

Automated Failure Analysis of Remote Edge Devices at Scale

Tyler Hoffman



Speaker

Tyler Hoffman, Co-founder @ Memfault



Previously a Firmware Engineer @ Pebble & Fitbit

Split time between writing & debugging firmware and building internal services to help monitor millions of devices.

Can find my thoughts and content on Memfault's Interrupt blog: <u>interrupt.memfault.com</u>



What are we trying to do today?

- Talk about firmware on embedded MCU systems
- Speed from **development** to **testing** to **mass production** and **scaling**
- Find and fix (quickly) 1 in 10,000 hour bugs in production
- Prevent issues from being released to 100% of users
- Establish measurable stats to track progress on goals



One of my favorite quotes

"These techniques are not necessary because experienced firmware engineers do not introduce bugs in their code"

- A (real) Firmware Engineer



4

Agenda









State of the union for firmware



Firmware is *pervasive*, but difficult

- In Q4 2020, 4.4 billion Cortex M's shipped
- Tons and tons of firmware running on these devices
- Very few tools to help debug these devices
- Software is becoming more complex leads to bugs
- Software issues lead to:
 - Bricked devices, RMA's, and security exploits

We need proper tools to help us build quality products





Increasingly complex software

- Firmware no longer is complete isolated
 - Must communicate to mobile phones & gateways
 - Protocols, devices, and security issues constantly changing
 - Libraries are *massively complicated* (e.g. mbedTLS)
- Expectations are growing
 - Everything connects to Internet
 - Must enable firmware updates
 - Competition is growing







Rudimentary Debugging Tools

***** MPU FAULT *****
Instruction Access Violation
***** Hardware exception *****
Current thread ID = 0x20000074
Faulting instruction address = 0xe0000000
Fatal fault in thread 0x20000074! Aborting.

[I][1605794400]	Wi-Fi	connected	
[W][1605796200]	Wi-Fi	disconnecte	d, reason: -2
[I][1605796500]	Wi-Fi	connected	
[W][1605797100]	Wi-Fi	disconnecte	d, reason: -2
[I][1605798000]	Batter	y status: 6	7%, 3574 mV

Tasks												×
ID	Priority	Name	Stack Info (Used / Siz			re)	e) Status				Run Count	
40	113	IP_Tas	k	564/1	.280 @ 0x2000	280 @ 0x20006708 Waiting			ting (eve	ng (event), with Timeout		2943
⇒ 198	108	GUI_T	ask	1324 /	4096 @ 0x200	0060F0 Running			ning			4170
97	106	USBMS	DTask	408 / 4	096 @ 0x20008C20			Susp	Suspended, with Timeout			468
140	105	IP_We	bServe	er 672/2	2048 @ 0x20006E78			Wai	Waiting (event object)			2
244	101 Idle	IP_FTF	Serv	Call Sta	Call Stack							
	Tue Tue		Function	Function		Stack Info			PC	Ret	urn Address	
		Rac	RadialMenu		Q	2000 61E8		0000 27A4	0000 2A0E			
		GUIDEMO	GUIDEMO_RadialMenu		Q	2000 61E8		0000 2A0A	0000 1162			
value		_Main	Main		e	2000 61F0		0000 1160	0000 1508			
± aItemInfo		GUIDEMO	GUIDEMO Main		Q	200	0 61F8	0000 1504	00	00 B3B0		
Para		OS_Start	OS_StartTask		Ø	2000 6238		0000 B3AE	N/.	A		
hDraw	hDraw 37		_	2000 61C4		-	4	4 long			_	
+ pPara 0x80			2000 61BC				struct PARA*					
xSizeWin 160			2000 61D0		4		int					
vSizeWin 4 995		5	2000 61CC		4		int					
xSize 518 2		248	8 2000 61D8		4 i		int					
vSize 5			2000 61D4		4 int		int	nt				
xPos 0			2000 61C8			4 int						
i	i -859 0		045 884	2000 61DC			4	4 int				

The existing tools are decent, but only when connected using JTAG

Debugging Data Exported to ???





Not much exists to help us

- Existing software solutions won't work
 - We aren't monitoring 10-100 servers, rather 10k 1m+ devices
 - We don't have MB's and GB's of RAM lying around, we have KB's
 - Assembly, C, C++, Rust (?)
 - Symbol files (.elf's) required for typical debugging experiences
 - AWS IoT,
- Few existing solutions that can help embedded devices
- Firmware is not an open-source ecosystem or sharing community

Whatever we need, we have to build ourselves



Some issues will only happen in the field

"a third of all software faults take more than 5000 execution-years to manifest themselves."

*Source: http://www.ganssle.com/tem/tem417.html



Agenda











Analyzing the issues



Do we have issues in the field...

They're likely having issues



Fundamentally, we need to know if and how many times devices are experiencing issues.

Once we have this information

- Can assess whether further work is needed.
- Can make assumptions about whether a firmware update improved stability
- Can guide business & product decisions
- Leaders know whether it's time to fix bugs and tech debt or build new features

```
void buggy function(void) {
  *(uint32_t *)0xbadcafe = 0x0;
}
void HardFault Handler(void) {
   // ... fault handling code ...
   NVIC SystemReset();
}
```



ERROR: Connection lost - reason -7



- Sending & Parsing logs
 - Log important device events
 - Log resets and registers
 - Hard and doesn't scale well

- Tracking simple metrics
 - # device resets
 - # devices alive and well
 - Average uptime of devices
 - Scalable, but less information



Requirements:

- 1. Path to the Internet
- 2. Central repository for device data







Agenda





Analyzing the issues



Analyzing the Issues

We have issues. How do we fix them?



Debugging in Production

- We don't have access to our normal toolset.
 - No UART & serial logs
 - No JTAG and quick re-flashing
 - No IDE, GDB, or step-through debugging
 - RMA's take weeks. Devices must persist the data.
- We need tools specifically for debugging in production
 - Build these tools as early as possible
 - Having them early will pay off **10x** (Speaking from experience)
 - If you have no interest in building them, check out Memfault

Using Logs to Analyze Issues

[I][1605794400] Wi-Fi connected [W][1605796200] Wi-Fi disconnected, reason: -2 [I][1605796500] Wi-Fi connected [W][1605797100] Wi-Fi disconnected, reason: -2 [I][1605798000] Battery status: 67%, 3574 mV

2 Wi-Fi disconnect events

1 Battery level event

```
***** MPU FAULT *****
Instruction Access Violation
***** Hardware exception *****
Current thread ID = 0x20000074
Faulting instruction address = 0xe0000000
Fatal fault in thread 0x2000074! Aborting.
```

1 crash due to MPU violation

Program Counter: 0xe0000000

- What is a coredump
 - A snapshot of RAM regions at a specific moment in time.
 - Can be captured on command, or at time of a fault, such as a hardfault or ASSERT.
 - Helps with post-mortem debugging
 - Requires symbol file & clever parsing
- What RTOS's support coredumps natively?
 - ESP32, Zephyr RTOS, MyNewt RTOS
 - Any RTOS with Memfault





```
void HardFault_Handler(void) {
    // ... fault handling code ...
    NVIC_SystemReset();
    // Don't just crash...
    capture_coredump();
}
```





A coredump is comprised of raw memory regions

- From coredumps, one can recover:
 - Threads and stacktraces
 - Heap allocations and statistics
 - Global variables
 - Local variables and function arguments
 - System and peripheral registers
 - Anything else you want to collect (and parse)







Threads, backtraces, local variables, and registers



Analysis						
Memory management fault detected at 0×2000a3c0						
Memory management fault on a data access						
Fault Register	Value	Hex Value				
CFSR	130	0×0000082				
HFSR	0	0×00000000				
SHCSR	458753	0×00070001				

Parse and interpret system registers



heap Q Order by Memo	ory Location V	Find Address Q	Regions V
heap_sz = unsigned int 0	@ 0×20002378	0x20009f58 ec 31 00 20 00	000000 .1
_mbedtls_heap = unsigned char[28000] {}	@ 0x200031ec	0x20009168 00000000 00	000000
<pre>wheap = buffer_alloc_ctx {}</pre>	@ 0x20009f4c	0x20009f70 00000000 000	00000
buf = unsigned char* {}	@ 0x20009f4c	0x20009f78 0000000 00	00000
len = size_t 28000	@ 0x20009f50	0x20009f80 0000000 00	00000
<pre>first = memory_header* {}</pre>	@ 0x20009f54	0x20009f88 01000f00 000	00000
<pre>first_free = memory_header* {}</pre>	@ 0x20009f58	0x20009f90 000000000000	000000
<pre>* = memory_header {}</pre>	@ 0x200031ec	0x20009198 00000000 000	~2 02 08
magic1 = size_t 4278233685	@ 0x200031ec	0x20009fa8 00000000 00	0 0 0 0 0
size = size_t 27968	@ 0x200031f0	0x20009fb0 0000000000000	00000
alloc = size_t 0	@ 0x200031f4	0x20009fb8 0000000 00	00000
prev = memory_header* {}	@ 0x200031f8	0x20009fc0 0000000 48	e20020 H
next = memory_header* {}	@ 0x200031fc	0x20009fc8 0000000 00	00000
prev_free = memory_header* {}	@ 0x20003200	0x20009fd0 ad c2 02 08 000	000000
next_free = memory_header* {}	@ 0x20003204	0x200091d8 00000000 00	546665 idlo
magic2 = size_t 3994130790	@ 0x20003208	0x20009fe8 00000000 000	000000 Iute
verify = int 0	@ 0x20009f5c	0x20009ff0 00000000 000	00000
heap_sem = sys_sem {}	@ 0×2001166c	0x20009ff8 0000000 00	00000

View all global variables at time of coredump



Coredumps – What can I debug?

- Crashes and other items worth extra investigation
 - Hardfaults, Memfaults, Asserts, etc.
 - Deadlocks, system stalls
 - Memory corruption issues
 - Heap out-of-memory
 - Stack overflows
 - Undefined behavior ("SHOULD NEVER HAPPEN" ^(c))



https://embeddedartistry.com/blog/2021/01/11/hard-fault-debugging/



When Coredumps Aren't Enough

- Coredumps aren't the perfect solution for all issues
 - Behavioral bugs require more of a timeline than a snapshot (if this & this, then that)
 - Devices with storage or bandwidth constraints
 - Issues in interpreted languages (MicroPython, JerryScript)
 - Determining trends
- When you need context before a crash...
 - Use logging or simple tracing
 - Store log lines in a RAM buffer, send it with coredump!

St	State Logs							
S	Search							
Lvl	LvI Message							
Т	Initializing Accel Subsystem							
Т	Launching sleep tracking app							
Т	Wifi Connected.							
Т	Analyzing Raw Algo Data (50B)							
1	Wifi Unavailable. Retrying in 5 min							
1	sleep tracking app closed							
1	Launching sleep tracking app							
E	Temp Sensor I2C transaction timeout, rv=-8							
	Wifi Connected.							
	Analyzing Raw Algo Data (100B)							

Agenda





Analyzing the issues



Preventing Issues

We fixed the issues. How do we don't introduce more?



Internal Testing

- You must be testing the devices internally
 - Catch and fix as many bugs before production as possible
 - Firmware team, QA, engineers, and other employees
 - Ask biggest fans of your product to be beta testers
 - Encourage users to report bugs
 - Less concerns around privacy & data collection
- Build and deploy nightly or regular firmware updates to devices
 - Constantly be fixing and verifying fixes
 - Perform experiments!

Your device just
reset. Please submit
a support ticket!

Unfortunately, some bugs will only be caught in production.

Staged Rollouts

- Send out firmware updates in a controlled manner
 - Send it out to 5%, then 10%, then 20%, etc.
 - At each step, assess quality and performance
 - New issues? pull the firmware update!
- How to implement staged rollouts?
 - Hosted firmware update solutions have this feature
 - At the very least, can set up devices to poll once every 24 hours (random interval) for new firmware update





Metrics

Uncover trends on anything that is numerical

- Common metrics are:
 - Task runtimes
 - Heap information
 - Connectivity statistics
 - Frequency of errors
 - Peripheral utilization and power states

Example App Booting							
Metric	ID:	1		Value:	15000	<pre># elapsed ticks</pre>	
Metric	ID:	2		Value:	8531	# Main Task ticks	
Metric	ID:	3		Value:	8223	# Timer Task ticks	
Metric	ID:	4		Value:	14	# Timer Task count	
Metric	ID:	5		Value:	3593	# "sensor on" ticks	
Metric	ID:	6		Value:	4696	<pre># heap-free low watermark</pre>	

Pinpoint power, connectivity, and performance issues and regressions

Metrics

```
void flash_sector_erase(uint16_t sector) {
    ...
    device_metrics_incr(kDeviceMetric_FlashSectorErases);
}
void flash_write_bytes(uint32_t addr, void *buf, size_t n) {
    ...
    device_metrics_incr_by(kDeviceMetric_FlashWriteBytes, n);
}
```

https://interrupt.memfault.com/blog/device-heartbeat-metrics



Metrics

Per Device Counts Averages Min/Max

Entire Fleet Counts

Averages Min/Max



D/2022



We can compare metrics between firmware versions





Agenda







Quantifying current issues





Analyzing the issues

Best Practices My Favorite Practices

Sharing my tips & tricks monitoring millions of devices with proper tooling in place



Life With Debugging Infrastructure

If you have logs & coredumps being sent from devices, you can employ what I call **offensive programming**.

This is the proper way to debug "in production"

- During runtime, raise errors immediately!
 - Bring bugs front and center
 - Undefined behavior reset to a known state
 - Assume that if bugs are raised early and often, they can be fixed before shipping
 - With a handful of devices, the 1 in 10,000 hour bugs can be caught



BAD BAD BAD BAD

https://interrupt.memfault.com/blog/defensive-and-offensive-programming

Argument validation





Assert on heap exhaustion





Easily capture and diagnose timing issues





State machine transition violations



Erase after free



Zero out allocations after free

Now standard in iOS 16.1 and macOS Ventura**



The End

- Happy to connect!
- https://www.linkedin.com/in/tyhoff/
- <u>https://memfault.com/webinars/</u>
- <u>https://interrupt.memfault.com</u>
- Email: tyler@memfault.com



Thank You Danke Gracias Grazie 谢谢 ありがとう Asante Merci 감사합니다 धन्यवाद **Kiitos** شكرًا ধন্যবাদ תודה



DEV/SUMMIT

The Arm trademarks featured in this presentation are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. All other marks featured may be trademarks of their respective owners.

www.arm.com/company/policies/trademarks

