structure. For instance, the value 0x1, 0x3, 0x7 are also Root nodes, which have a btree\_info\_t structure at the end.

Table 1.9: btree\_info\_t

Offset	Size	Name	Description
0xfd8	0x4	bt_flags	The B-tree's flags
0xfd0	0x4	bt_node_size	The B-tree's node size
0xfe0	0x4	bt_key_size	The B-tree's key size
0xfe4	0x4	bt_val_size	The B-tree's value size
0xfe8	0x4	bt_longest_key	The B-tree's longest key ever stored
0xfc0	0x4	bt_longest_val	The B-tree's longest value ever stored
0xff0	0x8	bt_key_count	Number of keys in the B-tree
0xff8	0x8	bt_node_count	Number of nodes in the B-tree

It is important to interpret the key-value offsets in the TOC in a special way. Since the offsets, in this case, have a static length (the BTNODE\_FIXED\_KV\_SIZE flag was set), the key and value are represented as 16-bit offsets. If the key value is not

fixed, the key and the value both use 32 bits (the first 16 bits is the offset, and the following 16 bits is the length). The key-value offsets occur in the TOC directly after the headers, but this can be deviated by the btn\_table\_space field.

When we read the offset for the key, we need to consider that the offset is relative to the start of the key area. However, when we read the offset for the value, it is relative from the end of the value area and backwards in the direction of the free space area. This is also why the key area in the illustration in Fig. 1.11 show Keys 1,2,3,..., while the Values are listed as ..., 3,2,1.

The last structure of the B-tree Root node is the btree\_info\_t, which we can use the Table 1.9 to interpret. The bt\_flags are 0x12, which consist of 0x10 (BTREE\_PHYSICAL) and 0x2 (BTREE\_SEQUENTIAL\_INSERT) which means avoiding splitting nodes in half during sequential inserts, avoiding a lot of half-full nodes [4, p. 130]. The bt\_node\_size (node size) is 0x1000 (4096) bytes. The bt\_key\_size (key size), bt\_value\_size (value size), bt\_longest\_key (longest key size), and the bt\_longest\_val (longest value size) are all 0x10 (16) bytes. Both the key and value sizes are necessary to know when using fixed sizes for keys and values. The bt\_key\_count (keys in this B-tree) is 0x5, and bt\_node\_count (nodes in this B-tree) is 0x1. The value 1 for node count means that this node is both the root and the leaf node.

## Finding the Volume

In Fig. 1.13 we have the object header in light yellow, followed by the node header in grey background, interpreted in the previous section.

The first record key-pair is highlighted in blue from offset 0x38, the first 16 bits key value 0x20 points to the key, while the second 16 bits value 0x30 points to the value. We know the key area starts at offset 0x1f8, and the value area starts before the btree\_info\_t structure, and if we count 0x30 backwards we find the start of the value that is in this case 0x10 bytes long, also highlighted in blue. The key can be found counting 0x20 from the start of the key area, and it is also in this case 0x10 in size. The key is also highlighted in blue. It is important to read the keys and values as 0x10 (16 bytes) each, as described in the btree\_info\_t structure.

In this case the key OID (the first 8 bytes) is 0x402 (LE), and the key XID (the next 8 bytes) is 0x58810e (LE), and its corresponding value physical OID address (the last 8 bytes in the value) is 0x1d562f2 (LE). The same can be done for all the other volumes in this node, highlighted using different highlight colours. We have finished this mapping in Table 1.10.

## Showing the Volume (APSB)

A volume in APFS uses the magic key APSB, which is a superblock that describes one volume. This is similar to a VBR in traditional file systems. It contains an FS,

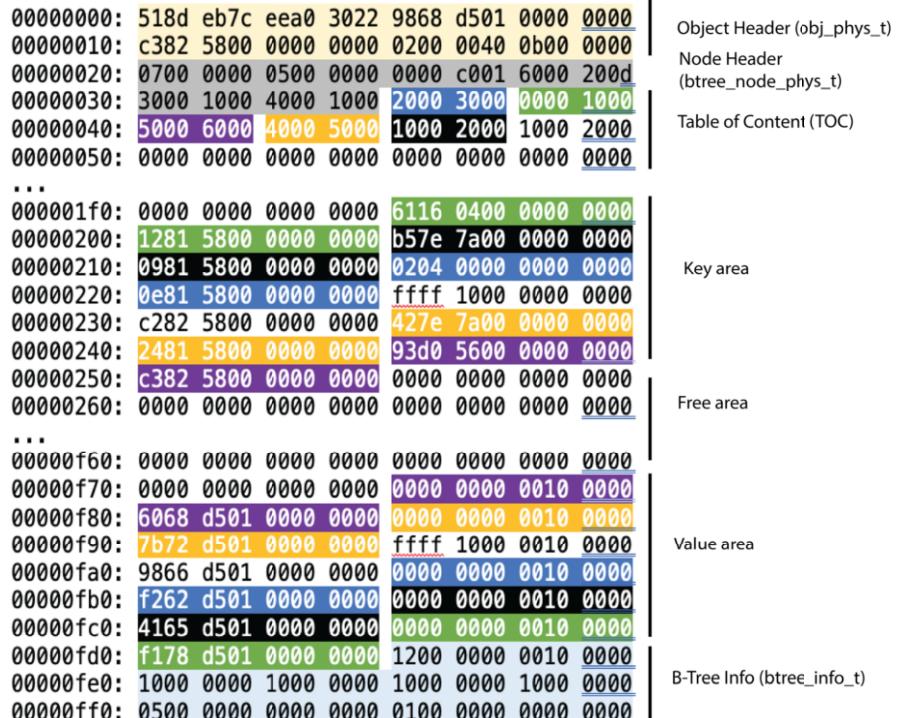


Fig. 1.13: Mapping of FS Volumes, from virtual to physical block.

Table 1.10: FS Volume mapping

Virtual OID	XID	Physical Address
0x402	0x58810e	0x1d562f2
0x41661	0x588112	0x1d578f1
0x56d093	0x5882c3	0x1d56860
0x7a7e42	0x588124	0x1d5727b
0x7a7eb5	0x588109	0x1d56541

files, metadata, and object map. In this section, we will be using the file system example for the virtual object id 0x402, physical address 0x1d562f2.

We can use Table 1.11 to interpret the values of the APFS volume superblock (APSB) shown in Fig. 1.14. We have not included all fields in this structure, only fields that we assume are important. Please refer to the complete `apfs_superblock_t` structure as found in the Apple developer documentation [4].

The magic of this apfs volume superblock has the signature APSB is ASCII values, found in the `apfs_magic` field. This value can be used when searching for APFS volumes in a corrupted file system. In the checkpoint superblock in Fig. 1.7

00000000:	c1e4	a178	de20	3929	0204	0000	0000	0000	.....x.	9).....
00000010:	0e81	5800	0000	0000	0d00	0000	0000	0000	..X.....	.....
00000020:	4150	5342	0000	0000	0200	0000	0000	0000	APSB	.....
00000030:	0000	0000	0000	0000	0000	0000	0000	0000	.....	.....
00000040:	f4fe	ca02	cf1	7b16	0000	0000	0000	0000	.....{	.....
00000050:	0000	0000	0000	0000	81e8	1300	0000	0000	.....	.....
00000060:	0500	0000	0000	0000	0600	0000	5900	460e	.....Y.F.	.....
00000070:	0100	0000	0200	0000	0200	0040	0200	0040	.....@. @	.....
00000080:	8f65	d501	0000	0000	8b87	7a00	0000	0000	.e.....z....	.....
00000090:	0765	d501	0000	0000	7065	d501	0000	0000	.e.....pe....	.....
000000a0:	0000	0000	0000	0000	0000	0000	0000	0000	.....	.....
000000b0:	0ed8	0b00	0100	0000	e645	0200	0000	0000	.....E....	.....
000000c0:	f90e	0100	0000	0000	ac06	0000	0000	0000	.....	.....
000000d0:	0000	0000	0000	0000	0100	0000	0000	0000	.....	.....
000000e0:	6b5d	a200	0000	0000	9f97	8f00	0000	0000	k].....	.....
000000f0:	939b	63ec	6dc4	3946	8910	5e68	39bf	0530	...c.m.9F.^h9..0	.....
00000100:	2dda	75c9	94c9	e315	0100	0000	0000	0000	-u.....	.....
00000110:	6866	735f	636f	6e76	6572	7420	2861	7066	hfs_convert (apf	.....
00000120:	732d	3234	392e	3630	2e32	3029	0000	0000	s-249.60.20)	....
...										
000002c0:	436f	7272	7931	3443	3932	2e44	3130	4431	Corry14C92.D10D1	.....
000002d0:	3031	4f53	0000	0000	0000	0000	0000	0000	010S.....	.....

Fig. 1.14: The APFS Superblock for the volume.

we had an array of file system object identifiers, and we can decide which of these file systems this volume belongs to by reading the apfs\_index value, here it is 0, meaning the first file system.

One very important field is the apfs\_unmount\_time, which must be interpreted as nanoseconds since January 1, 1970 at 0:00 UTC, not including leap seconds. In this case is the unmount time, 0x167bd1cf02cafef4, or Tuesday, 4 May 2021 09:06:18. We converted to decimal and divided with  $10^9$  to get the value as seconds. Then we used a UNIX epoch converter to translate it to a human-readable time format.

### ! Attention

The apfs\_unmount\_time is updated whenever unmounted, and this will normally mean at the reboot of the device, especially for the system volume. This also means if an investigator reboots the device, then this value is modified.

---

There is also another timestamp in the APFS volume superblock, in the field 0x100 (apfs\_last\_mod\_time) we find the 64 bit value 0x0x15E3C994B2AF9600, or Thursday, 26 December 2019 02:05:35. This is when the volume was last modified, and as it seems, modifying the metadata does not update this timestamp. One hypothesis is that only modifying file content or metadata related to files will update this timestamp.

Table 1.11: apfs\_superblock\_t

Offset	Size	Name	Description
0x0	0x20	apfs_o	The object header
0x20	0x4	apfs_magic	The magic signature for an APSB
0x24	0x4	apfs_index	The FS index in the container list of file systems
0x40	0x8	apfs_unmount_time	Last time this FS was unmounted (last reboot?)
0x58	0x8	apfs_alloc_count	Blocks allocated to this FS
0x60	0x14	apfs_meta_crypto	Information about encryption
0x74	0x4	apfs_root_tree_type	Root tree type
0x78	0x4	apfs_extentsref_tree_type	Type of the extent-reference tree
0x7C	0x4	apfs_snap_meta_tree_type	Type of the snapshot metadata tree
0x80	0x8	apfs_omap_oid	Object id of the object map
0x88	0x8	apfs_root_tree_oid	Virtual Object id of the root file system tree
0x90	0x8	apfs_extentsref_tree_oid	Object id of the extent-reference tree
0xB8	0x8	apfs_num_files	Number of regular files in the volume
0xC0	0x8	apfs_num_directories	Number of directories in the volume
0xD8	0x8	apfs_num_snapshots	Number of snapshots in the volume
0xE0	0x8	apfs_tot_blocks_alloced	Total number of blocks that have been allocated by this volume
0xF0	0x10	apfs_vol_uuid	Universal Unique identifier for the volume
0x100	0x8	apfs_last_mod_time	The time when this volume was last modified
0x108	0x8	apfs_fs_flags	Volume flags
0x110	0x8	apfs_formatted_by	Software created this volume
0x2c0	0x100	apfs_volumename	Name of the volume
0x3c0	0x4	apfs_next_doc_id	Todo
0x3c4	0x2	apfs_role	Role
0x3c6	0x2	reserved	Reserved

### ! Attention

The apfs\_last\_mod\_time is updated when the volume is modified. If what is relevant for the investigation is after this time, it may not be worth analysing this volume.

The apfs\_alloc\_count field yields the currently allocated blocks for this file system, 0x13e881 blocks, or 4.98 GiB of the storage is used. Later, another field describes the total number of blocks ever allocated, which increase for every new block allocated but do not decrease when a block is freed. This field is named apfs\_tot\_blocks\_alloced. Do not use this field when computing the currently used size of a volume.

Then there is some information about encryption in the field apfs\_meta\_crypto, which is for encryption purposes. We need to check in the apfs\_fs\_flags to see if the fs utilise encryption. The field apfs\_root\_tree\_type describes the type of tree. Here it is a B-tree. This is followed by the field apfs\_extentsref\_tree\_type, which is described as a physical B-tree. The next field is the apfs\_snap\_meta\_tree\_type, and it is also described as a physical B-Tree.

We need an object map because we need to map virtual object ids to physical object ids. We can find the object map object identifier in the apfs\_omap\_oid, and we depend on this field and the next field (apfs\_root\_tree\_oid) to find objects for

files and directories. The virtual object id for the root tree can be found in the field apfs\_root\_tree\_oid, and since this is a virtual object id, we need to map it to the physical object id using the apfs\_omap\_oid, see sect. 1.2.5.

The apfs\_num\_files describes the number of files in the volume, here 0x245e6 (148966) files. We can find the number of directories in the field apfs\_num\_directories, here 0x10ef9 (69369) directories. This volume contains 1 snapshot, described in the field apfs\_num\_snapshots. The apfs\_vol\_uuid is a unique identifier for this volume, in this case it is EC639B93-C4D6-4639-8910-5E6839BF0530. At offset 0x108, we find the volume Flags (apfs\_fs\_flags), and here it is 0x01, meaning this volume is not encrypted.

### > Important

Not all volumes on an iOS are encrypted.

The first entry starting from offset 0x110 is the name of the software that formatted this APFS volume, and it can be seen that it has been converted from HFS (hfs\_convert(apfs-249.60.20)). Finally, we have included the name of this volume, with grey background starting at offset 0x2c0. Here it is Corry14C92.D10D101OS.

## Volume Object mapping

The object mapping block (apfs\_omap\_oid) is shown in Fig. 1.15, and can be interpreted using the Table 1.12.

00000000:	46a0	4771	8aaf	87c8	8f65	d501	0000	0000	F.Gq.....e.....
00000010:	0d81	5800	0000	0000	0b00	0040	0000	0000	..X.....@.....
00000020:	0000	0000	0100	0000	0200	0040	0200	0040	.....@. @.....
00000030:	a365	d501	0000	0000	ea62	d501	0000	0000	.e.....b.....
00000040:	f500	5800	0000	0000	0000	0000	0000	0000	..X.....
00000050:	0000	0000	0000	0000	0000	0000	0000	0000	.....

Fig. 1.15: Physical mapping of FS B-trees.

As usual, we can find the object header in the first 0x20 (32) bytes, which describes this as a physical object map. There are no flags used (om\_flags), and there is one snapshot in this object map (om\_snap\_count). The root object tree is a physical B-tree (om\_tree\_type), and the same is the snapshot tree. We will focus on the current object map B-tree (om\_tree\_oid), where we find its physical object id, here 0x1d565a3. In this case, we want to map the virtual object id for the FS root B-tree, which is 8b87 7a00 0000 0000 (0x7a878b), to its physical object id.

At block 0x1d565a3 we find the top of the object map B-tree. However, in order to find the virtual object id we are searching for, we need to parse the B-tree, and in

Table 1.12: omap\_phys\_t

Offset	Size	Name	Description
0x0	0x20	om_o	The object header
0x20	0x4	om_flags	Flags used by omap
0x24	0x4	om_snap_count	Number of snapshots within omap
0x28	0x4	om_tree_type	Type of B-tree
0x2C	0x4	om_snapshot_tree_type	Type of snapshot B-tree
0x30	0x8	om_tree_oid	Object id of current B-tree
0x38	0x8	om_snapshot_tree_oid	Object id of snapshot tree
0x40	0x8	om_most_recent_snap	Transaction id of most recent snapshot
0x48	0x8	om_pending_revert_min	Smallest transaction id for a in-progress revert
0x50	0x8	om_pending_revert_max	Largest transaction id for a in-progress revert

our case, by using the 7th record number, we found the block that will contain the virtual object we are searching for, and now it is time to ask: Why did we select this record?

In Fig. 1.16 we show the hex dump of the object map root B-tree. The 7th key, and its value is highlighted in dark blue. The 8th key and value is highlighted in red. The 7th key is 0x6b9a81 (OID part), and the 8th key is 0x7ab7d6 (OID part). The key we are searching for (0x7a878b) is between these two keys. Therefore the virtual object id (key) we are searching for can be found by focusing on the former key (7th). This 7th value is found in the value area 0x18 bytes before the end of the value area, and it has the physical object id address 0x1d7309f. This is the physical OID address of the child node that should contain the virtual object id (0x7a878b) in one of the node's key or a key of a sub-node that we are searching for. Before we continue to this object, we need to check the object node header of the current object. Its node header is interpreted below.

- btn\_flags: 0x5, BTNODE\_ROOT, BTNODE\_FIXED\_KV\_SIZE
- btn\_level: 0x2, There are two levels of child nodes under this one.
- btn\_nkeys: 0x0a, there are 0xa (10) records.
- btn\_table\_space: 0x00 offset, 0x240, meaning TOC starts after the node header at 0x38, and is 0x240 in length. This also means the key are starts at 0x278, directly after the TOC.
- The shared free space starts at 0x120 in the key area, meaning  $0x278 + 0x120 = 0x398$ , and it is 0xbb0 in length, meaning it end at  $0x398+0xbb0=0xf48$ , where it meet the last part of the value area.

This means we are located on the top of the tree, the Root node, and we can expect two levels of child nodes under this root level. We continue to the physical location 0x1d7309f, as previously mapped, here we find the next level in this B-tree, shown in Fig. 1.17. This object node header is interpreted below.

- btn\_flags: 0x4, BTNODE\_FIXED\_KV\_SIZE, must be an index node.
- btn\_level: 0x1, There is one level of child nodes under this one.

```

00000000: 4242 7b31 71cc 7f5c a365 d501 0000 0000 BB{1q..`e.....
00000010: 0d81 5800 0000 0000 0200 0040 0b00 0000 ..X.....@....
00000020: 0500 0200 0a00 0000 0000 4002 2001 b00b .....@.....
00000030: 3000 8000 2000 4000 0000 0800 7000 4000 0... .@... p@...
00000040: 5000 3000 6000 3800 b000 6000 8000 4800 P.0.`8...`H...
00000050: 2000 1800 f000 8000 a000 5800 d000 7000 .....X..p...
00000060: d000 7000 d000 7000 d000 7000 d000 7000 ..p...p...p...
00000070: d000 7000 d000 7000 d000 7000 d000 7000 ..p...p...p...
00000080: 0000 0000 0000 0000 0000 0000 0000 0000 .....p...
...
000000270: 0000 0000 0000 0000 5304 0000 0000 0000 .....S....
000000280: 6fd2 1900 0000 0000 4000 1000 0000 0000 o.....@....
000000290: 409b 5700 0000 0000 819a 6b00 0000 0000 @.W.....k...
0000002a0: 830d 3c00 0000 0000 1001 1000 0000 0000 ..<.....
0000002b0: e900 5800 0000 0000 9000 1000 0000 0000 ..X.....
0000002c0: 174a 2100 0000 0000 8851 4700 0000 0000 .J!.....QG...
0000002d0: a95e 0800 0000 0000 07f6 5600 0000 0000 .^.....V...
0000002e0: 6fd2 1900 0000 0000 cf02 0200 0000 0000 o.....
0000002f0: 8b61 3500 0000 0000 54b0 6100 0000 0000 .a5.....T.a...
000000300: 38fe 2b00 0000 0000 ffff 1000 0000 0000 8.+.....
000000310: f300 5800 0000 0000 b6bf 7a00 0000 0000 ..X.....z...
000000320: e900 5800 0000 0000 6f3c 5b00 0000 0000 ..X.....o<[...
000000330: 2d4a 2100 0000 0000 1000 1000 0000 0000 .J!.....
000000340: f300 5800 0000 0000 d6fd 7a00 0000 0000 ..X.....z...
000000350: e900 5800 0000 0000 0001 1000 0000 0000 ..X.....
000000360: 4c9b 5700 0000 0000 d6b7 7a00 0000 0000 L.W.....z...
000000370: e900 5800 0000 0000 c000 1000 0000 0000 ..X.....
000000380: 459b 5700 0000 0000 e000 1000 0000 0000 E.W.....
000000390: 4c9b 5700 0000 0000 0000 0000 0000 0000 L.W.....
0000003a0: 0000 0000 0000 0000 0000 0000 0000 0000 .....p...
...
000000f40: 0000 0000 0000 0000 7800 0800 0000 0000 .....x...
000000f50: 6800 0800 0000 0000 d627 d701 0000 0000 h.....`...
000000f60: 8800 0800 0000 0000 f618 d701 0000 0000 .....`...
000000f70: 1000 0800 0000 0000 ea2f d701 0000 0000 .....`/...
000000f80: d527 d701 0000 0000 ffff 0800 0000 0000 .`.....
000000f90: 4830 d701 0000 0000 792e d701 0000 0000 H0.....y...
000000fa0: 7c2f d701 0000 0000 2b2f d701 0000 0000 |/.....+...
000000fb0: 5000 0800 0000 0000 9000 0800 0000 0000 P.....
000000fc0: 9f30 d701 0000 0000 2800 0800 0000 0000 .0.....(...
000000fd0: 267b d501 0000 0000 1200 0000 0010 0000 &{.....
000000fe0: 1000 0000 1000 0000 1000 0000 1000 0000 .....p...
000000ff0: 0e1e 0100 0000 0000 8b03 0000 0000 0000 .....p...

```

Fig. 1.16: The root of the object map B-Tree.

- btn\_nkeys: 0x8e, there are 0x8e (142) records.
- btn\_table\_space: 0x00 offset, 0x240, meaning TOC starts after the node header at 0x38, and is 0x240 in length. This also means the key area starts at 0x278, directly after the TOC.
- The shared free space starts at 0x900 in the key area, meaning 0x278 + 0x900 = 0xB78, and it is 0x8 in length, meaning it end at 0xB78+0x8=0xB80, where it meet the last part of the value area. The value area is not shown completely in Fig. 1.17.

```

00000000: 26d9 048b 62f3 a170 9f30 d701 0000 0000 &...b..p.0.....
00000010: f500 5800 0000 0000 0300 0040 0b00 0000 ..X.....@....
00000020: 0400 0100 8e00 0000 0000 4002 0009 0800 .....@.....
00000030: 1001 2000 9000 1000 1004 1002 d006 7003 ... ....p....
00000040: e002 7801 d002 7001 f001 0001 c002 6801 ..x..p....h...
00000050: f002 8001 0003 8801 b000 6000 b002 6001 .....`.....
00000060: 6002 3801 5002 3001 4001 a800 4002 2801 `8.P.0.@.(...
00000070: 7001 c000 6001 b800 3002 2001 a002 5801 p...`0...X...
00000080: 2001 9800 c001 e800 0001 8800 e000 7800 .....x.....
00000090: d000 7000 c000 6800 a000 5800 9000 5000 ..p...h...X...P...
000000a0: 8000 4800 7000 4000 5000 3000 3000 2000 ..H.p.@.P.0.0...
000000b0: 0000 0800 1000 1000 a006 5803 7008 4004 .....X.p.@...
000000c0: c005 e802 0002 0801 1006 1003 b005 e002 .....`.....
...
00000270: d004 7002 d004 7002 5133 7200 0000 0000 ..p...p.Q3r....
00000280: 3b4a 4900 0000 0000 8b87 7a00 0000 0000 ;JI.....z....
00000290: e900 5800 0000 0000 76a5 7a00 0000 0000 ..X....v.z.....
000002a0: e900 5800 0000 0000 1133 7200 0000 0000 ..X....3r.....
00000fb0: b030 d701 0000 0000 6bb6 d501 0000 0000 .0.....k.....
00000fc0: 8ab6 d501 0000 0000 235c d601 0000 0000 .....#`.....
00000fd0: 95b6 d501 0000 0000 ccaf d501 0000 0000 .....`.....
00000fe0: b9b6 d501 0000 0000 b11f d601 0000 0000 .....`.....
00000ff0: b230 d701 0000 0000 b130 d701 0000 0000 .0.....0.....

```

Fig. 1.17: An index node in the object map B-Tree containing the virtual object id searched for.

The object map index node can be seen in Fig. 1.17, where the 32nd key show offset 0x10 for key and 0x10 for value. In the key area, which starts at offset 0x278, we find this key 0x10 bytes further down. Here we find the virtual object id we are looking for (0x7a878b) and its transaction id 0x5800e9. However, since this is the only key with this OID value in the node, we know we have the latest one. The physical object id to which it is connected can be found 0x10 bytes from the end of this node, and where the value is 0x1d730b2. This is the physical address to the leaf node in this omap B-tree. Reading this object id (0x1d730b2), we can see in Fig. 1.18 the object map leaf node. As usual, we need to interpret the object node header.

- `btn_flags`: 0x6, `BTNODE_LEAF`, `BTNODE_FIXED_KV_SIZE`.
- `btn_level`: 0x0, There are no levels of child nodes under this one.
- `btn_nkeys`: 0x47, there are 0x47 (71) records.
- `btn_table_space`: 0x00 offset, 0x1c0, meaning TOC starts after the node header at 0x38, and is 0x1c0 in length. This also means the key are starts at 0x1F8, directly after the TOC.
- The shared free space starts at 0x700 in the key area, meaning  $0x1F8 + 0x700 = 0x8F8$ , and it is 0x8 in length, meaning it end at  $0x8F8 + 0x8 = 0x900$ , where it meet the last part of the value area. The value area is not shown completely in Fig. 1.18.

The first record points to the root virtual object id 0x7a878b, which can be found at physical object id 0x1d66e7f (last 8 byte of value).

```

00000000: 8dad a1cc 04e2 545b b230 d701 0000 0000 .....T[.0.....
00000010: f500 5800 0000 0000 0300 0040 0b00 0000 .....X.....@....
00000020: 0600 0000 4700 0000 0000 c001 0007 0800 .....G.....
00000030: 8001 9002 9001 9002 7005 8005 b002 c002 .....p.....
00000040: d004 e004 7001 8001 5006 6006 8004 9004 .....p..P`.....
00000050: a002 b002 d005 e005 3000 4000 8003 9003 .....0.@.....
...
00000760: b300 5800 0000 0000 8b87 7a00 0000 0000 .....X.....z.....
00000770: e900 5800 0000 0000 af87 7a00 0000 0000 .....X.....z.....
00000780: e900 5800 0000 0000 1003 1000 0000 0000 .....X.....z.....
...
00000a80: 0000 0000 0010 0000 7f6e d601 0000 0000 .....n.....
00000a90: f002 1000 0010 0000 1669 d501 0000 0000 .....i.....
...
00000fe0: 0000 0000 0010 0000 8d63 d501 0000 0000 .....c.....
00000ff0: 0000 0000 0010 0000 7d63 d501 0000 0000 .....}c.....

```

Fig. 1.18: Omap Leaf Node where we found the virtual object we searched for, which can be mapped to physical address 0x1d66e7f.

The physical address 0x1d66e7f, found in Fig. 1.18 is the physical address to the apfs\_root\_tree\_oid (the file system B-Tree) with virtual object id 0x7a878b. We discuss this B-Tree more in sect. 1.3.

### 1.3 APFS - Metadata Category

We have already found the FS Root B-tree, and now we will start explaining structures that are related to files and directories.

In Fig. 1.19 we show the content of the physical object id block 0x1d66e7f, and this block has the virtual object id 0x7a878b, found in byte 8 in the object header. This is typical for a virtual object id when we read it from its physical address. The virtual object id will still be stored from byte 8 in the object header. When we read the object node header, we find the following information.

- btn\_flags: 0x1, BTNODE\_ROOT,
- btn\_level: 0x3, There are three levels of child nodes under this one.
- btn\_nkeys: 0x4, there are 4 records.
- btn\_table\_space: 0x00 offset, 0x40, meaning TOC starts after the node header at 0x38, and is 0x40 in length. This also means the key area starts at 0x78, directly after the TOC.
- The shared free space starts at 0x5b in the key area, meaning  $0x78 + 0x5b = 0xD3$ , and it is 0xee5 in length, meaning it end at  $0xd3+0xee5=0xfb8$ , where it meet the last part of the value area.

We have listed all the four entries in this block using different background colors.

- OBJ\_ID\_MASK (0xffffffffffff)

00000000:	a2dc	b144	c152	e64a	8b87	7a00	0000	0000	.....D.R.J..z.....
00000010:	e900	5800	0000	0000	0200	0000	0e00	0000	...X.....
00000020:	0100	0300	0400	0000	0000	4000	5b00	e50e	.....@.[.....
00000030:	ffff	0000	ffff	0000	0000	1600	0800	0800	.....
00000040:	1600	2100	1000	0800	3700	1c00	1800	0800	...!.....7.....
00000050:	5300	0800	2000	0800	0000	0000	0000	0000	S.....
00000060:	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000070:	0000	0000	0000	0000	0100	0000	0000	0090	.....
00000080:	0c00	7072	6976	6174	652d	6469	7200	d4f1	...private-dir...
00000090:	0900	0100	0040	1700	636f	6d2e	6170	706c	.....@..com.appl
000000a0:	652e	5265	736f	7572	6365	466f	726b	00c2	e.ResourceFork..
000000b0:	f20a	0001	0000	4012	0063	6f6d	2e61	7070	.....@..com.app
000000c0:	6c65	2e64	6563	6d70	6673	0021	8f0b	0001	le.decmpfs.!....
000000d0:	0000	3000	0000	0000	0000	0000	0000	0000	...0.....
...									
00000fa0:	0000	0000	0000	0000	0000	0000	0000	0000	.....
00000fb0:	0000	0000	0000	0000	1efc	7a00	0000	0000	.....z.....
00000fc0:	b9d8	7a00	0000	0000	e5b2	7a00	0000	0000	.....z.....z.....
00000fd0:	e6b2	7a00	0000	0000	4200	0000	0010	0000	...z.....B.....
00000fe0:	0000	0000	0000	0000	7600	0000	de0e	0000	.....v.....
00000ff0:	0d1f	0a00	0000	0000	d183	0000	0000	0000	.....

Fig. 1.19: File System Root B-Tree.

- OBJ\_TYPE\_MASK (0xf000000000000000)
- OBJ\_TYPE\_SHIFT (60)

$$\text{Object}Id = \text{obj\_id\_and\_type} \& \text{OBJ\_ID\_MASK}$$

$$\text{ObjectType} = \text{obj\_id\_and\_type} \& \text{OBJ\_TYPE\_MASK} >> \text{OBJ\_TYPE\_SHIFT}$$

The first 8 bytes of the first record key is 0x9000000000000001, and when computing the object id we get 0x1. When computing the object type we get 0x9, which is a APFS\_TYPE\_DIR\_REC (found in Table 1.13). This means this record is a directory record.

### ! Attention

One FS object may have several records describing the object, and therefore there could be multiple records with the same object id.

It seems like the value field of these four records only contains virtual object ids. For the first record, this is file id 0x7ab2e6. This means we need to look up in the apfs\_omap\_id (0x1d565a3) again to map this virtual address to the physical address. The two highlighted keys in Fig. 1.16, show that the virtual object id we are searching for is between them, and therefore we select the first record, and we continue down the B-Tree until we find the correct physical address, which was 0x1d6558f.

Table 1.13: j\_obj\_types

Enum Name	Enum Value
APFS_TYPE_ANY	0
APFS_TYPE_SNAP_METADATA	1
APFS_TYPE_EXTENT	2
APFS_TYPE_INODE	3
APFS_TYPE_XATTR	4
APFS_TYPE_SIBLING_LINK	5
APFS_TYPE_DSTREAM_ID	6
APFS_TYPE_CRYPTO_STATE	7
APFS_TYPE_FILE_EXTENT	8
APFS_TYPE_DIR_REC	9
APFS_TYPE_DIR_STATS	10
APFS_TYPE_SNAP_NAME	11
APFS_TYPE_SIBLING_MAP	12
APFS_TYPE_FILE_INFO	13
APFS_TYPE_MAX_VALID	13
APFS_TYPE_MAX	15
APFS_TYPE_INVALID	15

Table 1.14: j\_drec\_key\_t

Offset	Size	Name	Description
0x0	0x8	hdr (objid and type)	The header of this record (type: j_key_t)
0x8	0x2	name_len	The length of the directory
0xA	name[name_len]	The name of this directory	

Table 1.15: j\_drec\_val\_t

Offset	Size	Name	Description
0x0	0x8	file_id	The node identifier
0x8	0x8	date_added	Timestamp describing when directory was moved/created here
0x10	0x2	flags	Flag describing inode file type (masked with DREC_TYPE_MASK)
0x12	var	xfields[]	Extended fields

Figure 1.20 shows that the directory with file name *private-dir* object id 0x1, and from its value field we need to go to virtual address 0x7a87e8, so again we need to look up in the volume object map to find the physical address, which was 0x1d563bc. Fig. 1.21 shows the same record for the filename *private-dir*, still object id 1, and we are now in the second index node, and we have 1 level of child nodes under this one, and we want to see the child node for this record, which can be found at virtual object id 0x7a878d. We look up in the volume object map, and we find that the physical address is 0x1d5632b.

```

00000000: a5e8 68e5 2eef 14d3 e6b2 7a00 0000 0000 ..h.....z.....
00000010: e900 5800 0000 0000 0300 0000 0e00 0000 ..X.....
00000020: 0000 0200 6600 0000 0000 4003 9c08 bc00 ....f....@....
00000030: ffff 0000 ffff 0000 0000 1600 0800 0800 .....
00000040: 1600 1800 1000 0800 2e00 2000 1800 0800 .....

...
000000370: 0000 0000 0000 0000 0100 0000 0000 0090 .....
000000380: 0c00 7072 6976 6174 652d 6469 7200 6402 ..private-dir.d.
000000390: 0000 0000 0090 0e00 436f 6465 5265 736f .....CodeReso
0000003a0: 7572 6365 7300 3609 0000 0000 0090 1600 urces.6.....
0000003b0: 4163 6365 7373 6962 696c 6974 792e 7374 Accessibility.st
0000003c0: 7269 6e67 7300 be17 0000 0000 0040 1200 rings.....@..
0000003d0: 636f 6d2e 6170 706c 652e 6465 636d 7066 com.apple.decmpf
0000003e0: 7300 7c25 0000 0000 0090 1600 4163 6365 s.|%.....Acce
0000003f0: 7373 6962 696c 6974 792e 7374 7269 6e67 ssibility.string
000000400: 7300 d731 0000 0000 0090 1300 4761 6d65 s..1.....Game
000000410: 4365 6e74 6572 2e73 7472 696e 6773 00c7 Center.strings..
000000420: 3e00 0000 0000 4012 0063 6f6d 2e61 7070 >.....@..com.app

...
000000f70: 908e 7a00 0000 0000 2c8e 7a00 0000 0000 ..z.....,z.....
000000f80: b98d 7a00 0000 0000 548d 7a00 0000 0000 ..z.....T,z.....
000000f90: ef8c 7a00 0000 0000 7b8c 7a00 0000 0000 ..z.....{.z.....
000000fa0: 078c 7a00 0000 0000 a88b 7a00 0000 0000 ..z.....z.....
000000fb0: 4b8b 7a00 0000 0000 f48a 7a00 0000 0000 K.z.....z.....
000000fc0: 898a 7a00 0000 0000 2a8a 7a00 0000 0000 ..z.....*.z.....
000000fd0: b689 7a00 0000 0000 4789 7a00 0000 0000 ..z.....G,z.....
000000fe0: d188 7a00 0000 0000 5c88 7a00 0000 0000 ..z.....\,z.....
000000ff0: e787 7a00 0000 0000 e887 7a00 0000 0000 ..z.....z.....

```

Fig. 1.20: File System B-Tree Index (level 1).

```

00000000: 687b bcf9 d5fe 63b5 e887 7a00 0000 0000 h{....c..z.....
00000010: e900 5800 0000 0000 0300 0000 0e00 0000 ..X.....
00000020: 0000 0100 5a00 0000 0000 0003 7c09 7c00 ....Z.....|..|..
00000030: ffff 0000 ffff 0000 0000 1600 0800 0800 .....
00000040: 1600 1d00 1000 0800 3300 2300 1800 0800 .....3.#.....

...
000000330: 0000 0000 0000 0000 0100 0000 0000 0090 .....
000000340: 0c00 7072 6976 6174 652d 6469 7200 1600 ..private-dir...
000000350: 0000 0000 0090 1300 4d65 6469 6153 7472 .....MediaStr
000000360: 6561 6d50 6c75 6769 6e73 0017 0000 0000 eamPlugins.....
000000370: 0000 9019 0041 6363 6573 736f 7279 4175 .....AccessoryAu
000000380: 6469 6f2e 6672 616d 6577 6f72 6b00 1700 dio.framework...
000000390: 0000 0000 2000 4170 706c 6542 6173 ..... AppleBas

...
000000fc0: 9387 7a00 0000 0000 9287 7a00 0000 0000 ..z.....z.....
000000fd0: 9187 7a00 0000 0000 9087 7a00 0000 0000 ..z.....z.....
000000fe0: 8f87 7a00 0000 0000 8e87 7a00 0000 0000 ..z.....z.....
000000ff0: 8c87 7a00 0000 0000 8d87 7a00 0000 0000 ..z.....z.....

```

Fig. 1.21: File System B-Tree Index (level 2).

Fig. 1.22 shows files from the root directory, where also the *private-dir* is located, highlighted in dark blue. When we interpret a directory key, we use the Table 1.14. In the 8 bytes before the file name, we find the object id 0x1 and object type (0x9,

000000000:	27a0	4940	acb8	9de1	8d87	7a00	0000	0000	'I@.....z....
00000010:	e900	5800	0000	0000	0300	0000	0e00	0000	..X.....
00000020:	0200	0000	4d00	0000	0000	8002	4006	2e00	...M.....@..
00000030:	ffff	0000	ffff	0000	0000	1600	1200	1200	.....
00000040:	1600	0f00	2400	1200	2500	0800	9000	6c00	....\$.%.l.
00000050:	2d00	1300	a200	1200	4000	0e00	b400	1200	-----@.....
...									
000002b0:	0000	0000	0000	0000	0100	0000	0000	0090	.....
000002c0:	0c00	7072	6976	6174	652d	6469	7200	0100	...private-dir...
000002d0:	0000	0000	0090	0500	726f	6f74	0002	0000	.....root...
000002e0:	0000	0000	302	0000	0000	0000	9009	002e	....0.....
000002f0:	5472	6173	6865	7300	0200	0000	0000	0090	Trashes.....
00000300:	0400	2e62	6100	0200	0000	0000	0090	0600	....ba.....
00000310:	2e66	696c	6500	0200	0000	0000	0090	0400	.file.....
00000320:	2e6d	6200	0200	0000	0000	0090	0d00	4170	.mb.....Ap
00000330:	706c	6963	6174	696f	6e73	0002	0000	0000	plications.....
00000340:	0000	900a	0044	6576	656c	6f70	6572	0002	....Developer..
00000350:	0000	0000	0000	9008	004c	6962	7261	7279	.....Library
00000360:	0002	0000	0000	0000	9007	0053	7973	7465	.....Syste
00000370:	6d00	0200	0000	0000	0090	0400	6269	6e00	m.....bin.
00000380:	0200	0000	0000	0090	0600	636f	7265	7300	.....cores.
00000390:	0200	0000	0000	0090	0400	6465	7600	0200	....dev...
000003a0:	0000	0000	0090	0400	6574	6300	0200	0000	....etc...
000003b0:	0000	0090	0800	7072	6976	6174	6500	0200	...private...
000003c0:	0000	0000	0090	0500	7362	696e	0002	0000	....sbin...
000003d0:	0000	0000	9004	0074	6d70	0002	0000	0000	....tmp...
000003e0:	0000	9004	0075	7372	0002	0000	0000	0000	....usr...
000003f0:	9004	0076	6172	0003	0000	0000	0000	3015	...var.....0.
00000400:	0000	0000	0000	3015	0000	0000	0000	9008	.....0.
00000410:	004c	6962	7261	7279	0016	0000	0000	0000	Library.....
...									
00000f60:	0400	0100	0000	12f5	ec46	ec09	4415	0800	.....F.D..
00000f70:	0100	0000	0000	0000	0200	0000	0000	0000	.....
00000f80:	00e4	9c33	c697	8d14	59ca	c360	8ced	5715	...3....Y.`..W.
00000f90:	59ca	c360	8ced	5715	00ac	12ac	7d98	8d14	Y.`..W....}
00000fa0:	0080	0000	0000	0000	1100	0000	0000	0000	.....
00000fb0:	2400	0000	0000	0000	0000	0000	5000	0000	\$......P..
00000fc0:	fd43	0000	0000	0000	0000	0000	0100	0800	.C.....
00000fd0:	0402	0500	726f	6f74	0000	0000	0200	0000	....root...
00000fe0:	0000	0000	28ca	b448	9673	c914	0400	0300	....(.H.s.....
00000ff0:	0000	0000	0000	b0dd	b448	9673	c914	0400	....H.s.....

Fig. 1.22: File System B-Tree Leaf, shows the files in the root directory.

APFS\_TYPE\_DIR\_REC). The next two bytes 0xc (12) describe the directory name's size. The next 12 bytes are the directory name *private-dir* + null terminator byte.

Then we interpret the value of this directory record, found from offset 0xfe0 and 0x12 (18) bytes, highlighted in dark blue. We use Table 1.15 to interpret the value. The node id (file\_id) is 0x3, and the directory was added at 0x14c9739648b4ddb0 (Monday, 19 June 2017 06:57:20 UTC). The flags field yield the inode file type, here 0x0004. We use flags & DREC\_TYPE\_MASK (0x000f), and we get file data type 4 (DT\_DIR). This means that this directory entry is describing a directory. This type can be found in Table 1.16.

Table 1.16: File Type Flags

Define Name	Define Value	Description
DT_UNKNOWN	0	An unknown directory entry
DT_FIFO	1	A named pipe
DT_CHR	2	A character-special file
DT_DIR	4	A directory
DT_BLK	6	A block-special file
DT_REG	8	A regular file
DT_LINK	10	A symbolic link
DT_SOCK	12	A socket
DT_WHT	14	A whiteout

We continue with the third record, from Fig. 1.22, highlighted in black. The key content is only 8 bytes, 0x3000000000000002, object id 2, and object type 3 (APFS\_TYPE\_INODE). Then we interpret the record value. The parent id is 0x01, private id is 0x2 (unique for this data stream, previously found to be describing the root filename (red highlight)). Then at 0xf80 we have four 8 byte fields that all describe timestamps; create\_time: 0x148D97C6339CE400 (Tuesday, 6 December 2016 06:45:30), mod\_time: 0x1557ED8C60C3CA59 (Wednesday, 26 September 2018 10:49:44), change\_time: 0x1557ED8C60C3CA59 (Wednesday, 26 September 2018 10:49:44), access\_time: 0x148D987DAC12AC00 (Tuesday, 6 December 2016 06:58:38). From offset 0xf0 we find the 8-byte internal flags 0x8000 INODE\_NO\_RSRC\_FORK, which means this inode does not have a resource fork. We find the number of directory entries in this directory in offset 0xfa8, and the value is 0x11 (17). This means we have 17 files or directories in the root directory. The owner of this file is owner id 0, and group id 0x50 (80). From offset 0xfd4 we find the name of this inode, which is root.

Table 1.17: j\_inode\_val\_t

Offset	Size	Name	Description
0x0	0x8	parent_id	The parent node id
0x8	0x8	private_id	This node id
0x10	0x8	create_time	Creation time
0x18	0x8	mod_time	Modification time
0x20	0x8	change_time	Change time
0x28	0x8	access_time	Access time
0x30	0x8	internal_flags	Internal flags
0x38	0x4	nchildren	Directory entries in this directory, or number sym links for a file
0x3C	0x4	default_protection_class	Default protection class <sup>4</sup>
0x40	0x4	write_generation_counter	A counter which increase <sup>5</sup> when node is modified
0x44	0x4	bsd_flags	Inode's BSD flags
0x48	0x4	owner	The user id
0x4c	0x4	group	Group id

## 1.4 APFS - File Name category

We refer to the APFS metadata category, since the file names are part of the parsing of the FS B-tree. The part of the key that contain the file names are related to this category. Make sure to notice that in the Fig. 1.22 we can see typically directory

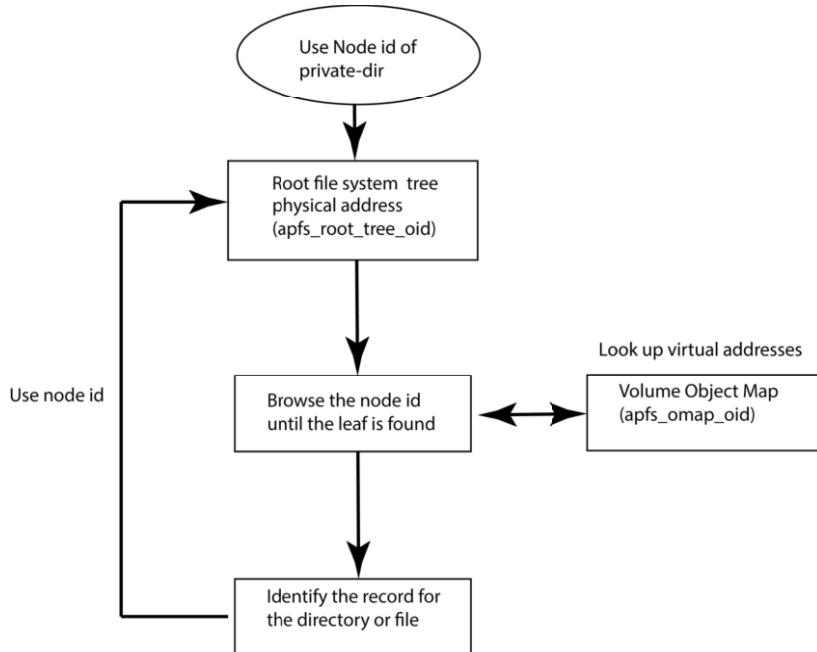


Fig. 1.23: How to browse for a file or directory in the File System Root B-Tree.

names found in the root directory; Trashes Applications, Developer, Library, System, bin, cores, dev, etc, private, sbin, tmp, usr, var.

One of the files we found was the *sbin* directory. However, we do not see the directory entries (files or directories) of this directory. We need to recognise the node identifier of the *sbin* directory (here 0x18c60), and then look it up using the File System Root B-Tree, and again we would need to use the volume object map. An overview of how to browse the File System Root B-Tree (physical address of apfs\_root\_tree\_oid) is shown in Fig. 1.23. In order to find the files in the root directory we already parsed the private-dir using the File System Root B-tree, and the resulting leaf node was shown in Fig. 1.22, which also includes the *sbin* directory we now want to focus on. In Fig. 1.24 we show the same leaf node, but have highlighted the *sbin* record. The key was found at offset 0x3be and is 0xf (15) bytes in size. The object id is 0x2 (meaning it belongs to the parent id 2(root)), and the type is 0x9 (APFS\_TYPE\_DIR\_REC). The size of the name is 0x5, and the name is *sbin*

```

00000000: 27a0 4940 acb8 9de1 8d87 7a00 0000 0000 .I@.....z.....
00000010: e900 5800 0000 0000 0300 0000 0e00 0000 ..X.....
00000020: 0200 0000 4d00 0000 0000 8002 4006 2e00 ...M.....@....
00000030: ffff 0000 ffff 0000 0000 1600 1200 1200 .....
00000040: 1600 0f00 2400 1200 2500 0800 9000 6c00 ....$...%....l.
00000050: 2d00 1300 a200 1200 4000 0e00 b400 1200 -.....@....
00000060: 4e00 1000 c600 1200 5e00 0e00 d800 1200 N.....^.
00000070: 6c00 1700 ea00 1200 8300 1400 fc00 1200 l.....l.
00000080: 9700 1200 0e01 1200 a900 1100 2001 1200 .....
00000090: ba00 0e00 3201 1200 c800 1000 4401 1200 ....2.....D.
000000a0: d800 0e00 5601 1200 e600 0e00 6801 1200 ....V.....h.
000000b0: f400 1200 7a01 1200 0601 0f00 8c01 1200 ....z.....
000000c0: 1501 0e00 9e01 1200 2301 0e00 b001 1200 .....#.....
...
000002b0: 0000 0000 0000 0000 0100 0000 0000 0090 .....
000002c0: 0c00 7072 6976 6174 652d 6469 7200 0100 ..private-dir...
000002d0: 0000 0000 0090 0500 726f 6f74 0002 0000 .....root...
000002e0: 0000 0000 3002 0000 0000 0000 0009 002e ....0.....
000002f0: 5472 6173 6865 7300 0200 0000 0000 0090 Trashes.....
00000300: 0400 2e62 6100 0200 0000 0000 0090 0600 ...ba.....
00000310: 2e66 696c 6500 0200 0000 0000 0090 0400 .file.....
00000320: 2e6d 6200 0200 0000 0000 0090 0d00 4170 .mb.....Ap
00000330: 706c 6963 6174 696f 6e73 0002 0000 0000 plications.....
00000340: 0000 900a 0044 6576 656c 6f70 6572 0002 ....Developer..
00000350: 0000 0000 0000 9008 004c 6962 7261 7279 .....Library
00000360: 0002 0000 0000 0000 9007 0053 7973 7465 .....Syste
00000370: 6d00 0200 0000 0000 0090 0400 6269 6e00 m.....bin.
00000380: 0200 0000 0000 0090 0600 636f 7265 7300 .....cores.
00000390: 0200 0000 0000 0090 0400 6465 7600 0200 .....dev...
000003a0: 0000 0000 0090 0400 6574 6300 0200 0000 .....etc...
000003b0: 0000 0090 0800 7072 6976 6174 6500 0200 .....private...
000003c0: 0000 0000 0090 0500 7362 696e 0002 0000 .....sbin...
000003d0: 0000 0000 9004 0074 6d70 0002 0000 0000 .....tmp.....
...
00000e70: 4415 0a00 608c 0100 0000 0000 00b6 3cc0 D...`.....<
00000e80: 4998 8d14 0400 a34d 0000 0000 0000 00ba I.....M.....

```

Fig. 1.24: File System Root B-Tree, with the sbin record.

+ null terminator. The key value was found at offset 0xe74 and is 0x12 (18) bytes in size. The node id is 0x18c60, the data added is 0x148d9849c03cb600 (Tuesday, 6 December 2016 06:54:55), and the flag is 0x4 (DT\_DIR). In order to identify the node id 0x18c60 for the *sbin* directory, we need to parse the Root File System B-Tree. After parsing the File System B-Tree and using the node id 0x18c60 that corresponds to the *sbin* directory, we found its physical address in 0x1d5ef8f, as shown in Fig. 1.25. Since we already have explained how to parse the FS B-Tree and the volume object map, we do not repeat this here.

```

00000000: 4b83 7bce f992 8f1c 8399 7a00 0000 0000 K.{.....z.....
00000010: e900 5800 0000 0000 0300 0000 0e00 0000 ..X.....
00000020: 0200 0000 2d00 0000 0000 8001 9603 1500 ....-.....
00000030: ffff 0000 ffff 0000 0000 2800 1200 1200 .....(.....
...
000000e0: d101 0800 9d08 6c00 d901 0f00 af08 1200 .....l.....
000000f0: e801 1400 c108 1200 fc01 1500 d308 1200 .....$.....
00000100: 1102 1300 e508 1200 2402 1500 f708 1200 .....$.....
...
00000380: 6d73 646f 732e 6673 0050 8c01 0000 0000 msdos.fs.`.....
00000390: 3060 8c01 0000 0000 9005 0066 7363 6b00 0`.....fsck.
000003a0: 608c 0100 0000 0090 0a00 6673 636b 5f61 `.....fsck_a
000003b0: 7066 7300 608c 0100 0000 0090 0b00 6673 pfs`.....fs
000003c0: 636b 5f65 7866 6174 0060 8c01 0000 0000 ck_exfat`.....
000003d0: 9009 0066 7363 6b5f 6866 7300 608c 0100 ...fsck_hfs`...
000003e0: 0000 0090 0b00 6673 636b 5f6d 7364 6f73 ...fsck_msdos
000003f0: 0060 8c01 0000 0000 9008 006c 6175 6e63 ...launc
00000400: 6864 0060 8c01 0000 0000 9006 006d 6f75 hd`.....mou
00000410: 6e74 0060 8c01 0000 0000 900b 006d 6f75 nt`.....mou
00000420: 6e74 5f61 7066 7300 608c 0100 0000 0090 nt_apfs`.....
00000430: 0a00 6d6f 756e 745f 6866 7300 608c 0100 ...mount_hfs`...
00000440: 0000 0090 0a00 6d6f 756e 745f 6e66 7300 .....mount_nfs.
...
00000710: 0000 0000 0000 0000 000a 001e 1702 0000 .....|.....
00000720: 0000 0000 0000 0000 000a 00c6 1602 .....|.....
00000730: 0000 0000 0000 0000 0000 000a 00c6 .....|.....
00000740: 4902 0001 0000 00f8 51d3 b6ca 73c9 140a I.....Q...s...
00000750: 0088 c309 0001 0000 00f0 beb9 73b5 b3e2 .....s.....
00000760: 1508 0002 0000 0000 0000 0060 8c01 0000 .....|.....
00000770: 0000 0000 1efa 4bde 4d81 14ba 915e ab0a .....K.....
00000780: b4e2 15ba 915e ab0a b4e2 1500 76ad e77d .....^.....v.}
00000790: 988d 1400 8000 0000 0000 000e 0000 0000 .....|.....
000007a0: 0000 0092 0000 0000 0000 0000 0000 0000 .....|.....
000007b0: 0000 00ed 4100 0000 0000 0000 0000 0001 .....A.....
000007c0: 0008 0004 0205 0073 6269 6e00 0000 0015 .....sbin.....

```

Fig. 1.25: File System Root B-Tree, showing some of the content of the sbin directory.

## 1.5 APFS - Content Category

Directory entries found in the FS B-Tree have many different types, and we have already scrutinised directories and inodes. However, a file needs somewhere to store its content. APFS uses extents for this. The data stream type `j_phys_ext_key_t` and `j_phys_ext_val_t` is normally used for this purpose. The private id (node id) from a file record found in the FS B-tree (`apfs_root_tree_oid`) is used as an identifier in the field `owning_obj_id` found in the structure `j_phys_ext_val_t`. If the file is fragmented, we need to browse the Extents B-Tree (`apfs_extentref_tree_oid`) in order to identify all the extents for the node id we are searching for.

There are also structures like `j_file_extent_key_t` and `j_file_extent_val_t` which describe an extent for a file, including the length (measured in bytes) of the extent and its physical block start address. However, sometimes a file can be compressed (especially system files that are part of the iOS system partition). These files utilise

compression using extended attributes describing the resource fork and the compression algorithm used. If compression is used, then the files will have an empty data fork [90]. Files created in the user data partition normally do not use compression, but it is possible. In our example iOS image, the volume was converted from HFS+, and contained compressed files. Some tools do not support reading these compressed files.

We will focus on the entry highlighted in red from Fig. 1.25. From the key, we can see that the object id is 0x18c60, which means its parent directory is the *sbin* directory. We can also see that it is a Directory Entry (0x9 - APFS\_TYPE\_DIR\_REC), which means it is a directory entry. The name of this directory entry is *fsck*. Moving to the corresponding value we can see that the node id (file\_id) is 0x10009c388, then we have the date added value 0x15e2b3b573b9bef0 (Sunday, 22 December 2019 13:13:31 UTC+0). Then we have the last two bytes in the value describing the file type, which is 0x8 ( DT\_REG ). DT\_REG is a regular file. We used Table 1.16 to interpret the file type flags. Now we know this directory entry is a regular file.

### > Important

Even if the key describes that it is a directory entry, this does not mean it is a directory. It is the last two bytes in the corresponding value field that yields the file type (0xA=symbolic link, 0x4=Directory, 0x8=Regular file). File types can be interpreted using Table 1.16.

---

We did not find more entries for this file in this node. Therefore, we should assume that we will find additional directory entries for the *fsck* by browsing through the FS B-tree for the specific object/node id. We used the FS B-tree (apfs\_root\_tree\_oid) and the Volume Object map to find the start block of the *fsck* file with node id (file\_id) 0x10009c388. We use this as an object id when browsing through the FS B-tree. This node id is less than the start of the second entry in the FS B-tree root. This means we will find it selecting the first entry.

Table 1.18: Addresses B-tree parsing

What	Virtual Address	Physical Address	Description
Volume OMAP		0x1d565a3	Using APFS superblock
FS B-Tree Root (L0)	0x7a878b	0x1d66e7f	Using APFS superblock and Volume OMAP
Index Node (L1)	0x7ab2e6	0x1d6558f	Using Volume OMAP
Index Node (L2)	0x7aaf3	0x01d64365	Using Volume OMAP
Leaf Node (L3)	0x7ab04c	0x1d645fb	Using Volume OMAP

Table 1.18 shows virtual addresses and the physical addresses that were found by browsing the Volume OMAP each time we had a virtual address.

000000000:	2406	5f14	e90a	9ad9	4cb0	7a00	0000	0000	\$.....L.z.....
000000010:	e900	5800	0000	0000	0300	0000	0e00	0000	_X.....
000000020:	0200	0000	3100	0000	0000	c001	6803	5400	....1.....h.T.
000000030:	ffff	0000	ffff	0000	0000	0800	7400	7400	.....t.t.
...									
0000001a0:	1303	0800	ec09	6c00	1b03	2100	200a	3400	.....l...!.4.
0000001b0:	3c03	1c00	340a	1400	5803	1000	4c0a	1800	<...4...X...L...
...									
000000500:	0000	8000	0000	0000	0000	88	c309	0001	.....
000000510:	0000	3088	c309	0001	0000	4017	0063	6f6d	..0.....@.com
000000520:	2e61	7070	6c65	2e52	6573	6f75	7263	6546	.apple.ResourceF
000000530:	6f72	6b00	88c3	0900	0100	0040	1200	636f	ork.....@.co
000000540:	6d2e	6170	706c	652e	6465	636d	7066	7300	m.apple.decmpfs.
000000550:	89c3	0900	0100	0080	0000	0000	0000	0000	.....
000000560:	0000	0000	0000	0000	0000	0000	0000	0000	.....
000000570:	0000	0000	0000	0000	0000	0000	0000	0000	.....
000000580:	0000	0000	0000	0000	0000	0000	0000	0000	.....
000000590:	0000	0000	0000	0000	0000	0000	0000	0000	.....
0000005a0:	0000	0000	0000	0000	0000	0000	0000	0000	.....
0000005b0:	0000	0000	0030	0000	0000	0000	d7c7	8701	....0
0000005c0:	0000	0000	0000	0000	0000	0200	1000		.....
0000005d0:	6670	6d63	0400	0000	10d2	0000	0000	0000	fpmc.....
0000005e0:	0100	3000	89c3	0900	0100	0000	d924	0000	..0.....\$..
0000005f0:	0000	0000	0030	0000	0000	0000	0000	0000	.....0.....
000000600:	0000	0000	d924	0000	0000	0000	0000	0000	.....\$.....
000000610:	0000	0000	608c	0100	0000	0000	88c3	0900	.....
000000620:	0100	0000	0056	c832	c1d4	8114	0056	c832	V.2.....V.2
000000630:	c1d4	8114	3940	ea74	b5b3	e215	0008	e578	....9@.t.....x
000000640:	0275	8b14	0840	0400	0000	0000	0100	0000	..u..@.....
000000650:	0400	0000	0200	0000	2000	0000	0000	0000	.....
000000660:	0000	0000	6d81	0000	10d2	0000	0000	0000	....m.....
000000670:	0100	0800	0402	0500	6673	636b	0000	0000	.....fsck....

Fig. 1.26: File System Root B-Tree, showing entries for the file *fsck* found in the *sbin* directory.

Fig. 1.26 shows more directory entries for the file *fsck*, for its inode and for extended attributes. In order to identify what is relevant of data content we first need to read the inode, highlighted in red. We can see from the key that it is for the inode 0x10009c388, and it is of the object type 3 (inode). The value field is described in Table 1.19, where we have included the first fields.

One of the most important inode fields for file content is the internal flags. In this case, it describes that this inode has a resource fork (INODE\_HAS\_RSRC\_FORK), which means we should find directory entries describing extended attributes.

From Fig. 1.26 at offset 0x513 highlighted in dark blue, we find the first directory entry for the extended attribute (APFS\_TYPE\_XATTR) for the *fsck* file, and it has the name *com.apple.ResourceFork*. The corresponding value can be seen highlighted in dark blue at offset 0x5e0.

From Table 1.20 we can see that the resource fork data stream points to another inode (file\_id), here 0x10009c389, which consists of 0x3000 bytes (3 blocks) and its real size is 0x24d9. We need to find the directory entry and use its extent in order to find the data belonging to this resource fork. At offset 0x534 highlighted in purple,

Table 1.19: Inode of fsck

Offset	Length	Field	Value description
0x0	0x8	parent_id	0x18c60 (sbin)
0x8	0x8	private_id	0x10009c388 (fsck)
0x10	0x8	create_time	0x1481d4c132c85600 (Friday, 28 October 2016 23:07:59)
0x18	0x8	mod_time	0x1481d4c132c85600 (Friday, 28 October 2016 23:07:59)
0x20	0x8	change_time	0x15E2B3B574EA4039 (Sunday, 22 December 2019 13:13:31)
0x28	0x8	access_time	0x148B750278E50800 (Tuesday, 29 November 2016 07:45:56)
0x30	0x8	internal_flags	0x44008 (Uncompressed size, <b>resource fork</b> , explicit protection class)
0x38	0x4	nlink	0x1 (number of hardlinks)

Table 1.20: xattr\_val\_t

Offset	Length	Field	Value description
0x0	0x2	flags	0x1 (XATTR_DATA_STREAM)
0x2	0x2	xdata_len	0x30
0x4	8	xdata	<b>0x10009c389</b> (first 8 bytes)
0xC	8	xdata	0x24d9 (size)
0x14	8	xdata	0x3000 (allocated size)
0x1C	8	xdata	0x0 (crypto id)
0x24	8	xdata	0x24d9 (total bytes written)
0x24	8	xdata	0x0 (total bytes read)

we have the second extended attribute (APFS\_TYPE\_XATTR) for the fsck file, with the name com.apple.decmpfs. This has to do with data compression.

From the corresponding value, we can see flags are 2 (XATTR\_DATA\_EMBEDDED), which means the data is embedded into this value field. The data length is 0x10 (16) bytes. The data starts with the fpmc name (cmpf when read as LE), which is the 4-byte magic signature (compression\_magic). The next 4 bytes is the compression type used, here 0x4 (unknown, type 1 is uncompressed). Then the next 8 bytes are the uncompressed size, here 0xd210 (53776) bytes. There is no extra data, and we assume that this means that the resource fork data stream 0x10009c389 is compressed.

At offset 0x550 highlighted in yellow, we have a directory entry for a file extent for the inode id 0x10009c389, and this extent starts at logical address 0 (the start of the file). Please note that this is the node number after the one assigned to fsck, and is the same as the one identified as the data stream of the resource fork belonging to fsck. We assume it is compressed.

The corresponding value starts at offset 0xb4, and is highlighted in yellow. The first 8 bytes is the field len\_and\_flags, and we can see the length is 0x3000 (number of bytes in the assigned blocks, here 3 blocks), and flags are not in use. The next 8 bytes are the physical block number this extent starts, here 0x187c7d7. The last 8

bytes are the encryption key or tweak used for the extent. Here it is 0x0 which means encryption is not used.

Extracting the file content is just extracting the three blocks starting from physical address 0x187c7d7. Then we will have an extracted file that is recognized as an Apple HFS/HFS+ resource fork. Since we extract a compressed resource fork, we will need to extract and decompress its data. However, to decompress the resource fork data content, we need to extract the compressed data from the resource fork and then decompress the data using the appropriate algorithm.

In most cases, we will be extracting non-compressed extents, and therefore it is out of scope to describe the resource fork format in this chapter.

## 1.6 APFS - Application Category

The APFS does not use a journal, instead it uses a feature called Atomic Safe-Save (ASS) to ensure that an FS operation is either completed, or it does not happen. This is implemented by using Copy on Write (COW), and the use of checkpoints.

## 1.7 Comparing our results with a commercial tool

We selected EnCase 8 as the commercial tool to compare our results with, and also to verify the accuracy and reliability of EnCase APFS support.

Name	Logical Size	File Created	Entry Modified	Last Written
ck-Resource	9,433	10/29/16 01:07:59 AM	10/29/16 01:07:59 AM	12/22/19 02:13:31 PM
cd-Resource	95,186	10/29/16 02:02:52 AM	10/29/16 02:02:52 AM	12/22/19 02:13:31 PM
un-Resource	15,395	10/29/16 01:07:58 AM	10/29/16 01:07:58 AM	12/22/19 02:13:31 PM
unchn-Reso...	151,686	10/29/16 04:37:58 AM	10/29/16 04:37:58 AM	12/22/19 02:13:31 PM
unt_nfs-Res...	20,461	12/05/19 06:23:35 AM	12/05/19 06:23:35 AM	12/22/19 02:19:37 PM
ount-Resou...	8,391	12/05/19 06:23:35 AM	12/05/19 06:23:35 AM	12/22/19 02:19:37 PM

Fig. 1.27: EnCase v8.08.00.140

In Fig. 1.27 we can see that the directories and the file names are missing the first two characters in the *ot* folder, which we assume should be the *root* directory, and the *stem* directory should be the *System*, etc. EnCase used the correct volume name, but the GUID for the Container is not exactly correct. The third section shows B348, but should be 0348. Other than that the GUID from the superblock is correct.

The file we extracted in the content section had the name *fsck* and was in the *sbin* directory. EnCase had this file in the *in* directory (should be *sbin*, but misses the first two characters). Our file was 0x24d9 (9433) bytes in size, and it corresponds to the Logical Size in the selected file in the right table view in the Fig. 1.27. EnCase uses the selected timezone when showing the timestamps. File Creation (create\_time), the Entry Modified (mod\_time), Last Written (change\_time), and Last Access (access\_time) are identical to our results when taking the used timezone (UTC+1) into consideration.

EnCase shows the file with a filename *ck-Resource*, and the first two characters are missing. The "-Resource" is something EnCase have added to the file, which is not a part of the real file name. It may be their approach to show that this is a resource fork.

We can not validate this version of EnCase when it comes to APFS support, especially because it does not show the accurate directory and file names. This can be fixed in an updated version.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



# Chapter 2

## Ext4



Rune Nordvik

**Abstract** The Ext4 file system is often used by Android cell phones and by Linux distributions. As a mobile forensic expert, it is necessary to understand the structures of this file system to recover data, verify tool results, and detect anti-forensics techniques that may be present in the file system. In this chapter, we will have a deep dive into topics important for an investigation. Many digital forensic tools do not recover much from the Ext4 file system [52], and therefore we show some of the most useful Ext4 recovery techniques proposed by current research.

The Ext4 file system is often used by Android<sup>1</sup> operating systems, and also by Linux desktop distributions [14], and this file system is open source. The Ext4 file system replaces the Ext2 and Ext3, but it is mostly backwards compatible. Carrier described Ext2 and Ext3 in his File System forensic analysis book [10], which includes information also relevant for Ext4. Fairbanks describes the Ext4 file system at a low level and from a Digital Forensics perspective. This chapter will describe file system information important for mobile forensic investigators and other digital forensic experts.

### 2.1 Introduction

This chapter will give in-depth knowledge about the Ext4 file system. We assume the readers know how to use a hex editor and how to interpret multi-byte fields in a structure. This includes how to read raw data based on the used endianness.

Even an open-source file system needs explanation because the source code is not necessarily easy to understand and does not highlight what is important for

---

Rune Nordvik

The Norwegian Police University College (Politihøgskolen), Slemdalsveien 5, 0369 Oslo, Norway,  
e-mail: [rune.nordvik@phs.no](mailto:rune.nordvik@phs.no)

<sup>1</sup> Android is an operating system developed by Google, which is based on the Linux operating system

investigations. We will show how important structures look on disk and explain how these structures should be interpreted. This chapter scrutinises the Ext4 file system of a Samsung S8 and focuses on the system partition. Detailed information about the Samsung S8 can be found at GSMArena [31]. The Samsung S8 system partition contains an Ext4 file system and should not contain user data. The phone was reset using the phone menus before acquisition to avoid including any personal data. However, using the phone's reset system is no guarantee for a complete wipe of previous data. Therefore, the data set will not be shared to ensure the anonymity of the device owner and that there is no possibility to identify any personal data from any partition. We will see that the system partition has not been reformatted during the reset. We do not know if this is also true for the user data partition because it is encrypted. The use of encrypted user data partition is often mandatory [1], and on the Samsung S8, it was enabled as default.

### ! Attention

It is difficult to verify if a device is fully wiped without full access to the file system. If the phones own reset process is very fast, you should assume the data is not wiped. Even if the data was encrypted before resetting the device, it may be decrypted by existing or future methods.

## 2.2 Ext4 - File system category

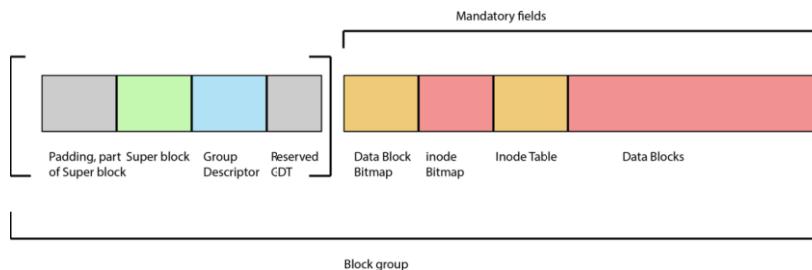


Fig. 2.1: Illustration of the elements of a block group.

The Ext4 contains multiple block groups, which have the same structure as shown in Figure 2.1. The first part of the block group is 1024 bytes reserved that can be used for boot code and form part of the first superblock [10, p. 402]. It is not a requirement that each block group has a superblock . The feature flag *sparse\_super* is set by default [10, p. 400], and it will store superblocks and group descriptors in block group 0, and then in block group 3<sup>x</sup>, 5<sup>x</sup>, and 7<sup>x</sup>. Another feature flag that

organises superblocks is the feature flag *sparse\_super2*. If set, the file system will only contain two superblock backups. If none of the *sparse\_super* flags is set, the superblock and group descriptor can be found in each block group.

### ! Tip: File system feature flags

Document the supported features of the file system under investigation. Do not assume all Digital Forensic tools support all features.

---

Since the superblock describes the overall file system; any superblock can be used to recover the file system. After the superblock we can find the group descriptor, which describes the block group. It has location addresses, statistics, and checksums about other mandatory elements in the block group. These other elements are the data block bitmap (allocation status of data blocks in the block group), the inode bitmap (allocation status of inodes in the inode table), the inode table, and finally, the data blocks. Data blocks can be assigned to file or directory content. A file can have any content, and a directory has directory entries. In order to include all metadata related to a file, it is necessary to connect the directory entry describing the file and its inode.

## 2.3 Superblock

The file system uses superblocks in order to describe important structures of the file system. This includes information such as the number of total inodes and blocks, how many inodes and blocks that are free, the size of inodes and blocks, information about file system checks, which OS the file system was created on, features the file system supports, a unique UUID (Universally Unique Identifier) for the file system volume, etc.

### 2.3.1 Temporary data about the File system

The superblock contains its temporary information, such as when it was created, or last mounted, or last written to. In the superblock, we can also find the first time an error found place and even the time of the last error. All timestamps found in the superblock are described as seconds since 1970 (Unix Epoch) and are defined as 32-bit fields that must be interpreted as little-endian.

From a mobile forensic expert view, it will be important to know when the file system was created since we can expect to find user allocated files created between the file system creation time and before the file system last written time. If we find

	Offset (h)	00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F	UTF-8
0000	0000788400	E0 2C 04 00 58 AF 10 00 00 00 00 DA 39 00 00	P
0010	0000788410	CC 0F 04 00 00 00 00 00 00 02 00 00 00 02 00 00	
0020	0000788420	00 00 00 00 00 00 00 00 00 70 1F 00 00 E7 0E BD 5E	A P A
0030	0000788430	E7 0E BD 5E 18 00 18 00 53 EF 01 00 03 00 00 00	I N
0040	0000788440	F0 88 58 49 00 4E ED 00 00 00 00 00 01 00 00 2C 00 00	B { , 9 C[,
0050	0000788450	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	Y78BBQ system
0060	0000788460	42 02 00 00 78 00 00 00 F1 CD 2A 39 FB 43 5B 2C	/system
0070	0000788470	95 32 59 3F 38 42 51 C8 73 79 73 74 65 6D 00 00	
0080	0000788480	00 00 00 00 00 00 00 00 00 2F 73 79 73 74 65 6D 00	
0090	0000788490	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00A0	00007884A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00B0	00007884B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00C0	00007884C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00D0	00007884D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00E0	00007884E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 AD D9 4F AB	
00F0	00007884F0	00 DE 53 3D 98 36 DC F3 C3 02 E3 01 01 00	L = 6
0100	0000788500	4C 00 00 00 00 00 00 00 F0 88 58 49 00 F3 01 00	I @
0110	0000788510	04 00 00 00 00 00 00 00 00 00 00 00 00 00 40 00 00	
0120	0000788520	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0130	0000788530	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0140	0000788540	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 04	
0150	0000788550	00 00 00 00 00 00 00 00 00 00 00 00 00 00 28 00 28	
0160	0000788560	01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0170	0000788570	00 00 00 00 00 00 00 00 00 49 17 09 00 00 00 00 00	I
0180	0000788580	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0190	0000788590	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	First time for error (s.first_error_time) N/A
01A0	00007885A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01B0	00007885B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	Time of most recent error (s.last_error_time) N/A
01C0	00007885C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01D0	00007885D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01E0	00007885E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01F0	00007885F0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	

Fig. 2.2: Timestamps found in the Ext4 superblock.

allocated files outside this time range, then this can be explained by one or more theories (hypotheses):

- Files that are part of the OS system installation process may keep one or more of their original timestamps.
- Apps may keep the timestamps when extracting container files, depending on how they extract the files.
- The cell-phone has lost its date/time due to power failure.
- The user could have reset the cell-phone date/time manually.
- Someone has manipulated the timestamps using a tool.

It is also possible to find previous files from before the file system was created, unallocated, from a previous file system. All these different theories about the reasoning for why we can find timestamps out of range is not complete and which theory (hypothesis) is the most likely should be tested.

### Tip: Use Experiments

Scientifically testing theories (hypotheses) is part of Digital Forensics. When it comes to file systems this can be done by performing experiments [35]. Do not base your investigation on assumptions!

When we reset the Samsung S8 device using the menus available, the system partition file system was not re-formatted. Since the partition was encrypted, we could not be certain if the user data partition file system was re-formatted. Figure 2.2 demonstrates

the creation date is from 2008. This is 9 years before this device was available on the market [31].

### 2.3.2 Supported features

The superblock defines the features supported in three different 32 bit fields;

- 0x5C s\_feature\_compat
- 0x60 s\_feature\_incompat
- 0x64 s\_feature\_ro\_compat

All the features found in these fields are supported by the file system driver version that created the file system. How the file system will be mounted depends on these three different fields. If the feature that is unrecognisable is found in the s\_feature\_compat, then the file system can be mounted with reading and writing support. If the feature not recognised is found in the s\_feature\_incompat, then it should not mount the file system. If the feature not recognised is found in the s\_feature\_ro\_compat, then the file system can be mounted as read-only.

	Offset (h)	00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F	UTF-8
0000	000D780400	E0 2C 04 00 50 AF 10 00 00 00 00 00 DA 39 00 00	P
0010	000D780410	CC 0F 04 00 00 00 00 00 02 00 00 00 02 00 00 00	p
0020	000D780420	00 80 00 00 00 80 00 00 70 1F 00 00 E7 0E BD 5E	A S
0030	000D780430	E7 0E BD 5E 18 00 18 00 53 EF 01 00 03 00 00 00	
0040	000D780440	F0 88 5B 49 00 4E ED 00 00 00 00 00 01 00 00 00	I N
0050	000D780450	00 00 00 00 00 00 00 00 00 01 00 00 00 00 00 00	
0060	000D780460	42 02 00 00 7B 00 00 00 F1 CD 2A 39 FB 43 5B 2C	B { ' 9 Cl,
0070	000D780470	98 32 59 3F 38 42 51 C0 73 79 73 74 65 60 00 00	2Y78BQ system
0080	000D780480	00 00 00 00 00 00 00 00 2F 73 79 73 74 65 60 00	/system
0090	000D780490	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00A0	000D7804A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00B0	000D7804B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00C0	000D7804C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00D0	000D7804D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00E0	000D7804E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 AD D9 4F AB	
00F0	000D7804F0	8D DE 53 3D 9B 36 DC F3 00 C3 02 E3 01 01 00 00	= 6 I
0100	000D780500	4C 04 00 00 00 00 00 00 F0 88 5B 49 0A F3 01 00	@
0110	000D780510	00 00 00 00 00 00 00 00 00 00 00 00 00 00 40 00 00	
0120	000D780520	00 80 07 00 00 00 00 00 00 00 00 00 00 00 00 00	
0130	000D780530	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0140	000D780540	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 04	

Compatible features (s\_feature\_compatible):  
0x2C = 0x20 Indexed directories, 0x8 Support extended attributes, 0x4 Journal

Incompatible features  
(s\_feature\_incompat) 242 = 0x200 Flexible block groups, 0x40 Files uses extents, 0x2 Directory entries record file type.

Read only compatible features  
(s\_feature\_ro\_compatible) 79 = 0x40 Large inodes, 0x20 Ext3 32000 subdirectory limit no longer applies, 0x10 Group descriptors have checksums, 0x8 Files space usage is stored in units of inode block sizes (huge files), 0x2 Allow storing files larger than 2 GiB (large files), 0x1 sparse superblocks.

Fig. 2.3: Feature flags in the Ext4 superblock.

#### Tip: Features have an impact

The features described in the superblock impacts how the inodes and directory entries should be interpreted.

## Compatible features

Even if the kernel does not understand one of the flags in this 32 bit field, it will mount the file system with read and write support. Table 2.1 can be used to interpret this field, and Figure 2.3 demonstrates an example of interpretation.

Table 2.1: Compatible features

Value	Description
0x1	Directory preallocation (COMPAT_DIR_PREALLOC)
0x2	Could mean the fs supports AFS magic directories. (COMPAT_IMAGIC_INODES)
0x4	Has a journal (COMPAT_HAS_JOURNAL)
0x8	Supports extended attributes (COMPAT_EXT_ATTR)
0x10	Has reserved GDT blocks for filesystem expansion (COMPAT_RESIZE_INODE)
0x20	Has directory indexes (COMPAT_DIR_INDEX)
0x40	“Lazy BG”. Not in Linux kernel (COMPAT_LAZY_BG)
0x80	“Exclude inode”. Not used. (COMPAT_EXCLUDE_INODE)
0x100	“Exclude bitmap”. Not used (COMPAT_EXCLUDE_BITMAP)
0x200	Sparse Super Block, v2 (COMPAT_SPARSE_SUPER2)

From an investigator’s perspective, not every compatible feature is relevant, however, the flag COMPAT\_SPARSE\_SUPER2 is especially important when locating the backup superblocks, in case the main one is partly corrupted or manipulated. If the COMPAT\_SPARSE\_SUPER2 flag is set, the super block field *s\_backup\_bgs*, found from superblock byte offset 0x24C, points to the two block groups that contain backup superblocks. It may seem strange that one field can point to two blocks, but this is because the field is an array of two 32 bits elements. In our example in Figure 2.3 the flag was not set. Another important flag is the COMPAT\_HAS\_JOURNAL. If this flag is set, recovery of data from the journal should be possible. Also note that when the journal is full, it will start writing transactions from the beginning of the journal file, effectively overwriting previous transactions [14]. More details about the Ext4 journal can be found in the sect. 2.7 at page 68. The feature COMPAT\_EXT\_ATTR is important since it allows extended attributes to be saved within the inode. This allows the user or programs to add extra information to individual files.

The flag COMPAT\_RESIZE\_INODE does not have a descriptive name, since it describes the number of blocks reserved for the extra Group Descriptor Table (GDT). These blocks are reserved for future file system expansion. These are important because the mandatory fields in the block group will be found after the blocks reserved GDT as illustrated in Figure 2.1, and this knowledge can be used for manual recovery of the file system.

Some of the flags are supported by Linux, but not necessarily used. For instance the COMPAT\_DIR\_PREALLOC which allows for pre-allocating a specific number of blocks to directories, defined in field *s\_prealloc\_dir\_blocks* at superblock byte offset 0xCD. The field is currently not used by the Linux kernel [41].

## Incompatible features

If the kernel does not understand one of the flags in this 32 bit field, it should not mount or repair the file system. Table 2.2 can be used to interpret this field, and Figure 2.3 demonstrates an example of interpretation.

Table 2.2: Incompatible features

Value	Description
0x1	Compression. Not implemented. (INCOMPAT_COMPRESSION)
0x2	Directory entries record the file type (INCOMPAT_FILETYPE)
0x4	Filesystem needs journal recovery. (INCOMPAT_RECOVER)
0x8	Filesystem has a separate journal device. (INCOMPAT_JOURNAL_DEV)
0x10	Meta block groups. See the earlier discussion of this feature. (INCOMPAT_META_BG)
0x40	Files in this filesystem use extents. (INCOMPAT_EXTENTS)
0x80	Enable a filesystem size over $2^{32}$ blocks. (INCOMPAT_64BIT)
0x100	Multiple mount protection. Prevent multiple hosts from mounting the filesystem concurrently (INCOMPAT_MMP)
0x200	Flexible block groups (INCOMPAT_FLEX_BG)
0x400	Inodes can be used to store large extended attribute values (INCOMPAT_EA_INODE)
0x1000	Data in directory entry. Feature still in development (INCOMPAT_DIRDATA)
0x2000	Metadata checksum seed is stored in the superblock (INCOMPAT_CSUM_SEED)
0x4000	Large directory >2GB or 3-level htree (INCOMPAT_LARGEDIR)
0x8000	Data in inode. Small files or directories are stored directly in the inode i_blocks and/or xattr space. (INCOMPAT_INLINE_DATA)
0x10000	Encrypted inodes are present on the filesystem (INCOMPAT_ENCRYPT)

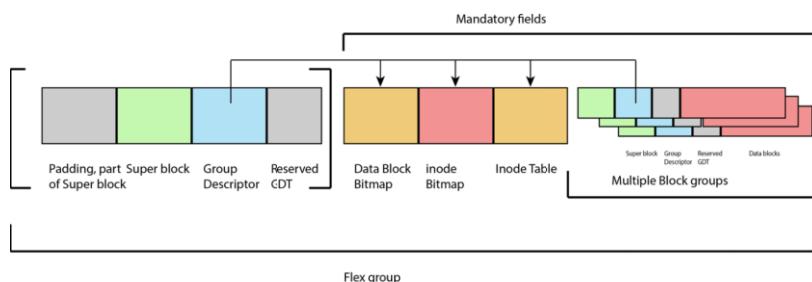


Fig. 2.4: Illustration of the elements of a flex group.

Flexible block groups are a unique way of organizing block groups into a set of flex groups. The first block of a flex group will include the bitmaps and the inode table for all groups within all the flex groups, and the other groups may contain super blocks

and group descriptors depending on the sparse superblock feature, and will include data blocks, as shown in Figure 2.4. The group descriptor is used to define where the bitmaps and inode table should be located, which enables flex groups, meaning they all point to the same bitmaps and inode locations as the first group descriptor. This is important to understand, because it deviates from the standard organisation of a block group as shown in Figure 2.5.

### Tip: Metadata location

The investigator should assume and test that the data block bitmap, inode bitmap, and inode table will exist in the first block group when flex groups are being used. In this case the metadata is near co-located in the beginning of the file system.

When flex groups are not used, it will be necessary to parse all superblocks and group descriptors in order to identify all bitmaps and the complete inode table.

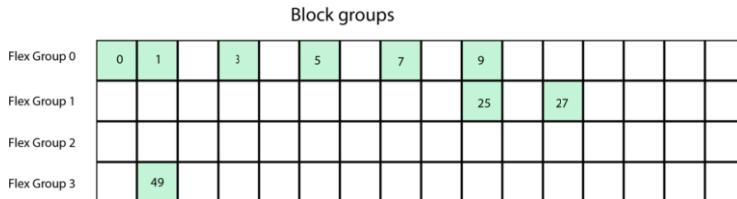


Fig. 2.5: Illustration of superblocks and group descriptors in flex groups or not flex groups when also the RO\_COMPAT\_SPARSE\_SUPER is in use. Based on [5].

### Read only compatible features

If the kernel does not understand one of the flags in this 32-bit field, it will mount the file system as read-only. Table 2.3 demonstrates that if the file system has large files, the superblock will use options like *RO\_COMPAT\_LARGE\_FILE* (exist files larger than 2 GiB), *RO\_COMPAT\_BIGALLOC* (extents are using clusters instead of blocks) and *RO\_COMPAT\_HUGE\_FILE* (file size is shown in logical blocks instead of sectors). Large files can be of investigative value since they may contain videos, file system containers, encrypted files, etc. It is important that these files are investigated.

### Tip: Large Files

Use tools that flag large files, or sort a file listing by file size.

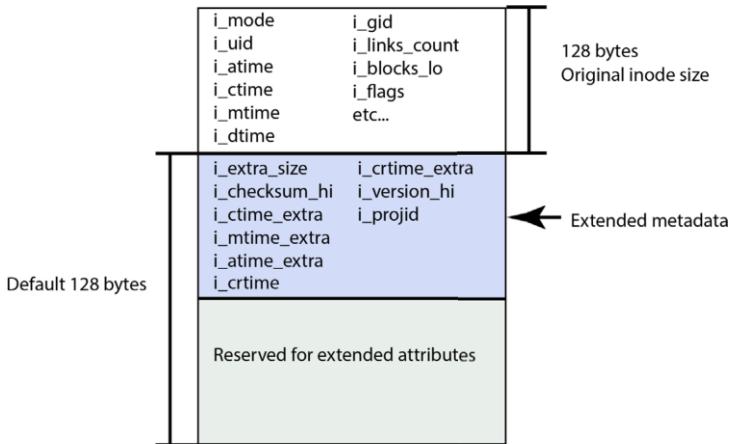


Fig. 2.6: Illustration of usage of extended metadata in inodes.

The default inode size was 128 bytes in Ext2 and Ext3, but in Ext4 it was typically 256 bytes. The option `RO_COMPAT_EXTRA_ISIZE` means that extended metadata is utilised, which will allow Ext4 features such as nano second parts of timestamps, the creation timestamp etc. The part after the extra inode size is still reserved for extended attributes. Figure 2.6 demonstrates that the extra metadata is found directly after the first 128 bytes of an inode, and that the extended attributes follow after the extended metadata. If extended metadata is not in use, then this area will be reserved for extended attributes. In our example Ext4 file system extended attributes was supported, see Figure 2.3 at page 45.

Using checksums of metadata is a measure to protect the metadata from being used if corrupted or manipulated. However, it does not protect metadata change if the checksum is updated and verifies the manipulated metadata.

Another important feature for the investigator is to check if the file system contains snapshots (a previous state of the file system). However, currently Ext4 in Linux/Android does not support snapshots.

The verity feature (`RO_COMPAT_VERITY`) may be interesting for the investigator, which means verity inodes may exist on the file system. These inodes have content that is read-only, and can be verified using a Merkle tree-based hash. A Merkle tree-based hash means that the file is divided into blocks that are hashed. Then these hashes are concatenated and represent larger blocks of data and re-hashed. This continues until there is one large block left, representing the complete file, which is hashed. It is this final hash the read-only file is verified against. The Figure 2.7 illustrates the Merkle tree hashes and is also explained by Merkle [50] who describes that in order to verify the public key, you only need the hashes of the first and second half of the public key and that you can compute half of the public key by knowing their quarters.

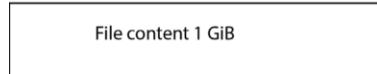
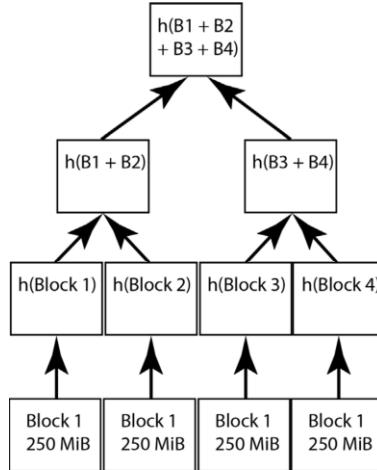


Fig. 2.7: Illustration of Merkle tree hash.

	Offset (h)	00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F	UTF-8
0000	000D780400	E0 2C 04 00 50 AF 10 00 00 00 00 00 DA 39 00 00	P
0010	000D780410	CC 0F 04 00 00 00 00 00 02 00 00 00 00 00 00 00	
0020	000D780420	00 00 00 00 00 00 00 00 70 1F 00 00 E7 0E BD 5E	D ^ S
0030	000D780430	E7 0E BD 5E 18 00 18 00 53 EF 01 00 03 00 00 00	A N
0040	000D780440	F0 88 55 49 00 4E ED 00 00 00 00 01 00 00 00 00	
0050	000D780450	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	B ( 9 Cl,
0060	000D780460	42 02 00 00 78 00 00 00 F1 CD 2A 39 FB 43 5B 20	ZY788Q system
0070	000D780470	78 32 59 3F 38 42 51 C0 73 79 73 74 65 6D 00 00	/system
0080	000D780480	00 00 00 00 00 00 00 00 F2 73 79 73 74 65 6D 00 00	
0090	000D780490	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00A0	000D7804A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00B0	000D7804B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00C0	000D7804C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00D0	000D7804D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
00E0	000D7804E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 AD D9 4F AB	
00F0	000D7804F0	BD DE 53 98 36 DC F3 F0 C3 02 E3 01 00 00 00 00	= 6 L I @
0100	000D780500	AC 04 00 00 00 00 00 00 F8 88 5B 49 0A F3 01 00	
0110	000D780510	04 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0120	000D780520	00 00 07 00 00 00 00 00 00 00 00 00 00 00 00 00	
0130	000D780530	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0140	000D780540	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0150	000D780550	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0160	000D780560	01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0170	000D780570	00 00 00 00 04 00 00 00 00 00 00 00 00 00 00 00	I
0180	000D780580	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
0190	000D780590	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01A0	000D7805A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01B0	000D7805B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01C0	000D7805C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01D0	000D7805D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01E0	000D7805E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
01F0	000D7805F0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	

Annotations for the Ext4 superblock:

- Block size (s\_log\_block\_size)  $2 = 2^{(10+2)}$  = 4096
- Cluster size (s\_log\_cluster\_size)  $2 = 2^{(10+2)} = 4096$
- Size of inode structure (s\_inode\_size), 100h=256 in decimal
- Block group number (s\_block\_group\_nr) 0=First block group
- Universally unique identifier for the volume (s\_uuid), here 392ACDF1-43FB-2C5B-98-32-593F384251C0
- Volume label (s\_volume\_name): system
- Directory last used as a mount point (s\_last\_mounted): /system
- Number of reserved Group Descriptor Blocks (s\_reserved\_gdt\_blocks): 0

Fig. 2.8: Information about blocks in the Ext4 superblock.

The feature flag `RO_COMPAT_READONLY` means that this file system should only be mounted as read only. Most implementation of Ext4 file system drivers complies

Table 2.3: Read only compatible features

Value	Description
0x1	Sparse superblocks. See the earlier discussion of this feature (RO_COMPAT_SPARSE_SUPER)
0x2	This filesystem has been used to store a file greater than 2GiB (RO_COMPAT_LARGE_FILE)
0x4	Not used in kernel or e2fsprogs (RO_COMPAT_BTREE_DIR)
0x8	This filesystem has files whose sizes are represented in units of logical blocks, not 512-byte sectors (RO_COMPAT_HUGE_FILE)
0x10	Group descriptors have checksums (RO_COMPAT_GDT_CSUM)
0x20	Indicates that the old ext3 32,000 subdirectory limit no longer applies (RO_COMPAT_DIR_NLINK)
0x40	Indicates that large inodes exist on this filesystem (RO_COMPAT_EXTRA_ISIZE)
0x80	This filesystem has a snapshot (RO_COMPAT_HAS_SNAPSHOT)
0x100	Quota (RO_COMPAT_QUOTA)
0x200	This filesystem supports “bigalloc”, extents are tracked in units of clusters (of blocks)(RO_COMPAT_BIGALLOC)
0x400	This filesystem supports metadata checksumming. (RO_COMPAT_METADATA_CSUM)
0x800	Filesystem supports replicas. This feature is neither in the kernel nor e2fsprogs (RO_COMPAT_REPLICA)
0x1000	Read-only filesystem image; the kernel will not mount this image read-write and most tools will refuse to write to the image (RO_COMPAT_READONLY)
0x2000	Filesystem tracks project quotas (RO_COMPAT_PROJECT)
0x8000	Verity inodes may be present on the filesystem (RO_COMPAT_VERTITY)

with this setting, but there may exist driver implementations or tools who allow writing to the file system even if it is set to read only.

#### Tip: Test if a file system is read only

The investigator should perform experiments to test if it is possible to write to an identical copy of the read only file system using the same driver or tools found on the device under investigation.

### 2.3.3 The group descriptor

The group descriptor describes information about a particular group [14], for instance, the locations of the block bitmap, inode bitmap, and the inode table. In order to find the group descriptor, we need to know the block size, as shown in Figure 2.8. The value in this field is 2, and the formula we need to use is  $10^{(10+s\_log\_block\_size)}$ . We can find the group descriptor in the block following the superblock. In order to

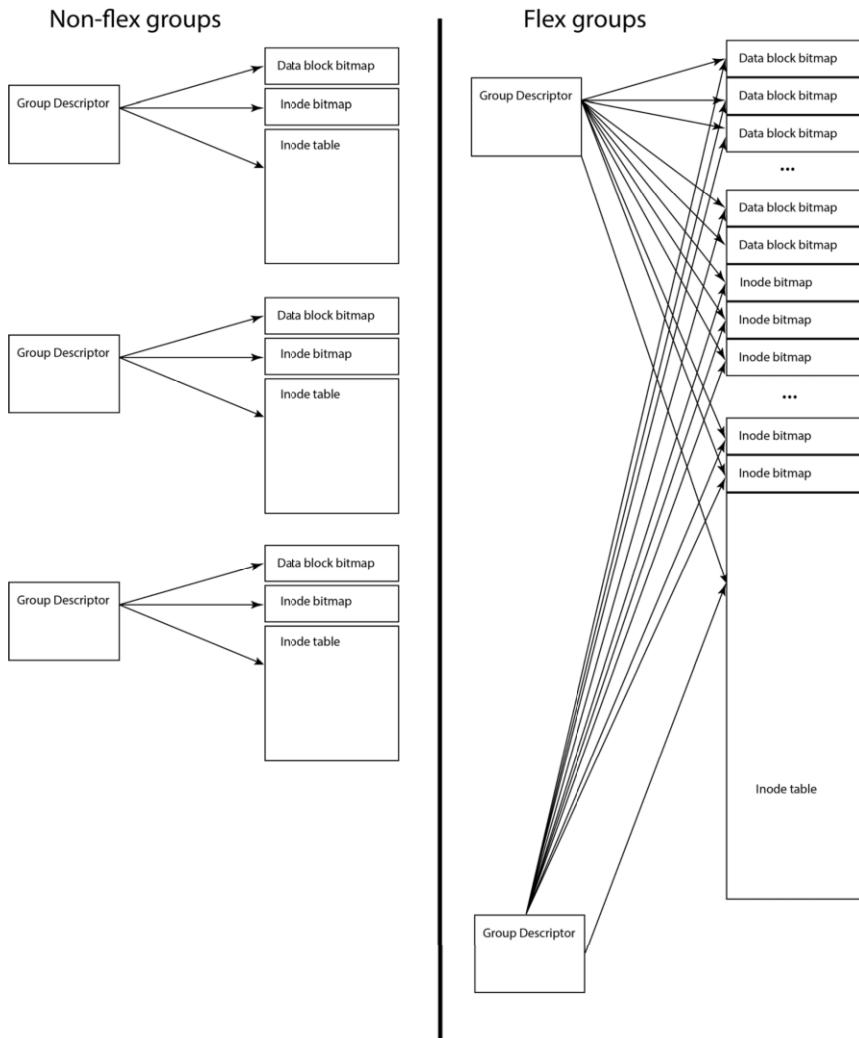


Fig. 2.9: Different designs for Group Descriptors

find the group descriptor, we, in this case, move 4096 bytes, one block, forward from the start of the superblock, from byte offset 0, not from 1024.

If the Ext4 has the 64-bit feature (INCOMPAT\_64BIT) enabled, then the location of the bitmaps and the inodes table has two fields each. The first fields should describe the lower bits for the location, while the last describes the upper bits. These fields should describe the block location of the block bitmap. In our example, the 64-bit feature was not enabled, and therefore each group descriptor is only 32 bytes. The inode table can be found in the block defined in `bg_inode_table_lo` at group descriptor byte offset 0x08. The locations are relative to the start of the superblock.

	Offset (h)	00 01 02 03 04 05 06 07	08 09 0A 0B 0C 0D 0E 0F	UTF-8	
0000	000D781800	02 00 00 00 12 00 00 00	22 00 00 00 00 00 F5 08	"	
0010	000D781810	7D 02 04 00 00 00 00 00	00 00 00 00 F5 08 18 E7 }	f .	Location to data block bitmap (0x2)
0020	000D781820	03 00 00 00 13 00 00 00	19 02 00 00 00 00 66 1F	f	Location to inode bitmap (0x12 or 18 in decimal)
0030	000D781830	0A 00 00 00 00 00 00 00	00 00 00 00 66 1F 2E BB	p	Location to inode table (0x22 or decimal 34)
0040	000D781840	04 00 00 00 14 00 00 00	18 04 00 00 00 00 70 1F	p	
0050	000D781850	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 8C 88	p	
0060	000D781860	05 00 00 00 15 00 00 00	07 00 00 00 00 00 70 1F	p	
0070	000D781870	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 70 C3	p p	
0080	000D781880	06 00 00 00 16 00 00 00	FE 07 00 00 00 00 70 1F	p	
0090	000D781890	00 00 05 00 00 00 00 00	00 00 00 00 70 1F CF 80	p π	
00A0	000D7818A0	07 00 00 00 17 00 00 00	F5 09 00 00 00 00 70 1F	p	
00B0	000D7818B0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 22 5F	p " -	
00C0	000D7818C0	08 00 00 00 18 00 00 00	EC 0B 00 00 00 00 70 1F	p	
00D0	000D7818D0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 67 5D	p gJ	
00E0	000D7818E0	09 00 00 00 19 00 00 00	E3 D0 00 00 00 00 70 1F	p	
00F0	000D7818F0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 46 39	p F9	
0100	000D781100	0A 00 00 00 1A 00 00 00	DA 0F 00 00 00 00 70 1F	p	
0110	000D781110	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 04 FD	p	
0120	000D781120	0B 00 00 00 18 00 00 00	D1 11 00 00 00 00 70 1F	p	
0130	000D781130	00 00 05 00 00 00 00 00	00 00 00 00 70 1F FB B3	p	
0140	000D781140	0C 00 00 00 1C 00 00 00	C8 13 00 00 00 00 70 1F	p	
0150	000D781150	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 7A CE	p z	
0160	000D781160	0D 00 00 00 1D 00 00 00	BF 15 00 00 00 00 70 1F	p	
0170	000D781170	00 00 05 00 00 00 00 00	00 00 00 00 70 1F A5 2F	p /	
0180	000D781180	0E 00 00 00 1E 00 00 00	B6 17 00 00 00 00 70 1F	p	
0190	000D781190	00 00 05 00 00 00 00 00	00 00 00 00 70 1F D6 14	p	
01A0	000D7811A0	0F 00 00 00 1F 00 00 00	AD 19 00 00 00 00 70 1F	p	
01B0	000D7811B0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F 2A C6	p *	
01C0	000D7811C0	10 00 00 00 20 00 00 00	A4 1B 00 00 00 00 70 1F	p	
01D0	000D7811D0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F DD 89	p ^	
01E0	000D7811E0	11 00 00 00 21 00 00 00	9B 1D 00 00 00 00 70 1F	!	p
01F0	000D7811F0	00 00 05 00 00 00 00 00	00 00 00 00 70 1F CF FA		

Fig. 2.10: Group descriptors in a flex group

The data block location is defined in field `bg_block_bitmap_lo` at offset 0x0, and the inode bitmap is defined in field `bg_inode_bitmap_lo` at offset 0x4. All these values are 32 bits and must be interpreted as little-endian, as shown in Figure 2.10. However, in Figure 2.10 we see that there are multiple 32-byte units, where each of them is a group descriptor, one for each block group in the flex group. A very similar copy of this group descriptor block is found in all other group descriptor blocks. However, `bg_flags` values may deviate. It is important to understand that not all block groups have superblocks or group descriptors if either the superblock RO\_COMPAT\_SPARSE\_SUPER or the COMPAT\_SPARE\_SUPER2 feature flag is set. The field `bg_flags` can have any combination of these values :

- 0x1 Inode table and bitmap are not initialized
- 0x2 Block bitmap is not initialized
- 0x4 Inode table is zeroed (on initialisation)

In Figure 2.10 the flags value is 0x4 for block group 0, 1. While it is 0x5 for the rest of the block group descriptors in this descriptor block, which means that these block groups have not initialized their inode table or inode bitmaps. We could verify that there were only initialized inode tables in the first two locations (block 0x22 and 0x219).

Table 2.4: Group descriptor

Offset	Size	Name	Description
0x0	0x4	bg_block_bitmap_lo	Location to data block bitmap
0x4	0x4	bg_inode_bitmap_lo	Location to inode block bitmap
0x8	0x4	bg_inode_table_lo	Location to the inode table
0xC	0x2	bg_free_blocks_count_lo	Free blocks in block group
0xE	0x2	bg_free_inodes_count_lo	Free inodes in block group
0x10	0x2	bg_used_dirs_count_lo	Used directories in block group
0x12	0x2	bg_flags	Important for bitmaps and inode tables
0x14	0x4	bg_exclude_bitmap_lo	Location of snapshot exclusion bitmap
0x18	0x2	bg_block_bitmap_csum_lo	Data block bitmap checksum
0x1A	0x2	bg_inode_bitmap_csum_lo	Inode bitmap checksum
0x1C	0x2	bg_itable_unused_lo	Unused inodes in group

## Universal Unique Identifier

In the superblock the field, *s\_uuid*, assigns a unique identifier for the file system volume. This should be unique for every instance of a volume created, however, if we flash a partition, the target may be assigned the same UUID for its file system as the original source.

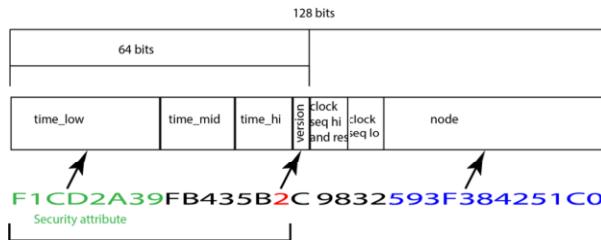


Fig. 2.11: Structure of the UUID v.2

The structures of the UUIDs are defined in RFC4122 [44], and the one used here is version 2 as shown in Figure 2.11. It uses a 60-bit timestamp (in which the four least significant bytes are overwritten with a security attribute) with an Epoch from 15th of October 1582, and a node identifier (MAC address) at the last 6 bytes. How is this important? Assuming the vendor is following the standard, it can approximate the file system creation and be connected to a MAC address. The MAC address in this example is globally unique and is a multicast address. However, the author is not convinced that the vendor follows the standard for the following reasons (1) the timestamp does not reflect the time of file system creation (2) the MAC address organisational part (OUI) is not recognised as a known organisation/vendor.

## 2.4 Ext4 - Metadata Category

Here we describe the inodes, inode bitmap, extended attributes.

### 2.4.1 The inode

The index node (inode) is defined in the structure *ext4\_inode*, which defines most of the metadata related to a file, except its file name. Previous versions of Ext used a 128-byte size inode, while the Ext4 standard uses 256 bytes. However, the first 128 bytes are backwards compatible with previous versions of Ext. The information in this section is based on the Ext4 source code and the interpretation found at Kernel.org [41].

Table 2.5: Inode offset table

Offset	Size	Name	Description
0x00	0x2	i_mode	User privileges and type of file
0x02	0x2	i_uid	Lower 16 bits of the owner id
0x04	0x4	i_size_lo	Lower 32 bits of the file size
0x08	0x4	i_atime	Last access time
0x0C	0x4	i_ctime	Last inode change time
0x10	0x4	i_mtime	Last data modification time
0x14	0x4	i_dtime	Deletion time
0x18	0x2	i_gid	Lower 16 bits of group id
0x1A	0x2	i_links_count	Number of hard links pointing to this file
0x1C	0x4	i_blocks_lo	Lower 32 bits of 512 byte blocks this file uses
0x20	0x4	i_flags	Inode flags
0x24	0x4	i_osd1	For Linux this is the inode version
0x28	0x3C	i_block[]	Block map or Extent tree.
0x64	0x4	i_generation	File version for NFS
0x68	0x4	i_file_acl_lo	Lower 32 bit address of extended attribute block
0x6C	0x4	i_size_high	Higher 32 bit address of file size
0x70	0x4	i_obso_faddr	Obsolete fragment address
0x74	0xC	i_osd2	OS descriptor 2
0x80	0x2	i_extra_isize	Size of the used area of inode - 128
0x82	0x2	i_checksum_hi	Upper 16-bits of the inode checksum
0x84	0x4	i_ctime_extra	Extra change time bits
0x88	0x4	i_mtime_extra	Extra modification time bits
0x8C	0x4	i_atime_extra	Extra access time bits
0x90	0x4	i_crttime	File creation time, in seconds since the Unix Epoch
0x94	0x4	i_crttime_extra	Extra file creation time bits
0x98	0x4	i_version_hi	Upper 32-bits for version number
0x9C	0x4	i_projid	Project ID

## 2.4.2 User privileges and type of file

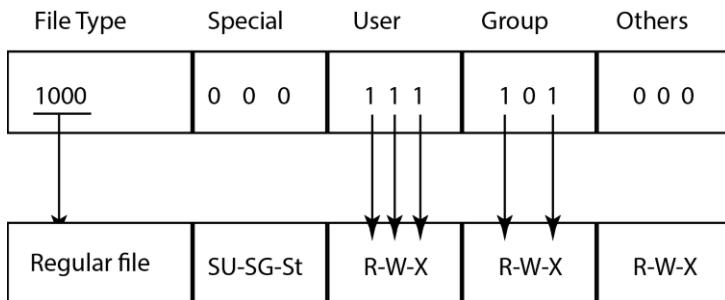


Fig. 2.12: File type and privileges.

As illustrated in Figure 2.12 the *i\_mode* field name 12 least significant bits are used for user privileges. These privileges are important when investigating a file or directory since it explains ownership and user privileges. However, it is also important to understand that these privileges may be changed if the user has the privileges to do so. Let us assume they are 000111101000. The 3 least significant bits describe all others and are 000, which corresponds with r (read)-w(write)-x(execute). In this case, none of them are set, which means that users not defined as the owner of the file or not within the filegroup will not have privileges for this file. The second 3 least significant bits are 101, and they describe the group. In this case, read and execute is set, while the write is not. The third 3 least significant bits have the value 111, and they describe the owner. The 3 most significant bits are special privileges. Here they are 000, and it means the Set-UID, Set-GID, and the Sticky bit are not set. Set-UID makes sure an executable uses the owner as the user executing the file and not the actual user executing it.

Similarly, it is possible to force using the defined group id for this file as the executable group instead of the real user group assigned to the user executing the file. The least significant bit of the special bits is for the Sticky bit, which affects directories. If it is set, it means that all files within this directory can only be modified by the owner. The remainder 4 bits of the *i\_mode* field are for describing the type of file. An inode can describe a regular file, a directory, a device (character-based or block-based), a symbolic link, a named pipe (FIFO) or a socket, as shown in Table 2.6. Knowing the type of the file tells the investigator what kind of inode is under investigation. This can give insight into if an inode describes a communication socket (two-ways communication) or FIFO (one-way communication), or if it is just a pointer (symbolic link) to another inode, or if the inode is used to access a storage device (for instance a sd\_card). All kinds of devices can be accessed through an inode describing a device. There are two main types of devices: block and character. A block device stores data in predefined blocks that may be randomly accessed. A character

device can be read from and written to and accessed as a sequential stream of bytes. A file system is a block device, and most devices could also be character devices. Hard disks could have interfaces for both block devices and character devices [55].

The difference between block device and character device is that the former is described data in predefined blocks, and these blocks may be randomly accessed. A character device is accessed through a stream of data in sequence, for instance a network card (REF).

Table 2.6: Inode file types

4 MSb	Meaning
0001	Special FIFO file (named pipe)
0010	Character device
0100	Directory
0110	Block device
1000	Regular file
1010	Symbolic link
1100	Socket

### 2.4.3 Temporary metadata describing inodes

Almost every inode has fields describing important timestamps. For backward compatibility, these are located from hex offset 0x08 from the start of the inode, and are 32-bit integers describing the number of seconds since 1970 (Unix Epoch). However, extra 32 bit fields in the inode use the least significant 2 bits to expand the timestamp to 34 bits. The remainder of the 30 bits is used for nanoseconds granularity.

We can only trust the timestamps if there is no malware installed on the device or any other tools to manipulate the inode metadata. The mobile device clock also needs to be accurate. The following section explains one method of manipulation.

When a file is created, the current time is set for all the timestamps in the inode. This means that if all the timestamps are the same, the file has not been changed after creation and it has not been accessed at a later time as long as the flags do not contain the flag 0x80 (bit 7 is set, counting from bit 0), which means the file system does not update the access date. If this flag is not set, the access date will update when a user or a program access the file. The investigator should always check if access times are close in time, indicating a program has accessed multiple files in the session. For instance, an anti-virus program may have opened each of the files without resetting the access time. A digital forensic logical extraction of selected files will update the accessed timestamp if the accessed timestamps are updated, assuming the tools requesting these files are using the operating system.

## 2.4.4 Temporary metadata manipulations

### ! Attention

It has been reported that it is possible to use the nano seconds part of a timestamp to hide information in Ext4 [18].

---

It is difficult to detect manipulations of the least significant parts of an Ext4 timestamp because most current listing tools do not show timestamps with the nanoseconds granularity, and even if they do, it is difficult to detect these manipulations by the user. The data hiding in the nanosecond part of a timestamp can easily get corrupted if all timestamps fields are used for hiding data. Timestamps such as *i\_ctime* and *i\_mtime* can be changed by user activity. However, the created timestamp (*i\_crtime*) will not change since it defines the creation of a file, and a delete operation will not affect such a date [18]. Although, a deleted inode gets unallocated in the inode bitmap and can therefore be reallocated by new inodes. This reallocation will destroy parts of the hidden data, which requires error measures in order to recover hidden data [18]. To preserve secrecy, the user can utilize cryptography. [18] describe that they used symmetric string cyphers in their proof of concept tool. They also repaired the inode checksums for each manipulated inode. Therefore, the detection of manipulated inodes is difficult to detect.

### Tip: Detect Manipulation

Document the Apps, tools, or malware installed on a mobile device. Try to identify their abilities from trusted sources. Investigate tools that have abilities to manipulate metadata.

---

Some tools may have timestamp manipulation or steganography abilities. This is one of the reasons why digital forensic experts should document the Apps, tools, or malware installed on a mobile device.

Fortunately, modern mobile devices have protection mechanisms to avoid installing software that is not approved by the device provider. Apple uses the App-Store, while Google uses the Google Play protect functionality. However, the latter can be easily disabled by the user. In addition, devices can be jailbroken on iOS or rooted on Android, allowing users to install anything.

### ! Attention

Malware needs to survive a reboot, and therefore it will try to stay hidden in the file system. Data hiding within file system metadata is a known approach [37].

---

## 2.4.5 Links count

The field *i\_links\_count* shows the number of directory entries referring to this inode. The directory entry has the inode number to which inode it points to. Multiple directory entries could be pointing to the same inode that indicate hard links. When the last directory entry pointing to an inode is deleted, this inode is marked as unallocated in the inode bitmap [10, p. 426]. This is not the same with soft links. Adding a soft link will not increase the links count of the inode it points to. Instead, it will create a symbolic inode. This symbolic inode points to a file path (directory entry), not to an inode [10, p. 426].

### Blocks used by a file

The number of 512 byte blocks (sectors) used by a file is defined in the *i\_blocks\_lo* field. However, if the *inode i\_flags* has the EXT4\_HUGE\_FILE\_FL file option set and the superblock has the huge file feature enabled then the field *i\_blocks\_hi* needs to be added using this formula.

$$(i\_blocks\_lo + i\_blocks\_hi \ll 32)$$

If the *i\_flags* has the EXT4\_HUGE\_FILE\_FL inode but file system does not have the huge file feature then field *i\_blocks\_hi* needs to be added using this formula.

$$i\_blocks\_lo + (i\_blocks\_hi \ll 32)$$

### Inode flags

This field has several options describing special properties for the file. A few flags that could be important for the investigation:

- 0x10 File is immutable, which means the file can not be changed.
- 0x20 File can only be appended.
- 0x80 Does not update access time. This is important because we know this timestamp is no longer updated on access.
- 0x800 Encrypted inode, which means the file content is encrypted.
- 0x4000 File data must be written through the journal. This is important since the previous content from this file may be found in the journal as long as the journal has not been overwritten with new transactions. This also depends on what is being recorded to the journal.
- 0x40000 This is a huge file, which has special meaning when computing the block size of a file.
- 0x80000 The file uses extents, which we explain in sect. 2.4.5. If this is not set it may use direct or indirect block pointers.

- 0x10000000 The inode contains inline data.
- 0x20000000 Create children with the same Project ID.

## Block map, Extent tree or inline data

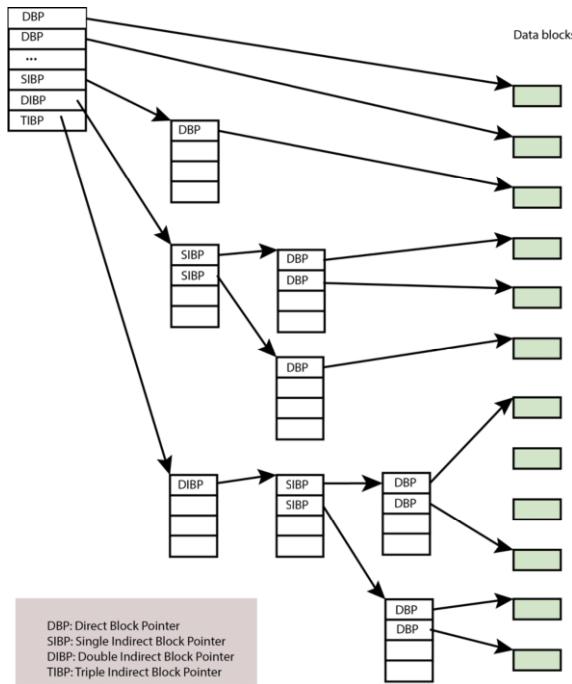


Fig. 2.13: Direct or indirect block pointers. Primarily used by previous Ext versions, but still supported in Ext4.

As illustrated in Figure 2.13, previous versions of Ext used block maps (direct or indirect blocks). However, Ext4 can still support it. Block maps are inefficient when a file uses many blocks since a maximum of 15 block pointers is available. The first 12 should be direct block pointers, while the last three could be single, double, or triple indirect block pointers [10].

For Ext4 it is more usual to find the use of extents, and they have their own structure. Table 2.7 demonstrates the extent header. Figure 2.14 illustrate how extents may be organised in an extent tree.

The 0xF30A (value interpreted as ) can be used as a signature to find inodes, which can be used to recover files and metadata. However, using the extents magic to carve for inodes will not deviate between extent headers found within the inode at inode offset 0x28 or extent headers found in data blocks used by an extent tree.

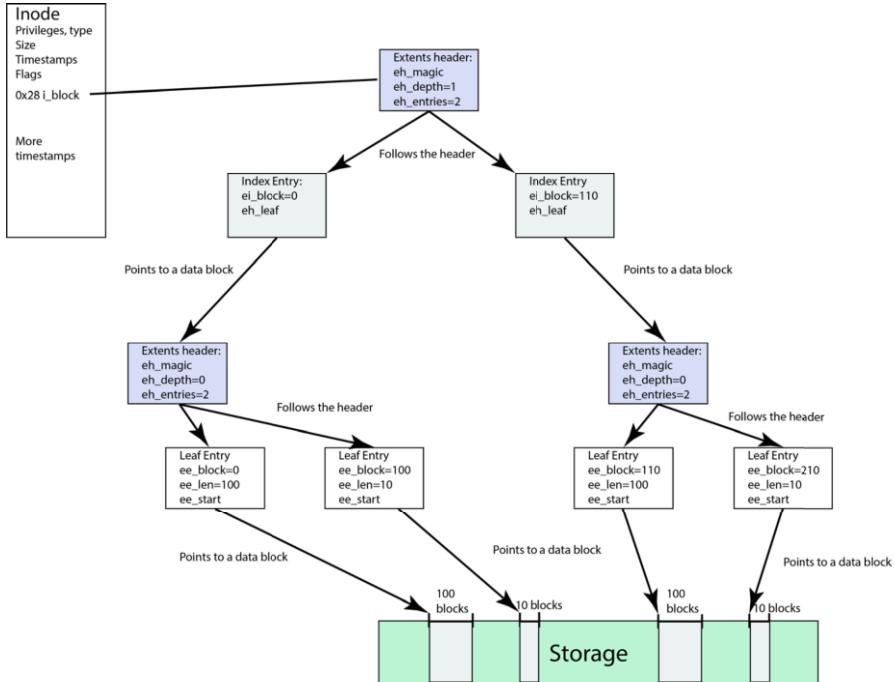


Fig. 2.14: Illustration of Extents in Ext4. Extent trees are not really needed with only 4 leaf extents, but more than 4 leaf extents will need a tree. We do not need two index entries, since the leafs could easily be included in one of the available 340 extents in one 4096 byte block. Each extent will describe one fragment. Ext4 tries to avoid fragments whenever possible, therefore, it is unusual to have many levels in the tree.

Table 2.7: Extent header

Offset	Size	Name	Description
0x0	0x2	eh_magic	Value 0xF30A
0x2	0x2	eh_entries	Number of extent entries
0x4	0x2	eh_max	Max number of entries
0x6	0x2	eh_depth	Depth of this extent node, 0= points to data.
0x08	0x4	eh_generation	Generation of the tree

This kind of recovery will not identify the file's name since the name is not included in the inode structure. However, some techniques can be used to connect the names found in directory entries and inodes found in the inode table [52].

After the header, the extent entries will follow. These are either extents pointing to new extent indexes or to the /data blocks containing the file content itself. If the extent header *eh\_depth* is larger than zero, the extent index entries will point

to blocks containing other extent entries. If the *eh\_depth* is equal to zero, then the extent describes and points to the blocks containing the file content.

An Extent will define a contiguous number of blocks, and if the file is fragmented, it will contain multiple extents. An extent header entry (*ext4\_extent\_idx*) have the following structure, which is necessary to parse in order to find all blocks that a file is using eventually.

Table 2.8: Extent index entry

Offset	Size	Name	Description
0x0	0x4	ei_block	Covers file blocks from block forward
0x4	0x4	eh_leaf_lo	Lower 32 bits of the block containing next extent node in the tree
0x8	0x2	eh_leaf_hi	Higher 16 bits of the block containing next extent node in the tree
0xA	0x2	eh_unused	Not in use

If the extent header has the depth 0, it will contain the leaf extent node (*ext4\_extent*), which describe the blocks used for file content.

Table 2.9: Extent leaf entry

Offset	Size	Name	Description
0x0	0x4	ee_block	First logical file block of this extent
0x4	0x2	ee_len	The length of the extent in blocks
0x6	0x2	ee_start_hi	Higher 16 bits of the extent physical start block
0x8	0x4	ee_start_lo	Lower 32 bits of the extent physical start block

The first block of a file will always start from logical block 0, but have a completely different physical disk location. The logical start block (*ee\_block*) is defined to the extent, which is necessary in order to organize the fragments correctly. The start of the first physical block, where this extent starts can be found in *ee\_start\_hi* and *ee\_start\_lo*. The extent contains the blocks from this location and contains the number of contiguous blocks defined in the length field (*ee\_len*).

After the last extent entry, there is a checksum named *eb\_checksum* which is computed by using the file system uuid (from the superblock)+inum (from the directory entry)+igeneration (from the inode)+extent block (not including the checksum). This checksum is not necessary since the inode is already checksummed [41]. This checksum can be computed using the crc32c algorithm to identify manipulation attempts [41]. A crc32c library can be used to test this checksum [19].

The 60 byte *i\_block* can also contain inline data, as long as the file system supports this and that the inode flag has defined that inline data is used. To create an Ext4 file system that supports inline data, it has to be formatted with the *mke2fs -O inline\_data*. This area can also be used to store small extended attributes added by using the *xattr*

tool. Large extended attributes will have a pointer to them (either direct or indirect block pointer, or an extent).

	Offset (h)	00 01 02 03 04 05 06 07	08 09 0A 0B 0C 0D 0E 0F	UTF-8	
0000	000D7A2100	ED 41 00 00 00 10 00 00	F0 88 5B 49 DB 0E BD 5E	I ^	User privileges, and type of file (here directory)
0010	000D7A2118	DB 0E BD 5E 00 00 00 00	00 00 1B 00 00 00 00 00		eh_magic (when extents are used)
0020	000D7A2120	00 00 00 00 04 00 00 00	0A F3 01 00 04 00 00 00		Number of extents (0x1)
0030	000D7A2130	00 00 00 00 00 00 00 00	01 00 00 00 02 1F 00 00		Depth of node (0x0 = points to data)
0040	000D7A2140	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		Lower 32 bits data block address (0x1F92)
0050	000D7A2150	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		
0060	000D7A2160	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		
0070	000D7A2170	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		
0080	000D7A2180	20 00 00 00 DC DB 26 B6	DC DB 26 B6 00 00 00 00	I	
0090	000D7A2190	F0 88 5B 49 00 00 00 00	00 00 00 00 00 00 00 00		
00A0	000D7A21A0	00 00 02 EA 07 06 40 00	00 00 00 00 1A 00 00 00	@	
00B0	000D7A21B0	00 00 00 00 73 65 6C 69	6E 75 78 00 00 00 00 00	selinux	
00C0	000D7A21C0	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		
00D0	000D7A21D0	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00		
00E0	000D7A21E0	00 00 00 00 75 3A 6F 62	6A 65 63 74 5F 72 3A 73	u:object_r:s0	
00F0	000D7A21F0	79 73 74 65 6D 5F 66 69	6C 65 3A 73 30 00 00 00	ystem_file:s0	

Fig. 2.15: Example of a directory inode.

An example of an inode is shown in Figure 2.15. This inode is the second element in the inode table. This means we are looking at inode number 2 (the root directory) [10, p413]. If we look at the first two bytes of this inode, we can see the value 0x41ED (LE). The four most significant bits are 0100, which means it is a directory (see Table 2.6). The extent tree starts in inode byte offset 0x28, starting with a header. The header includes the 0xF30A magic for extents, it contains one extent, and this extent is a leaf node (depth is 0). The generation field is not in use. Directly after the extent header, we find the only extent. It starts from logical block 0, has a length of 1 block, the location of the block is 0x1F92 (we do not need to think of the higher 16 bits address in the ee\_start\_hi since it is 0). Since every block is 0x1000 (4096 bytes long), we know the block can be found at byte 0x1F92000 relative to the start of the Ext4 partition. This is demonstrated in 2.5.

## File version

The *i\_generation* is meant to be used for NFS (Network File System) and a random value will be created for every new file created. This is described in the function *\_ext4\_new\_inode* found in the *alloc.c* file. Note that the Ext2 and Ext3 file system uses another approach where the generation is set on mount, and then it is increased with 1 for every file created using the *ext2\_new\_inode* or *ext3\_new\_inode* function.

The value is not guaranteed to be unique, and the author observes that multiple inodes may contain a zero value. If the value is not zero, and the creation date of the two inodes are equal and found in the same file system, then they are both describing the same inode. This must be considered when comparing inodes found outside the inode table, for instance, in the journal. Correlating and interpreting different fields in this way can be used to find all previous instances of the same inode, assuming they are not overwritten. If we have found an inode that obviously is deleted in the inode table, finding previous versions of this inode can give us the previous extents

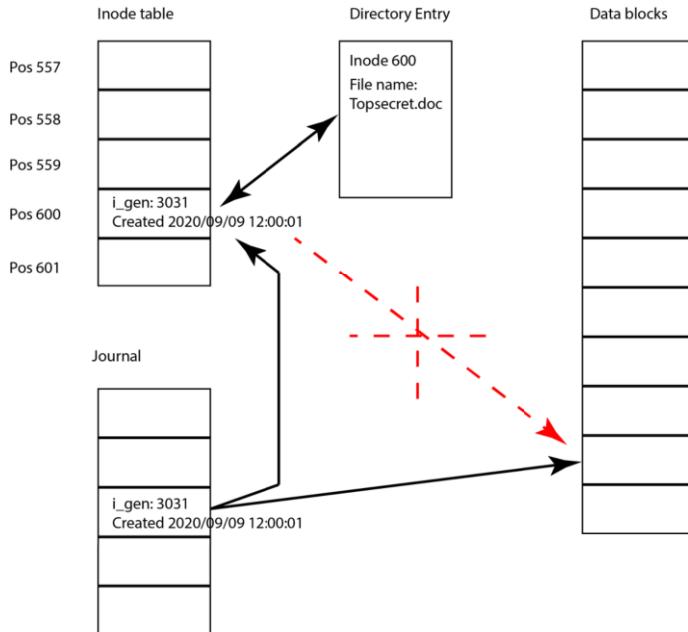


Fig. 2.16: Recovery of a file using a previous state of an inode.

in a non deleted state, allowing recovery of a previous state of this file. Knowing the position of the inode in the inode table allows us to search for the directory entry by parsing every directory. If the directory entry is found, it gives us the name and location in the file system directory tree, assuming that the directory entry is not overwritten. This recovery methodology for Ext4 is illustrated in Figure 2.16.

## Operating System Descriptor 2

The `osd_2` field has different content based on the OS used to create it, we describe this when Linux (Android) is used as the OS.

The operating system descriptors must be used together with similar fields described earlier in the inode. For instance, the higher value of the user id must be used together with the lower value of the user id, and so on. Most forensic tools show the owner or group of a file, but this can also be verified manually.

## Project ID

The field `i_projid` is used for creating children with the same project id. This can be used to define size quotas for a group of files; for example, setting how much space

Table 2.10: Os Descriptor 2 (Linux)

Offset	Size	Name	Description
0x0	0x2	<code>l_i_blocks_high</code>	Upper 16-bits of the block count
0x2	0x2	<code>l_i_file_acl_high</code>	Upper 16-bits of the extended attribute block
0x4	0x2	<code>l_i_uid_high</code>	Upper 16-bits of the Owner UID
0x6	0x2	<code>l_i_gid_high</code>	Upper 16-bits of the GID
0x8	0x2	<code>l_i_checksum_lo</code>	Lower 16-bits of the inode checksum
0xA	0x2	<code>l_i_reserved</code>	Unused

a user can save in the user directory. This requires that the superblock supports this feature (RO\_COMPAT\_PROJECT).

## 2.5 Ext4 - File Name category

Directory entries are important since they include the name of a file or directory, and contain the inode number of the file or directory. It is easy to find the byte location in the inode table by multiplying the inode number with 256. This requires that the investigator knows where the inode table starts, which we have shown can be located by scrutinizing the first group descriptor. The directory entry depends on one of the incompatible features for recording file types in directory entries defined in the superblock . Figure 2.3 on page 45 demonstrated directory entries record file types. Therefore, we need to use the following structure as defined in Table 2.11.

Table 2.11: Directory Entry

Offset	Size	Name	Description
0x0	0x4	<code>inode</code>	The inode this entry points to
0x4	0x2	<code>rec_len</code>	Length of this entry
0x6	0x1	<code>name_len</code>	Length of name
0x7	0x1	<code>file_type</code>	The file type of this entry
0x8	Var	<code>name[name_len]</code>	ASCII name of file

If the superblock does not define the use of recording file names in directory entries, then we use an almost identical structure. The only difference is that the `name_len` and the `file_type` is merged into a 2 byte field that describes the `name_len`.

Figure 2.17 depicts all the directory entries found in the root directory. The location can be found by scrutinizing the second inode in the inode table. From this location, we can find every allocated file and directory in the file system, and the deleted files in their directory entries have, if not overwritten. This is why some

	Offset (h)	00 01 02 03 04 05 06 07	08 09 0A 0B 0C 0D 0E 0F	UTF-8
0000	000F712000	02 00 00 00 00 00 01 02	25 00 00 00 02 00 00 00	.
0010	000F712018	0C 00 02 02 2E 00 00	0B 00 00 00 14 00 0A 02	..
0020	000F712020	6C 6F 73 74 2B 66 6F 75	6E 64 00 00 00 0C 00 00 00	lost+found
0030	000F712038	0C 00 03 02 61 70 70 00	6A 01 00 00 0C 00 00 03 02	app j
0040	000F712048	62 69 6E 00 03 02 00 00	14 00 0A 01 52 75 69 6C	bin buil
0050	000F712058	64 2E 70 72 6F 70 00 00	D4 02 00 00 04 14 00 0A 02	d.prop
0060	000F712060	63 61 6D 65 72 61 64 61	74 61 00 00 E3 02 00 00	cameradata
0070	000F712070	28 00 18 01 63 6F 70 70	61 74 69 62 59 6C 69 74	compatibilit
0080	000F712080	79 5F 6D 61 74 72 69 78	2E 78 6D 6C E4 02 00 00	y_matrix.xml
0090	000F712090	14 00 09 02 63 6F 6E 74	61 69 6E 65 72 00 00 00	container
00A0	000F7120A0	E7 02 00 00 0C 00 03 02	65 74 63 00 09 05 00 00	etc
00B0	000F7120B0	14 00 09 02 66 61 6B 65	2D 6C 69 62 73 00 00 00	fake-libs
00C0	000F7120C0	DB 05 00 00 14 00 0B 02	66 61 6B 65 2D 6C 69 62	fake-lib
00D0	000F7120D0	73 36 34 00 DD 05 00 00	18 00 05 02 66 6F 6E 74	s64 font
00E0	000F7120E0	73 00 00 00 12 07 00 00	14 00 09 02 66 72 61 6D	s fram
00F0	000F7120F0	65 77 6F 72 6B 00 00 00	Af 09 00 00 18 00 00 02	ework
0100	000F712100	68 69 64 64 65 6E 00 00	B4 09 00 00 14 00 0A 01	hidden
0110	000F712110	69 6E 66 6F 2E 65 78 74	72 61 00 00 35 09 00 00	info.extra
0120	000F712120	0C 00 03 02 6C 69 62 00	94 0C 00 00 10 00 05 02	lib lib
0130	000F712130	6C 69 62 36 34 00 00 00	BA 0F 00 00 14 00 00 01	lib64 lkm_sec_info
0140	000F712140	66 6B 6D 5F 73 65 63 5F	69 6E 66 6F 3B 0F 00 00	media F
0150	000F712150	18 00 05 02 6D 65 64 69	61 00 00 00 46 10 00 00	preload K
0160	000F712160	18 00 07 02 70 72 65 6C	6F 61 64 46 6F 74 61 4F	preloadota0
0170	000F712170	18 00 0F 02 70 72 65 6C	6F 61 64 46 6F 74 61 4F	nly N priv
0180	000F712180	6E 6C 79 04 10 00 00	18 00 08 02 70 72 67 76	-app prod
0190	000F712190	2D 61 70 70 09 12 00 00	18 00 07 02 70 72 6F 64	uct reco
01A0	000F7121A0	75 63 74 00 0B 12 00 00	1C 00 14 01 72 65 63 6F	very-from-boot.p
01B0	000F7121B0	76 65 72 79 2D 66 72 6F	6D 2D 62 6F 6F 74 2E 70	saivo
01C0	000F7121C0	0C 12 00 00 0C 00 04 02	73 61 69 76 4F 12 00 00	tima_measure
01D0	000F7121D0	28 00 15 01 74 69 6D 61	5F 6D 65 61 73 75 72 65	ment_info P
01E0	000F7121E0	6D 65 6E 74 5F 69 6E 66	6F 00 00 00 50 12 00 00	tts
01F0	000F7121F0	0C 00 03 02 74 74 00 00	83 12 00 00 0C 00 00 02	
0200	000F712200	75 73 72 00 58 13 00 00	18 00 06 02 76 65 6E 64	usr X vend
0210	000F712210	6F 72 00 00 48 16 00 00	18 00 10 02 76 6F 69 63	or @ voic
0220	000F712220	65 62 61 72 67 65 69 6E	64 61 74 61 6F 16 00 00	ebargeindatao
0230	000F712230	0C 00 04 02 78 62 69 6E	01 F7 01 00 0C 00 03 02	xbin
0240	000F712240	6F 6D 63 00 71 16 00 00	BC 0D 0C 07 63 73 63 5F	omc q csc_
0250	000F712250	63 6F 6E 74 65 6E 74 73	00 00 00 00 00 00 00 00	contents

Fig. 2.17: Content of root directory.

digital forensic tools show that a specific file is deleted, but the content may be harder to recover, discussed in sect. 2.6.

## 2.6 Ext4 - Content Category

The content of a file is pointed to by the inodes, as long as the file is allocated and not deleted.

### 2.6.1 Recovery of files

When a file gets deleted, the file extents header zero out the number of extents and depth of the tree. However, most of the extent entries may also be zeroed out. [14] shows an extent index that is not zeroed out after deletion, and he shows that extent leaves are zeroed out (except for the logical ei\_block). This means that recovery of an inode is most likely to succeed if the inode uses extent trees since it is possible

to parse down the tree to the leaf extent(s) that describes the block addresses used for data content [14]. The deleted timestamp is set to the time of deletion, and that many deleted files have been modified and changed equally to the deleted timestamp (Ext2, Ext3) [10, p. 420]. [14] also shows that on Ext4 the deleted inode's accessed, changed, and the modification time is set equal to the deletion time. On deletion, the file size, link count, and the number of blocks used by the file is zeroed out. However, it is possible to carve for file content only.

### Inode Carving using extent magic signature

Since some of the extent information is wiped, recovery of deleted data in Ext4 is more difficult than previous versions of ExtX. However, it is possible by performing carving for previous inodes or other metadata structures. An inode does not have a special static signature, even though it is possible to search for the *eh\_magic* if the inode uses extents [11]. [11] describe that the magic signature persist even for deleted files, and that they used the type field (*i\_mode*) 4 most significant bits to identify false positives not corresponding to one of the 7 different file types, as seen in Table 2.6 at page 57. Dewald and Seufert (2017) [11] also describes that it is possible to combine further pattern testing; for instance, specific timestamp intervals, or a set of access rights, in order to filter out even more false positives. However, this will not identify inodes that do not use extents. This approach for inode metadata carving is very well suited to identify and recover Ext4 inodes from a re-formatted partition (for instance, an Ext4 file system re-formatted to NTFS). [11] did not manage to connect the file names to the carved inodes when using the inode carving method.

### 2.6.2 Generic metadata time carving

Another approach for metadata carving is thinking that the timestamps near co-located could act as a dynamic signature based on equality. File systems have structures describing their files' metadata, and they usually have temporal information near co-located. We can use equality to identify a set of timestamps based on their known granularity. The scientific paper selected this approach Generic Metadata Time Carving [52]. Using this technique, it is possible to find all inodes that match the equality pattern, not only inodes using extents.

### 2.6.3 Additional file content

Even though different file types are part of the content category, they are not filesystem-specific, and therefore not included in this chapter.

## 2.7 Ext4 - Application Category

The Ext4 journal is used for recovery purposes when the file system becomes out of sync. Modern file systems often use journals. Depending on the flags in the superblock it can be used, but it does not need to. The journal is described as an application-level feature [10, p. 437]

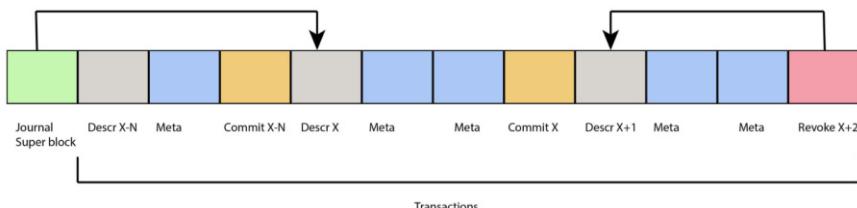


Fig. 2.18: Ext4 Journal transaction overview, based on an illustration from Carriers file system forensic analysis book [10, p. 438]

The first interesting part of the journal is the journal superblock, which contains a pointer to the first descriptor in the journal [14, 10]. The first descriptor may not be in the beginning of the journal because of the circular writing of transactions. When a transaction is stored at the end of the journal, the next transaction will be written at the start of the journal (overwriting previous transactions) [14, 10]. Every descriptor is followed by a set of metadata and/or data transactions and, finally, end with a commit block. If a file system crashes before the commit block is written, then the commit block is missing on the next mount of the file system, then a revoke block is created. This points to the previous descriptor, effectively undoing any of its transactions. Then the file system will be in a consistent state.

The journal is normally found in inode number 8, but can be placed in any other inode defined in the superblock. If the INCOMPAT\_JOURNAL\_DEV is set, the journal can be located on another device described by its UUID, defined in the superblock.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



# Chapter 3

## The Flash-Friendly File System (F2FS)



Chris Currier

**Abstract** The Flash-Friendly File System (F2FS) is used not just by removable media but also by mobile devices and more. In this chapter, we look under the hood to better understand the structure of and recognize this file system. From a forensic perspective, we look for deleted files to see if we can retrieve them.

### 3.1 Introduction

The Flash-Friendly File System (F2FS) is a Linux system specifically designed for NAND Flash memory. This type of memory is common in removable storage devices and mobile devices. Samsung Electronics developed the system in 2012. One thing to mention is that the resources for F2FS are minimal. Many books and other resources, even regarding forensic examination, barely mention F2FS. Due to its increasing importance in the field of mobile forensics, we want to address file system information important to investigators.

#### 3.1.1 NAND (Not And) Flash Memory

Universal Serial Bus (USB) flash drives (thumb drives), Solid State Drives (internal/external storage), SD Cards, and even mobile devices use NAND Flash memory. For a physical extraction, the chipset (flash memory) is what we are trying to get access to and obtain data from. This flash memory also contains a processor.

Depending on the internal geometry or flash memory management, aka Flash Translation Layer (FTL), NAND-based storage devices display different characteristics where parameters are added for configuring on-disk layout, allocation selection

---

Chris Currier

MSAB, Hornsbruksgatan 28 SE-117 34 Stockholm Sweden e-mail: [chris.currier@msab.com](mailto:chris.currier@msab.com)

and algorithms for cleaning [42]. NAND is not an acronym. It stands for ‘NOT AND’. It is a Boolean operator and logic gate. Both NAND and NOR gates are depicted in Fig. 3.1.



Fig. 3.1: NAND and NOR Gates

The NAND  $(!(A \wedge B))$  logic yields FALSE when both input values (A and B) are True and yields TRUE if any input value is False. In contrast, the NOR  $(!(A \vee B))$  operator yields TRUE if both input values (A and B) are False and yields FALSE if any input value is True (see Table 3.1 below).

Table 3.1: NAND and NOR Gate Logic

A	B	$!(A \wedge B)$	$!(A \vee B)$
T	T	F	F
T	F	T	F
F	T	T	F
F	F	T	T

We can further compare NAND and NOR flash memory [88]. The differences are as follows:

### NAND flash memory

- contains an integrated circuit that uses NAND gates to store data in memory cells.
- devices write and erase data faster.
- devices store more data than NOR flash memory of the same physical size.

### NOR flash memory

- uses NOR gates to store data in memory cells.
- devices write data slower.
- devices read data faster.
- data storage is less efficient.

Memory cells of flash memory can store more than one bit per cell across different voltages have a significantly limited lifetime of around 10000 write cycles. This necessitates an even distribution (wear-levelling) of the write operations over the

entire flash memory. That is why flash mass storage is given an abstraction layer by its controller, the Flash Translation Layer (FTL).

### 3.1.2 Flash Translation Layer (FTL)

A flash translation layer is located in the controller of flash memory. It is responsible for the actual use of the memory. In doing so, it has to master a whole range of tasks:

*“Unlike jffs2 and logfs, f2fs is not targeted at raw flash devices, but rather at specific hardware that is commonly available to consumers – SSDs, eMMC, SD cards, and other flash storage with FTL (flash translation layer) already built in.” [8]*

This includes relying on the FTL for the wear levelling. Meaning that writing to the storage media is done evenly and not just to the first cells. Constant writing and rewrites to just the first cells of this flash memory would eventually corrupt the media. The FTL is a combination of hardware and software that can perform a number of central tasks for memory use through this interaction. It essentially ensures that writes are distributed evenly across the memory. This significantly increases the lifespan of an SD card. However, the FTL offers a conventional block device interface. It does not care about the erase-before-write property of a NAND flash device. Flash write-only can write zeros. And Flash erase can write the ones. Flash erase sets all bits to 1, so the flash write can leave the bit alone or switch it to 0. Because of this, in addition to FTL, special file systems developed for flash memory such as JFFS, Yaffs and Log FS are used further to increase the memory cells’ lifetime and better address the erase-write problem.

## 3.2 Flash Filesystems

After taking a brief look at the hardware basics, we will now turn to the actual file system. File systems Log FS take the special properties of SD cards into account and operate as log-structured file systems. They write data sequentially to the flash memory, similar to a cyclic logbook, thus ensuring that all cells are used evenly. However, these file systems are exotic because they have an unfavourable side effect: Data and metadata end up sequentially in multiple versions on the storage medium. It is the task of an elaborate and relatively slow garbage collection to remove obsolete and deleted data from the log.

The F2FS file system addresses this problem. A compromise is made: It structures data for write operations like a log-structured file system in sequential series that are as long as possible but leaves it to the flash translation layer to eliminate the redundancies.

### 3.2.1 The Log-Structured File System (LSFS) or (LFS)

As already discussed, NAND flash devices can have different characteristics depending on their internal geometry, and the flash management scheme (FTL) used. In order to meet these, the new file system has several parameters with which it can be optimally adjusted to the respective memory.

F2FS is based on the Log-Structured File System (LSFS). This structure uses a block structure. The blocks are then written with data (files). The block/data is mapped using index nodes referred to as *inodes*. When data in a block is updated, the inode needs to be updated. Responsible for holding the location of the inodes is the imap. The imap will have a 4-byte entry(or pointer) for each inode. The imap will be written at the end. In the simplified example (see Fig. 3.2 below), we can see Block 0 and Block 1. The pointer for the files here is in INODE A. New Data is added to Block 2, and an INODE B is created to point to that data. The imap is updated and at the end to point to INODE A and B.

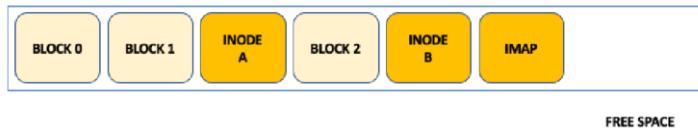


Fig. 3.2: Simple representation of the Blocks in a Log-Structured File System (LSFS)

If we have a single file split up into multiple sections spanning across different areas in a file system, this is fragmentation. Fragmentation causes issues with speed and more. We want the file to be complete in the same physical area. After data has been deleted, the Log-Structured File System takes the live data and brings them together in sections, updating the inode(s) and imap:

*“A log structured file system writes all modifications to disk sequentially in a log-like structure, thereby speeding up both file writing and crash recovery. The log is the only structure on disk; it contains indexing information so that files can be read back from the log efficiently.”* [75]

### 3.2.2 Flash-Friendly File System (F2FS): Enter F2FS

The F2FS file system is a bit more complex than the basic diagram (Fig. 3.2) and information in sect. 3.2.1 with regards to the Log-Structured File System described above. Similarly to other filesystems, F2FS is comprised of blocks; each block is 4K in size. Although, “the code implicitly links the block size with the system page size.” [8]

F2FS block addresses are 32 bits. [8] records that “*the total number of addressable bytes in the file system is at most  $2^{(32+12)}$  bytes or 16 terabytes*”. The author acknowledges that this is unlikely to be a limitation for current flash hardware.

The name ‘Flash-Friendly File System’ derives from its design i.e., a filesystem that is designed for the NAND flash memory-based storage. [42] documents that a log structure file system approach was adopted while adapting to newer forms of storage. In addition, the filesystem was designed to fix some issues of the aging log-structured file system, such as the snowball effect of the wandering tree and the heavy cleaning workload.

### 3.2.3 Wandering Tree Problem

The Wandering Tree issue for the Log-Structured File System is that there are so many pieces when updating. When file data in LFS is updated and written to the end of log, there are several things to consider:

- its location has changed, and its direct pointer block must be updated.
- as a consequence of the change to the direct pointer block, its indirect pointer block must also be updated.
- upper index structures must recursively be updated (e.g. inode, inode map, and checkpoint block).

Bitutskiy cited in [3] describes this as the wandering tree problem and “*in order to enhance the performance, it should eliminate or relax the update propagation as much as possible*”.

## 3.3 On-Disk Layout of F2FS

A classic hard disk stores information through remanence (remaining magnetization), unlike flash memory. Rotating circular disks are used for storage. In order to locate a memory cell on such a medium, it is divided into different areas. The geometry of a hard disk is the division of the hard disk into tracks and sectors. It is essential first to introduce terms such as *sector* and *partition*. The former has to do with the geometry of a block device. The second term is aimed more at the logical management of a hard disk.

### Sector

The term “Sector” refers to the physical sector or location on a physical disk. If you think back to the mechanical hard drives that had a platter it was divided into sectors or physical blocks (see Fig. 3.3). The starting position is Sector 0.

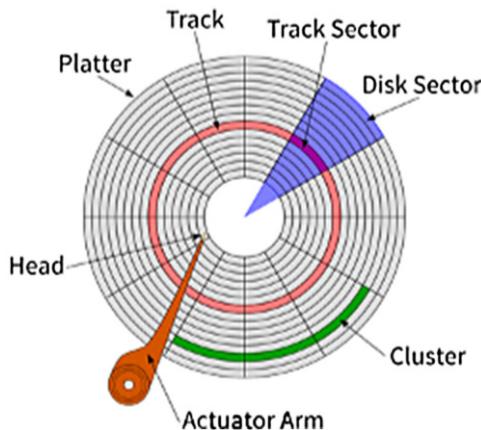


Fig. 3.3: A Breakdown Representation Inside a Hard Disk Drive (HDD) [53]

## Partitions

When you format a hard disk drive or flash memory; you prepare it with a file system. In doing so you may have a single partition or create multiple partitions such as a system partition, recovery partition, and/or a user partition. You may also create additional partitions for different file systems i.e. HFS, NTFS, EXT, etc. Fig. 3.4 is an example from Microsoft's Disk Management showing multiple partitions for a single hard disk drive.

Disk 0	300 MB Healthy (Recovery Partition)	512 MB Healthy (EFI System Partition)	OSDisk (C:) 475.35 GB NTFS (BitLocker Encrypted) Healthy (Boot, Page File, Crash Dump, Basic Data Partition)
Basic 476.81 GB Online			

Fig. 3.4: Disk Partitioned and on the C: Volume formatted with the New Technology File System (NTFS)

### Important

NAND flash memory-based storage devices have different characteristic according to their internal geometry and compared to a traditional hard disk, which stores data on rotating disks using a magnetic record.

### 3.3.1 Creation of F2FS partitions with `Mkfs.f2fs`

F2FS file systems are usually created with a special tool called *Mkfs.f2fs*. It can be used to create a F2FS file system (usually in a disk partition).

*“The `mkfs.f2fs` is for the use of formatting a partition as the f2fs filesystem, which builds a basic on-disk layout.”* [89]

It is normally operated from the command line. The most important parameters are summarized in Table 3.2. If you prefer a graphical user interface, you can alternatively use the gparted program under Linux to create an F2FS partition (see Fig. 3.5).

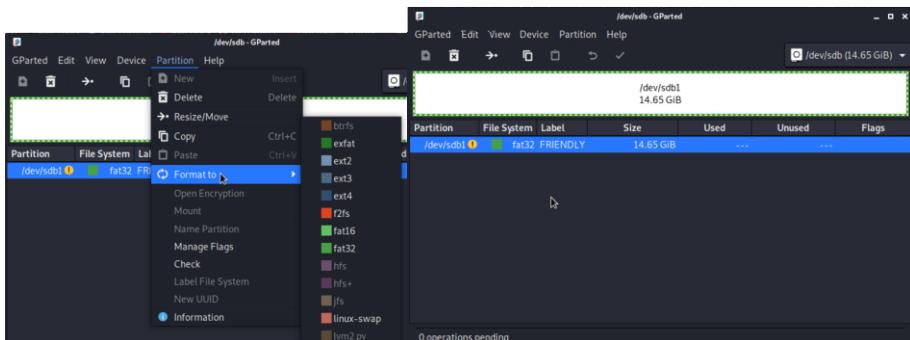


Fig. 3.5: Using GParted in Kali Linux to Format the USB Flash Drive to F2Fs

Table 3.2: `Mkfs.f2fs` Command Options

Command Option	Description
-l [label]	Give a volume label, up to 512 unicode name.
-a [0 or 1]	Split start location of each area for heap-based allocation. 1 is set by default, which performs this.
-o [int]	Set overprovision ratio in percent over volume size. 5 is set by default.
-s [int]	Set the number of segments per section. 1 is set by default.
-z [int]	Set the number of sections per zone. 1 is set by default.
-e [str]	Set basic extension list. e.g. “mp3,gif,mov”
-t [0 or 1]	Disable discard command or not. 1 is set by default, which conducts discard.

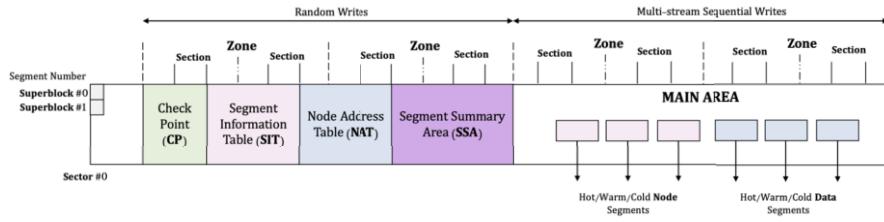


Fig. 3.6: Flash-Friendly File System Representation of How it appears Physically on the Disk [45]

### 3.3.2 F2FS on Disk

The F2FS is split into blocks that are 4K in size. Blocks are collected into *segments*. A segment is 512 blocks or 2MB in size. Each *section* is comprised of several consecutive segments. A *zone* is comprised of a series or set of sections. An *area* is comprised of multiple sections. The default size when using the `mkfs` utility is  $2^0$ . Hence, there is one segment per section. A *volume* is comprised of six areas. The structure is depicted in Fig. 3.6. As mentioned, F2FS is split into six areas in total. Each is briefly described below and further discussed in this chapter:

- **Superblock (SB)**: holds the partition information and F2FS parameters; it is unchangeable.
- **Check Point (CP)**: represents the file system status; bitmaps for SIT and NAT; orphan inode list; summary entries of the active segment.
- **Segment Information Table (SIT)**: contains the valid segments and bitmap information in the Main Area.
- **Node Address Table (NAT)**: a block address table.
- **Segment Summary Area (SSA)**: summary entries representing the owner information including parent inode number and node/data offset.
- **Main Area**: node blocks store indices of data blocks; a data block contains directory or user file data.

As pointed out, F2FS divides the drive into six different, consecutive areas. At the beginning of the partition, we find the Superblock (SB). This is followed by a second copy of the Superblock. It is used if the first Superblock becomes corrupt. The Checkpoint (CP) region follows this. This region contains, among other things, information on the active segments and orphaned or expired nodes. Next comes the segment information table (SIT). It provides information about the blocks stored in the main area and their status (active or inactive). It is in turn followed by the Node Address Table (NAT), which can query the addresses of the respective active nodes. The following Segment Summary Area (SSA) provides information about which node owns which blocks. The first five blocks thus represent the metadata of the partition. The Main Area (MA), as the sixth region, contains the actual data blocks. Next comes the segment information table (SIT). It provides information

about the blocks stored in the main area and their status (active or inactive). It is in turn followed by the Node Address Table (NAT), which can query the addresses of the respective active nodes. The following Segment Summary Area (SSA) provides information about which node owns which blocks. The first five blocks thus represent the metadata of the partition. The Main Area as the sixth region contains the actual data blocks with the files and directories.

## Superblock

For the F2FS file system, the start of the logical partition is the Superblock. Fig. 3.6 identifies that there is superblock 0 and 1 that is in place as a redundancy in case there is a failure. Like other terms (inode, dentry, etc.) in this chapter, Superblock is based on Unix and not unique to F2FS. The Superblock (SB) [45]:

- is located at the beginning of the partitions.
- two copies exist, as redundancy for failure.
- includes basic partition information.
- includes several default parameters of F2FS.

The most important data fields of the Superblock are shown in Fig. 3.7. Like many binary formats, the Superblock starts with a Magic (Header). An example taken from a *Huawei P9* is given in Fig. 3.8. In this case, it is Hex 1020F5F2.

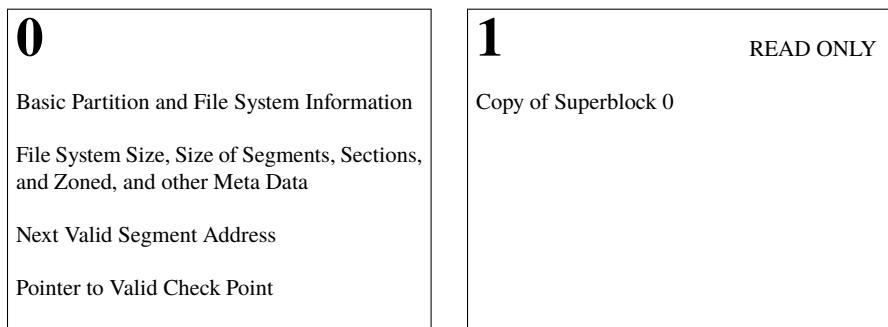


Fig. 3.7: F2FS Superblock the Starting Point and Backup Copy

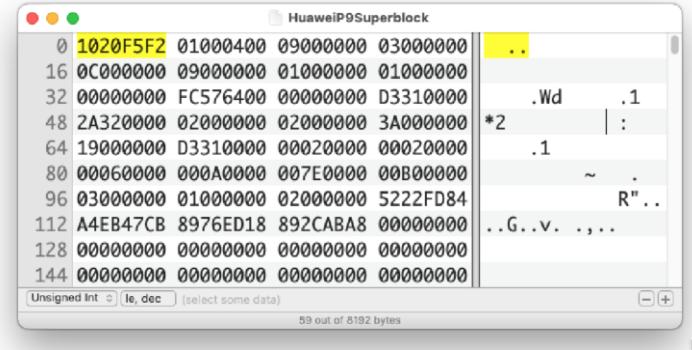


Fig. 3.8: Huawei P9 Extraction Showing the Start of a Superblock

## Zone

A zone contains several sections for easier management. The default number of sections in a zone is 1, but there may be any number of sections in a zone. The purpose of zones is to separate into different parts of the device the six open sections in the device. Flash drives are often made from a number of sub-devices. Each sub-device can process Input/Output (IO) requests. These requests can be processed independently and hence processed in parallel. Therefore, the six open sections can process requests and write in parallel [8]. One of the issues with NAND Flash memory is writing to an erase block first. F2FS uses zones, and each zone has its own erase block.

## Section and Segment

Fragmentation is an issue any file system wants to avoid. F2FS uses sections to organize and keep blocks together in segments. 512 blocks make up a segment (2MB). These segments contain such things as Checkpoint, Tables, and the Main Area.

## Check Point (CP)

If you have used and installed Microsoft Windows, then you may have seen the term restore point before. The idea is you can have a restore point in case an application install goes wrong, or there is some other issue. F2FS also has a built-in feature to manage this called the Checkpoint. The Checkpoint also has NAT and SIT Journaling, which will be discussed in the Cleaning section (see 3.4.6).

0	LATEST STABLE VERSION	1	LAST STABLE VERSION
	File System Status BitMaps: Valid NAT & Valid SIT Orphaned inode Lists Active Segments		File System Status BitMaps: Valid NAT & Valid SIT Orphaned inode Lists Active Segments

Fig. 3.9: F2FS Checkpoint Current and Last Stable Versions

### Segment Information Table (SIT)

The Segment Info Table (SIT) assists in identifying blocks that are in use “Valid” and those that are “Invalid” i.e. containing deleted data and may be cleaned. The SIT also tracks when a segment is empty of valid blocks and can be reassigned with live data.

### Node Address Table (NAT)

The Node Address Table is for addressing the Main Area node blocks. The structure of the NAT contains the latest version, the inode number and the block address. There are three types of nodes: *inode*, *direct node*, and *indirect node*. Table 3.3 depicts the concept of the Node Address Table (NAT), which is used to read from the device. Each unique node is assigned a node ID (see sect. 3.4.1), which is recorded in the table, along with the physical on-disk location (block address).

Table 3.3: NAT Example Table

node ID		block address
0		addr0
...		addr...
, N		addrN

#### ➤ Important

The term *inode* stands for *index node*. This forms the basic data structure for managing file systems with Unix-like operating systems. Each node is uniquely identified within a file system by its *inode number*. Each name entry in a directory refers to precisely one inode. This contains the file’s metadata and refers to the data of the file or the file list of the directory.

## Segment Summary Area (SSA)

Like the NAT, the Segment Summary Area (SSA) is concerned with the Main Area portion. The area deals with “Valid” Blocks of data in the Main Area. As was mentioned with the Log-Structured File System, when data can be removed, i.e. cleaned, this will probably cause fragmentation.

The valid blocks that may now be fragmented from each other can be copied and moved so they are all together. Speeding up the process to get to the data as it is all together and not spread across the drive in different locations.

*“The Segment Summary Area (SSA) stores summary entries representing the owner information of all blocks in the Main area, such as parent inode number and its node/data offsets. The SSA entries identify parent node blocks before migrating valid blocks during cleaning”* [45].

## Updates to the SIT and NAT

When data is updated it is not until a new check point is created that the changes are made to the Node Address Table(NAT) and Segment Info Table (SIT). Until this occurs, the updated information is;

- held in memory.
- if only a few updates, they can be written into Segment Summary blocks.
- updated info is written into the Checkpoint block for when the checkpoint is created.

## Shadow Copy

If you have been doing computer forensics, you are probably aware of, or at least heard, the term shadow copy or Volume Shadow Copy. This saves data by creating a snapshot as a safety net. F2FS does look for and use the last valid checkpoint. There are two Checkpoints. One is for identifying the most recent live or valid data. The second one is the shadow copy. Both the NAT and SIT also use shadow copies.

## Main Area

The Main Area is where the blocks that contain file data are located. As F2FS uses different sections, it allows for the data (e.g. directory or file content) to be kept separate from the node (e.g. the indexing information) [8]. The six active logs in the Main Area are managed using the following temperature-based categorisations, which are based on several strategies (Fig 3.4, according to [45]):

Each block in the Main Area is 4KB, and each is allocated by its type: data or node. In the Main Area, there are three data blocks and three node blocks. Data

blocks contain either a directory or user file data whereas a node block contains either an inode or indices of the data blocks. Data/node blocks cannot be stored in sections at the same time [45]. F2FS implements a search functionality i.e., a file *lookup operation* using the following set of steps outlined by [45], which assumes the file `/dir/file`. Fig. 3.10 and Table 3.5 identify the steps for the F2FS’s lookup operation.

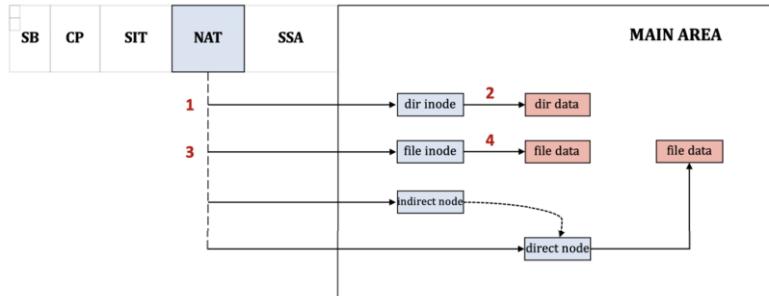


Fig. 3.10: F2FS File Lookup

Thus, we traverse through the file system tree with every file request. The starting point of our search is always Node Address Table.

## 3.4 File Structure of F2FS

### 3.4.1 Node Structure

File systems like the Log-Structured FS (LFS) use the index node (to identify the physical location of nodes), one large log, and updating direct and indirect nodes caused issues, such as the Wandering Tree (see sect. 3.2.3).

F2FS uses the Node Allocation Table (NAT) for finding the physical location of the nodes. The node blocks themselves have a Node ID. Following is a breakdown of the three types of node blocks, as recorded by The Linux Kernel Archives [89]:

Table 3.4: Temperature-based Categorisations of Main Area Blocks

Temperature	Node Block	Data Block
Hot	Direct node blocks of directories	Dentry blocks
Warm	Direct node blocks except those allocated as ‘Hot’	Data blocks except those allocated as ‘Hot’ and ‘Cold’
Cold	Indirect node blocks	Multimedia data or migrated data blocks

Table 3.5: F2FS File Lookup Operation

Step	Description
1	A block is read to obtain root inode. The block location is collected from the Node Address Table (NAT).
2	Searches for a directory entry- <code>dir</code> - inside the root inode block from the data blocks. The corresponding inode number for the directory is obtained.
3	The inode number is translated into a physical location. This location is obtained using the Node Address Table (NAT).
4	The inode named <code>dir</code> is collected by reading the corresponding block.
5	The directory entry named <code>file</code> is identified in the <code>dir</code> inode. The inode for file is translated into a physical location and the corresponding block is read to obtain the inode of <code>file</code> . The data stored in the Main Area, and various indices from the corresponding file structure, can then be retrieved.

- **inode:** 4KB assigned to each inode block. Each comprises of 923 data block indices.
- **direct node:** There are two direct node pointers.
- **indirect node:** There are two indirect node pointers and one double indirect node pointer.

Whether it is a direct or indirect node: In both cases, these contain references to 1018 data blocks [89]. The NAT is used by F2FS to map all node blocks using translation. The pointer-based file indexing system, which uses both direct and indirect node blocks in addition to the Node Address Table, is considered to prevent the ‘wandering tree’ problem [45, 89]. Unlike traditional LFS design, F2FS avoids the problem by updating a single direct node block and its corresponding entry in the Node Address Table. This update process prevents the “the propagation of node updates caused by leaf data writes” [89]. This is dissimilar to the traditional LFS design where both direct and indirect pointer blocks are updated recursively and cause a snowball/chain (i.e. wandering tree) effect, which is inefficient [45]. The comparisons are documented in Table 3.6.

Table 3.6: Comparison of an Updated File Between LSFS and F2FS

Description	LSFS	F2FS
Data is Updated	Direct and Indirect pointer blocks are updated recursively.	Only updates one direct node block and its NAT entry.
If the file is larger than 4GB	Updates one more pointer block for a total of three.	Still updates only one.

The **Index Node (inode) Block** does not have the physical address for a file. Instead, it has the points to direct or indirect pointers with the node number. Fig. 3.11

demonstrates the structure of the inode block. The figure also depicts the use of several pointers in an inode block:

- direct pointers to the file's data blocks.
- two single indirect pointers.
- two double-indirect pointers.
- one triple-indirect pointer.

F2FS also reserves 200 bytes in an inode block to store extended attributes. If a file is very small, it can be saved directly in the inode. This procedure is also called in-lining. In this case, however, the file size must be less than 3,692 bytes. We can also inline extended attributes [45].

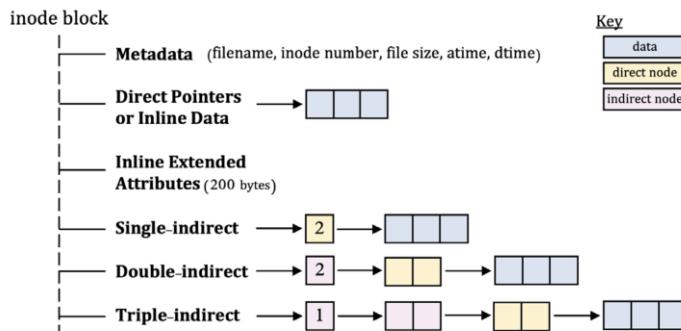


Fig. 3.11: Index Node (inode) Block

A **Direct Node Block** has the physical block address of the file and is updated when a file is updated. This Direct Node address is updated in the Node Address Table (NAT). When a file is updated, the **Indirect Node Block** is not, since it does not have the physical address. Instead, indirect node blocks hold identifiers (node IDs) that locate another node block, following the pointer-based structure.

### 3.4.2 File Creation and Management

File systems are different and use different ways of managing file locations. In FAT32, a file name is altered when deleted, replacing the first character with a hex *E5* character. The filename itself is not part of the file but rather stored as a new directory entry. At which the front of a real Library, you go to the card to look up the book name or author. The card points you to where you will find the book. F2FS has directory entries also, and these are called *dentries*. Like the Library analogy and other file systems, the system has the file information, including the inode number.

## Directory Structure

Directory Entry or *Dentry* keeps track of the index nodes (inodes) and occupies 11 bytes. A dentry contains a bitmap and two arrays of slots and names (see Fig. 3.12). A bitmap entry identifies if a slot is “Valid”. Each slot includes the (1) hash value of the file name (1 byte), (2) inode number (4 bytes), (3) length of the file name (4 bytes) and (4) file type (directory, symlink, regular file . . . - 1 byte).

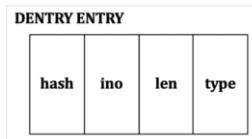


Fig. 3.12: F2FS Dentry Entry Structure

To manage a large volume of dentries and improve efficacy, multi-level hash tables are utilised by the F2FS file system. Each level has a hash table with a dedicated number of hash buckets. Several steps occur when F2FS looks up file names in a directory (see Table 3.7).

Table 3.7: F2FS Multi-level Hash Tables

Step	Description
1	Calculates the hash value of the file name.
2	Traverses the hash tables incrementally starting at level 0 until it reaches the maximum level which has been allocated and recorded in the inode.
3	Scans one bucket at each level (level 0 through to level $N$ ), resulting in an $O(\log(\# \text{ of } \text{dentries}))$ complexity.
4	For speed and efficacy, F2FS compare the bitmap, hash value and file name to find a dentry.

In addition, for example, there is a requirement for larger directories in server environments. The F2FS file system can be configured in the first instance to allocate space for many dentries [45]. “With a larger hash table at low levels, F2FS reaches to a target dentry more quickly.” [45] A bucket (see Fig. 3.13) consists of two or four dentry blocks.

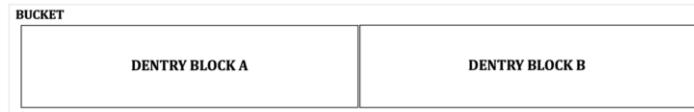


Fig. 3.13: F2FS Bucket Structure

A dentry block consists of 214 dentry slots and file names. In order to determine whether a dentry is valid, a bitmap is used again. Due to the described properties, a dentry block is always exactly 4 KB in size. This value is determined as follows:

```
Dentry Block (4 K) =
bitmap (27 bytes) + reserved (3 bytes) +
dentries (11 * 214 bytes) + file name (8 * 214 bytes).
```

Fig. 3.14 depicts the structure of a dentry block. To clarify, deleted directories and entries can be recognised because they are marked as invalid in the bitmap. The dentry concerned is thus free and can be used otherwise.

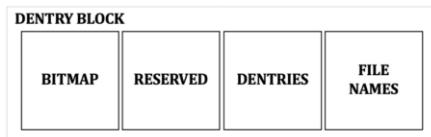


Fig. 3.14: F2FS Dentry Block Structure

### 3.4.3 Fsck.f2fs Identifying Files

In Kali Linux, using Fsck.f2fs on the USB Flash Drive, we could see the folder and files that were not deleted. Note the inode identifiers in bold. Any deleted folder or deleted files are not displayed. You should start to see terms that you are now familiar with. In the example above, the file *pngpicture.png* was found with the inode id 0x6. The Fsck.f2fs tool automatically scans for file system errors and corrects them [94]. The following listing shows an example of the output of the programme for a flash stick:

```
(ccurrier@kali)$ sudo fsck.f2fs -t /dev/sdb1
Info: [/dev/sdb1] Disk Model: Flash Disk
Info: Segments per section = 1
Info: Sections per zone = 1
Info: sector size = 512
Info: total sectors = 30717952 (14999 MB)
Info: MKFS version
      "Linux version 5.10.0-kali3-amd64 (devel@kali.org) (gcc-10 (Debian 10.2.1-6)
          10.2.1-20210110"
```

```

...
Info: superblock features = 0 :
Info: superblock encrypt level = 0, salt = 000000000000000000000000000000000000000000000000000000000000000
Info: total FS sectors = 30717952 (14999 MB)
Info: CKPT version = 373c4953
Info: Checked valid nat_bits in checkpoint
Info: checkpoint state = c5 : nat_bits crc compacted_summary unmount
|-- folder <ino = 0x4>, <encrypted (0)>
|  |-- pngpicture.png <ino = 0x6>, <encrypted (0)>
|  '-- textstays <ino = 0xa>, <encrypted (0)>
|-- dump_sit <ino = 0x5>, <encrypted (0)>
[FSCK] Unreachable nat entries [Ok..] [0x0]
[FSCK] SIT valid block bitmap checking [Ok..]
[FSCK] Hard link checking for regular file [Ok..] [0x0]
[FSCK] valid_block_count matching with CP [Ok..] [0xa0]
[FSCK] valid_node_count matching with CP (de lookup) [Ok..] [0x5]
[FSCK] valid_node_count matching with CP (nat lookup) [Ok..] [0x5]
[FSCK] valid_inode_count matched with CP [Ok..] [0x5]
[FSCK] free segment_count matched with CP [Ok..] [0x1d0f]
[FSCK] next block offset is free [Ok..]
[FSCK] fixing SIT types
[FSCK] other corrupted bugs [Ok..]
Done: 2.829615 secs

```

### 3.4.4 Metadata

The term metadata should be familiar to forensic examiners, usually referred to simply as data about data. Consider the properties of a file. In F2FS, three types of nodes are used that hold information about actual files.

There are index nodes referred to as inodes , direct nodes and indirect nodes. An inode consists of forensically important information, such as file size, allocated blocks, ownership (e.g., UID and GID) and Modified, Accessed and Changed (MAC) times [95]. These three timestamps are important and can tell us the following information when the file or directory is:

- Modified: updated when the file or directory is written.
- Accessed: updated when the file or directory is read.
- Changed: updated when the inode is modified.

The MAC-timestamps are all time specifications in *ms* (Unix epoch) and can easily be converted to a readable format with an appropriate converter. In the below .png file example, which we have already seen briefly in the last section, you can see some of the metadata for the file to include timestamps, file name, and file size.

```

[print_node_info: 353] Node ID [0x6:6] is inode
i_mode          [0x     81a4 : 33188]
i_links         [0x      1 : 1]
i_size          [0x    3c41f : 246815]
i_blocks        [0x      3e : 62]
i_atime         [0x606db122 : 1617801506]
i_atime_nsec   [0x24990b38 : 614009656]
i_ctime         [0x606db114 : 1617801492]
i_ctime_nsec   [0x2bc016f9 : 734009081]

```

i_mtime	[0x606b7024 : 1617653796]
i_mtime_nsec	[0x ad1c7fa : 181520378]
i_generation	[0xd0856627 : 3498403367]
i_namelen	[0x e : 14]
i_name	[pngpicture.png]

### 3.4.5 Multi-Head Logging

What is Multi-Head logging, and what does it have to do with Hot/Warm/Cold Data? The historical origin of the terms "hot" and "cold" goes back to the different data storage devices and their vibration. Hot data was located near the heat of spinning drives and CPUs, while cold data was stored on a tape or drive far from the data centre floor.

Wait, what is Hot/Warm/Cold Data? This involves the frequency of writes to the Main Area to both the Node and Data segments (see Fig.3.15). In the case of F2FS, the blocks that are updated particularly frequently are designated as warm. The direct blocks reference a block by its actual physical address on the disk. On the other hand, indirect blocks are assigned to the "cold" category. This is because they have only logical NodeID. This logical address must be adjusted or changed much less frequently during updates. Again, we distinguish between node blocks (containing linking and meta information) and data blocks (the actual data content). The Log-Structured File System uses a single log area, while F2FS splits this into six types of logs in the Main area split between node segments and data. Refer back to Table 3.4 for descriptions of Hot, Warm and Cold Node and Data Segments. See Table 3.8 for further descriptions.

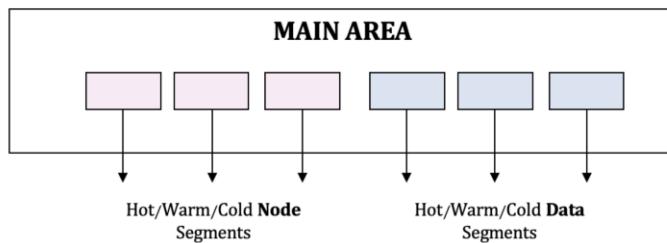


Fig. 3.15: Node and Data Segments in the Main Area [45]

In looking at F2FS resources, you will find benchmark testing that shows 2, 4, and 6 logs utilized for writing. You may be thinking of a standard log file, but this refers to multiple writing schemes. Six logs are the default value.

Table 3.8: Additional description of why Block Types are assigned a certain temperature [45]

Description		LSFS	F2FS
Direct Node Blocks	Hot	Updated frequently as they have the physical address.	
Indirect Node Blocks	Cold	Has only the node number and are only created or updated when a dedicated node block is added or removed.	
Directory Direct Node and Data Blocks	Hot	Different write patterns compared to blocks for regular files.	
Data Blocks	Cold	Valid for extended period of time.	

### 3.4.6 Cleaning

Deleted data is never meant to be stored permanently, no matter the operating system. This data that is once allocated goes unallocated. This data is available to be overwritten.

However, the F2FS file system allows the system to be cleaned. The idea behind this is performance. So that when data is deleted, the system can clean and defragment the live or valid data. Cleaning could be initiated by the user in the device settings, initiated due to a lack of free sections, or part of a regular background cleaning. The cleaning process is divided into three steps (see Table 3.9). First of all, the section to clean must be selected. Different selection strategies are used for this purpose. The selection of the target section is followed by identifying the invalid blocks. For this, the bitmap can inside the SIT is used. The cleanup process ends with the creation of a new checkpoint to free the blocks that have been released in the meantime for reallocation.

### Adaptive Logging

You just read about the two, four, and six multi-head logging options to handle different types of data based on their frequency of writes. The Adaptive logging has to do with where the writes will occur and possibly involve the cleaning process in F2FS. The Log Structured File System used two logging features:

- Normal: Uses clean segments, and the data is written in order.
- Threaded: Looks for invalid data areas to write data to.

So F2FS is using a dynamic approach with Adaptive Logging. There are two strategies based on the presented circumstances. See Table 3.10 and the illustration (Fig. 3.16) below.

Table 3.9: Cleaning Process: Three Steps [45]

Step	Title	Description
1	Victim Selection	Identify a victim section among non-empty sections. Two policies: <ul style="list-style-type: none"> <li>• Greedy Policy – foreground cleaning to minimize the latency visible to applications. Selecting section with the smallest number of valid blocks. Migrating valid blocks.</li> <li>• Cost Benefit Policy – Selects victim section not only on its utilization but also its “age”*.</li> </ul>
2	Valid Block Identification & Migration	Victim selected. Need to identify valid blocks in the section quickly. A validity bitmap per segment is in the SIT. F2FS retrieves parent node blocks containing their indices from the SSA information. If the blocks are valid, F2FS migrates them to other free logs. For background cleaning, F2FS does not issue actual I/Os to migrate valid blocks. Instead, F2FS loads the blocks into page cache and marks them as dirty. Then, F2FS just leaves them in the page cache for the kernel worker thread to flush them to the storage later.**
3	Post-Cleaning Process	All valid blocks have been migrated. Victim selection registered as a new free section, an F2FS ‘pre-free’ section. Checkpointing occurs where the section is made a free section and it can be reallocated. This process is used to manage any loss of data referencing by checkpoints due to for example, unexpected power outages that may result in pre-free sections being reused before checkpointing.

As long as the data medium still has sufficient storage space, the append logging known from LSF is used. In this mode, only clean segments that are not yet occupied are written to. Remember: Sequential writes are always preferable to random writes in flash memory. In Threaded Logging mode, the whole thing is reversed. Now the remaining invalid blocks in dirty segments are collected, written. So now the remaining gaps are being used. F2FS automatically switches between the two modes depending on the remaining free memory. This is an attempt to achieve the best possible write performance.

### Roll-Back Recovery

F2FS uses checkpoints to ensure integrity in the event of power failures or other disruptions. A Checkpoint uses not only the checkpoint but also both the Node Address Table (NAT) and the Segment Info Table (SIT). There is NAT and SIT journaling and bitmap addressing the valid NAT and SIT within the checkpoint.

Table 3.10: F2FS Adaptive Logging [45]

Title	Description
Append Logging	Writing to clean segments. Need cleaning operations if there is no free segment. Cleaning causes mostly random read and sequential writes. <i>Node is always written with append logging policy.</i>
Thread Logging	Lower than %5 of total sections by default. Writing to Dirty segments. Reuse invalid blocks in dirty segments. No need to clean. Cause random writes. <i>F2FS gracefully gives up normal logging and turns to threaded logging for higher sustained performance.</i>

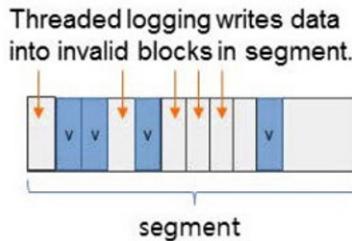


Fig. 3.16: Threaded Logging is tactical in going after invalid blocks in a segment

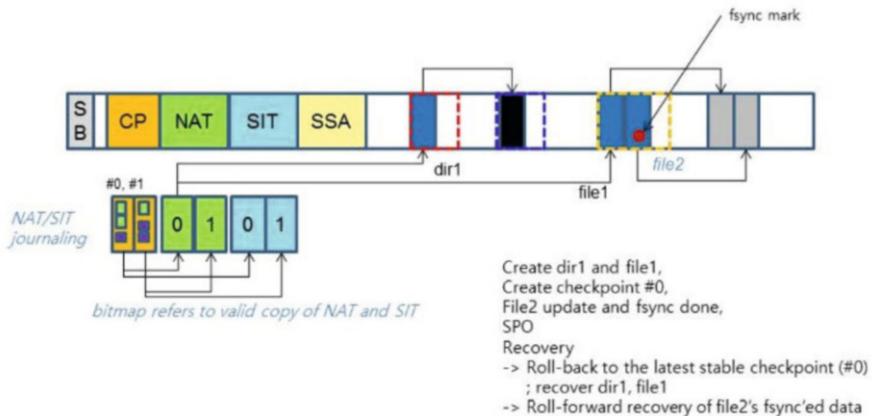


Fig. 3.17: F2FS Index Structure (Source: [46])

If there is a disruption of a power issue then F2FS can do a rollback recovery using the most recent valid Checkpoint (see Fig. 3.17). Earlier checkpoints were discussed, and F2FS has two to maintain a stable one. The header and footer identify the stable

Checkpoint is the same. There is also a version number if both checkpoint headers and footers match. In this instance, the file system will choose the most recent Checkpoint. F2FS does not always need to rely on the Checkpoint for recovery when `fsync()` is involved. F2FS focuses on the data blocks and direct node blocks (which are marked). The Roll-Forward Recovery Procedure consists of the following steps according to Lee et al. [45]:

1. Search marked Direct Node Blocks.
2. Per marked Node Block, identify old and new Data Blocks by checking the difference between the current and previous node block.
3. Update SIT: Invalidate old Data Blocks.
4. Replay new data writes; update NAT and SIT accordingly.
5. Create Checkpoint.

### > Important

Under Linux/Unix, the C function `fsync()` from the standard library transfers changed buffered data from the working memory to the file on the physical device. The call blocks until the device have reported that the transfer is complete. Of course, the buffer should only be flushed occasionally to avoid stress for the hardware.

## 3.5 Forensic Analysis

The Flash-Friendly File System (F2FS) is made for removable media and used with some Android devices, or maybe changed over from another Android (Linux) based File System such as EXT. You will find plenty of videos on how to change the File System on an Android device and speed comparisons between F2FS and other file systems. In this section, we want to turn to how to read and acquire F2FS-formatted memory sticks or even SSDs. This section will talk about forensic analysis of the different F2FS formatted devices.

### 3.5.1 F2FS Sample Dataset

The examples discussed in this chapter with the associated binaries can be found at [github.com](https://github.com/Xamnr/F2FS) under the following link: <https://github.com/Xamnr/F2FS>. For the examples, three different drives were analysed: 1) the memory of a Huawei P9, 2) the content of a USB memory stick, and 3) an SSD. All three volumes were formatted with F2FS. The volumes were filled with text documents, image and video files. Afterwards, some of the files were deleted. The first example comes up with a binary file of a HuaweiP9 Superblock <File:HuaweiP9Superblock> and Checkpoint <File:HuaweiP9Checkpoint.zip>.

Beyond this, for USB flash drive three extracts had been made:

1. BASE: Formatted F2FS
2. ADDED: Two Folders created. One folder and four test files created. 2 png files and 2 text files.
3. DELETED: 1 png (Moved) and 1 text file (copied). 1 png deleted and deleted test folder.

The second example contains a sample of a USB flash drive. A total of 2 folders and four files were added to the drive:

```
<ADDED EXTRACTION>
| -- folder
| -- Test
|   | -- pngpicture.png
|   | -- pngtodelete.png
|   | -- textstays
|   | -- texttodeletediconderoga
```

The third example contains the formerly added files, but this time some of them had been deleted:

```
<DELETED EXTRACTION>
| -- folder
|   | -- pngpicture.png <moved>
|   | -- textstays <copied>

deleted:
| -- pngtodelete.png
| -- textstays
| -- texttodeletediconderoga
```

The example is completed by three dumps from the NAT, SIT and SSA region of the flash drive.

### 3.5.2 F2FS and Windows

The main issue for forensics as it relates to F2FS is that Microsoft Windows does not recognize F2FS formatted devices. Windows sees the USB Device itself but does not recognize the partition.

Fortunately, at least forensic tools should recognize the device (as shown below with MSAB's XRY) and be able to do extractions. Best practice is to use a write blocker for the media. Once the drive has been detected, the analysis process with XRY is quite straightforward.

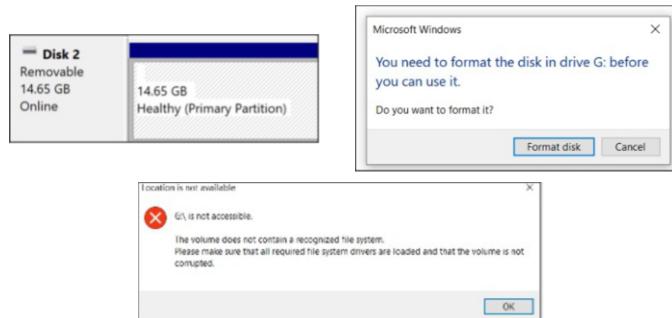


Fig. 3.18: F2FS formatted USB drive connected to a Microsoft Windows 10 computer

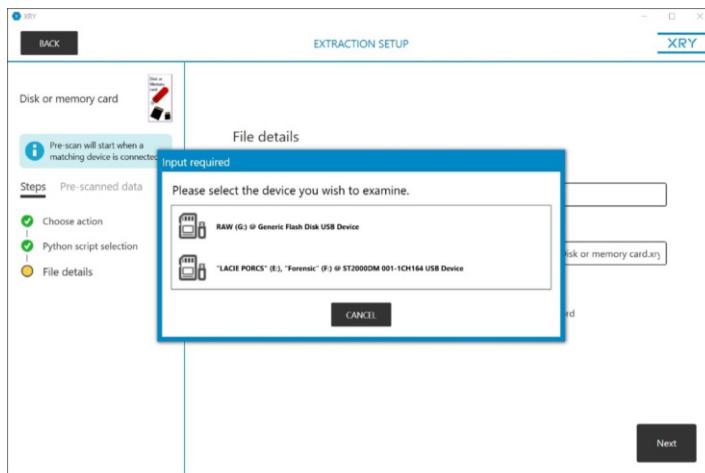


Fig. 3.19: MSAB's XRY 9.3 showing the Disk Connection Options

### 3.5.3 Data-Extraction with XRY

The first step is, of course, to establish a connection to the device and the storage medium. For Android devices it will depend on whether the mobile device is supported by model or by a generic one (see Fig. 3.19). Mobile Device support may not include all options such as Logical and Physical extraction, but rather only one. Remember Micro SD Cards inserted into the mobile device should be removed at some point and done separately.

The forensic tool(s) should be able to image the mobile device and/or the removable media without an issue. Since F2FS is made for Flash memory it could be found on Solid State Drives as well. After we have successfully read in the flash drive, the next step is to analyse the data it stores. For an example, we read the internal memory of a Huawei P9 that was formatted with F2FS. The actual extraction was

carried out using the software XAMN. Fig. 3.20 shows the result for of our example drive. Obviously, in addition to some regular files, deleted file artefacts were also found during the scan process. The F2FS formatted Flash drive also had deleted files that were retrieved as well. Remember recovering deleted files relies on a lot of variables, but at least we know it is possible with the F2FS file system.

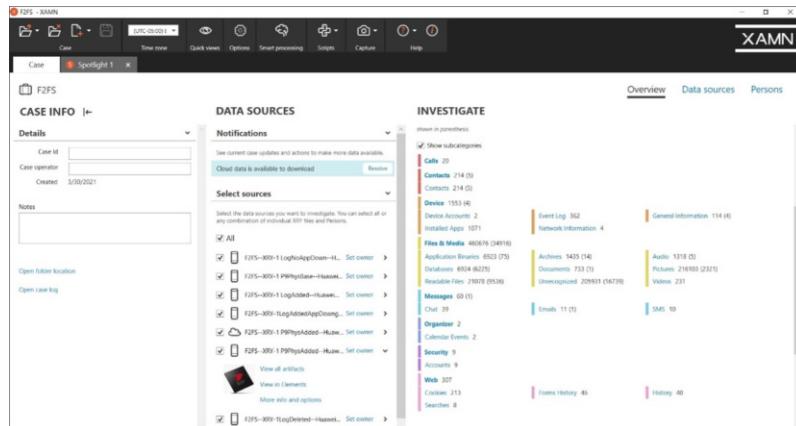


Fig. 3.20: MSAB’s XAMN Showing that the extraction did decode data from the Huawei P9 cell phone Including deleted files

The header of the checkpoint area is initially less interesting from a forensic point of view (see Table 3.12). It only contains information about the number of available segments.

### 3.5.4 Superblock Examination

After successfully making a forensically sound copy of the medium to be examined, we can begin the actual analysis. Our investigation should also begin here since an F2FS partition starts with the Superblock(SB). In this region, important information about the structure of the rest of the file system can be found.

As examiners we commonly see the file system as raw data with our tools as seen in Fig. 3.21. Can we make sense of this data? In researching F2FS, a resource that describes the Superblock data structure was found <sup>1</sup>. This data was used to create Tables 3.11 and 3.12 for the Superblock and the Checkpoint, respectively<sup>2</sup>.

<sup>1</sup> [www.programmersought.com/article/49182049693/](http://www.programmersought.com/article/49182049693/)

<sup>2</sup> \* Offset 3204+ Kernel Information and more: Linux version 5.10.0-kali3-amd64 (devel@kali.org) (gcc-10 (Debian 10.2.1-6) 10.2.1 20210110, GNU ld (GNU Binutils for Debian) 2.35.1) #1 SMP Debian 5.10.13-1kali1 (2021-02-08)

Superblock Area															
Address	Hex	ASCII													
.....0000	10 20 F5 F2	01	00	0E	00	09	00	00	03	00	00	00	00	00	00
.....0016	0C 00 00 00 09 00 00 00 01 00 00 00 01 00 00 00	ó	ó	ó	ó	ó	ó	ó	ó	ó	ó	ó	ó	ó	ó
.....0032	00 00 00 00 00 97 3A 00 00 00 00 00 15 1D 00 00	:	“	”	”	”	”	”	”	”	”	”	”	”	”
.....0048	4A 1D 00 00 02 00 00 00 02 00 00 00 22 00 00 00	J	“	”	”	”	”	”	”	”	”	”	”	”	”
.....0064	0F 00 00 00 15 1D 00 00 00 02 00 00 00 02 00 00	N	I	“	”	”	”	”	”	”	”	”	”	”	”
.....0080	00 06 00 00 00 0A 00 00 00 4E 00 00 00 6C 00 00	«\	J	”	”	”	”	”	”	”	”	”	”	”	”
.....0096	03 00 00 00 01 00 00 02 00 00 00 AB 5C 97 4A	ú	E	ó	Y	G	L	”	”	”	”	”	”	”	”
.....0112	FA 45 4F 19 B0 6F DD 47 4C 22 C4 3E 00 00 00 00	ú	E	ó	Y	G	L	”	”	”	”	”	”	”	”
.....0128	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	ú	E	ó	Y	G	L	”	”	”	”	”	”	”	”

Fig. 3.21: Viewing the USB Flash Drive’s F2FS Superblock with MSAB’s XAMN

The Superblock (SB) provides us with important information about the structure of the partition and the exact location of NAT, SIT, SSA and MAIN area. The values given are to be understood as multiples of the block size. In our example, the SIT starts at block no. 1536. The correct address is obtained when we multiply this number by the block size. In our case, the SIT thus starts at byte 6,291,456 (1536\*4096). The total size of the SIT can again be determined via the 4-byte value at offset 56. In our case, this is *hx02000000*. Since this is a little-endian (LE) value with the least significant bit on the far left, the size is 2. Incidentally, all other header fields are also LE values. The node ID of the root directory can also be determined in the superblock (offset 96).

If your forensic tool supports viewing the raw data then that is a start. Hopefully the tool has options to translate the code such as showing the bit options i.e. 16 or 32, Little or Big Endian, and others such as the GUID. If your tool does not have these options then you can use the HxD Hex Editor/Viewer with the Inspector feature. HxD can be found here [mh-nexus.de/](http://mh-nexus.de/) or for English [mh-nexus.de/en/](http://mh-nexus.de/en/).

### ! Remember

HxD is an editor, so work off of a copy of the file.

### 3.5.5 Examine NAT, SIT & SSA with Linux

Since F2FS was developed specifically for use with Linux or Android, it is obvious to conduct an investigation with this system as well. This section will show the forensic analysis of F2FS with open-source digital forensic tools if you are inclined to use Linux, i.e. Kali, Santoku, or another forensic type, and command line. Then you may

Source: [www.programmersought.com/article/49182049693/](http://www.programmersought.com/article/49182049693/)

Source: [www.programmersought.com/article/37962049663](http://www.programmersought.com/article/37962049663)

Table 3.11: Super Block (USB Flash Drive) Example Values

Offset (deci- mal)	Bytes	Description	Hex	ASCII	Value	Format
0	4	Magic Number	10 20 F5 F2	10 20 F5 F2	1020F5F2	N/A
4	2	Major Version	01 00	..	1	Int 16 LE
6	2	Minor Version	0E 00	..	14	Int 16 LE
8	4	Log 2 Sector size in bytes	09 00 00 00	....	9	Int 32 LE
12	4	Log 2 Sectors per block	03 00 00 00	....	3	Int 32 LE
16	4	Log 2 Block Size in bytes	0C 00 00 00	....	12	Int 32 LE
16	4	Log 2 Block Size in bytes	0C 00 00 00	....	12	Int 32 LE
20	4	Log 2 Blocks per Segment	09 00 00 00	....	9	Int 32 LE
24	4	Segments per Sector	01 00 00 00	....	1	Int 32 LE
28	4	Sections per Zone	01 00 00 00	....	1	Int 32 LE
32	4	Checksum offset inside super block	00 00 00 00	....	0	Int 32 LE
36	8	Total # of User Blocks	00 97 3A 00 00 00 00 00	..:....	3839744 <sup>3</sup>	Int 64 LE
44	4	Total # of Sections	15 1D 00 00	....	7445	Int 32 LE
48	4	Total # of Segments	4A 1D 00 00	J...	7489	Int 32 LE
52	4	Segments for Checkpoint	02 00 00 00	....	2	Int32 LE
56	4	# of Segments for SIT	02 00 00 00	....	2	Int32 LE
60	4	# of Segments for NAT	22 00 00 00	“...”	34	Int 32 LE
64	4	# of Segments for SSA	0F 00 00 00	....	15	Int 32 LE
68	4	# of Segments for Main	15 1D 00 00	....	7445	Int 32 LE
72	4	Start Block address of Segment 0	02 00 00 00	....	2	Int 32 LE
76	4	Start of block address for Check- point	00 02 00 00	....	512	Int 32 LE
80	4	Start block address of SIT	00 06 00 00	....	1536	Int 32 LE
84	4	Start block address of NAT	00 0A 00 00	....	2560	Int 32 LE
88	4	Start block address of SSA	00 4E 00 00	.N..	19968	Int 32 LE
92	4	Start block address Main	00 6C 00 00	.l..	27648	Int 32 LE
96	4	Root inode number	03 00 00 00	....	3	Int 32 LE

want to (write-protected, of course) gather additional system information about the removable media (USB Flash Drive in this case). You could mount the dd(bin) image as well. You will find F2FS Tools in the official Linux kernel GitHub repository <sup>4</sup>. Alternatively, you can install the necessary tools via package management. Under Ubuntu, for example, the following command can be used for this purpose:

```
# sudo apt-get install f2fs-tools
```

File system information can be dumped from the device using the F2FS Tools dump function. The primary interest is obtaining data from the NAT, SIT, and SSA. However, we can also use it to obtain a file if we want to. The dump.f2fs shows the on-disk inode information using the inode number and a dump of all SSA and SIT entries to file recognised by dump\_ssa and sump\_sit. On the web page of the

---

<sup>4</sup> <http://git.kernel.org/pub/scm/linux/kernel/git/jaegeuk/f2fs-tools.git>

Table 3.12: Checkpoint (USB Flash Drive) Example Values

Offset (decimal)	Bytes	Description	Hex	ASCII	Value	Format
0	8	CP version for comparing old and new	41 49 3C 37 00 00 00 00	AI<7	926697793	Int 64 LE
8	8	# of User Blocks	00 3A 38 00 00 00 00 00	:8	3684864	Int 64 LE
16	8	# of Valid blocks in Main Area	42 00 00 00 00 00 00 00	B	66	Int 64 LE
24	4	# of Reserved segments for garbage cleaning (GC)	81 00 00 00	.....	129	Int 32 LE
28	4	# of overprovision segments	F8 00 00 00	Ø	248	Int 32 LE
32	4	# of free segments in Main area	0F 1D 00 00	.....	7439	Int 32 LE

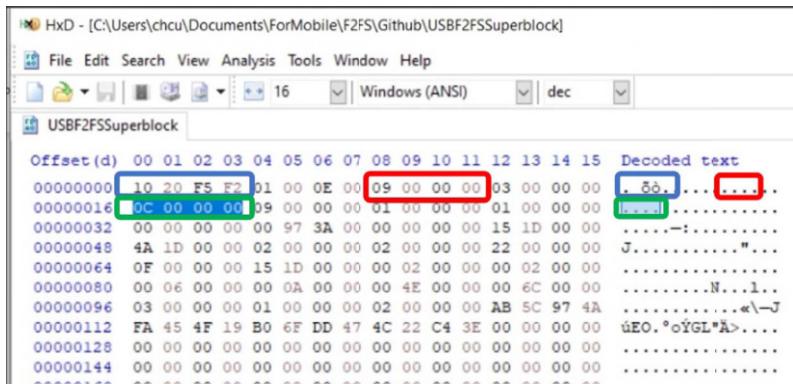


Fig. 3.22: HxD with color coded selections of Superblock to match table data 3.11

the Linux Kernel Organisation [89] all options can be retrieved, with which the command can be executed (see Table 3.13).

Table 3.13: Usage: dump.f2fs command options

Command Option	Description
-d	debug level [default:0]
-i	inode no (hex)
-n	NAT dump nid from #1 #2 (decimal), for all 0 -1
-s	SIT dump segno from #1 #2 (decimal), for all 0 -1
-S	Sparse_mode
-a	SSA dump segno from #1 #2 (decimal), for all 0 -1
-b	Blk_addr (in 4 KB) Block Address
-V	Print the version number and exit

The tool can also be run without special parameters. In this way, information about the size of the disk and the sector size can be determined first of all. A typical output of the `dump.f2fs` command results without any options where `/dev/sdb1` is the USB Flash Drive looks like this:

```
#sudo dump.f2fs /dev/sdb1
Info: [/dev/sdb1] Disk Model: Flash Disk
Info: Segments per section = 1
Info: Sections per zone = 1
Info: sector size = 512
Info: total sectors = 30717952 (14999 MB)
Info: MKFS version
"Linux version 5.10.0-kali3-amd64 ..."
...
Info: superblock features = 0 :
Info: superblock encrypt level = 0, salt = 00000000000000000000000000000000
Info: total FS sectors = 30717952 (14999 MB)
Info: CKPT version = 373c4953
Info: checkpoint state = c5 : nat_bits crc compacted_summary unmount
Done: 0.177078 secs
```

### Node Allocation Table (NAT) Data

In section 3.3.2 we have already learned about the task and function of the Node Allocation Table. With `dump.f2fs` we can output the contents of the table for our disk. The output of command `sudo dump.f2fs -n 0~1 /dev/sdb1` for our example is shown in Table 3.14. Remember: All the node blocks are mapped by NAT. Hence, the position of each node is translated by the NAT table. Apparently there are exactly 4 valid direct accounts on the device. In addition to the logical node ID, the block address is also specified. Since the default block size is 4 KB by default, we can thus determine the exact physical on-disk location of the node.

Table 3.14: Example for a Node Allocation Table (NAT)

nid: 3	ino: 3	offset:0	blkaddr: 27666	pack:1
nid: 4	ino: 4	offset:0	blkaddr: 27661	pack:1
nid: 6	ino: 6	offset:0	blkaddr: 28167	pack:1
nid: 10	ino: 10	offset:0	blkaddr: 28167	pack:1

### Show the Segment Info Table (SIT) Data

The second important data structure besides the NAT is the Segment Information Table. The term segment here means a contiguous lump of disk blocks. Normally 512

continuous Blocks grouped into one segment. As already discussed in section XX, a segment is assigned a sector by default. For our example USB stick, the command would look like this: `sudo dump.f2fs -s 0~1 /dev/sdb1`. The (shortened) output in this case looks like this:

```
segment_type(0:HD, 1:WD, 2:CD, 3:HN, 4:WN, 5:CN)

segno:0          vblocks:2          seg_type:3          sit_pack:1
00 00 02 80 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
segno:1          vblocks:3          seg_type:4          sit_pack:1
02 50 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
...
segno:7444       vblocks:0          seg_type:0          sit_pack:1

valid_blocks:[0xa0]      valid_segs:5 free_segs:7440
```

The result is a bitmap for valid blocks inside the different segments. The segment summary contains 512 entries, which is the 2MB segment size. A summary entry for a 4KB-sized block in a segment contains the. The first segment with number 0 seems to have 2 valid block. It is a hot data block with directory entries inside (`seg_type=3`). In contrast, the second segment contains Data blocks (`seg_type=4`). There are a total of 160(`hxA0`) valid blocks on the F2FS partition that are stored in 5 valid segments.

### Look inside the Segment Summary Area (SSA) Data

A look at the SSA allows us to find out who exactly owns each block. For example, to get an insight into the segment with the number 0, we can use the following command: `sudo dump.f2fs -a 0~1 /dev/sdb1`

```
segno: 0, Current Node
[ 0: 3][ 1: 3][ 2: 4][ 3: 5][ 4: 5]
[ 5: 5][ 6: 5][ 7: 5][ 8: 4][ 9: 4]
[10: 5][11: 5][12: 4][13: 4][14: 5]
[15: 5][16: 3][17: 5][18: 3][19: 3]
[20: 4][21: 3][22: 4][23: 3][24: 3]
[25: 0][26: 0][27: 0][28: 0][29: 0]
...
```

As clearly seen, the blocks with the numbers 0 and 1 from segment no 0 both belong to the node with the node ID 3. The third block belongs to the node with the number 4 and so on.

## Obtain a file by it's node ID

If you recall earlier the linux command `sudo fsck.f2fs -t /dev/sdb1` obtained some data including the files and their node identifiers.

```
|-- folder <ino = 0x4>, <encrypted (0)>
|   |-- pngpicture.png <ino = 0x6>, <encrypted (0)>
|   '-- textstays <ino = 0xa>, <encrypted (0)>
```

In order to get the file, again we can to use the linux command line with the `F2fs.dump` command. Your forensic tools or even Linux should do this for you, but just to reinforce how this F2FS directory works with nodes you could use the command:

```
#sudo dump.f2fs -i 0x6 /dev/sdb1

Info: [/dev/sdb1] Disk Model: Flash Disk
Info: Segments per section = 1
Info: Sections per zone = 1
Info: sector size = 512
Info: total sectors = 30717952 (14999 MB)
...
[print_node_info: 353] Node ID [0x6:6] is inode
i_mode          [0x     81a4 : 33188]
i_advise         [0x      3 : 3]
i_uid            [0x      0 : 0]
i_gid            [0x      0 : 0]
i_links          [0x      1 : 1]
i_size           [0x    3c41f : 246815]
i_blocks          [0x      3e : 62]
i_atime          [0x606db122 : 1617801506]
i_atime_nsec     [0x24990b38 : 614009656]
i_ctime          [0x606db114 : 1617801492]
i_ctime_nsec     [0x2bc016f9 : 734009081]
i_mtime          [0x606b7024 : 1617653796]
i_mtime_nsec     [0x ad1c7fa : 181520378]
i_generation     [0xd0856627 : 3498403367]
i_current_depth  [0x      0 : 0]
i_xattr_nid      [0x      0 : 0]
i_flags           [0x      0 : 0]
i_inline          [0x      1 : 1]
i_pino            [0x      4 : 4]
i_dir_level       [0x      0 : 0]
i_namelen         [0x      e : 14]
i_name            [pngpicture.png]
i_ext: foofs:0 blkaddr:ef400 len:3d
i_addr[0x0]        [0x    ef400 : 979968]
i_addr[0x1]        [0x    ef401 : 979969]
i_addr[0x2]        [0x    ef402 : 979970]
i_addr[0x3]        [0x    ef403 : 979971]
...
i_addr[0x3b]        [0x    ef43b : 980027]
i_addr[0x3c]        [0x    ef43c : 980028]
```

```
i_nid[0]          [0x      0 : 0]
i_nid[1]          [0x      0 : 0]
i_nid[2]          [0x      0 : 0]
i_nid[3]          [0x      0 : 0]
i_nid[4]          [0x      0 : 0]
```

```
Do you want to dump this file to ./lost_found/? [Y/N] Y
Info: checkpoint state = c5 :
nat_bits crc compacted_summary unmount
Done: 3.409981 secs
```

We can thus query important meta-information about the file. In addition to file names (*pngpicture.png*) and size, the MAC timestamps are displayed. The number and concrete address of the blocks and the size of the file are printed out. The file occupies a total of 63 blocks. With a block size of 4096 bytes, this corresponds to 258.048 bytes. The actual size of the file is somewhat smaller, with 246.815 bytes.

### 3.5.6 Carving for artefacts with XAMN

XAMN is an intuitive tool that helps you find and analyse data faster and easier, we can even find a formerly deleted audio file (Fig. 3.23). For the file content search, all three samples of the dataset were first loaded for analysis with XAMN.

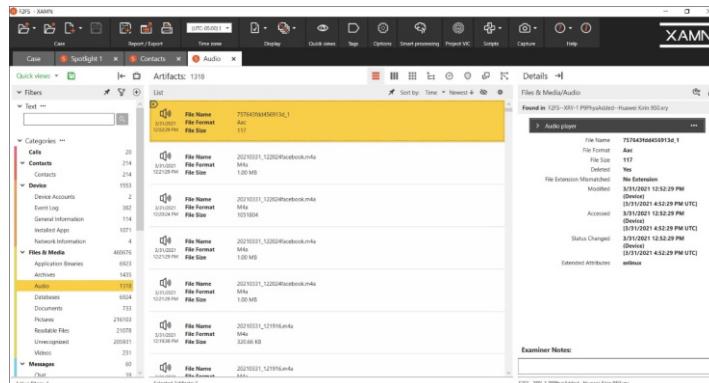
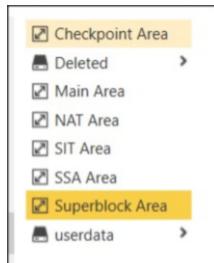


Fig. 3.23: MSAB's XAMN Spotlight showing a deleted audio file

First of all, it can be noted, that the program is able to detect F2FS partitions and their regions. Each region can be called individually and its content can be examined separately.



Looking at the screen capture on the left you should recognize the F2FS Areas. Or more importantly that you are dealing with a Flash-Friendly File System. Both the Superblock and Checkpoint Area are highlighted.

The files themselves and their meta information must be located in the Main Area of the disk. Accordingly, we checked to see if there were any references to the files added or deleted in the examples. A search for the PNG picture file name resulted in two hits (see Fig. 3.24). A simple search with a hex editor should give the same result. But also the names of the image data could be found (see Fig. 3.25). This appears to be the contents of the directory. In addition, a dot '.' and then two dots '..' are recognizable in the dump. Presumably these stand, as usual in Linux for the current directory or for the parent directory entry.

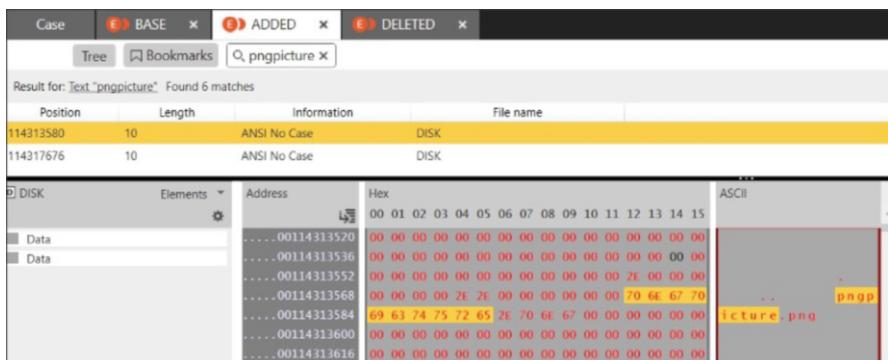


Fig. 3.24: MSAB's XAMN Elements showing the “pngpicture” search results

We are then able to find all four file names together in both the ADDED and DELETED extractions.

Even more, all three samples of the USB flash device could be successfully imported into XAMN. A look at the file tree shows that all file artefacts were found for the three scenarios (1..BASE,2..ADDED,3..DELETED). The correct directory structure could also be reconstructed. In addition to the regular files, the moved or deleted files are displayed as well (Fig. 3.26).

..00114325840	00 00 00 00 00 00 00 00 00 00 00 00 2E 00 00 00 00	.
..00114325856	00 00 00 00 2E 2E 00 00 00 00 00 00 70 6E 67 70	.. pngp
..00114325872	69 63 74 75 72 65 2E 70 6E 67 00 00 70 6E 67 74	icture.png odelete.png
..00114325888	6F 64 65 6C 65 74 65 2E 70 6E 67 00 74 65 78 74	text stays
..00114325904	73 74 61 79 73 00 00 00 00 00 00 00 74 65 78 74	text todeleteticonderoga
..00114325920	74 6F 64 65 6C 65 74 65 74 69 63 6F 6E 64 65 72	
..00114325936	6F 67 61 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
..00114325952	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	

Fig. 3.25: Found all four file names together in the extraction notice the “..” and “.” Before the data



Fig. 3.26: File Tree showing the three USB Flash Drive Extractions

## PNG File Signature Analysis

In the normal logging, blocks are written to clean unused segments. Thus there is a good chance that blocks occupied by a file are also written together on the volume. In our example, among other things, two .png image files were uploaded to the disk and then deleted again. This raises the question of whether we can recover the file contents using carving? We want to see and know about the ability to carve the PNG Files. Keep in mind the Cleaning ability of F2FS would most likely cause additional changes.

It is often assumed by laymen that once files have been deleted they cannot be recovered - neither in whole nor in part. This (careless) assumption is found especially often with respect to Unix. The possibilities of finding data fragments at different places of the file system or the hard disk can be used in an investigation. Often these individual fragments can be put together to a file with the method called file *Carving* or at least essential information can be extracted. In our example we want to search for the two image files in *png* format. Binary formats of image files often start with a

magic number by which we can recognize the file. Referring to Gary Kessler's File Signature website: [www.garykessler.net/library/file\\_sigs.html](http://www.garykessler.net/library/file_sigs.html), a search for "PNG" resulted in the following hits (Trailer means the end of the file):

Table 3.15: PNG File Signature as shown from Gary Kessler's website

Description	Result
File Header (Hex)	89 50 4E 47 0D 0A 1A 0A
File Header (ASCII)	%oPNG...
File Description	Portable Network Graphics file
File Extension	PNG
File Trailer or Footer (Hex)	49 45 4E 44 AE 42 60 82
File Trailer or Footer (ASCII)	IEND@B^...

Knowing both the Header and footer gave us a few options. We decided to use a Global Regular Expression (GREP) to see if we could recover the complete file from file header to footer (or trailer). To verify the file header and footer is correct. We can look at the extraction and pick a PNG file and view the raw file data. As you can see below (see Fig. 3.27 and 3.28) both the header and the footer look good.

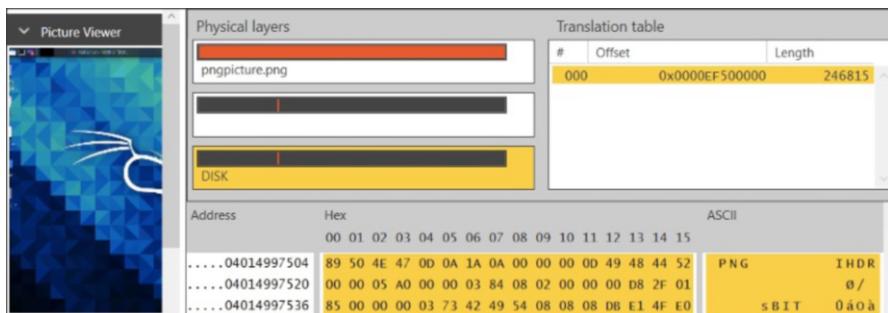


Fig. 3.27: MSAB's XAMN Source Mode showing the PNG File Header (Hex): 89 50 4E 47 0D 0A 1A 0A

Excellent, now we can carve for the deleted png pictures to see if we can find them. A search for file headers in both the Added data USB and the Deleted USB resulted in 101 png file header hits for both of them. Note the deleted png file carved below. The forensic Analysis Tool used below (see Fig. 3.30 and 3.31) was able to find the deleted files and present them.

A search for Test did reveal a hit with the word Folder and Test (the two folders on the USB Flash Drive). Note the red characters are showing that these are different by comparison with the other extraction or extractions.

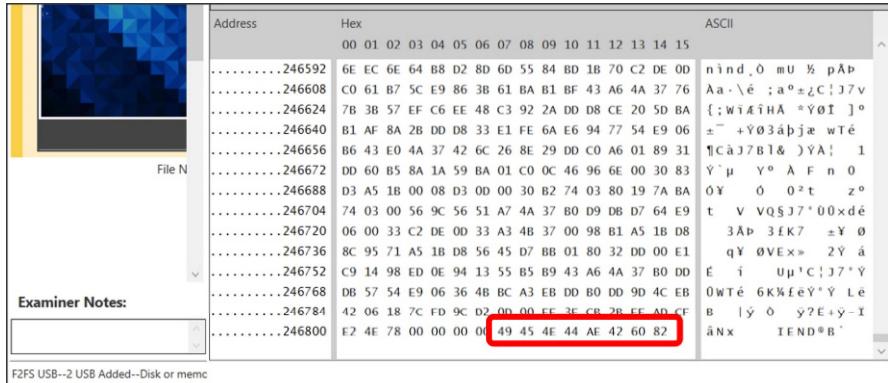


Fig. 3.28: MSAB's XAMN Source Mode showing the PNG File Footer (Hex): 49 45 4E 44 AE 42 60 82

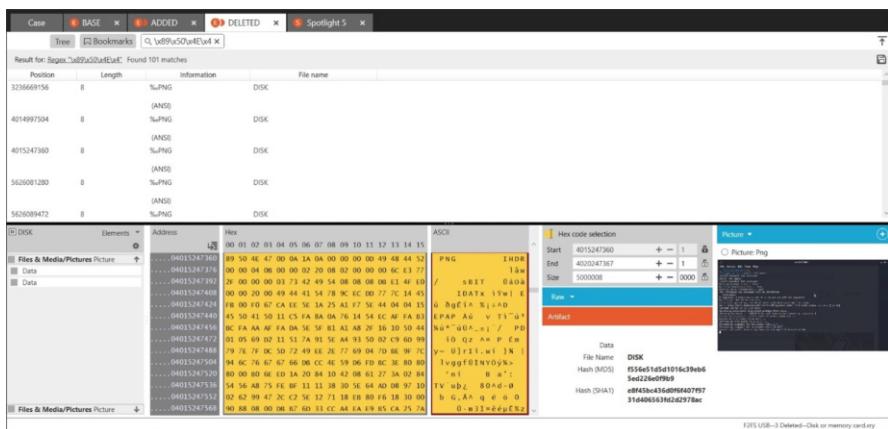


Fig. 3.29: MSAB's XAMN Elements showing the PNG file header search and the deleted PNG File

### 3.5.7 Node Allocation Table (NAT) Comparisons

The changes made to samples 2 (ADDED) and 3 (DELETED) must also have led to a change in the Node Access Table (NAT). Blocks must have become invalid or new blocks must have been occupied on the data carrier as a result of deleting or moving files. For the examination with XAMN we must first open in the NAT area. Then we can compare the two areas with each other. The results are shown in Fig. 3.32 for the ADDED sample and in Fig. 3.33 for the DELETED sample. The only differences in the NAT area was found between the USB Flash Drive, with Data Added and when data was moved/deleted; these are shown in red.

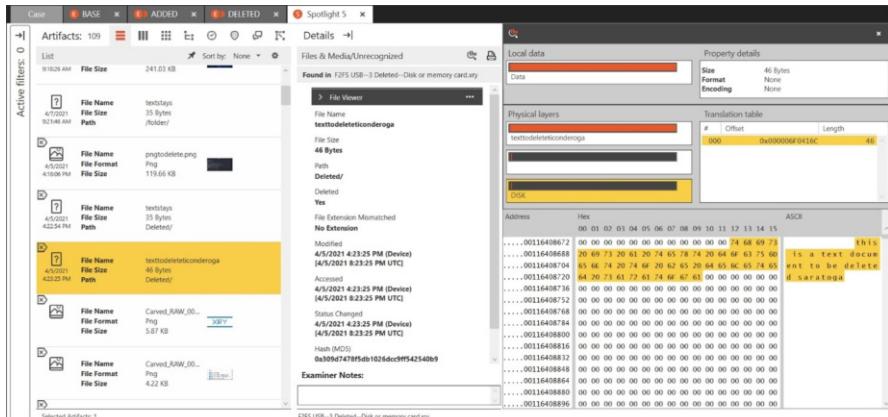


Fig. 3.30: MSAB's XAMN showing the Deleted Text File

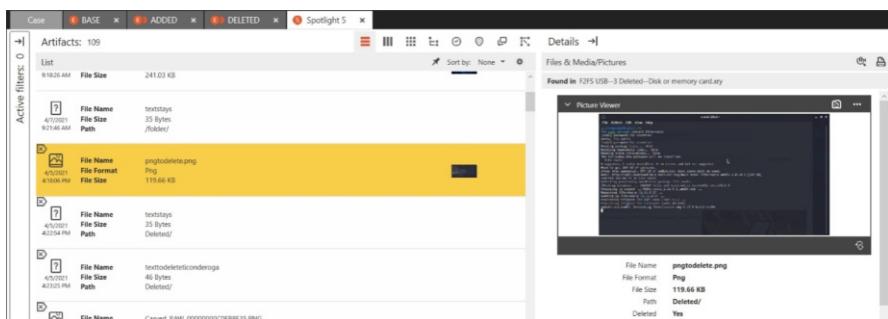


Fig. 3.31: MSAB's XAMN showing the Deleted PNG File

Address	Hex	ASCII
.....02097104	00 01 02 03 04 05 06 07 08 09 00 00 00 00 00 00	
.....02097120	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
.....02097136	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
.....02097152	00 00 00 00 00 00 00 00 00 00 01 00 00 00 01 00	
.....02097168	00 00 02 00 00 00 01 00 00 00 03 00 00 00 00 00	
.....02097184	01 6C 00 00 04 00 00 02 6C 00 00 00 05 00 1 1	
.....02097200	00 00 07 6C 00 00 06 00 00 00 01 6E 00 00 00 00 1 n	
.....02097216	07 00 00 00 02 6E 00 00 08 00 00 00 03 6E 00 n n	
.....02097232	00 00 09 00 00 00 04 6E 00 00 00 00 00 00 00 00 n n	
.....02097248	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	
.....02097264	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	

Fig. 3.32: USB Flash Drive NAT bytes in red are different when compared with the DELETED sample's NAT



Fig. 3.33: USB Flash Drive NAT bytes in red are different when compared with the ADDED sample's NAT

### Additional Data Structure

In addition to the works in this chapter, a F2FS forensic paper [96], in Chinese, was found and the following data structure tables (Table 3.16 and 3.17) are included as they may be deemed useful. For convenience a copy of the the original and the translated version using Google Translate (accuracy cannot be guaranteed) can be found here: <http://github.com/Xamnr/F2FS>.

Table 3.16: Common file metadata information data structure [96]

Intra-block Offset	Byte Length	Content Description
0x10	8	File size
0x20	24	Timestamp
0x5C	255	File name
0x168	3692	923 group index address
0x0FE8	20	File identification number id
0x0FEC	20	File node number ino

Table 3.17: Data structure of catalog file metadata information block [96]

Intra-block Offset	Byte Length	Content Description
0x20	24	Timestamp
0x58	4	Byte length of the directory name
0x5C	255	Directory name
0x168	3692	Directory subfile information
0x0FE8	20	Catalog file identification number nid
0xFEC	20	Catalog file node number ino

### 3.6 F2FS - Application fields

One of the largest Android Manufacturers Samsung and the original creator of the Flash-Friendly File System (F2FS) is using F2FS in combination with UFS in some of their devices over using EXT, such as the Galaxy Note 10 and Galaxy Tab S6. Early on Motorola and Google used F2FS. Huawei and ZTE have also used F2FS on some of their devices.

An interesting, albeit dated article that is of interest entitled “Drone Forensic Analysis Using Open Source Forensic Analysis Using Open Source Tools” [94]. Drones use flash media, many have removable media, so it is not a surprise to see the connection. What is interesting is the mounting of the DD image to use the F2FS Tools. Working off a copy of the dd image would be advised. There are certainly a lot of developments in F2FS, including last year with Linux 5.11 and encryption.

### 3.7 Conclusion

The Flash-Friendly File System (F2FS) has been around for some time and as you can see still may be used. You have seen that it was specific for Flash memory and that this includes some mobile devices. With regards to the Android mobile devices the user may elect to use F2FS over EXT4 if that is an option. Forensic Tools should be able to handle the Flash-Friendly File System, so test them to be sure. The issue is will we find the data that has been deleted. As you saw recovering deleted files is a possibility, however, not a certainty as there are so many variables involved.

**Acknowledgements** Many thanks to Changman Lee, Dongho Sim, Joo-Young Hwang, and Sangyeun Cho, Samsung Electronics Co., Ltd. for their documentation and the 2015 USENIX Conference presentation by Joo-Young Hwang entitled: “F2FS: A New File System for Flash Storage” [45] and to Neil Brown for his 2012 article “An f2fs teardown” [8].

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter’s Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter’s Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



# Chapter 4

## QNX6



Conrad Meyer

**Abstract** The QNX6 filesystem is present in Smartphones delivered by Blackberry (e.g. Devices that are using Blackberry 10) and modern vehicle infotainment systems that use QNX as their operating system. In 2015 QNX as an OS was used in over 50 million vehicles [6] and can hence be considered as one of the most important operating systems in the automotive world. Today's digital forensics tools don't recover a lot from this filesystem, have difficulties with different block sizes, or even don't support the filesystem at all. So it's crucial for the forensic examiner to understand the principles of this filesystem used. This chapter gives an overview of how the filesystem generally stores the files and metadata to give the examiner the chance to get the most information out of the evidence.

### 4.1 Introduction

This chapter gives an insight into the different structures and principles of the QNX6 filesystem developed by QNX. The filesystem was first introduced within QNX Neutrino 6.4 real-time operating system, which today is owned and developed by Blackberry. It is a power-safe file system [7] and can withstand a sudden loss of power without corrupting or losing data. This property is especially useful for the forensic examiner, as it can easily happen that evidence (e.g. a vehicle or smartphone) loses its power supply due to a battery pack running empty.

---

Conrad Meyer

Central Office for Information Technology in the Security Sector (ZITiS), Zamdorfer Straße 88, Munich, Bavaria e-mail: [conrad.meyer@zitis.bund.de](mailto:conrad.meyer@zitis.bund.de)

Table 4.1: Standard Parameters of the QNX6 Filesystem

Parameter	Value	Remark
Max physical Size	2 TB 2	
Supported Standard Logical Blocksizes	512, 1024, 2048, 4096 Bytes	
Max Filename Length	510 bytes	UTF-8

Table 4.1 shows the standard values that are regularly used when formatting a volume with the QNX6 filesystem. Note, that especially in-car infotainment systems, those values can be different (e.g. larger blocksize). All the addressing inside the filesystem is based on the blocksize, extracted out of the superblock.

The following sections will give the reader an insight into the binary structures of the most important parts of the filesystem, like a superblock or inode and some basic knowledge about the mechanism when files are deleted.

## 4.2 QNX6 Filesystem Structure

To understand the principle behaviour and main functions of the QNX6 filesystem, the following chapter shows the structure of a volume and how files, directories and metadata are linked. Volumes can be formatted in QNX6 in little-endian or big-endian style. All the examples in the following show a QNX6 Volume formatted with little endianness. Fig. 4.1 shows the main parts of a QNX6 filesystem and their standard size and addresses. The system area contains the Bitmap of the allocated



Fig. 4.1: Layout of a QNX6 filesystem volume

and unallocated Blocks of the Filesystem. Each bit represents a Block. Suppose the volume is formatted in the standard way. In that case, the volume will start with a volume boot record, which contains standard ASCII coded bootloader messages (Fig. 4.2), already giving a hint that the Volume is formatted with QNX.

### ! Attention

Sometimes, on non standard volumes a partition directly starts with the Superblock.

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
00000000	EB	10	90	00	20	60	D2	00	10	00	00	00	D8	3F	06	00	é	`ò Ø?
00000010	00	80	FA	31	C0	8E	D0	BC	00	20	B8	C0	07	50	B8	36	€úlAžðn	,À P,6
00000020	01	50	CB	00	00	00	00	00	00	00	00	00	00	00	00	00	PE	
00000030	00	00	66	90	00	00	00	00	00	00	00	00	8D	B4	00	00	f	
00000040	10	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00000050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00000060	FF	FF	00	00	00	93	00	00	FF	FF	00	00	00	93	00	00	ÿý	" ÿý "
00000070	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00000080	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	I	ÿý >Í
00000090	18	00	90	7C	00	00	00	00	FF	FF	00	00	9B	CF	00		ÿý "Í QNX v1	
000000A0	FF	FF	00	00	00	93	CF	00	0D	0A	51	4E	58	20	76	31	.	.2b Boot Loader
000000B0	2E	32	62	20	42	6F	6F	74	20	4C	6F	61	64	65	72	00	Unsupported BI	
000000C0	0D	0A	55	6E	73	75	70	70	6F	72	74	65	64	20	42	49	OS	RAM Error
000000D0	4F	53	00	0D	0A	52	41	4D	20	45	72	72	6F	72	00	0D	Disk Read Error	
000000E0	0A	44	69	73	6B	20	52	65	61	64	20	45	72	72	6F	72	Missing OS Im	
000000F0	00	0D	0A	4D	69	73	73	69	6E	67	20	4F	53	20	49	6D	age Invalid OS	
00000100	61	67	65	00	0D	0A	49	6E	76	61	6C	69	64	20	4F	53	Image Unsuppo	
00000110	20	49	6D	61	67	65	00	0D	0A	55	6E	73	75	70	70	6F	rted Multi-Boot	
00000120	72	74	65	64	20	4D	75	6C	74	69	2D	42	6F	6F	74	00	:	~ Güö
00000130	3A	20	00	0D	0A	00	0E	1F	88	16	11	00	FB	FC	F6	06	t è? ö	u *
00000140	03	00	02	74	03	E8	3F	00	F6	06	03	00	01	75	06	BE	" è< »*U'Áí r! ü	
00000150	A8	00	E8	3C	00	BB	AA	55	B4	41	CD	13	72	21	81	FB	U <u>ü</u> öÁ t ,	PžA,
00000160	55	AA	75	1B	F6	C1	01	74	16	B8	00	02	50	8E	C0	B8	PřlÁtç»	é@ È%
00000170	00	02	50	66	31	C0	89	C7	BB	08	00	E8	40	00	CB	BE	À è	ëS' í fà í
00000180	C0	00	E8	C0	00	EB	53	B4	0F	CD	10	83	E0	7F	CD	10	À- Át »	í èöÁ
00000190	C3	AC	08	C0	74	09	BB	07	00	B4	0E	CD	10	EB	F2	C3		
000001A0	66	03	06	04	00	BE	40	00	89	5C	02	89	7C	04	6C	44	f	¾@ t\ %  ÇD
000001B0	06	66	89	44	08	8A	16	11	00	B4	42	CD	13	C3	56	E8	fwd S	'Bí ÁVé
000001C0	DE	FF	72	02	5E	C3	BE	DF	00	E8	C5	FF	EB	0C	E8	CF	þyr ^ÃwB	éÄye ēï
000001D0	FF	73	06	F6	C4	10	75	EE	F9	C3	F4	EB	FD	00	00	00	ÿs öÄ uiùÅéy	
000001E0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000001F0	00	00	00	00	00	00	00	00	00	00	00	00	00	55	AA	U*		

Fig. 4.2: Sector 0 of a QNX6 Partition/Volume

In the following, we will have a closer look at all the structures above. We will follow those structures to construct a file and its metadata out of the filesystem information. The example filesystem is in little-endian mode.

### 4.2.1 Superblock

The filesystem maintains two Superblocks or global root blocks. One of those blocks, called the working Superblock, manages the modified data, while the other one, the stable Superblock, consists of the original version of all the blocks. Which Superblock is the active one is determined by the 64-bit long serial number. The Superblock with the higher serial is the active one. After all, active write operations are done, and the integrity is checked, the former working superblock becomes the new stable one by updating the serial number (old superblock serial +1).

The superblock contains the global information of the filesystem. Table 4.2 contains the offset address of the main features of the Superblock.

Table 4.2: Main Features and their Offset in the QNX6 superblock

Parameter	Offset in Superblock	Size (bytes)
Serialnumber	0x8	8
creation timestamp	0x10	8
last access timestamp	0x14	8
Volume ID	0x20	16
Blocksize	0x30	4
Root Inode Inodes	0x48	array 16 x 4 bytes
Root Inode bitmap	0x98	array 16 x 4 bytes
Root Inode longfilenames	0xE8	array 16 x 4 bytes

### ! Attention

When used with the standard driver issued by Blackberry and the default settings, you can determine the last access to the filesystem by selecting the stable superblock (highest serial) and checking the access timestamp (assuming that system time is used was valid). However, some non-standard drivers don't touch this timestamp, so for reliable results, you have to test the drivers from the System where the image originated in each case!

The superblock contains three root inodes that point to the main parts of the filesystem. The first array root inode contains the pointers to the inodes that contain the data (files, directories, data). The second one contains the pointers to the bitmap of the allocated blocks, and the third one is the pointers to the long filenames (filenames > 27 utf8 characters, up to 510 characters). The data inside those root inodes is shown in Table 4.3. Those root inodes contain pointers to the corresponding filesystem parts. If the level parameter is zero, the root inode has 16 direct pointers. By adding another level, indirect pointers are added, as shown in Fig. 4.4. Each indirect pointer then points to a block containing inodes or indirect 32-bit pointers, depending on the defined number of levels. The actual data is always at the lowest level of the tree. Given the value of blocks that such a tree can address is  $16 * (\text{block size in bytes} / 4)^{\text{level}}$ . So, for example, with a level value of 2, and a block size of 1024 bytes, already 1,048,576 blocks can be addressed.

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
00002000	22	11	19	68	46	DA	79	9A	23	00	00	00	00	00	00	" hFÜÿ\$	
00002010	1E	00	00	00	43	94	6C	60	00	01	00	00	04	00	03	'C'1`	
00002020	54	08	BE	35	56	35	4F	2B	8C	24	B2	EB	CB	2A	42	90	
00002030	00	10	00	00	00	19	00	00	A7	16	00	00	F8	C7	00	00	
00002040	7E	7F	00	00	01	00	00	00	00	80	00	00	00	00	00	~ €	
00002050	CD	00	00	00	FF	í											
00002060	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002070	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002080	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002090	01	01	00	00	00	00	00	00	FF	18	00	00	00	00	00	ÿ	
000020A0	00	00	00	00	01	00	00	00	FF	ÿÿÿÿÿÿÿÿ							
000020B0	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
000020C0	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
000020D0	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
000020E0	00	01	00	00	00	00	00	00	00	B0	03	00	00	00	00	*	
000020F0	73	7F	00	00	FF	s 9999999999999999											
00002100	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002110	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002120	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002130	01	01	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
00002140	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002150	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002160	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002170	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00002180	00	01	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
00002190	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021A0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021B0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021C0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021D0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021E0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	
000021F0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿÿ	

Offset	Title	Value
2000	Magic	21 11 19 68
2004	Checksum	46 DA 79 9A
2008	Serial	23 00 00 00 00 00 00 00
2010	CTime	01.01.1970
2014	ATime	06.04.2021
2018	Flags	00 01 00 00
201C	Version1	04 00
201E	Version2	03 00
2020	Volumeld	94 08 BE 35 56 35 4F 2B 8C 24 B2 EB CB 2A 42 90
2030	BlockSize	00 10 00 00
2034	Number of INodes	00 19 00 00
2038	Free INodes	A7 16 00 00
203C	Number of Blocks	F8 C7 00 00
2040	Free Blocks	7E 7F 00 00
2044	Allocation groups	01 00 00 00

Root Node		
2048	size	00 80 0C 00 00 00 00 00
2050	Pointer	CD 00 00 FF
2090	Levels	01
2091	Mode	01
2092	Spare	00 00 00 00 00 00 00

Fig. 4.3: An example of a QNX6 superblock.

Table 4.3: Structure of the root inodes

Parameter	Offset in root inode	Size (bytes)
Size	0x0	8
Pointer	0x8	array 16 x 4 bytes
Levels	0x48	1
Mode	0x49	1

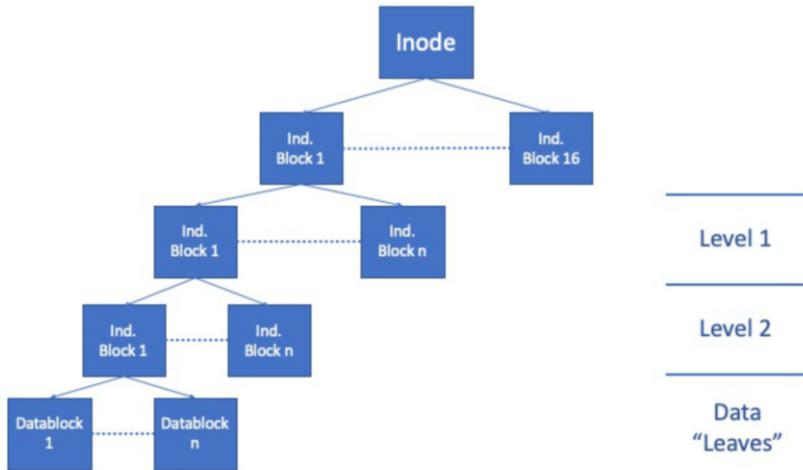


Fig. 4.4: Illustration of inode levels, here a level value of 3

#### 4.2.2 Bitmap

The Bitmap block is used to determine whether a block in the filesystem is used or not. Each bit in the bitmap represents a block. A value of 0 means the Block is unused, 1 means that the Block is allocated. If the volume size is smaller than the bits available in the Bitmap Block, the unused bits are stuffed with ones. The bitmap incorporates two parts. First, system area 1 is split into two halves, where the upper half is used by superblock 1, and the lower half is used by superblock 2. This bitmap area contains the bitmap, inode and indirect addressing blocks of those structures. Second, the bitmap of the blocks that are not used for the filesystem structure (bitmap and inodes). The preallocation of the first system area block leads to the effect that each superblock always works on its own filesystem structure, and to the point that there is always a non-corrupted structure, even in the case of a sudden power loss (a superblock is just becoming the stable one, if all write operations are done, see sect. 4.2.1).

Fig. 4.5 depicts the end of the used space of the bitmap pointed to in the example superblock from Fig. 4.3. The bitmap comprises two blocks, starting at 0x3000, and the volume contains a total of 0xC7F8 blocks. In Fig 4.5, the stuffing of the unused space with ones therefore starts at 0x48FF: Bitmap starting address: 0x3000 + number of blocks 0xC8f8 divided by 8 (each Block represented by 1 bit).

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
000037F0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003800	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003810	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003820	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003830	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003840	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003850	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003860	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003870	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003880	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00003890	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038A0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038B0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038C0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038D0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038E0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
000038F0	00	00	00	F0	FF	8YYYYYYYYYYYYYY												
00003900	FF	YYYYYYYYYYYYYYYY																
00003910	FF	YYYYYYYYYYYYYYYY																
00003920	FF	YYYYYYYYYYYYYYYY																
00003930	FF	YYYYYYYYYYYYYYYY																
00003940	FF	YYYYYYYYYYYYYYYY																
00003950	FF	YYYYYYYYYYYYYYYY																
00003960	FF	YYYYYYYYYYYYYYYY																
00003970	FF	YYYYYYYYYYYYYYYY																
00003980	FF	YYYYYYYYYYYYYYYY																
00003990	FF	YYYYYYYYYYYYYYYY																
000039A0	FF	YYYYYYYYYYYYYYYY																

Fig. 4.5: An example of a QNX6 Bitmap

#### 4.2.3 Inode

On the lowest level of the root inode tree, in the "leaves", the direct inode data is found. Depending on the level defined, also those inodes can address other indirect inode addressing blocks. An inode contains a vast amount of data useful for the forensic examiner, e.g. permissions, access time, change time, and modification time. Table 4.4 shows the offsets and the size of the various parameters in an inode.

Table 4.4: Structure of an inode

Parameter	Offset	Size (bytes)
size	0x0	8
uid	0x8	4
gid	0xC	4
ftime	0x10	4
mtime	0x14	4
atime	0x18	4
ctime	0x1C	4
mode	0x20	2
blockpointer	0x24	array 16 x 4 bytes
Levels	0x54	1
status	0x49	1 (see table 4.5 )

Table 4.5: inode status byte

Value	Status
0x1	directory
0x2	deleted
0x3	normal

As QNX OS is in line with the POSIX standards; also the timestamps are. The epoch is the standard POSIX (or UNIX) epoch, the 01.01.1970, 00:00 UTC. From that epoch, the timestamps are counted in seconds. The modified timestamp (mtime) is the time of the last write operation on this specific file. The access timestamp (atime) tells the examiner the time the file was last read. The change timestamp (ctime) is changed when the permissions of a file are changed. So ctime can be changed without a change in atime. The timestamp ftime is not fully referenced in the POSIX standard. Like in many other filesystems, it is the timestamp when the file was created. The inode 1 always contains the root directory, and inode counting starts with 1.

### ! Attention

When it comes to timestamps, the forensic expert has to pay attention to the reliability of the timestamps given. This is especially true for QNX6. Not all timestamps are actualised on some systems, as with QNX with the standard QNX6 file-system driver. Whenever possible, tests with the system you are examining should be performed (e.g. changing permissions, modifying files, etc.)!

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
00006360	FF	FF	FF	FF	00	03	00	00	00	00	00	00	00	00	00	00	YYYY
00006370	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
00006380	00	10	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
00006390	F8	03	00	00	4D	93	6C	60	13	94	6C	60	4D	93	6C	60	o M"1" "1" M"1"
000063A0	ED	41	03	00	72	7E	00	00	FF	IA r YYYYYYYYYY							
000063B0	FF	YYYYYYYYYYYYYYYY															
000063C0	FF	YYYYYYYYYYYYYYYY															
000063D0	FF	YYYYYYYYYYYYYYYY															
000063E0	FF	FF	FF	FF	00	03	00	00	00	00	00	00	00	00	00	00	YYYY
000063F0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
00006400	00	10	00	00	00	00	00	00	00	00	00	00	00	00	00	00	

Offset	Title	Value
5880	Size	0010000000000000
5888	Uid	00000000
589C	Gid	00000000
5890	File time	01.01.1970 00:00:00
5894	Mod. time	06.04.2021 16:58:53
5898	Access time	06.04.2021 17:02:11
589C	Change time	06.04.2021 16:58:53
58A0	Mode	ED 41
58A2	ExtMode	05 00

Blockptr	
58A4	BlockPtr 0
58A8	BlockPtr 1
58AC	BlockPtr 2
58B0	BlockPtr 3
58B4	BlockPtr 4
58B8	BlockPtr 5
58BC	BlockPtr 6
58C0	BlockPtr 7
58C4	BlockPtr 8
58C8	BlockPtr 9
58CC	BlockPtr 10
58D0	BlockPtr 11
58D4	BlockPtr 12
58D8	BlockPtr 13
58DC	BlockPtr 14
58E0	BlockPtr 15

58E4	File levels	00
58E5	Status	03
58E6	Unknown	00 00
58E8	Zero	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

Fig. 4.6: An example of a QNX6 Inode.

#### 4.2.4 Directories

Inodes with the status 0x3 point to a directory file system object that contains sub-directories and file entries with names shorter than 27 UTF-8 characters. An entry starts with the inode number of that entry, where you can find the metadata like timestamps and the pointers to the Data or other directories, followed by a name length field and the actual name. A directory always contains a "." and a ".." entry. The "." entry contains the inode number of the directory inode, and the ".." entry

contains the inode number of the parent directory inode. In the example Fig. 4.7, those entries are both pointing to the same inode number because the directory shown is the root directory.

Table 4.6: Directory entry

Parameter	Offset	Size (bytes)
Inode number	0x0	4
Namelength	0x4	1
Name	0x5	up to 27

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
07F12F00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F12F01	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	.
07F12F02	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	..
07F13000	01	00	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	
07F13010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13020	01	00	00	00	00	02	2E	2E	00	00	00	00	00	00	00	00	
07F13030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13040	02	00	00	00	00	05	2E	62	6F	6F	74	00	00	00	00	00	.boot
07F13050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13060	03	00	00	00	00	03	62	69	6E	00	00	00	00	00	00	00	bin
07F13070	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13080	04	00	00	00	00	03	65	74	63	00	00	00	00	00	00	00	etc
07F13090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F130A0	05	00	00	00	00	04	69	6E	66	6F	00	00	00	00	00	00	info
07F130B0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F130C0	06	00	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	lib
07F130D0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F130E0	07	00	00	00	00	03	6F	70	74	00	00	00	00	00	00	00	opt
07F130F0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13100	08	00	00	00	03	75	73	72	00	00	00	00	00	00	00	00	usr
07F13110	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13120	1C	00	00	00	08	66	6C	61	73	68	2E	73	68	00	00	00	flash.sh
07F13130	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F13140	1D	00	00	00	13	66	6F	72	6D	61	74	41	70	70	43	68	formatAppCh
07F13150	6B	50	65	72	73	2E	73	68	00	00	00	00	00	00	00	00	kPers.sh
07F13160	1E	00	00	00	0E	66	6F	72	6D	61	74	42	6F	6C	6F	31	formatBolol
07F13170	2E	73	68	00	00	00	00	00	00	00	00	00	00	00	00	00	.sh
07F13180	21	00	00	00	0E	66	6F	72	6D	61	74	42	6F	6C	6F	32	! formatBolo2
07F13190	2E	73	68	00	00	00	00	00	00	00	00	00	00	00	00	00	
07F131A0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	

Fig. 4.7: An example of a QNX6 directory. Here, the root directory is shown.

A long directory entry has a different structure (Table 4.7). It includes the Inode, in which the timestamps and pointers to the data are. Furthermore, the long filenames inode Number, where the entry's name is found, is noted in this structure. An example of a long filename/directory entry is displayed in Fig. 4.8.

Table 4.7: Long Directory entry

Parameter	Offset	Size (bytes)
Inode number	0x0	4
size	0x4	1
Long Filenames Inode Number	0x8	4
checksum	0x12	checksum

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
07F75000	08	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	00	.	
07F75010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F75020	01	00	00	00	02	2E	2E	00	00	00	00	00	00	00	00	00	..	
07F75030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F75040	2B	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	00	+	lib
07F75050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F75060	58	02	00	00	18	66	69	6C	65	66	6F	72	6D	61	74	68	X	fileformat
07F75070	61	6E	64	62	6F	6F	6B	2E	61	73	63	69	69	00	00	00		andbook.ascii
07F75080	59	02	00	00	FF	00	00	00	2B	00	00	00	99	D8	6D	5B	Y	ÿ + %m[
07F75090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F750A0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		

Fig. 4.8: An example of a QNX6 inode entry of a long filename

#### 4.2.5 Long Filenames Inode

If a file or directories length is longer than 27 UTF-8 characters, the name is stored in the long filenames node. Long filenames Inodes start counting with zero. The structure is shown in Table 4.8, an example is Fig. 4.9.

Table 4.8: Long Filenames Inode

Parameter	Offset	Size (bytes)
filename length	0x0	2
filename	0x2	up to 510 bytes

#### 4.3 Example: Construction of a file

To understand how a file can be retrieved from the filesystem data, we will manually find the file /usr/fileformatandbook.ascii with its content and metadata by using the

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
03221000	24	00	66	69	6C	65	66	6F	72	6D	61	74	68	61	6E	64	\$ fileformatand	
03221010	62	6F	6F	6B	76	65	72	79	6C	6F	6E	67	6E	61	6D	65	bookverylongname	
03221020	2E	61	73	63	69	69	00	00	00	00	00	00	00	00	00	00	.ascii	
03221030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		

Fig. 4.9: An example QNX6 long filenames entry

filesystem information. We will begin the reconstruction from the root directory. As already mentioned in the previous chapter, inode 1 contains the root directory. From there, we will start finding the file in the filesystem structure. The first step is to determine the valid stable superblock by the serial number. The superblocks inode root block is shown in Fig. 4.10

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
00002000	22	11	19	68	46	DA	79	9A	23	00	00	00	00	00	00	00	" hFÚyš#	
00002010	1E	00	00	00	43	94	6C	60	00	01	00	00	04	00	03	00	C"1"	
00002020	94	08	BE	35	56	35	4F	2B	8C	24	B2	EB	CB	2A	42	90	" %5V5O+€\$=€€*B	
00002030	00	10	00	00	00	19	00	00	A7	16	00	00	F8	C7	00	00	S øç	
00002040	7E	7F	00	00	01	00	00	00	00	80	0C	00	00	00	00	00	~ €	
00002050	CD	00	00	00	FF	Í yyyyyyyyyyyyyy												
00002060	FF	yyyyyyyyyyyyyyyy																
00002070	FF	yyyyyyyyyyyyyyyy																
00002080	FF	yyyyyyyyyyyyyyyy																
00002090	01	01	00	00	00	00	00	00	FF	18	00	00	00	00	00	00	ÿ	

Fig. 4.10: Inode Root block used in the file reconstruction example

The root block tree has one level, meaning that we go on with the indirect inode block in the next step. The formula can easily calculate the physical address of those blocks:

$$\text{blockaddress} = \text{blocknumber} * \text{blocksize} + \text{offset}$$

On standard QNX6 Volumes, the offset is the superblock size + the offset of the beginning of the superblock. Thus, the first indirect inode block is located at  $0xCD * 0x1000 + 0x3000 = 0xD0000$ , where  $0xCD$  is the block number,  $0x1000$  the blocksize and  $0x3000$  the global offset due to the superblock with size  $0x1000$  and start at  $0x2000$ . From the indirect inode (Fig. 4.11), we can retrieve the number  $0x03$ , and by this, the address of the first inode block, which is located at  $0x6000$ .

The first inode in this block is the root inode. If we take the first block pointer,  $0x7F10$ , of this inode, we get the address of the root directory:  $0x7F13000$ . This root directory, Fig. 4.13 is already familiar to us, as the second version of it is shown in Fig. 4.7, but this time, it is the root directory maintained by the first superblock.

In the root directory, we take the inode number for the /usr directory,  $0x08$ . With this number, we go back to the first Inode Block, where the inode 8 is located at  $0x6380$  ( $0x6000$ , where inode 1 is located plus  $7 * 0x80$  offset, for the preceding inodes). From that inode (Fig. 4.14) we can then calculate the /usr directory offset in the way we already did for the root directory. The /usr directory is defined at block  $0x7F72$

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
000D0000	03	00	00	00	CF	00	00	00	D0	00	00	00	D1	00	00	00	Í	Ð Ñ
000D0010	D2	00	00	00	D3	00	00	00	D4	00	00	00	D5	00	00	00	Ó	Ó Ó Ó
000D0020	0B	00	00	00	D7	00	00	00	0D	00	00	00	OE	00	00	00	*	
000D0030	DA	00	00	00	DB	00	00	00	DC	00	00	00	DD	00	00	00	Ú	Ú Ú Ý
000D0040	13	00	00	00	DF	00	00	00	E0	00	00	00	16	00	00	00	ß	à
000D0050	17	00	00	00	18	00	00	00	19	00	00	00	1A	00	00	00		
000D0060	1B	00	00	00	1C	00	00	00	1D	00	00	00	1E	00	00	00		

Fig. 4.11: Indirect inode block

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
00006000	00	10	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00006010	1E	00	00	00	CC	43	6D	38	10	94	6C	60	0C	44	6D	38	íCm8 "1"	Dm8
00006020	FD	41	09	00	10	7F	00	00	FF	ÿA	ÿÿÿÿÿÿÿÿ							
00006030	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿ																
00006040	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿ																
00006050	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿ																
00006060	FF	ÿÿÿÿ																
00006070	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00006080	00	10	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
00006090	1E	00	00	00	1E	00	00	00	1E	00	00	00	0C	44	6D	38		Dm8

Fig. 4.12: inode 1 which contains the pointers to the root directory

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
07F13000	01	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	00	.	
07F13010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13020	01	00	00	00	02	2E	2E	00	00	00	00	00	00	00	00	00	..	
07F13030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13040	02	00	00	00	05	2E	62	6F	6F	74	00	00	00	00	00	00	.boot	
07F13050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13060	03	00	00	00	03	62	69	6E	00	00	00	00	00	00	00	00	bin	
07F13070	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13080	04	00	00	00	00	03	65	74	63	00	00	00	00	00	00	00	etc	
07F13090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F130A0	05	00	00	00	04	69	6E	66	6F	00	00	00	00	00	00	00	info	
07F130B0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F130C0	06	00	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	lib	
07F130D0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F130E0	07	00	00	00	03	6F	70	74	00	00	00	00	00	00	00	00	opt	
07F130F0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13100	08	00	00	00	03	75	73	72	00	00	00	00	00	00	00	00	usr	
07F13110	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
07F13120	1C	00	00	00	08	66	6C	61	73	68	2E	73	68	00	00	00	flash.sh	
07F13130	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		

Fig. 4.13: Root Directory

which is at offset 0x7F5000. Here we see now our filename and the corresponding inode Number, where the metadata and pointer to the file content is.

We see that the *fileformatandhandbook.ascii* file has the inode number 0x258. Knowing this, we have to find the offset where this inode is defined. With a block size of 0x1000 and an inode size of 0x80, each inode block contains 0x20 inodes, so the inode we are looking for is the 24th inode in inode block number 19. Going back to Fig. 4.11, the 19 inode block is at physical block 0xE0, calculated address 0xE3000

Offset	Title	Value
6380	Size	00 10 00 00 00 00 00 00
6388	Uid	00 00 00 00
638C	Gid	00 00 00 00
6390	File time	01.01.1970 00:16:56
6394	Mod. time	06.04.2021 16:58:53
6398	Access time	06.04.2021 17:02:11
639C	Change time	06.04.2021 16:58:53
63A0	Mode	ED 41
63A2	ExtMode	03 00
00006370		00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00006380		00 10 00 00 00 00 00 00 00 00 00 00 00 00 00 00
00006390		F8 03 00 00 4D 93 6C 60 13 94 6C 60 4D 93 6C 60
000063A0		ED 41 03 00 72 7F 00 00 FF FF FF FF FF FF FF FF FF
000063B0		FF
000063C0		FF
000063D0		FF
000063E0		FF FF FF FF 00 03 00 00 00 00 00 00 00 00 00 00 00
000063F0		00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

Fig. 4.14: Inode 8, which has the pointer to the /usr directory in our example

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
07E75000	08	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	00	.
07E75010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07E75020	01	00	00	00	02	2E	00	00	00	00	00	00	00	00	00	00	..
07E75030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07E75040	2B	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	00	+ lib
07E75050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07E75060	58	02	00	00	18	66	69	6C	65	66	6F	72	6D	61	74	68	X fileformat
07E75070	61	EE	64	62	6F	6F	6B	2E	61	73	63	69	69	00	00	00	andbook.ascii
07E75080	59	02	00	00	FF	00	00	00	2B	00	00	00	99	D8	6D	5B	Y ÿ + "m[
07E75090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	

Fig. 4.15: /usr directory with the entry of the file we are looking for

+ 0xB80 (24th inode in Block). In this inode, depicted in Fig. 4.16 we find all the relevant filesystem metadata for this file and the pointers to the filesystem content.

Following now the pointers to the content, beginning with 0x19D, we can retrieve the file block by block (Fig. 4.17).

After demonstrating the retrieval of the example file from the file system data, it is easy to understand the next section, which shows the possibilities to reconstruct deleted files.

## 4.4 Deleted Files

There are some possibilities to recover deleted files in a QNX6 Volume, depending, when the file or directory was deleted and what happened with the filesystem in the meanwhile. Deleting an entry (directory or file) in QNX6 means that the Status in

Offset	Title	Value
E3B00	Size	9C 16 00 00 00 00 00 00
E3B88	Uid	00 00 00 00
E3B8C	Gid	00 00 00 00
E3B90	File time	06.04.2021 16:57:31
E3B94	Mod. time	06.04.2021 17:02:39
E3B98	Access time	06.04.2021 17:02:56
E3B9C	Change time	06.04.2021 17:02:39
E3BA0	Mode	FD 81
E3BA2	ExtMode	01 00

Blockptr	
E3BA4	BlockPtr 0
E3BA8	BlockPtr 1
E3BAC	BlockPtr 2
E3BB0	BlockPtr 3

000E3B60	FF FF FF FF 00 03 00 00 00 00 00 00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	YYYY
000E3B70	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	00
000E3B80	9C 16 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	0'1`/`1`@`1`/`1`
000E3B90	FB 92 6C 60 2F 94 6C 60 40 94 6C 60 2F 94 6C 60	40 94 6C 60 2F 94 6C 60 FF FF FF FF FF FF FF FF	Ý 2 YYY
000E3BA0	FD 81 01 00 9D 01 00 00 1C 32 00 00 FF FF FF FF FF	FF	YYYYYYYYYYYYYYYY
000E3BB0	FF	FF	YYYYYYYYYYYYYYYY
000E3BC0	FF	FF	YYYYYYYYYYYYYYYY
000E3BD0	FF	FF	YYYYYYYYYYYYYYYY
000E3BE0	FF FF FF FF 00 03 00 00 00 00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	YYYY
000E3BF0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00	YYYY

Fig. 4.16: Inode entry of our example file

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI	ASCII
0019FFD0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
0019FFE0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
0019FFF0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00		
001A0000	4	68	69	73	20	69	73	20	61	20	54	65	73	74	66	69	This is a Testfi	
001A0010	6C	65	20	66	6F	72	20	74	68	65	20	46	69	6C	65	20	le for the File	
001A0020	46	6F	72	6D	61	74	20	68	61	6E	64	62	6F	6F	6B	2E	Format handbook.	
001A0030	20	54	68	69	73	20	54	65	73	74	66	69	6C	65	20	6A	This Testfile j	
001A0040	75	73	74	20	72	65	70	65	61	74	73	20	74	68	65	20	ust repeats the	
001A0050	73	61	6D	65	20	74	65	78	74	20	6F	76	65	72	20	61	same text over a	
001A0060	6E	64	20	6F	76	65	72	20	61	67	61	69	6E	2E	20	54	nd over again. T	
001A0070	68	69	73	20	69	73	20	61	20	54	65	73	74	66	69	6C	his is a Testfil	
001A0080	65	20	66	6F	72	20	74	68	65	20	46	69	6C	65	20	46	e for the File F	
001A0090	6F	72	6D	61	74	20	68	61	6E	64	62	6F	6F	6B	2E	20	ormat handbook.	
001A00A0	54	68	69	73	20	54	65	73	74	66	69	6C	65	20	77	61	This Testfile wa	

Fig. 4.17: Content of our example file

an Inode switches to "deleted" (see Table 4.5) and that the entries inode number is deleted from the directory as shown in Fig. 4.18. By this, it is not possible to recover a file by its name, because there is no link anymore between the filename and the inode containing the metadata and the pointers to the file content. If a directory is updated after a file was deleted (e.g. a new file is added), the filesystem driver moves the directory to another block. The filename is "lost" from the regular filesystem

directory tree. Also, the blocks, which contain the content of the files are set to unused in the bitmap, which means, they are free to be overwritten by other data. Knowing this, there are still some possibilities to recover files, with and without their respective names.

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
07FD7FFF0	FF	YYYYYYYYYYYYYYYYYY															
07FD8000	08	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	00	.
07FD8010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8020	01	00	00	00	00	02	2E	2E	00	00	00	00	00	00	00	00	..
07FD8030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8040	2B	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	00	+ lib
07FD8050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8060	00	00	00	00	18	66	69	6C	65	66	6F	72	6D	61	74	68	X fileformat
07FD8070	61	6E	64	62	6F	6F	6B	2E	61	73	63	69	69	00	00	00	andbook.ascii
07FD8080	00	00	00	00	FF	00	00	00	2B	00	00	00	99	D8	6D	5B	ÿ + %m[
07FD8090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	

Sector 261.823 of 409.568	Offset:	7FD7FFF	= 255   Block:														
07FD8000	08	00	00	00	01	2E	00	00	00	00	00	00	00	00	00	00	.
07FD8010	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8020	01	00	00	00	00	02	2E	2E	00	00	00	00	00	00	00	00	..
07FD8030	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8040	2B	00	00	00	03	6C	69	62	00	00	00	00	00	00	00	00	+ lib
07FD8050	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8060	58	02	00	00	18	66	69	6C	65	66	6F	72	6D	61	74	68	X fileformat
07FD8070	61	6E	64	62	6F	6F	6B	2E	61	73	63	69	69	00	00	00	andbook.ascii
07FD8080	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD8090	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD80A0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	
07FD80B0	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	

Fig. 4.18: Directory entry before (bottom) and after (top) deletion

The first possibility, if the file was just deleted recently, it may still be present in the non-active filesystem structure of the second superblock. If this is the case, the file can normally be fully recovered, even with its content (still, it is possible that the content is not original).

Second, you can parse the inodes to recover files with their metadata without the associated filename. This fact is quite problematic because the Blocks do not necessarily still contain the files original data.

In conclusion, we see that the reconstruction of files is sometimes possible. However, compared to some other filesystems (e.g. NTFS), there is a smaller possibility to recover deleted files from the filesystem information. In some special cases where you can prove the integrity of a file in another way (e.g. some packed/zipped files), it is still helpful to take advantage of the inode structure and the possibility to put together fragmented files from the pointers inside the inode.

## 4.5 Forensic Tools supporting QNX6 filesystems

The Linux kernel includes a read-only driver for QNX6 (and QNX4) file systems. Also, some mobile forensic tools like UFED physical analyzer support this file system to a certain degree. Until today, those tools just support volumes formatted with the standard values shown in Table 4.1. Lately, there have been some projects in the Autopsy / Sleuthkit community to support QNX6, but until today, none of the projects has come to an end.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



## **Part II**

# **Mobile File Formats**

File format analysis examines one specific file. An App or a program typically interprets the data contained in a file. Files can contain user data as well as configuration data, caches or any data. One aim of file format analysis is to recover corrupted files or restore deleted entries from files that only mark entries as deleted but do not overwrite the deleted data.

In this part of the book, the general design of five common file formats in mobile devices is described, and different analysis techniques are presented. This part abstractly approaches the topic and is not limited to how a specific tool analyzes a file format.

File formats are used by Apps or programs and provide mechanisms to store data in a structured way. File formats can organize metadata as well as data such that the specific App can use this. The described file formats in this part are typically used in mobile devices.

# Chapter 5

## SQLite



Dirk Pawlaszczyk

**Abstract** SQLite is, without doubt, the most widely used database system worldwide at the moment. The single file database system is used, among other things, in operating systems for cell phones, such as Android, iOS or Symbian OS. On a typical smartphone, we usually find several hundred SQLite databases used by a wide variety of apps. Due to its widespread use, the database format is of particular importance in mobile forensics. It is not uncommon for the suspect to try to cover his tracks by deleting database content. Recovering deleted records from a database presents a special challenge. In this chapter, the on-disk database format of the SQLite database system is highlighted. Therefore, we take a closer look at the database header as well as record structure on a binary level. We first examine the structure of the data. Recovery options for erased records are discussed as well. Special attention is paid to the slack areas within the database: unallocated space, Freelist as well as free blocks. In this context, we discuss basic techniques for carving and acquisition of deleted data artefacts. Despite the main database format and recovery options, temporary file types like write-ahead logs and rollback journals are analyzed as well.

### 5.1 Introduction

A large amount of data is being stored and processed in relational databases. The most widely used database system in the world is undoubtedly SQLite since it is the default solution for the Android and iOS operating systems. So it is not surprising, that web browsers, messenger services and mobile applications employ the free and serverless database solution as their storage format of choice [61],[60]. At the moment, there are more than a trillion SQLite instances in active use [81]. In the vast majority of criminal investigations involving information technology, one task is to make information stored in such databases accessible. Evidence acquisition for

---

University of Applied Sciences (Hochschule Mittweida), Technikumplatz 17, 09648 Mittweida, Germany, e-mail: [pawlaszc@hs-mittweida.de](mailto:pawlaszc@hs-mittweida.de)

databases is traditionally made with SQL, a powerful query language. Also, SQLite supports most of the SQL language commands. In this way, the data can be accessed with one of the freely available viewers. Unfortunately, this form of analysis usually does not allow access to deleted records or temporary data content such as recently added but not committed entries. This creates the need for alternative ways to analyze such databases forensically.

## 5.2 The SQLite File Structure

SQLite is a single-file database engine, i.e., all tables are managed in only one file on disk. There is no intermediary server process; an application has to communicate with first, for storing data. It does not work this way. Instead, the database can be integrated directly into an application. Therefore, it provides a library and an easy to use programming interface. This fact has significantly contributed to the current spread and popularity of the program. We will discuss the basic structure of a database before turning to the details of carving for data records.

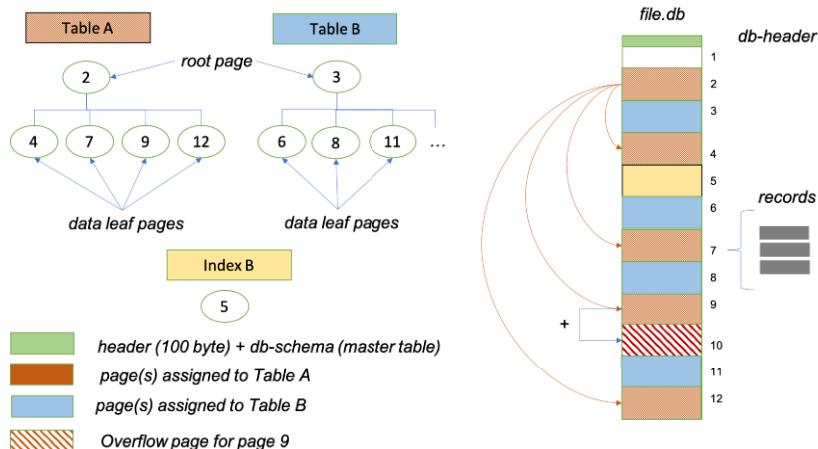


Fig. 5.1: Schematic structure of a SQLite database

Like most structured binary formats, the database file starts with a header part [80]. Its size is exactly 100 bytes. Beyond this, the database file is divided into pages of equal size. The file size is thus always a multiple of the page size. A page number uniquely identifies a single page, whereas the first page has the number one. The default page size usually is 4096 bytes. However, it can be adjusted if necessary to a minimum of 512 bytes and a maximum of 64KB [32]. Of course, the header is part of the first page. In a relational database system, all data is stored in tables. This is also the case with SQLite. In turn, a table is distributed over one or more pages of

the database on the binary level (see Fig. 5.1). Each data page again contains one or more records, for precisely one table. To access and acquire all records of a particular table, we must first determine which pages of the database are associated with this table. This information can again be taken from the first page of the database. Besides the header string, this page contains one more piece of information - the database schema. Necessary information such as the root page numbers, column names, and column types of the tables are stored here, in a data structure called *SQLite\_Master Table*. We will discuss the details of this table in sect. 5.2.3. To represent a table and its pages, SQLite uses a balanced tree data structure (B+tree) under the hood. In a B+tree, the raw data elements are stored exclusively in the leaf nodes, while the inner nodes contain only links. Since the maximum size of a page is limited from above, we can gain more space for links or branches in the inner nodes by moving the leaves' data records. Moreover, this limits the height of the tree. Since data elements are normally accessed via the tree's root, a lower height reduces the number of nodes to be traversed. Many relational database systems manage their records in this way.

Table 5.1: SQLite page types and byte flags

Page Type	1st Byte in Page
table b-tree interior page	0x05
table b-tree leaf page	0x0d
index b-tree interior page	0x02
index b-tree leaf page	0x0a
overflow page	0x00 (for db-size < 64GB)
freelist page	0x00 (first 8 bytes filled with zero-bytes)
pointer map	0x01 or 0x02 or 0x03 or 0x04 or 0x05
locking page	0x00 (only, if db-size > 1 GB)

A page with links to other pages only is called a *b-tree interior page* [80]. The record nodes are saved in *table b-tree leaf pages*. Beyond this, a table can have multiple indexes. An index contains links to normal table records to speed up searching and sorting by specific fields. Whenever we create an index, SQLite creates a B-tree structure to hold the index data as well. Similar to normal tables we can distinguish between *index b-tree interior pages* as well as *index b-tree leaf pages*. When a data record is too large for a single data leaf page, the excessive bytes are spilt onto so-called *overflow pages*. Several overflow pages are filled at once to store large amounts of data such as Binary Large OBjects (BLOBs). Together all overflow pages for one record form a linked list. To capture all the data associated with a record, we need to read all the pages. The payload for an record and the preceding pointer are combined to form a cell.

Despite the five data page types, SQLite knows three more page classes. A database file might contain one or more pages that are not in active use. Whenever the last record is deleted from a page, this page is released. The freed page will be reused when new pages are required and filled with new table contents. In the

meantime, all unallocated pages are stored in a so-called freelist (sect. 5.3). These freelist pages are of particular forensic value since most of the removed content can be found here.

A further not yet discussed page type are so-called *pointer maps*. A pointer map has the function of not losing track when pages are moved from one position in the database file to another. This page type is created whenever the database is reorganized or cleaned up. A pointer map provides a lookup table to quickly determine page types and their parents. However, this page type exists only in auto-vacuum databases. The *locking page* is the last page type in SQLite. The first page of this page class starts at byte offset  $2^{30}$  (1,073,741,824) and always remains unused. Conversely, this means that a locking page only appears when the database size is more extensive than 1 GB. Since it is empty, it has only a technical, but no forensic value and is therefore not considered further.

We can usually determine the type of page by looking at the page's first byte. The flag-byte at offset 0 indicates the page class. Table 5.4 lists all the page types discussed so far. However, not every database will include all of these types. With the page size and type information at hand, an investigator can walk through the database and identify all areas of interest.

### 5.2.1 The Database Header

Every forensic investigation starts with analysing the file header. The header contains important information that will help us to carve for deleted records. The fields of the header have a precisely defined size and position (see Fig. 5.2). The individual (multi-byte) fields are encoded as big endian (BE) values.

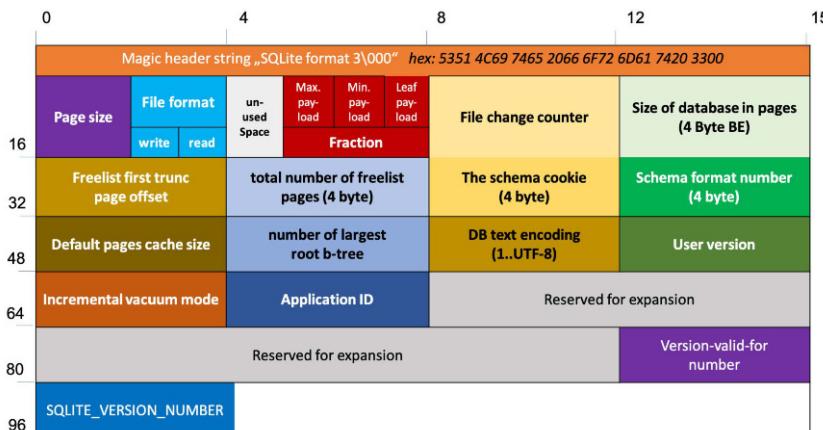


Fig. 5.2: The SQLite Database Header Format and fields

We will discuss the fields below and evaluate them in terms of their respective value for a forensic investigation [80]:

- Each database starts with the header string. The magic header value is always set to "SQLite format 3". We can use the header information to carve the beginning of a database file on the binary level. Offset 15 marks the end of the magic header string. It holds a special character, the null terminator (0x00).
- At offset 16, we can find a two-byte big-endian integer value representing the database's page size. The value in this field must be a power of two. The range of values is between 512 and 32768. There is one exception: The value 0x0001 is viewed as a big-endian 1. It represents the value 65,536 - the largest possible page size - since this number will not fit in a two-byte usually.
- The two flag bytes at offset 18 and 19 control the read and write permission for the database. The values should typically always be either a 1 or a 2. For the rollback journalling mode (sect. 5.4.2), both values are set to 1. In contrast, number 2 in both fields indicates a WAL journalling mode (sect. 5.4.3). If the write version has a value greater than 2, this database file must be accessed as read-only. These two fields' value can indicate whether other files (WAL file or journal file) are present.
- The 1-byte integer value at offset 20 of the header is used to apply for certain SQLite extensions. The number of bytes specified here reduces the usable area within the page. In this way, for example, special salt or nonce values can be stored for each page when using the cryptographic extension. This value is usually 0. The value can be odd.
- The bytes on offset 21 to 23 have fixed values per definition. Maximum and minimum payload fraction must be 64, 32. The byte for the leaf payload fraction always holds the value 32.
- With each transaction carried out on the database, the 4-byte big-endian integer at offset 24 is usually incremented by one. A process that wants to read data from the database can determine whether there has been a change since the last access.
- With the 4-byte integer on offset 28 stores the size of the in-header database in pages. However, this value may differ from the file's actual size when accessing a database before version 3.7.0. Alternatively, you can determine the actual file size and divide by the page size to infer this value.
- At offset 32, we can find a 4-byte big-endian integer which indicates the beginning of the so-called freelist. As already pointed out, unused pages in the database file are stored within this data structure. This field has a significant meaning, as it allows us to access pages of the database that are no longer visible. It holds the offset of the first page of the list. If the value is zero, the list is empty.
- At offset 36 represents the total number of entries on the freelist. Together with the start address, one can thus automatically iterate over the released pages.
- Each change to the database schema, such as adding or deleting a table or creating an index, automatically leads to an increment of the value at offset 40.

- The 4-byte value at offset 44 represents the format number. This field has a value between 1 and 4. For a SQLite database created with the latest version of the database, the value is always 4 and thus supports the more SQL commands. Databases created before November 2005 usually have a value of 3 or less.
- The value default pages cache size at offset 48 queries or sets the suggested maximum number of pages of disk cache for a database file.
- The 4-byte big-endian integer value at offset 52 is only used to manage pointer-maps for auto vacuum-databases. A non-zero value means that this database file contains pointer map-pages.
- All strings in the database are encoded with the same encoding. There are only 3 valid encodings: UFT8 (value 1), UTF16LE (value 2), UTF16BE (value 3). For the analysis of the database, this field value must always be read first.
- The integer at offset 64 is true for incremental\_vacuum and false for auto\_vacuum mode. A value is larger than 0 means that the database reclaims space after data has been deleted. An autovacuum database thus contains few deleted artefacts - if any. It is defragmented automatically.
- The Application ID at offset 68 can be set by the Application programmer. It is not used by SQLite.
- Offset 92 covers the value of the change counter. The integer at offset 92 indicates which transaction the version number is valid for.
- The 4-byte integer at offset 96 stores the SQLITE\_VERSION\_NUMBER value. The version number of the database library with which changes were last made to the database is noted here.

All remaining header bytes are reserved for future expansion. Consequently, we can ignore them.

### ► Important

As can be seen from what has been said, various header fields must be read and analyzed as the first step of every examination. Thus, the page size (offset 16) and the number of pages (offset 28) must always be determined, since we need to know the structure and size of the database. In order to interpret the strings correctly, the encoding must also be examined (offset 58). A look at the freelist entries at offset 32 and 36 tells us whether unused pages in the database exist. If we do not find any references to free pages, it may be an auto vacuum database (offset 64). Using the flags for transaction management at offset 18 or 19, we can also find out which additional SQLite files may exist. This is of particular interest because these files can also contain records of former transactions. Thus, old states of the database have been overwritten in the meantime could be made visible again. The header's remaining information is more technical and is, therefore, less interesting for the investigator.

## 5.2.2 Storage Classes, Serial Types and Varint-Encoding

In order to understand the binary format of records we first need to clarify what data types SQLite knows at the binary level and how they are encoded. Like most other databases, SQLite uses strict typing. Therefore, each value stored is mapped to one of the five storage classes (Table 5.2). The word storage class is just another term for a data type. However, the latter is more commonly used in connection with programming languages. SQLite supports storage classes for integers (INTEGER), floating-point numbers (REAL), strings (TEXT), binary objects (BLOB), and other numeric data such as dates (NUMERIC). The storage class thus determines how the binary data is to be interpreted. Conceptually, each column of a table is assigned with a specific affinity. The affinity denotes the preferred storage class for a column. The data type of a column defines what value the column can hold. However, the SQL standard knows several data type names for one SQLite storage class. For example, there exist more than ten different integer data types in SQL. For texts, there exist seven different types. Accordingly, each data type is mapped to exactly one storage class.

A second essential aspect is a length occupied by a cell value. An integer, for example, will consume a length between zero and a maximum of 8 Bytes. A floating-point number is mapped to a 64-bit field. A text can have an arbitrary length. SQLite uses the so-called serial types to map storage class and length. In simplified terms, this type is a number. The concrete value of the number provides information about the length of a cell value. At the same time, the storage class can be derived from the numerical value. Table 5.3 lists all possible serial types. For serial types 0, 8, 9, the value is zero bytes in length. The serial type is used whenever the type and length of a cell must be determined. Usually, each table row has a corresponding header that summarizes the serial bytes for each column. As a rule, a serial type occupies exactly one byte. Especially with texts or BLOBS, this principle is sometimes deviated from as soon as the numerical value's length exceeds 127. In this case, additional bytes may be added to map the serial type.

Table 5.2: Mapping from SQL types to SQLite storage classes [80]

SQL Data Type	Storage Class
INT, INTEGER, INTUNSIGNED, LONG, TINYINT, SMALLINT, MEDIUMINT, BIGINT, INT2, INT8	INTEGER
TEXT, CHARACTER, CLOB, VARCHAR, NCHAR, NATIVE CHARACTER, VARYINGCHARACTER	TEXT
REAL, DOUBLE, DOUBLEPRESICION, FLOAT	REAL
NUMERIC, DEZIMAL, BOOLEAN, DATE, DTIME	NUMERIC
BLOB ( no datatype specified )	BLOB

SQLite uses a particular encoding for storing serial types. The representation form used is a variable-length integer (varint). SQLite version 3 uses this simple byte-oriented encoding where each byte contains 7 bits of the integer being encoded. The most significant bit (MSB) is a flag bit, indicating more bytes to follow. Since most integers in a database have relatively small values, we can keep memory consumption low this way. Storing with a fixed-length integer will mostly generate unnecessarily many null bytes. Instead, SQLite uses a static Huffman encoding of 64-bit two's-complement integers that needs less space for small positive values. The serial type varints for large strings and BLOBs might extend up to nine-byte varints. The following illustration should once again make clear the storage principle of varint-values:

1 Byte	0XXXXXXXXX	..127
2 Bytes	1XXXXXXXXX 0XXXXXXXXX	..16384
3 Bytes	1XXXXXXXXX 1XXXXXXXXX 0XXXXXXXXX	..2097152
4 Bytes	1XXXXXXXXX 1XXXXXXXXX 1XXXXXXXXX 0XXXXXXXXX	..268435456

Since texts have a variable size, and the length calculation is performed by a formula. A numerical value above 12 or 13 can only occur with texts or BLOBs. An odd value will be correspondingly for texts. On the other hand, if the value is even, then it is the BLOB storage class. For example, to store the word *Test*, the value  $21(0x15)$  -  $2 * \text{text length} + 13$  - is stored as the length specification. A JPEG file with, let us say, the length of 109 Bytes would be encoded with the serial type number 230 since  $N * 2 + 12$  is what we need to calculate for a binary object. However, since we cannot map this value with 7 bits, we have to add a second byte for the varint:

```
decimal: 230 = 128 + 64 + 32 + 4 + 2
binary: 1110 0110
varint: 1000 0001 and 0110 0110 (2-Byte: 1X.. 0X..)
```

Thus, we must first calculate the respective length specification each time we need to know the exact length of a table cell. The serial values 8 and 9 are noteworthy features. They can be used to map the two values 0 or 1. An extra content byte is not necessary in this case. With the information presented, we are now able to decode the cells of a table row.

### 5.2.3 Decoding The SQLite\_Master Table

A database schema is a set of data definitions that define the structural design of a database. As already explained, the schema, or the master table, resides on the database's first page, just behind the header. Technically, it is a regular table [85]. Table 5.4 shows all columns and their meaning for the master table. The schema table contains all database objects in the database and the statement used to create each object. With the schema table's help, all table names, the corresponding column names and data types can be determined. Each table entry is opened by two additional fields: the *rowid* and the payload (see Fig. 5.3). Both values are only visible on the

Table 5.3: Serial Type Codes Of The Record Format [80]

Serial Type	Size	Meaning
0	0	Value is a NULL.
1	1	A 8-bit twos-complement integer.
2	2	A big-endian 16-bit twos-compl. integer.
3	3	A big-endian 24-bit twos-compl. integer.
4	4	A big-endian 32-bit twos-compl. integer.
5	6	A big-endian 48-bit twos-compl. integer.
6	8	A big-endian 64-bit twos-compl. integer.
7	8	A big-endian 64-bit floating point number.
8	0	integer 0 (schema format $\geq 4$ ).
9	0	integer 1 (schema format $\geq 4$ ).
10,11	variable	Reserved for internal use. Variable size.
$N \geq 12$ , even	$(N-12)/2$	Value is a BLOB with $(N-12)/2$ bytes length.
$N \geq 13$ , odd	$(N-13)/2$	Value is a string in the text encoding and $(N-13)/2$ bytes in length. The nul terminator is not stored.

binary level. Any row of the master table and therefore every database object is assigned to a unique, non-NULL, signed 64-bit integer - the *rowid*. This value is used as the access key for the data in the underlying B-tree. On the binary level, each table row starts with a rowid number greater than null. Most tables in a typical SQLite database schema are rowid tables. A *rowid table* is defined as any table in an SQLite schema that is not a virtual table and is not a WITHOUT ROWID table. The rowid is not part of the table definition. A payload field that stores the length of the record follows directly after the rowid.

Table 5.4: Structure of the sqlite\_master table [85]

Column Name	Description
type	type of database object (table, index etc.)
name	name of the database object
tblname	table that the database object is connected to
rootpage	root page
sql	SQL statement used to create the database object.

Interestingly, we can find descriptions for tables that have already been removed. If an object in the database is erased, the schema table's corresponding record is marked as removed. If a table is dropped, the rowid value for the line in question is set to 0x0000. The entry that is no longer needed is only overwritten when a new database object is added. In the meantime, the entry is still accessible. Figure 5.3 shows an example of a deleted entry for a table in hex mode. The table header and

all columns of the record are intact. Only the rowid value at Offset 3935 has been wiped with zero bytes.

In the example below, the signature 0x7461626C65 represents the object type of a table. The table name, i.e. "users", directly follows the type column. However, we must parse and analyse the corresponding SQL statement from the fifth column to get all column names and the corresponding type information.

Root page of this table	Row ID	Length of payload (i.e. 161 Bytes)	Type of database object	Object name
3920	42272046 4C4F4154 204E554C 4C0A2900 0000A117		B' FLOAT NULL ) .	
3940	17170182 1B746162 6C657573 65727375 73657273		, tableusersusers	
3960	02435245 41544520 5441424C 45207573 65727320		CREATE TABLE users	
3980	280A0927 69642720 494E5420 554E5349 474E4544		( `id` INT UNSIGNED	
4000	204E4F54 204E554C 4C2C0A09 276E6160 65272054		NOT NULL, `name` TEXT	
4020	45585420 4E4F5420 4E554C4C 2C0A0927 7375726E		EXT NOT NULL, `surname` TEXT NULL, `code` INT NULL, `cod	
4040	616D6527 20544558 54204E55 4C4C2C0A 0927636F		e` FLOAT NULL )	
4060	64654127 20494E54 204E554C 4C2C0A09 27636F64			
4080	65422720 464C4F41 54204E55 4C4C0A29 00000000			

SQL statement used to create the database object.  
Columns(0x171701821B):  
String(5 bytes), String(5 bytes), String(5 bytes), Integer (1 byte), String(135 bytes)

Fig. 5.3: Record of a dropped table from the sqlite\_master (example)

By analyzing the SQL statement, a storage class can be derived for each table column. For the five columns of the table <users>, the following columns can be identified: INT, TEXT, TEXT, INT, REAL. This type of vector can be considered as a kind of fingerprint. Sometimes, a found record could be recovered, but it is not clear to which table it belongs. With the help of the table's signature derived in this way, an assignment can still be made, even for a deleted record. Of course, this rule is not always 100% accurate. It is not excluded that two tables have the same signature. However, it can help us make an educated guess, which will be correct in most cases.

## 5.2.4 Page Structure

All records are stored on pages. Approaching the data of a table requires a leaf page scan. To access the data, we must understand the structure of a page. Each page starts with a header, with a total size of 8 bytes in the case of a data leaf page (see Fig. 5.4). All header bytes are big-endian values. The header starts with the page type at offset 0. In the case of a leaf page, the page starts with the value 0x0D. It can be classified from the other pages by reading this value. The 2-byte value at offset 1 marks the beginning of the first free block on the page. A free block is created whenever a record is deleted from the database. All free blocks are organized as a linked list,

whereas the first two bytes of the free block point to the offset of the following free block within the list. If the free block is the last on the chain, this value is zero. If we want to identify deleted records, our search should start right here in the free block list [80].

page header (8 bytes)				
page type (1 byte)	first free block on page (2 bytes)	number of cells on page (2 bytes)	start-offset of cell content area (2 bytes)	fragmented bytes (1 byte)

Fig. 5.4: Fields of a b-tree leaf page header

Another essential value is located directly behind the free block field at offset 3. The 16-bit two's-complement integer field is called *number of cells*. Its value indicates how many active cells exist within the current page. In SQLite, the serial type header and the values of a particular table row are combined into a structure called "cell". So if we want to access a record, we need to locate the matching cell. Fortunately, all cell offsets are stored in an array directly after the page header. Hence, to read a regular record of a table, we need to iterate through the cell pointer field.

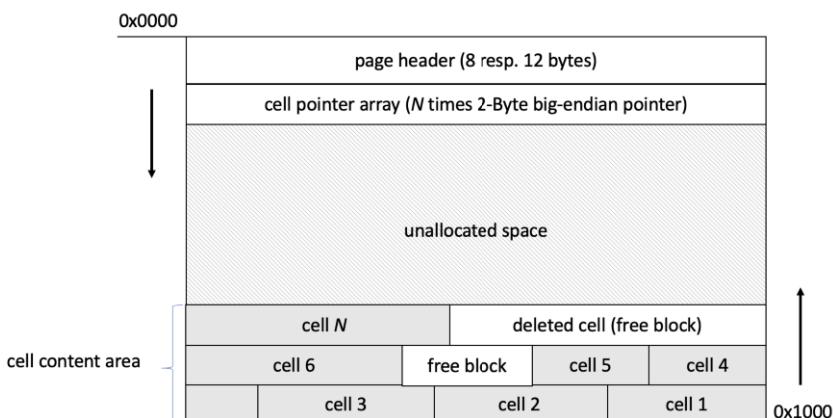


Fig. 5.5: Structure of a regular data leaf page (permanent and temporary)

The next header field at offset 5 provides the start-offset of the content area. A b-tree leaf page is divided into regions (see Fig. 5.5). The cell content area is always located at the bottom of the page. The header and the cell pointer array are always located at the beginning of the page. Between them resides the unallocated space. As the content area grows from the highest memory address towards the lower address,

overlapping the two mentioned regions is prevented. The concept is thus similar to the management of heap and stack areas within memory management. The last value in the header denotes the number of fragmented bytes. A free block requires at least 4 bytes of space. Areas between 1 to 3 bytes form a fragment and thus cannot hold any data records.

Figure 5.6 shows an example of the header of a page on a binary level. In addition to 15 cells, we can also find at least one free block of offset 3620(0x0E24). The content area in this example starts at 0x0DEC. The cell pointer array is highlighted in yellow. Interestingly, we can find five more cell pointers shown in red. The value of the surplus cell offsets corresponds to the start offset of the cell content area. From this, we can conclude that apparently, five other records must have existed on the page in the past. Nevertheless, they have been deleted in the meantime. Thus, in addition to the 15 regular records, there should be five more deleted records on the page. However, the deletion turned the cells into free blocks. So, to find and restore them, we need to examine each element of the free block list.

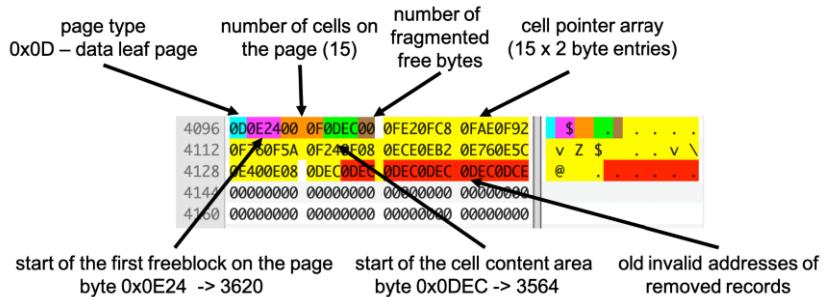


Fig. 5.6: Sample header and cell content array for a data leaf page

It is not always possible to find all deleted records by checking the free blocks. If a record is deleted that resides directly at the unallocated area, the offset value for the cell content area start is moved up in the direction to a higher address. Of course, this address denotes the cell pointer offset of the next regular record. The data set is thus moved to the unallocated area by changing the border. We must consider this case in our search since this record will never appear in the free block list.

However, it gets even worse. If a complete page is deleted, SQLite typically wipes the first 4 Bytes of the header with zeros. So, in this case, the offset for the first free block is erased. Thus, we do not know where precisely the list begins. What does this, in turn, mean for our search for hidden records? The best way to approach our search for slack areas is to use the exclusion principle. *Slack space* is the leftover storage that exists on a page when records do not need all the space which has been allocated. Slack areas are always created when records are deleted. Hence, the total amount of slack space can thus be calculated as shown in the equation below.

```
slack space with (possible) deleted content = page content
    - header (8 bytes)
    - N times 2-Byte cell pointer
    - fragmented bytes
    - N times cell
```

If we exclude the regular, well-known areas of the page, we automatically access the slack areas. Only the areas determined in this way can contain deleted data artefacts. In any case, we must always consider the unallocated space and the free block list when searching within the page. Fortunately, leaf pages are always structured the same. However, there is a second type of leaf page, the index leaf page. In its structure, this page corresponds to a regular data leaf page, except for one difference. The index leaf page starts with the value 0x0A at offset 0. However, what has been said so far also remains valid for the second type of page.

### 5.2.5 Recovering Data Records

Now that we know the location of the records, we can start reading them. This information can be derived from the cell offset array (see the last section). Every cell has the same structure (see Fig. 5.7). The cell header opens with a payload value. It indicates the total size of the cell in bytes. This value does not include the cell header itself. Normally, the payload field is followed by the rowid (see sect. 5.2.3). As already explained, the pseudo-column is usually generated automatically by SQLite. It is used to enable efficient access via the table tree. However, not all records have a rowid. For example, index records are created without this field. If the option "WITHOUT ROWID" is part of the CREATE TABLE statement, this field is also missing. Thus, the cell header has a minimum size of 1 byte for a mandatory payload value. The values in the cell header and all other header fields are varint values without a fixed size. So to read a record, we always have to read value by value. Skipping or omitting bytes is not possible because the fields do not have a fixed offset. The actual cell starts again with a header. This time, it is the header of the data record.

The *header size* field indicates how many bytes the header contains. Its value includes the actual header size byte. The individual serial types follow immediately. Column by column, we must first determine the storage class and space for each table cell. The header is followed directly by the actual data record. Since we operate on a binary level, the exact length of each field to be read and the data types can only be determined via the serial bytes in the header. However, it might be challenging to determine the exact beginning or end of the column cell values without this information. An intact header is, therefore, an essential prerequisite for successful data recovery.

## > Information

The recovery of deleted data depends on the data management policy used. This, of course, differs from application to application. We can distinguish three cases:

1. **Wipe with zeros.** The free block is completely overwritten with zero-bytes. Recovery of data is impossible even if the removed area is identified.
2. **Truncate or remove deleted area.** The second policy is made on a small size of data. It deletes the record itself, and there is no way even to trace the occurrence of deletion. Some iPhone system files are handled this way.
3. **Add to a free list.** The last policy is to mark the record or page as free. The data itself remains in the database. This procedure generates the least I/O-traffic compared to the other two strategies. It is therefore used as the default behaviour of SQLite.

In the case of a data record that has been deleted, it sometimes happens that the cell header and parts of the record header are replaced with new information [59]. These new data fields cover the free block's length in bytes and the address of the following free block. Since both pieces of information are mapped to a 16-bit fixed-length integer, a total of four bytes of the respective cell are overwritten. In total, we can discern six situations when dealing with a deleted record (see Table 5.5).

Many records are deleted without being marked or overwritten. As explained earlier, some records are deleted by merely moving the cell content area's border upwards. Thus, the records slip into the unallocated area of the page. When clearing the browser cache, for example, almost all entries are removed from a caching table. Instead of first marking each record as deleted, the links to the affected pages are deleted from the table tree. Anything else would be a time-consuming process. Instead, the page as a whole is skipped. In both cases, however, the deleted records remain intact. Complete reconstruction is, therefore, possible. Sometimes a record is removed from the middle of the content area of an active page. In this case, the record

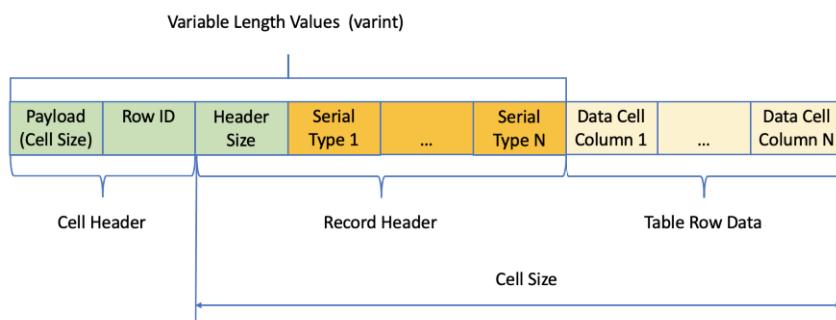


Fig. 5.7: Schematic structure of a data leaf cell

Table 5.5: Recovery Situations

Wiped Data	Recoverability
cell is intact (no wiped bytes)	yes
payload bytes	yes
payload bytes + rowid	yes
payload bytes + rowid + header length	yes
payload bytes + rowid + header length + 1st serial	partly
two or more serial type are wiped	no

is converted to a free block. Thus, the beginning is overwritten, at least partially. The previously occupied space will be released for reallocation. This, in turn, can result in different cases that influence the recoverability of the data record. Sometimes only the payload got wiped. In another case, the payload field, together with rowid, may be overwritten. We can mostly do without this field information. As long as the rest of the cell record remains intact, we can read the required column lengths and types and correctly interpret the data. Even a wiped header length field should not be a big problem. This field only holds the total length of the header. It can be reconstructed by summing the individual serial lengths. It gets tricky when columns are also overwritten. Without a valid column type and length specification for our first column, we cannot reconstruct the remaining columns correctly. However, the first column of a table is often an ID column with a numerical value. Knowing the length of the first column of a regular record on the same page can indirectly infer the first column's length for our destroyed record. Unfortunately, this rule does not work in every case. For example, if the first column contains a text with variable length, we will most likely not restore the record correctly. If more than one serial type has been overwritten, reconstruction seems unlikely. We then have too many possible lengths to consider. Strictly speaking, the number of possible lengths for a column grows exponentially with the number of overwritten length or type information in the header.

Figure 5.8 shows the content area of a data leaf page. There are a total of three records on the page. The cells are located at the end of the page. Remember, the cell content area always grows from higher towards the lower address. The record in the middle is deleted. The records before and after it are intact. Cell header, record header and all data are unaltered. Even without knowledge of the table, it can be deduced from the serial types alone that it is a table with apparently two columns. The first column can store integers (serial types 0x02 resp. 0x03). The second column is a string since the value is odd and greater than 13 (see sect. 5.2.2).

We can see that the second of the three data cells have been deleted because the first 4 bytes of the data set have been overwritten with the free block identifier. The identifier is 0x0000000C. The first two bytes have the value 0x0000. From this, we can conclude that it is the last free block within the page. The second half of the identifier tells us something about the length of the free block. It is exactly 12 bytes (0x000C). The free block is outlined in red in the illustration. As we can see, the actual

8120	00000000	00000000	00000000	00000000	00000000	00000000	00000000	00000000
8148	00000000	000A0303	0315019C	4052756E	65	000000	0C177530	.@Rune u8
8176	43687269	73	090103	02154E20	4469726B	Chris	N Dirk	

Payload Length	ROWID	Header Length	Serial Typ(es)	Column 1	Column 2
0x0A 10 Byte	0x03	0x03	03 15 INT STRING	0x019C40 -> 105536	0x52756E65 „Rune“
XX	XX	XX	XX 17 ?? STRING	0x7530 -> 30000	0x4368726973 „Chris“
0x09 9 Byte	0x01	0x03	02 15 INT STRING	0x4E20 -> 20000	0x4449726B „Dirk“

Fig. 5.8: Example data page with three records (one is wiped)

data fields of the deleted record are still intact. However, the PPL-field, ROWID, header length byte, and the first column's serial type are no longer accessible. The serial type of the second column is not wiped. From the length specification of the free block and the knowledge about the length of the second column, we can infer the length of the first column in this case. Accordingly, the first column of our data set can only be 2 bytes in size:

```
length of the first column field =
    12 byte (total free block length)
    - 5 (0x15 -13 / 2) (length of text column)
    - 4 (free block identifier)
    - 1 (serial type byte for 2nd column)
```

Thus, we can recover deleted content in many cases, even when parts of the header have been overwritten.

### 5.3 Accessing The Freelist

As soon as the last record on a page is deleted, it is transferred to the free list. At the same time, the link within the table tree is removed. From now, the page cannot be accessed from an active table. However, it can be assigned to a new table at any time. Meanwhile, the content of the page is still accessible. Usually, it is not wiped or replaced with random values. The pages are just sitting on the free list, waiting to be used again. Like the slack areas in the standard database pages, these unused pages may contain forensically exciting values such as chat protocols, short messages, or web pages visited [61].

The freelist is a simple linked list consisting of *trunk pages* 5.9. Each trunk page initially contains a 4-byte integer pointer referencing to the next trunk page in the list

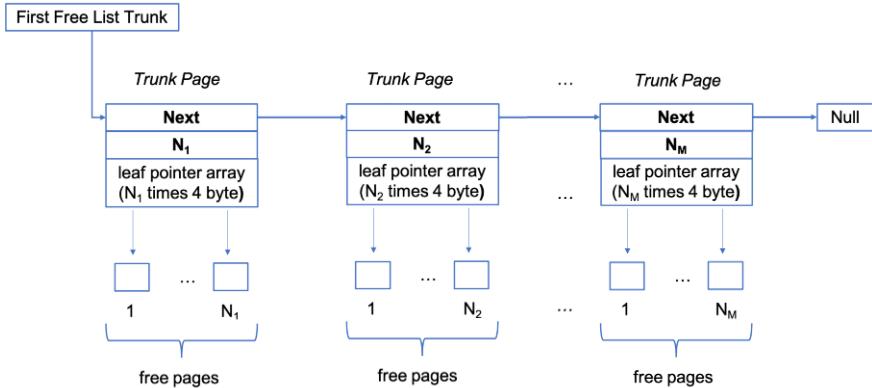


Fig. 5.9: Schematic principle of a freelist trunk list

[80]. A zero-byte value means that this is the last trunk page in the list, and the list ends here. The second 4-byte value in a trunk page contains the number of leaf page offsets. To analyse all the list pages, we must first visit each trunk page and query the offsets stored. Nevertheless, where do we have to start our search for freelist treasures?

The starting address for the freelist can be calculated very easily [59]. We must first determine the start offset of the first trunk page from the header at offset 32 of the database. Second, we need the page size. The latter can also be determined from the header. From these two values, we can calculate the actual offset of the first trunk page:

$$\text{offset of 1st trunk page} = (\text{trunk page number} - 1) * \text{page size}.$$

A trunk page consists of an array of 4-byte big-endian integers. As pointed out, the first 4 bytes of the trunk page header references the next trunk page within the list. The next, a four-byte big-endian integer holds the length of the leaf pointer array of the current page. With these two pieces of information at hand, we can quickly iterate over the array's entries.

The basic algorithm is shown in Listing 1. An example of a trunk page will illustrate what has been said so far (see Fig. 5.10). In addition to the reference to the next trunk page at offset 0, the number of page pointers to follow is visible (offset 4). The offset of the first free page can be found directly behind the two header integers at Offset 8. The second pointer is exactly 4 bytes behind. In the example, there are a total of 555 entries on the TrunkList page. The data size is therefore  $8 + 555 * 4 = 2228$  bytes. Thus, all unused pages can be found and accessed with linear time complexity with the described algorithm.

**Algorithm 1** Freelist Page Recovery

---

▷ Input: SQLITE db filepointer

```

1: read pagesize ← 4 byte BE on byte 0x10
2: read trunk ← for the first freelist trunk on byte 0x20
3: while trunk ≠ null do
4:   start = (4 Byte BE in offset - 1) * pagesize.
5:   db.seek(start)                                ▷ go to start of the trunk page
6:   read trunk ← for the next freelist trunk page (4 Byte BE)
7:   read length ← number of cell entries (4 Byte BE)
8:   for j = 0, 1, ..., length - 1 do           ▷ iterate over trunk page array
9:     db.seek(start + 8 + (4 * j))
10:    read freepage ← next free page number
11:    fpstart = (freepage - 1) * pagesize.
12:    db.seek(fpstart)                          ▷ go to start of next free page
13:    readPage()                                ▷ start analyzing the hidden page
14:   end for
15: end while

```

---

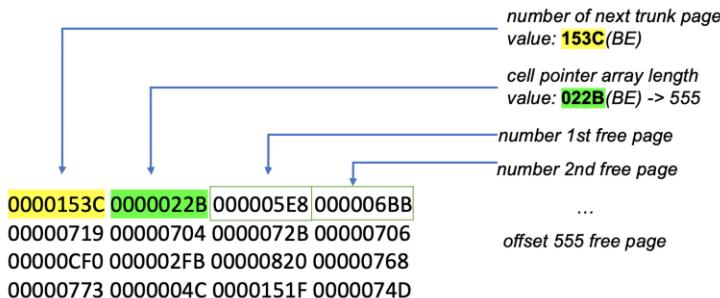


Fig. 5.10: start of a freelist trunk page (example)

## 5.4 More Artefacts

As explained earlier, SQLite manages all records in a single database file. However, access management, transaction handling, and integrity protection are performed with the help of primarily additional temporary files [84]. Despite the main database file, SQLite uses nine distinct types of temporary files (see Fig. 5.11). Below we will take a look at the other file types of SQLite. The focus is on searching for records no longer in the regular database but can still be found in one of those files.

### 5.4.1 Temporary File Types

SQLite creates several temporary files when managing the database. A *transient database*, for example, is a temporarily created file when the database is reorganized.

Data pages that are no longer required are removed. The whole process is comparable to the defragmentation of a hard disk. Pages are joined together, and gaps are closed. Then, the temporary file's content is copied back into the original database file, and the temporary file is deleted. However, this file type is generated only for databases for which the VACUUM property is activated. Since the database copy is deleted immediately afterwards, it is not easy to locate it on the disk. However, it might be possible to find old page versions of the database on the medium through carving. From time to time, SQLite makes use of *transient indices*. Each index is therefore stored in a separate temporary file. For example, if the ORDER-BY or GROUP-BY clause is used in an SQL statement, a corresponding index file is created to manage the intermediate results. The index is automatically deleted at the end of the statement that uses it.

In the case of complex SQL statements, partial queries are sometimes stored in a temporary file. In SQLite, this method is called "materializing" the subquery. This is the case, for example, with large SQL INNER JOIN statements. The query optimizer decides for which query a separate swap file is created.

Database users can create a temporary table using the "CREATE TEMP TABLE" command. Since this unique table is created only for a particular database connection and is not visible to other database users, it is swapped out to a separate file. Again, the temporary database file used to store temporary tables is removed automatically when the database connection is closed. When SQLite performs a transaction with multiple statements, a *Statement Journal File* can be used to undo individual steps. Assume that by executing a statement, 100 rows of a table are modified. After half of the records have been modified, the execution must be aborted due to an error. The rows of the database that have been modified so far are written back with the statement journal's help. All five of the temporary file formats discussed can contain data or temporary results of the database transactions. However, these data are highly volatile. In most cases, the temporarily stored results are already deleted when the statement is finished. Thus, it is not very likely for an investigator to come into contact with such artefacts. We will, therefore, not consider them further.

There are four remaining file types in SQLite. Unlike the formats discussed so far, these are files that are often encountered when examining a database. These files are *Rollback Journals*, *Write-ahead Logs*, *Shared-Memory Files* as well as *Super Journals*. They can usually be found in the same directory as the actual database file. Admittedly, the data stored in it is also classified only temporary within the official documentation of SQLite. However, the data stored in them is updated or overwritten much less frequently. We almost always find one of these file types. For this reason, these are also listed under the heading *other permanent files* in Fig. 5.11. Thus, the chance to acquire data from these files is much more likely. However, in some cases, the use of one file format excludes the use of the second. For example, the shared memory file and write-ahead log are usually found together. In contrast, the rollback journal is only found in a directory if the first-mentioned files are absent. Of the file formats mentioned above, super journals are relatively rare. The files are created only in transactions where multiple databases are updated simultaneously in an atomic transaction. Accordingly, without a super-journal in place, transaction commit on

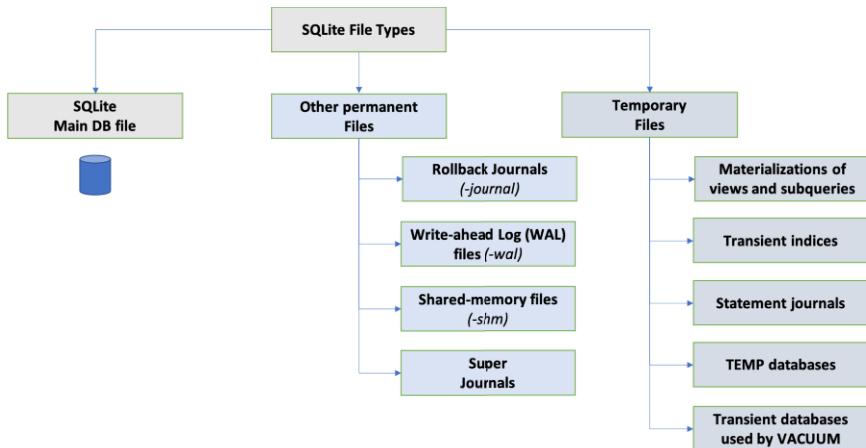


Fig. 5.11: The SQLite file types (permanent and temporary)

a multi-database transaction would be atomic for each database individually, but it would not be atomic across all databases. Due to the relatively low usage level, we will not take a closer look at this file format. Instead, we will focus on the two remaining journal formats. Besides availability and confidentiality, data integrity forms a central goal of every database system. SQLite is no exception. SQLite maintains its integrity by using journals and transactions. Below we will examine the two integrity protection techniques offered by SQLite in more detail: Write-ahead logs and Rollback Journals [80],[32].

#### 5.4.2 Rollback Journals

The idea behind the rollbacks is simple: If a database gets into an inconsistent state due to write access, it is reset to the last valid state. To implement atomic commit and rollback capabilities, SQLite offers a file called *rollback journal*. Rollback refers to resetting the individual processing steps of a database transaction [79]. The system is thus wholly returned to the state before the start of the transaction. In the case of SQLite, a copy is first created for all database pages possibly affected by the transaction and stored in the rollback journal. If something goes wrong during transaction processing, the database can always be reset to the last valid state if required. Note: SQLite permanently stores the entire page in the journal file, even if the transaction modifies only a single record.

A journal file is usually created when a new transaction is started and deleted after the transaction is completed. Although this is the default behaviour, in many cases, there is a deviation from this approach. For example, if the application developer activates the *exclusive locking mode* for a database, then the rollback journal is not

immediately deleted. An application can enable the exclusive locking mode by using the following pragma-statement:

```
PRAGMA locking_mode=EXCLUSIVE;
```

In this case, the journal file may be truncated, or the file's header may be wiped with zero bytes. Which behaviour of this occurs depends on the SQLite version used. However, the file is preserved in any case as long as the locking mode is activated. Fortunately, many applications that use rollback journals for transaction safety operate in this mode, reducing unnecessary IO operations. The same behaviour as is seen in EXCLUSIVE locking mode can also be reached by setting the *journal mode pragma* to PERSIST instead of DELETE which is the default behaviour in SQLite:

```
PRAGMA journal_mode=PERSIST;
```

No matter which of the two modes is activated, an investigator can restore the old execution states of the database. In this way, data records that may have been deleted in the meantime can be made visible again.

### ! Attention

The rollback journal file is always located in the same directory as the actual database. One can quickly identify the journal by the file name: It has the same name as the database but with the extension "**-journal**". Thus, the name of a journal file is precisely eight characters longer than the original name of the database [84].

---

A rollback journal is a binary format. Just like the main database file, it contains a small header. The header has a fixed size of a maximum of 28 bytes. The individual header fields and their meanings are shown in Table 5.6. Next to the Magic Header String, information about the total number of database pages stored in the journal. The header also records the original size of the database file. So if a change causes the database file to grow, we will still know the original size of the database. Unfortunately, the fields carried in the header are usually automatically overwritten after a COMMIT and wiped with null bytes. Thus, we will rarely be able to recover useful information from it. However, the header is usually preserved if a transaction cannot be completed due to a power down.

The journal file has a preset page size. The value can be determined via the offset 20 in the header. Even if this value can no longer be determined due to wiping, there is a way out. The default value of the first sector is 512. The remaining space of the first journal page is filled with zero bytes. Since the default page size is 512 bytes, the header is thus always followed by a padding area of zero bytes. After the header and padding area, zero or more page records will follow. Such a record contains a copy of precisely one database page. Additionally, each record is introduced by a one-field header. Only with this value, SQLite can reset the correct page in the database in case of a rollback. On offset four, the original content of the database

Table 5.6: Rollback Journal Header Format

Offset	Size	Description
0	8	Header string: 0xd9, 0xd5, 0x05, 0xf9, 0x20, 0xa1, 0x63, 0xd7
8	4	The "Page Count" - The number of pages in the next segment of the journal
12	4	A random nonce for the checksum
16	4	Initial size of the database in pages
20	4	Size of a disk sector.
24	4	Size of pages in this journal.

page follows. The journal page record ends again with a 4-byte big-endian value. It holds the checksum for this page. The value is used to guard against incomplete write operations.

Table 5.7: Rollback Journal Page Record Format [80]

Offset	Size	Description
0	4	The page number in the database file
4	N	Original content of the page prior to the start of the transaction
N+4	4	Checksum

Since the header is always reset for each new transaction, the page records directly following the header are always the most current. However, journal records of past transactions can still be stored in the same journal. For example, suppose a transaction changed ten database pages. The following transaction only rewrote five pages. In that case, the database subsequently contains the database's state before the last transaction plus five more pages from the previous. The following example shows the beginning of the second journal page of a rollback file:

```
|0x1200|61746506 BAC4E54E 0000000B 0D000000 |ate....N.....|
|0x1210|0B0E2C00 0F620F35 0FC20F96 0F0E0EFD |...,..b.5.....|
|0x1220|0ED10EB8 0E9A0E5A 0E2C0000 00000000 |.....Z.....|
0xBAC4E54E  -> Checksum of the 1st journal records
0x0000000B  -> page 11 in the database (start of the 2nd journal)
0xD0000000  -> start of a data leaf page (snapshot)
```

The start of 2nd journal record can be calculated as follows:

```
0x0200  1st sector (header + padding area) - 512 byte
+ 0x0004  page record page number (record start) - 4 byte
+ 0x1000  1st page in journal - 4096 byte
+ 0x0004  checksum of 1st journal page (record end) - 4 byte
-----
0x1208  start offset of the 2nd journal record
```

The example shows the end of the first journal page and the beginning of the second journal frame. While the green highlighted value at offset 0x1204 still belongs to the first journal page, the value at offset 0x1208 already initiates the next journal record. Generically, the address of each journal could be determined as follows:

$$\text{Record}_{\text{start}}(N+1) = \text{size of 1st sector} + N \times (\text{page size} + 8)$$

However, how can we determine whether the database's journal page belongs to the last transaction or is not perhaps older? A different random nonce is used each time a transaction is started to minimize the risk that unwritten sectors might by chance contain data from the same page that was a part of prior journals. The last nonce is a 4 Byte integer value and can be found at offset 12 in the journal header. By changing the nonce for each transaction, stale data will still generate an incorrect checksum. Since the entire page is always saved from the database, we can restore the actual data described in section [5.2.5](#).

### 5.4.3 Write-Ahead Logs

As pointed out in the last section, a copy of the data page to be changed is first created before writing directly into the database file in a classic rollback journal [86]. Version 3.7.0 of the SQLite database engine introduced an alternative concept for transaction management [84]. With *write-ahead logs (WAL)*, this procedure is reversed. The content of the original database file is not changed. Instead, every change is appended into a separate WAL file. It works like a roll-forward journal. All changes are first written to the WAL file. Even a COMMIT does not automatically update the database file [79]. If, for example, other reading database connections exist simultaneously, they can operate as usual on the original unaltered data. Meanwhile, a concurrently running write process stores its changes into the WAL file. Moving the WAL file transactions back into the database is called a *checkpoint*. Usually, SQLite does a checkpoint automatically. If the WAL size reaches a threshold size of 1000 pages, a checkpoint is triggered by default. As soon as we examine a database that works in WAL mode, we must also analyse the included WAL archive. Simultaneously, this also means that we may have different versions of the same database page in the main database and the WAL file. As long as no checkpoint has been carried out, the WAL file exclusively contains the latest changes. The database is, therefore, still in an old state. If we look at both files together, we can get a consistent view [86].

To access the content of a WAL file, all we have to do is open the corresponding database file. When opening a WAL mode database, the WAL file's content is automatically transferred back to the database. In other words, a checkpoint is executed. However, this procedure is usually not recommended for various reasons. With this approach, old artefacts that are evidentially valuable to the investigator could be overwritten and thus lost. Moreover, we would be violating a fundamental rule of any forensic investigation: Never change the evidence.

## > Important

It is best not to work with a standard database viewer when evaluating a database in WAL mode. Even by opening the database, one risks losing old data due to checkpointing.

But how should we proceed then? One possibility is the use of a special forensic database browser. An example would be the FQLite<sup>1</sup> browser. This program reads the database and the WAL file separately. Since access is read-only, all data is preserved.

## ! Attention

A particular database will use either a rollback journal or a write-ahead log. It is not possible to use both at the same time. The write-ahead log is always located in the same directory as the actual database. One can quickly identify the journal by the file name: It has the same name as the database but with the extension "**-wal**".

Let us now turn to the actual structure of the file. The WAL file starts with a header. Zero or more so-called WAL-frames follow it. Just as with the rollback journal, a frame represents the altered content of exactly one page of the database. The file header has a size of exactly 32 bytes. It starts with a 4 byte long Magic Number (see Table 5.8). At offset 4 follows the file format version. Again, this is a 4-byte unsigned integer value. The size of one page of the database is stored at offset 8. Using the field *checkpoint sequence number* at offset 12, we can again determine how many checkpoints have already been executed since their creation.

Table 5.8: WAL Header Format [86]

Offset	Size	Description
0	4	Magic number. 0x377f0682 or 0x377f0683
4	4	File format version. For example 3007000.
8	4	Database page size. Example: 1024
12	4	Checkpoint sequence number
16	4	Salt-1: random integer incremented with each checkpoint
20	4	Salt-2: a different random number for each checkpoint
24	4	Checksum-1: First part of a checksum on the first 24 bytes of header
28	4	Checksum-2: Second part of the checksum on the first 24 bytes of header

The last fields of the header form two salt values and two checksum values. Using these fields, we can determine which frames belong to the current checkpoint and

<sup>1</sup> <https://github.com/pawlaszczyk/fqlite>

have not yet been transferred to the database. Figure 5.12 shows an example of the header of a WAL archive in FQLite.

Each WAL frame also starts with a header [84]. The structure of the header with its fields is shown in Table 5.8. The header consists of exactly six big-endian values, each with a size of 4 bytes. The first object is the page number this frame is assigned. Using the page number, we can identify the place in the database where the change takes effect. The value at offset four can be used to determine whether a COMMIT was performed. A value other than 0 is a so-called *commit frame*. Let us remember that a COMMIT does not automatically update the database. Like the header of the WAL file, each frame header ends with two salt values and two checksums. The four big-endian 32-bit unsigned integer values are located from Offset 8 to 24.

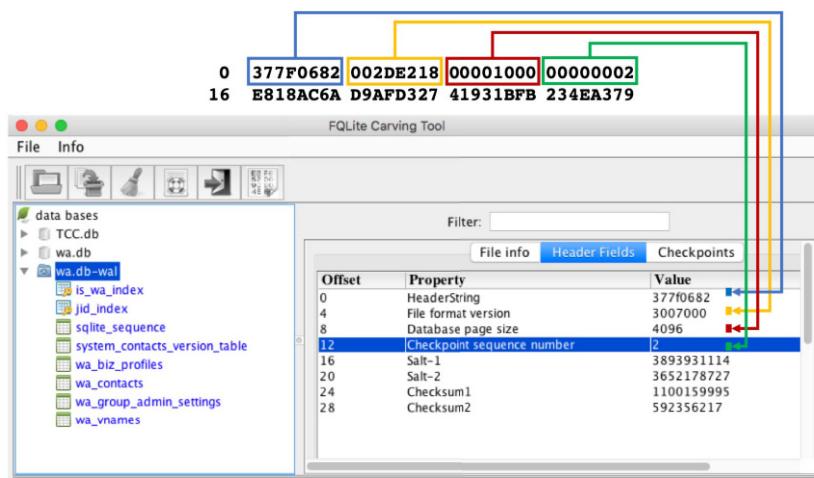


Fig. 5.12: View on a WhatsApp-DB WAL-Header with FQLite Carving Tool

A WAL archive always grows from the beginning. It can cause frames from different checkpoints to appear in the same file, whereas current ones are always at the file's beginning. Fortunately, we can use the mentioned salt values to determine relatively quickly whether the frame under consideration is valid or whether it belongs to an older state already transferred to the database. Whether a frame is valid can be determined as follows [86]:

1. The Salt-1 and Salt-2 values from the header must both match the values in the respective frame.
2. The 8-byte checksum in the frame must match the cumulative checksum over the first 24 bytes of the WAL header plus the first 8 bytes and the contents of all previous frames.

If a checkpoint was executed successfully, the WAL file is reset afterwards. In this case, the salt values are overwritten. The value of salt-1 is incremented, while a

Table 5.9: WAL Frame Header Format [86]

Offset	Size	Description
0	4	Page number
4	4	For commit records, the size of the database file in pages after the commit. For all other records, zero.
8	4	Salt-1 copied from the WAL header
12	4	Salt-2 copied from the WAL header
16	4	C checksum-1: Cumulative checksum up through and including this page
20	4	C checksum-2: Second half of the cumulative checksum.

new random value is assigned to salt-2. Previously valid frames are automatically discarded due to this procedure. However, the previous frames usually remain in the archive due to the I/O- operations when the file is truncated. Thus, there is an excellent chance to make past states of database pages visible again with the help of the WAL file.

Let us take a look at how write-ahead logs work. Figure 5.13 shows the frames list of a WAL file. Below the header field for the Salt-1 value, several frames are shown. All frames with matching salt values belong to the same checkpoint. The first seven frames thus form a unit. The remaining frames are part of an older checkpoint. As we can see, the salt in the header matches the salt in the first unit. Accordingly, the pages have not yet been transferred to the database. In other words, the WAL file contains the latest version of page 2,4,6,18. The pages within the database are out of date. The next checkpoint is usually executed when opening the database, and these data records are transferred to the database. Since WAL files always work at a page level, the complete database page is updated. Remember, the salt value changes for each checkpoint. Thus the Salt-1 field in the header is discarded afterwards.

Interestingly, page 6 has been updated three times. When a checkpoint occurs, each page will be written back to the database in the same order written to the WAL file. Pages are written from the start of the WAL file. Accordingly, the update order would be 2,6,4,12,6,18,6. This allows a timeline to be created, starting with the first to the last update step.

## 5.5 Conclusions

The SQLite database format has great importance in the field of mobile forensics. In this chapter, we have therefore tried to take a look behind the scenes. As quickly became apparent, the file format of SQLite has some similarities to a classic file system, where files are usually stored in blocks. Instead of blocks or clusters, data content in SQLite is managed in pages. As has been shown, even records are often recoverable after they have been deleted. Analogous to a file system, these are usually

Header	
Salt-1	123456
123456	2
123456	6
123456	4
123456	12
123456	6
123456	18
123456	6
111110	5
111110	2
111110	6
111110	17

Fig. 5.13: Frame list of a WAL file (example)

not wiped but merely marked as deleted. However, we do not manage files but data sets.

We further identified different slack spaces of an SQLite database. Besides free blocks and the unallocated space, we can find deleted records, especially in the freelist area of the database. The carving techniques discussed within this chapter can help make these data sets visible again in many cases.

Of the temporary file-formats considered, rollback journals and the WAL files are of particular interest to the investigator, as they may contain old or previously altered data. However, special care must be taken when acquiring data from these files. Thus, the data stored in a WAL file can be reconstructed manually or with specialized forensic tools. Using an ordinary SQLite reader, on the other hand, can lead to the loss of data.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.



# Chapter 6

## Property Lists



Christian Hummert and Georgina Louise Humphries

**Abstract** Property List files (\*.plist) are a widely used data storage format used by Apple software. Most of the system properties are stored in plists, but also, many apps store their configuration in plist-files. The data held within Property is regularly of high evidential value for forensic analysts, so understanding the format is essential for the forensic investigation of Apple mobile devices and computers. Not all of today's digital forensics tools recover plists properly. Especially for carved or damaged plists, the support is insufficient. So the forensic examiner must understand the principles of this file format. This chapter gives an overview of the plist structure to give the examiner the knowledge to get the most information out of the evidence possible.

### 6.1 Introduction

Property List files (\*.plist) are one of the widely used data storage formats used by Apple software . Most of the system properties are stored in plists (many of them are located in `/Library/Preferences/`), but many apps store their configuration in plist-files. Therefore, property lists can be found in various places on Apple systems. They sometimes can even be found on devices other than Apple operating systems (especially if other Apple software like Safari or iTunes is installed). The data held within Property is regularly of high evidential value for forensic analysts, so understanding the format is essential for the forensic investigation of Apple mobile devices and computers.

---

Christian Hummert  
Agentur für Innovation in der Cybersicherheit, Halle, Germany, e-mail: [hummert@cyberagentur.de](mailto:hummert@cyberagentur.de)

Georgina Louise Humphries  
Politihøgskolen - Norwegian Police University College Department of Postgraduate Studies, Oslo, Norway, e-mail: [georgina.louise.humphries@phs.no](mailto:georgina.louise.humphries@phs.no)

Property lists offer a structured and efficient way to represent and persist hierarchies of objects to disk. They are the standard way to save and load data between the internal representation within objects in Objective-C or Swift programs and disk files. The standard objects of the Cocoa framework have built-in methods to deal with plist files.

The first plists were developed with the NeXTSTEP operating system. NeXTSTEP is a discontinued operating system that merges the Mach kernel and the BSD kernel. It was developed by NeXT, a company founded by Steve Jobs in 1985, and released in 1989. NeXTSTEP had a text-based, human-readable format for plists, serialized to ASCII in a syntax somewhat like a programming language. Apple replaced the format with an XML-based format and also introduced the binary plist format. Since Mac OS X 10.7, in addition, JSON notation can be read and written. So, there are four different property list formats[9]:

- NeXTSTEP property lists (deprecated since OS X 10.0)
- XML property lists (introduced with OS X 10.0)
- Binary property lists (introduced with OS X 10.2)
- JSON property lists (introduced with OS X 10.7, but not 100% compatible)

The formats except the binary plist have the advantage of being human-readable. In contrast, the binary plist offers the most efficient representation on disk and fast serialization/deserialization. OS X offers the `plutil` utility (introduced in OS X 10.2) to check the syntax of property lists or convert a property list file from one format to another. It also supports converting plists to Objective-C or Swift object literals. Another tool that comes with OS X is `PlistBuddy` (it can be found at `/usr/libexec`). `PlistBuddy` allows to merge plists or edit their content.

Property Lists in XML or JSON notation can be easily edited and evaluated in any desired text editor. Apples IDE Xcode also contains a hierarchical viewer and editor for binary and XML plists. In addition, Apple offers an Apple Script interface to create, edit and write property lists (since OS X 10.5). Due to the human-readable notation of NeXTSTEP, XML and JSON property lists, they are not an obstacle for forensic investigations. Therefore, this chapter concentrates on the binary plists (`bplist`) format and will afterwards describe some of the interfaces to plist files.

## 6.2 Binary plist Structure

Apple disclosed the structure of the binary property list format; it is documented in the comments of the Apple-provided open-source `CFBinaryPList.c`<sup>1</sup> and declarations of the `ForFoundationOnly.h`<sup>2</sup>. Every binary plist file comprises four sections: a header, an object table, an offset table and a trailer (compare table 6.1).

Each `bplist` file begins with an 6-byte header, containing the magic `bplist` (Hex: `0x62706C697374`). The header is followed by a 2-byte version. The most common

---

<sup>1</sup> <https://opensource.apple.com/source/CF/CF-550/CFBinaryPList.c>

<sup>2</sup> <https://opensource.apple.com/source/CF/CF-550/ForFoundationOnly.h>

Table 6.1: Structure of a bplist file.

Offset	Size	Description
0x00	6	bplist header (0x62706C697374)
0x06	2	format version
0x08	LEN1	object table
0x08 + LEN1	LEN2	offset table
0x08 + LEN1 + LEN2	32	trailer

version on Apple devices is **00**, but there are at least two other versions of binary property lists, too; **bplist15** or **bplist16** occur. Unlike for **bplist00**, there is no documentation for either format. The **bplist15** format appears to be internal to CoreFoundation. The **bplist16** format is internal to Foundation, too, and is used almost exclusively in Objective-C remoting over XPC. The format of **bplist16** is similar, but not compatible with **bplist00**, noting the following differences: Files in format version 16 do not have a trailer, and the items start directly at the head of the property list, right after the **bplist16** magic, and are packed (not aligned). In addition, in **bplist16** there are more data types available.[47]

The bplist file ends with a 32-byte long trailer. The structure of the trailer is shown in table 6.2. The bytes 0 to 4 of the trailer are unused. Byte 5 contains the sort version. Byte 6 stores the information of the size in byte of each offset entry in the offset table. Similarly, byte 7 stores the information of the size of each object reference in a container. At offset 0x8, there is an 8-byte entry that saves the number of objects that are encoded inside the object table. The following 8 bytes save the offset of the first offset in the offset table (usually zero). The last 8 bytes of the trailer denotes the start of the offset table, counting from the start of the bplist.

Table 6.2: Structure of the bplist trailer.

Offset	Length	Description
0x0	5	unused
0x5	1	sort version
0x6	1	size per offset in offset table in bytes
0x7	1	size per object reference in a container
0x8	8	number of objects in object table (big endian)
0x10	8	offset of the first offset in the offset table (big endian)
0x18	8	offset of the offset table (big endian)

The second section in every bplist file is the object table. The object table contains all the data objects of the plist. All object types are identified by a single byte, also called a marker (compare Table 6.3). This byte encodes the type of an object and the size of the data.

Table 6.3: Format of object types.

Object	Marker	(Additional Info)	Description
null	0000 0000		
bool	0000 1000		false
bool	0000 1001		true
fill	0000 1111		fill byte
int	0001 nnnn ...		$2^n$ bytes (big endian)
real	0010 nnnn ...		$2^n$ bytes (big endian)
date	0011 0011 ...		8 byte float (big endian)
data	0100 nnnn [int] ...		nnnn bytes unless 1111 then [int] count followed by bytes
string	0101 nnnn [int] ...		nnnn chars unless 1111 then [int] count followed by bytes
string	0110 nnnn [int] ...		Unicode string, nnnn chars unless 1111 then [int] count followed by bytes
	0111 xxxx		unused
uid	1000 nnnn ...		nnnn+1 bytes
	1001 xxxx		unused
array	1010 nnnn [int] objref*		nnnn entries unless 1111 then [int] count followed by entries
	1011 xxxx		unused
set	1100 nnnn [int] objref*		unused
dict	1101 nnnn [int] keyref*		nnnn entries unless 1111 then [int] count followed by entries
	1110 xxxx		unused
	1111 xxxx		unused

The marker is the binary representation of a single byte. All other objects can be uniquely identified by the marker byte's 4 most significant bits (MSB). At the same time, the least significant bits (LSB) of the marker byte denotes sizing information. If the object size is small enough, the size is encoded immediately in the 4 right-most bits, and then the actual data values follow. If the object size is larger, the LSB matches 0xF (1111), denoting that the next bytes encode size information before the actual value bytes.

The size is encoded as follows: The MSB equals 0x1 (0001), and the LSB contains a value  $x$ . The size will be stored in the following  $2^x$  bytes in big-endian.

For example: let us assume the object table contains the sequence 0x5F 10 19. The first 0x5F is converted into its binary representation 0101 1111. The MSB is 0101, so the object denotes a string. The LSB matches 1111, so the size is encoded in the next byte. The next byte is 0x10 = 0001 0000. The MSB equals 0x1, and the LSB shows that  $2^0 = 1$  byte follows, which stores the size of the string. The next byte is 0x19 resulting in a 25-byte long string.

Markers corresponding to objects such as ints, real numbers, strings are immediately followed by a multibyte sequence representing their actual values. This is not always the case, though. In the case of object containers, such as arrays and

dictionaries, the marker byte is followed by object references that are simply offset to the offset table. The length of this offsets is determined in the bplist trailer, and are counted from the beginning of the offset table. Therefore, a container element is just a reference that points back to a position in the offset table, which points back to the object table and specifically to a marker corresponding to the individual object. This technique flat-maps the actual multi-level hierarchy and allows all objects to have fixed sizes.

The third section in bplists contains offsets to the object table and guides to the actual values of objects. The size of each offset is defined in the file trailer. All offsets are calculated from the beginning of the file (not the end of the header). The number of offsets stored in the offset table is also given in the trailer.[\[40\]](#)

### 6.3 Example

Given is the following plist (Table 6.4) from a MacBook Pro:

Table 6.4: Example plist (object table colored in blue, offset table colored in red, trailer colored in yellow).

62 70 6C 69 73 74 30 30 D2 01 02 03 04 5E 42 61	b p l i s t 0 0 “	^ B a
74 74 65 72 79 48 69 73 74 6F 72 79 5F 10 13 54	t t e r y H i s t o r y _ T	
6F 74 61 6C 4E 75 6D 62 65 72 6F 66 45 76 65 6E	o t a l N u m b e r O f E v e n	
74 73 09 10 0A 08 0D 1C 32 33 00 00 00 00 00 00	t s - 2 3	
01 01 00 00 00 00 00 00 00 05 00 00 00 00 00 00		
00 00 00 00 00 00 00 00 35		5

The given bplist is version 00 as the header states. To analyze the bplist in a first step the trailer is marked (here yellow) the trailer comprises the last 32 bytes of the file. Now the trailer can be decoded (result in Table 6.5):

Table 6.5: Decoded example bplist trailer.

Content	Offset	Length	Description
0x00	0x5	1	sort version
0x01	0x6	1	size per offset in offset table
0x01	0x7	1	size per object reference
0x0000000000000005	0x8	8	number of objects in object table
0x0000000000000000	0x10	8	offset of the first offset in offset table
0x0000000000000035	0x18	8	offset of the offset table

Now, it is clear that the bplist contains five objects in the object table, and the offset table starts at 0x35, whereas the first object-offset starts at  $0x35 + 0x00$ , and each

offset has the size of one single byte. In consequence, the offsets from the offset table are 0x08, 0x0D, 0x1C, 0x32 and 0x33, which leads to the following four objects from the object table:

1. 0xD2 01 02 03 04
2. 0x5E 42 61 74 74 65 ...
3. 0x5F 10 13 54 6F 74 ...
4. 0x09
5. 0x10 0A

The first object starts with 0xD2, which is 1101 0010 in binary. The MSB (1101) shows that the object-type is a dictionary. The LSB (0010 = 2) shows that the dictionary has two entries. The data 0x01 02 03 04 has to be interpreted as object references that are simply offsets to the offset table. Before interpreting the dictionary as an object container, the other four entries should be decoded.

The second object starts with 0x5E, which is 0101 1110 in binary. The MSB (0101) shows that the object type is a string. The LSB (1110 = 14) shows that the string is 14 chars long. The content of the string is "BatteryHistory".

The third object starts with 0x5F, which is 0101 1111 in binary. The MSAB (0101) shows that the object is another string. The LSB (1111) shows that the string is longer than 14 chars, and the size is encoded in the following bytes. The next byte is 0x10, which is 0001 0000 in binary. The MSB (0001) is defined as 0001, and the LSB (0000) shows that the following  $2^0 = 1$  bytes encode the length of the string size. So the next byte has to be decoded. It is 0x13, which states that the string has a length of 0x13 = 19 chars. The content of the string is "TotalNumberOfEvents".

The fourth object starts with 0x09, which is 0000 1001 in binary. The MSB (0000) indicates the object as bool, and the LSB (1001) indicates the content as "true".

The fifth object starts with 0x10, which is 0001 0000 in binary. The MSB (0001) indicates the object as an integer. The LSB indicates that the integer is  $2^0 = 1$  byte long. The content of the integer is 0x0A which is 10. Now the dictionary (the first object) can be decoded. The dictionary has two entries: the objects at offset 0x01 and 0x02 in the offset table, which are the two strings. The first entry is connected to the object on offset 0x03, which is the boolean. The second entry is connected to the object on offset 0x04, which is the integer. That gives the following result (Table 6.6):

Table 6.6: Decoded example plist.

dictionary (2 entries)			
	BatteryHistory	TRUE	
	TotalNumberOfEvents	10	

Fig. 6.1 shows the same plist file decoded with Apples XCode IDE and confirms the correct decoding.

Key	Type	Value
Root	Dictionary	(2 items)
BatteryHistory	Boolean	YES
TotalNumberOfEvents	Number	10

Fig. 6.1: Illustration of the elements of a block group.

## 6.4 Forensic Tools Supporting plists

There is quite a bunch of tools supporting the decoding of plist files. Most of the tools support binary plists as well as XML property lists. As the property lists are in an Apple format, the macOS universe gives the best support. If the given property list file is in XML format, it can be edited in any text editor. On the other hand, if the given property list file is in the binary format, it can be converted to XML first by running on the macOS shell:

```
plutil -convert xml1 file.plist
```

If an XML property list should be converted back this is possible with:

```
plutil -convert binary1 file.plist
```

A more convenient way to edit and browse plist files is to install the Apple development platform Xcode. The suite includes a graphical editor which is easy to use (compare Fig. 6.2).

For old versions (Xcode 4.2 and earlier) there was a separate application for editing property lists (/Developer/Applications/Utilities/Property List Editor.app).

Nevertheless, despite the Apple world, there are plenty of alternatives. One free plist parser for binary property lists is `binplist`<sup>3</sup> a parser module written in Python. For forensic use, it is possible to create an instance of the `BinaryPlist` class and then call the `Parse()` method, with a file-like object as an argument.

```
with open("myfile.plist", "rb") as fd:
    bplist = BinaryPlist(fd)
    top_level_object = bplist.Parse(fd)
```

The `Parse()` method returns the top-level object, just as `readPlist`. Once parsed, `BinaryPlist.is_corrupt` can be checked to recognize whether the plist had corrupt data. This allows to maybe decode the corrupted data manually and gather the maximum information on the plist even when they are corrupt.

---

<sup>3</sup> <https://github.com/google/binplist>

Key	Type	Value
Root	Dictionary	(17 items)
LastSuccessfulDate	Date	2021-05-14T12:48:29Z
AutomaticCheckEnabled	Boolean	0
LastAttemptSystemVersion	String	11.2.2 [20D80]
LastBackgroundCCDSuccessfulDate	Date	2016-09-30T13:05:07Z
AutomaticallyInstallMacOSUpdates	Boolean	0
LastUpdatesAvailable	Number	1
SkipLocalCDN	Boolean	0
LastRecommendedUpdatesAvailable	Number	1
LastAttemptBuildVersion	String	11.2.2 [20D80]
AutomaticDownload	Boolean	0
RecommendedUpdates	Array	(1 item)
Item 0	Dictionary	(5 items)
Identifier	String	MSU_UPDATE_20E241_patch_11.3.1
MobileSoftwareUpdate	Boolean	1
Display Name	String	macOS Big Sur 11.3.1
Product Key	String	MSU_UPDATE_20E241_patch_11.3.1
Display Version	String	11.3.1
LastFullSuccessfulDate	Date	2021-05-14T12:48:29Z
LastRecommendedMajorOSBundleIdentifier	String	
PrimaryLanguages	Array	(2 items)
Item 0	String	de
Item 1	String	de-DE
LastSessionSuccessful	Boolean	1
LastBackgroundSuccessfulDate	Date	2017-07-20T19:07:32Z
LastresultCode	Number	2

Fig. 6.2: View of an example plist in the Xcode editor.

Another Python module dealing with binary property lists is `ccl-bplist`<sup>4</sup>. A C-library to handle plists is `libplist`<sup>5</sup>.

Forensic Suites on Windows create a different picture. The MSAB Forensic suite<sup>6</sup> offers property list decoding as well as the Cellebrite<sup>7</sup> suite. Oxygen Forensics<sup>8</sup> offers the Oxygen Forensic Plist Viewer that automatically unpacks one plist file and displays its content as a folder tree (compare Fig. 6.3). The entries can be converted to different encodings. In addition, the binary content of cells, such as images, video and sound, can be decoded and visualized. The Forensic Toolkit (FTK)<sup>9</sup> does support plist decoding, whereas the Encase Forensic Software<sup>10</sup> offers a plugin script to decode these files.

<sup>4</sup> <https://code.google.com/archive/p/ccl-bplist/>

<sup>5</sup> <https://github.com/JonathanBeck/libplist>

<sup>6</sup> <https://www.msab.com/>

<sup>7</sup> <https://www.cellebrite.com/>

<sup>8</sup> [www.oxygen-forensic.com](http://www.oxygen-forensic.com)

<sup>9</sup> <https://accessdata.com/products-services/forensic-toolkit-ftk>

<sup>10</sup> <https://security.opentext.com/encase-forensic>

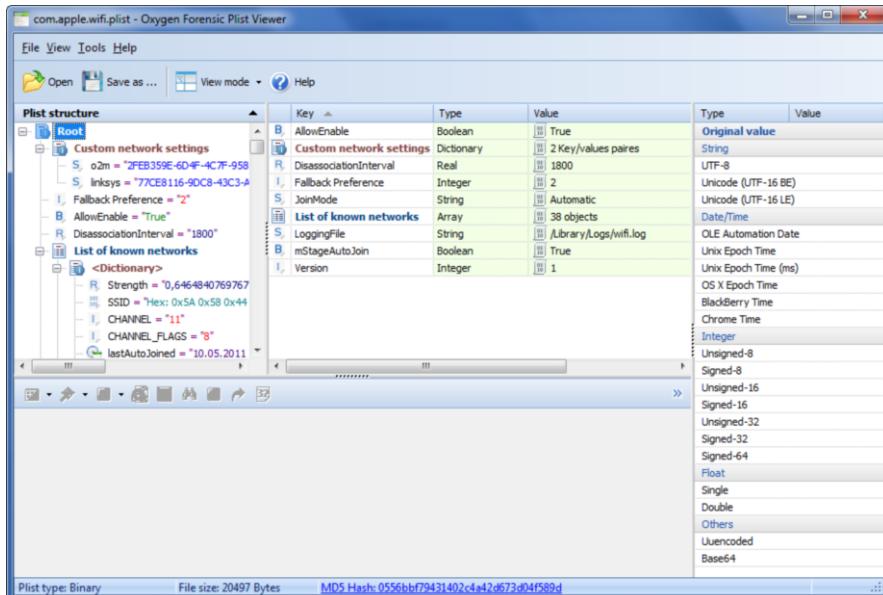


Fig. 6.3: View of an example plist in the Oxygen Forensic Plist Viewer.

## 6.5 Conclusions

Property Lists have great importance in the field of mobile forensics. In this chapter, a look behind the scenes was given. There are four different types of property lists. The chapter concentrated on the binary property list format because it is not intuitive to decode and has a certain prevalence. The binary property list format contains a header, a trailer and an object table. The type of the single objects can be found in a look-up table.

There are plenty of tools that support the decoding of plist files. But in a forensic manner, information from corrupted files (especially parts of files carved from the unallocated space) must be revealed. To fulfill this challenge, deep knowledge about file structures is crucial.

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

