

Pentesting Hardware - A Practical Handbook

Mark Carney

November 22, 2018

Contents

Introduction

Hardware hacking is dope, yo...

Acknowledgements

I am deeply indebted, to ... my directors, producers, make-up artists... whatever.

Copyright

For the current releases, this material is released under Creative Commons v3.0 - quote me all you like, and reference my work no problem, print copies for yourselves, but just leave my name on it.

Part I

Hardware Skillz

First off, we need to talk about what hardware is, so that we can talk about the threats in the next part.

This part will be more preparatory stuff - What software is best to get ahold of? What hardware testing kit is best? How does Protocol X work in a little more detail?

This is not designed to be an ‘undergraduate textbook’, but rather a high-level overview for the interested with links to proper technical descriptions where required.

Chapter 1

PCB Eye-Spy

So, first off, we should discuss what each part of a Printed Circuit Board, or PCB, is - what they look like, what they do, and how they operate.

1.1 Form Factors

There are a few things to note about form factors on PCB's - when researching what a part is, or just working out what it could be.

1.1.1 Through-the-hole

When you get your electronics kits from Maplins/RadioShack/wherever, they are generally through the hole kits - that means that the components poke through the PCB, and you solder on the underside and snip off the extra wire. These are very common for hobbyist kits, but much rarer on production PCB's. As such, I'm going to mention them and move on.

1.1.2 SMC

Surface Mount Components was a bit of a revolution in PCB manufacture. The idea is that you can just make a PCB, put on a few thin blobs of solder paste, use a robot arm to place the components on top of the board, and then bake in an oven until it's all bonded. SMC PCB design was instrumental in the full automation of PCB population and testing. It made life stupendously easier for manufacturers.

SMC's are packaged on reels of plastic that resembles bubble-wrap. Each plastic well contains one component, and the Pick-n-place machines (which, let's be honest, do what they say in the name) can easily load a component into a robot arm and then gently drop it onto the PCB in the precise right place.

Many Integrated Circuits (ICs) are in an SMC format called 'Small Outline Integrated Circuit' or SOIC for short. These allow for easy access to

pins with test clips (discussed later on) and allow for easy identification of components and traces between them.

Another common IC package type is Ball Grid Array, or BGA. It allows a chip to be much smaller, but have many more active pins in the smaller packaging, that are resistant to the leg-shortcircuits that can occur when soldering very tiny pins coming from an IC - it is very easy to have short circuits between them, owing to the sticky nature of hot solder on tinned metal surfaces. Indeed, the prevalence of BGA's actually encourages one of my favourite things to find on a PCB - JTAG debug capability.¹

In the next few sections we'll assume that your PCB is SMC for the most part.

1.2 Basic Components

1.2.1 Resistors

Resistors essentially provide some level of resistance. They are basically essential energy wasters. *Essential??* Oh yes - we need to talk about short circuits for this.

A short circuit is when the power finds a path of such little resistance, that all the energy in our power supply (battery, PSU, perpetual motion machine², etc.) will race to the *path of least resistance* - it's not a cliché, it's how electricity flows. So, resistors let us manage this flow by setting different resistances.

You recall that Ohm's law that you heard about? Well, in the land of hardware hacking, it's everywhere and influences nearly everything you will be doing, so understanding what it's about is key.

Resistance lets us do many things, but chief amongst them is control the current through a device. By Ohms law - $V = I \times R$ - we know that if we fix the voltage and the resistance, then we fix the current at $I = \frac{V}{R}$. This is a very useful property - you'll see later how we can actually use this to measure current.

The greatest resistance we need consider is through air - hence traces can get very very small and close together, as our voltages are so low, we don't really come into contact with the high energy stuff. Building in a resistor lets us use components without worry of shorting our power supply (lots of fuses and trip switches are in place to prevent massive surges caused by a short circuit - mainly because things go 'bang' when that happens).

To measure resistance, use a voltmeter with a resistance measuring setting. I'm not going to talk about colour bands on resistors, as so few of these are ever found in the wild, it's not worth knowing that stuff anymore.³

¹See relevant section in the testing chapters.

²To any physicists - IT'S A JOKE! Tell your face...

³It really used to be a thing, but now with SMC's, it's not really worth our while

1.2.2 Capacitors

Storage tanks that use an electric field to hold voltage. ...and in normal words? Put simply, a capacitor is two plates that are very very very close together but don't touch. The various parameters of the body of a capacitor dictate how much energy it can store (measure in Farads, but a Farad is huge, so most capacitors you'll meet are in pico-Farads (pF) or microFarads (μF - that's a Greek 'mu' letter, but very common to see in its place is 'uF'). The electric field between these plates is what stores the energy inside a capacitor.

Capacitors are used for all sorts, but the most common usage is on power lines and in filters - because they store energy, you can discharge that energy when you need some backup. What do I mean? Well, powerlines are notoriously fickle and unreliable. As such, a capacitor can store charge and then 'top up' a power flow to help smooth it out. Most modern electronics are very sensitive to power fluctuations (see the section on glitching), so this is something you'll spot very easily.

These are the things that whistle and whine when a TV power supply is dying, and the larger ones (look like little cylinders) have 3 or 4 arm crosses on top. These are so that when the bit between the plates - the *dielectric*, meaning, 'between the electric bits' - fails (which it can do) you can see it, because this plate will 'pop-out' and the component will just look really different from the others. Now you can go open a business fixing TV's.⁴

1.2.3 Inductors

The 'other' electrical energy storage device - these store energy in a magnetic (as opposed to an electric) field. These are commonly used in analogue electronics and signal processing circuits (read: radios such as those in WiFi modules, BLE, etc. etc.) and in power blocks (both as one component of a transformer and in filtering).

They have the ability to block AC current but allow DC current to pass - as such, they are found in switch-mode power supplies - when used in conjunction with capacitors, you can eliminate lots of noise and 'ripples' in the power output.

We're not going to see a huge amount of these initially, so don't worry for now, but it's important to know that they are there, and that they are a thing.

discussing it.

⁴Don't. There's a lot of high-energy parts in old TV's that pack enough of a punch to seriously harm, or even kill you. If you're daft enough to touch those components, then don't put yourself on my conscience.

1.2.4 Transistors

Electronic switches. They have three pins - collector, emitter, base. Collector you can consider to be where 'stuff goes in', emitter is then self explanatory, and the 'base' is the switch bit. When base is 0V, then nothing flows from collector to emitter. When the base is raised over a particular voltage, then a signal can flow from collector to emitter. The simplest use is as an amplifier - a weak audio signal is sent to base, and this makes the higher voltage from the collector flow according to the base signal. Digitally, this means that you could replace slow valves in computer circuits - valves being electronic switches that rely on a heating element to give out electrons and are thereby slow and also prone to needing to be replaced.

The advent of the transistor was only met by the advent of the LED in the history of electronic engineering - it still may be the most significant breakthrough in this field during the 20th C. regardless.

1.2.5 Integrated Circuits

IC's can be anything. **ANYTHING!** IC's can have moving parts⁵, capacitors, transistors, inductors - literally, anything built into them. To be useful, capacitors and resistors and inductors tend to be external (as very small versions don't have the same effectiveness). They're amazing little doo-dads, and we're going to spend most of this book working out what's going on inside them. There are many tricks and tips for getting knowledge about IC's, but in general, this is the whole aim of hacking hardware - accessing what's going on inside IC's that we can't see.

1.2.6 Other Common Components

Before we get to some more interesting stuff, let's cover some other common components:

Power Regulators

- these take in one voltage, and either step it up or step it down to another. They're very reliable, and often get very warm (so they could well have a heat sink on them - don't remove it!). They have three pins, generally, and are usually next to or on a power line.

How to identify a power line? To keep the resistance down, these are often the thickest traces on a PCB.

⁵For a fun example of why these are important, see the documented issues with SiTime Micro-Electro-Mechanical Systems, or MEMS, clock oscillators, and how being exposed to Helium caused them to stop working... No, really! It was something that only affected iPhones, as they are the only ones to use them. Here's the article on Hackaday: <https://hackaday.com/2018/10/31/helium-can-stop-your-iphone-maybe-other-mems-too/>

1.2.7 Miscellaneous

[I'm sure I'll find some misc stuff to put here...]

1.3 IC Identification

So, this section is intended to give you an overview of ‘usual packages’ that these functioning IC’s are found in. That said, there’s something that should be noted explicitly:

The only way to give a certain ID of a chip is via the model number on the top.

Granted, this is not always possible (cf. the ‘sandwich’ method used to put the RAM on top of the SoC on PiZero boards, or the placement of heat-sinks that take light munitions to remove). Nevertheless, it is a good idea to have an overview of the function of each chip, with some notes about how things are achieved/connected and what to spot.

ProTip: If you can’t read the text on an IC, then rub some chalk over it, and lightly brush it off with a cosmetics brush. This will leave the chalk dust in the pits of the etched surface, and will stand out against the black epoxy - this makes is superbly easy to read with no effort.

1.3.1 SoC

The ‘System on Chip’ was a recent revolution in IC manufacturing. The first SoC was found in a digital watch, and this marked the start of a steady flow of such IC’s that would steadily condense greater and greater functionality into smaller and smaller footprints.

The core idea is to have one IC that does all the main functionality. Nowadays, they generally run a specific firmware loaded off a nearby flash-based ROM chip. They generally do a lot of stuff for the developers these days - they tend to implement USB stacks (saves you big-banging your own), as well as Wi-Fi and Bluetooth/Bluetooth Low Energy (BLE) all in one package with a MIPS/ARM (or some other RISC-like CPU) all built in. They can have their own built in RAM, or they can defer to other chips for this. They can have their own in-built secure storage, to act as a sort of ‘bootloader’ or just a place to put juicy stuff like private keys and such.

SoC’s run the show, and they fetch and load the firmware developed to perform the tasks at hand for a particular device. Thankfully, many many SoC’s run firmware that is essentially a pared down linux installation. No GUI, but they can drive screens, interface with protocols and other chips, and crucially, have a standard set of debug/development protocols that you can use to look under the hood of a device’s working.

1.3.2 Flash ROM

A ROM (Read Only Memory) chip is effectively a solid state storage chip. Much like you get in USB sticks, phones, SD cards, SSD's, etc. How do they work? Well, they have an array of ‘floating gate’ transistors that are arranged as if they are in a NAND circuit (a well known layout if you use transistors a lot). A floating gate has a captured electron in an isolated piece of substrate that influences the field effect of the collector. But how does that electron get there? Well, quantum tunneling.⁶

ROM Chips, aside from this modern physics wizardry, are very straightforward devices to look at. They tend to use SPI. (see protocols chapter)

1.3.3 RAM Chips

Random Access Memory (RAM) is a well known computer concept, but here you don’t get RAM on little snap-in boards, but rather long, flat IC’s with lots of connections. These tend to stand out as the traces have to be of precisely the same length - as such, serpent-like ‘wiggly’ traces are common going to/from these chips to the SoC - these keep the trace lengths the same and stop the signals going out of sync.

1.3.4 Interface Chips

USB, Ethernet, Bluetooth/BLE, Wi-Fi - these are all quite specialised protocols, and people long ago coded these into silicon. As such, SoC’s will potentially have an interface chip for any protocols they don’t natively support, or to act as latency-free support chips (if the on-board system is doing a lot of processing, the on-board ethernet may well suffer, so an external ethernet controller IC may be brought in - very common thing on CPE’s).

Common interface-specific IC’s you will find include:

- Bluetooth/BLE controllers
- Ethernet controllers - referred to as ‘PHY’ IC’s, as they correspond to the physical layer of the OSI model
- DECT phone controllers
- Wi-Fi controllers - these are usually shielded on many PCB’s so may have to be identified from a bootlog
- USB Controllers - a USB stack takes up valuable firmware space, so for full USB support, small controller chips are used instead

⁶Specifically Fowler-Nordheim, or FM-tunneling

1.4 Electricity and Electronics basics

[This section will be a reference for some basic formulae relating to Electricity and electronics in general. Included as a reference just in case it's ever needed.]

Chapter 2

Toolbox

There's a lot that you need to know about when doing hardware, and some specific skills include:

- soldering
- connecting to hardware interfaces (like UART, SPI, etc.)
- Reverse Engineering PCB's

2.1 Hardware Hacking Toolkit

So, the protocols we will cover in the next chapter are all well and good, but how do you connect to a UART port? How do you talk to a chip over SPI? How do you get your computer to talk JTAG? None of these actions are particularly hard, but you will need some kit.

Here I will present my software and hardware toolkits - I've put these together so that you, the reader, has a 'shopping list' of tools and equipment to get/download in one place. Some of this will be repeated in later chapters, along with the protocols in play for a full picture.

2.1.1 Software

2.1.2 Connection Software

First, here's some software that I find useful for connecting to things like a bus pirate (see below):

minicom

A nice interface, reasonably intuitive, however sometimes prone to just not letting you type input to a device. Minicom is very reliable for logging and connecting to devices, and doesn't kick up a fuss nor crash out easily.

screen

‘screen’ is amazing! All your key-shortcuts you had to learn to leave running processes on a server work here (Ctrl+a+d for example). It works as follows:

```
|| screen /dev/ttyUSB0 115200
```

Here, 115200 is the baud rate (see the walkthrough), and /dev/ttyUSB0 is the device you want to connect to. It’s really simple, and will also log your output with the -l argument. Swiss-Army tool indeed!

flashrom

Flashrom can be installed using yum or apt-get, and essentially is designed for flashing BIOS ROM IC’s. However, it has a config for going over buspirate, and is very handy for downloading firmware images (including bootloaders) directly off the flash IC’s themselves. Here’s an example command:

```
||>$ flashrom -p buspirate_spi:dev=/dev/ttyUSB0,spispeed=1M -o ./firmware-dump.dat
```

This will dump the firmware down in no time! (Well, up to 20 minutes, but hey, there are harder ways of doing it!)

2.1.3 Analysis Software**binwalk**

Created by DevTTYS0, binwalk (<https://github.com/devttys0/binwalk>) is a superb firmware analysis tool. Replete with entropy analysis capability, it parses through a binary firmware file and identifies important offsets as well as compression/filesystem types and settings. If something is LZMA compressed SquashFS in Little endian mode, binwalk will probably see it and tell you.

The -e extract tool is monumentally useful, too - if it sees something it knows, it will parse out that data and decompress it, presenting it to you in a new subfolder for you to inspect. I will be using this in the firmware analysis walkthroughs, so I can thoroughly recommend getting ahold of this.

Jefferson

Where binwalk fails, others pick up the mantle - one such example is Jefferson (<https://github.com/sviehb/jefferson>) Created by sviehb, Jefferson is a JFFS2 compatible firmware analysis tool. This is very useful when analysing firmwares from vendors like Huawei, who have a tendency to use such formats.

‘Logic’ by Saleae

The Logic package from Saleae (<https://www.saleae.com/downloads>) is a free-to-use logic analysis software. Interfaces with many analysers, and provides a nice suite of on-the-fly analysis tools. Simply set you sample rate top left, connect the logic analyser (see below), press record, then power on the device/signal you want to record.

Analysis includes the ability to play with baud rates to see which work when assessing a UART connection, through to sniffing SPI traffic (if the sample rate is high enough). Supporting Saleae by buying their stuff is worth this software alone - they could charge a lot for it!

OpenOCD

The Open On-Chip Debugger software is an open-source debug protocol suite that interfaces with the latest bus-pirate firmware, as well as many other devices like a GoodFET. To install for use with a bus pirate on a Debian baesd system:

```
|| sudo apt-get install libtool autoconf texinfo libusb-dev
    libftdi-dev
git clone git://git.code.sf.net/p/openocd/code
cd code
./bootstrap
./configure --enable-maintainer-mode --disable-werror --enable-
    buspirate
make
sudo make install
```

You’ll then need to add a `MyConfig.cfg` file with the following contents:

```
|| source [find interface/buspirate.cfg]

buspirate_vreg 0
buspirate_mode open-drain
buspirate_pullup 1

buspirate_port /dev/ttyUSB0
```

Alternatively, if you’re running with a Shikra or other FT232H board, you’ll need a config like this:

```
|| interface ftdi
ftdi_vid_pid 0x0403 0x6014
ftdi_layout_init 0x0c08 0x0f1b
adapter_khz 1000
```

In either case, copy this file to `/usr/local/share/openocd/scripts/interface`, or wherever your openOCD installation locates its script files. You can then run OpenOCD with the following command:

```
|| sudo openocd -f interface/MyConfig.cfg -f target/<target-device
    >.cfg
```

In this command, <target-device> refers to the config file for your target. We'll cover more detail on how these files work in the JTAG chapter.

You can now access the OpenOCD running instance by going to `telnet localhost 4444` and interacting with the device.

Using GDB will be discussed in the advanced section, but once started, GDB can connect once loaded by typing

```
|| (gdb) target remote localhost:3333
```

We'll cover much more detail later on about how JTAG/SWD works and how to effectively use them for debugging hardware. This quickstart guide, however, should be enough to get you going if you are having some issues.

2.1.4 Hardware

Arduino (any flavour)

The Arduino broke free the microcontroller (MCU) world to hobbyists and hackers alike! [I will likely do a full section on this later] From implementing a custom protocol, to using them for general purpose hackery, they are very versatile and useful devices. Having a few of them hanging around has saved my life on more than one occasion; my most common use for them is bit-banging¹ some protocol or other in order to get the target IC to believe that some event/input is taking place.

FT232H Board - the Shikra and Beyond!

Possibly one of the most versatile IC's for hardware hacking to emerge in the last few years, the FTDI FT232H multiprotocol IC to have emerged in recent years! The list of protocols this things supports natively is impressive:

- UART
- SPI
- JTAG/SWD
- I²C
- GPIO pins controllable over USB

A very versatile chip, they first came to prominence in the Shikra board from Xipiter systems, however comparable clones have emerged on AdaFruit and AliExpress. **Note** - make sure you get a board with 5V and 3.3V output. You'll thank me later!

¹We'll discuss this later, but bit-banging is essentially sending 0's and 1's corresponding to a protocol over a digital output pin by writing binary strings to the pin in order to simulate some functionality.

Bus Pirate

These little beauties are ace! They speak many protocols - the main things to do are as follows. First off, get ahold of a version 3 hardware - the earlier versions are great, but don't support the more advanced firmware.

Second, get ahold of the firmware that supports OpenOCD - Open On-Chip Debugger is a protocol that lets you do debugging in the chip itself by means of the debug protocols like JTAG and SWD. This may seem a little advanced, but it's worth entertaining as an idea - it will give you valuable information, as well as another avenue with which you can dump the firmware.

[I'll put a section here showing how to do the upgrade if it's of interest.]

JTAGulator

Though one of the more upmarket items, these things are *invaluable* and comes equipped with 30 test points and a hot-pink PCB to boot. What do they do? They let you find debug ports such as JTAG, hence the name! A much updated version of the JTAGEnum firmware for Arduino, the JTAGulator does all of the combinatorics of working combination of some mess of wires soldered onto test points will give you access to a JTAG debug port on a target board. We'll be discussing in-depth workings later in the book.

Available from Grand Designs, having been designed by the amazing Joe Grand, they can be bought assembled from Parallax. For details, see here: <http://www.grandideastudio.com/jtagulator/>

IC Connectors

There are two main ways that we will connect to IC's.

Spring-Loaded Taps - The first are using spring-loaded 'leg taps' that come attached to 10-line cables, usually with female pin connectors at the other end.² These are little connectors with small metal grabbing arms - when you press down, the arms extend and open out. You can then place them over the device pin you want to investigate, and release slowly. The arms will grab onto the IC leg, and you can then access the voltage going across that pin by hooking the other end of the connector to test equipment.

[Photo of these taps]

8-pin and 16-pin SOIC Clips - SOIC clips are *really* useful, as they are designed for debugging SOIC package IC's. They usually have a marked cable that is supposed to map to the marked leg on the IC. The clip then just sits on top of the IC and lets you analyse the signals from a pinout you connect to the other end of the clip cable (provided with the clip in most cases).

²I usually refer to them as 'bird's nests', as that's what it feels like you're untangling sometimes.

[Photo of SOIC clip and of a SOIC clip connecting to a flash ROM chip]

Both of these connector types are relatively inexpensive, and will serve you very, very well!

UART to USB adapter

Usually marketed as ‘CP2102’ chips, this IC is designed to do USB-to-TTY UART conversion, mapping directly to /dev/ttyUSB0. Connect the device, start screen or minicom, and then start up the target device and you’re all set!

Multimeter

These things are simply invaluable! Checking resistances, voltages, continuity - the average multimeter is a versatile and indispensable piece of equipment. A decent one only costs ~20, but they last years and will be referenced throughout the book. I will be assuming you have access to one, they’re that important.

If you want to go all out - get a Fluke multimeter, and never worry about accurate, reliable measurements again (so long as you keep it calibrated).

Logic Analyser

These pieces of equipment are very useful - they are essentially multi-channel digital oscilloscopes. They don’t measure fine-grained voltages for tracing the actual voltages on a chip, but so long as it breaks a threshold (usually 0.5V) these will record and send you data that you can then do analysis on. We will cover the precise methodologies for doing this later on, but suffice to say, these are worth every penny!

Saleae (<https://www.saleae.com/>) make some of the best consumer-grade ones around, but many oscilloscopes come with some implementation of a logic analyser built in these days. Saleae’s ‘Logic’ software is also free to download, and compatible with many of the ‘imitation’ devices found on the market.

Oscilloscope (*Optional - but only to a point... You will need one sooner or later.*)

Although only really required when doing work with glitching or more advanced testing/reverse engineering, an oscilloscope can be a very handy tool. Essentially, most modern desktop oscilloscopes are self-contained units that allow you to capture and analyse analogue and digital signals on the fly. They have many many options, and proper usage of the most advanced models can yield significant results!

That said, many cheaper devices connect wirelessly or over USB and can be particularly easy to use. Beware, however, as many of the cheaper ones don't isolate their ground signals, as such, creating a ground loop³ can be a significant risk!!

For starters, consider a picoscope 2206B from PicoTech:

<https://www.picotech.com/oscilloscope/2000/picoscope-2000-overview>

Bits and Bobs

So, there are things that are just worth having around, the list including:

- Wire - just generic plastic coated wire, and some thin copper enamelled wire is so handy to have. Tin the ends (see soldering) and you've got something you can connect to a PCB in minutes.
- Breadboards - also called 'prototyping boards', these are white plastic boards you can plug into. I'll discuss below how they work and how to use them, but they are a worthwhile investment.
- Spare Component Selections - I recommend having a stock of the following:
 - Resistors - a range between 10Ω and $100k\Omega$ is best.
 - Capacitors - a range of small ceramic ones and larger cylindrical ones is best. From a few pF to a few μF is a good idea. Handy for setting up aux power supplies, if nothing else.
 - Transistors - in case you need them for a specific circuit
 - LED's - if only to give a visual signal that something is happening, LED's are low-resistance sources of confirmational light in my circuits.
- Breakout Units - Arduino 'starter kits' come with 1,001 little breakout boards. Sensors, joysticks, stepper motors, LED digital output blocks - there are all sorts in these kits, and they are fairly inexpensive. If nothing, it'll give you something to learn Arduino programming on.

2.2 Soldering

AKA 'How to play around with molten tin' - it's actually a lot of fun. Fundamentally there are three parts to soldering:

³These occur when a signal has multiple paths to ground - through induction and other actions, they can gather energy and start oscillating. Search online for the kinds of noise/signals that you can generate, and then remember that such powerful signal loops are usually well beyond the tolerances for most IC's.

1. soldering iron - ideally a soldering station that is temperature controlled, and has a heat-gun for reflow work.
2. solder - leaded or without lead, good quality with a rosin core and at least 2 sizes; thicker 3-4mm and thin 0.5-1mm diameter solder.
3. flux - either a flux pen or flux gel.

For best results, we highly recommend a temperature controlled iron, and if you can stretch to it, one that comes on a reflow station with a hot-air reflow gun to boot. One of these stations will mean you will never want for soldering capability!

2.2.1 Soldering Basics

Solder is an alloy mostly consisting of tin. Tin's melting point is 231.9 °C which is very low compared to most metals - it'll melt at the highest temperature of a domestic oven, and as such, is widely used owing to the fact there's quite a bit of it hanging around the planet, and it's relatively safe to use - no special handling is required, melting point is low meaning no high temperature equipment is needed, and it's also quite soft and easy to work with. We use it to seal our copper piping, connect up our copper contacts on PCB's. Tin in the form of solder is a useful material.

THAT SAID - 231.9 °C will give you a nasty burn, and as always, take great care when soldering. This isn't a toy, it isn't a game, and if you don't think you should put a hot soldering iron there, then you're right. So don't!

Most solder alloys mix tin with silver and copper (and sometimes a fourth element such as manganesees - as in Sn-Ag-Cu-Mn solders used in BGA reflow work) which pulls the melting point down to around 183 °C, or thereabouts, so it can be worth a quick check what the melting point of your particular solder is.

Soldering Technique

Here is a basic solid technique for doing soldering:

NB - I'm assuming you're using a soldering iron tip that is suitable, and solder wire that has the flux built into it. 'Fluxed solder' is much easier to use, so if you're new to all this, I recommend starting with that.

1. Make sure the iron is at full temperature - if that means waiting a while, wait. When you start heating things up, it will pass a lot of energy on to the things you want to work on, so make sure it's good and hot. Got make a coffee or pet a cat or something...
2. 'tin' the soldering iron tip - no matter what size soldering tip you are using, dab solder on the end, and wipe it off on a wet sponge. This will

leave a shiny surface that solder loves to stick to - *Warning* - it will oxidise and brown very quickly, so this will need to be done regularly.

3. tin all the surfaces you want to solder - tin sticks to tin really well. So, if you're soldering wires it is best to tin them first. If you're soldering to a circular contact test pad, then prepare this first (see below for both)
4. Once you're ready, make sure the PCB is well anchored/secure - if things fall over you'll damage some component, or even yourself, so avoid leaving things hanging precariously.
5. Once the components are in place and ready to be soldered, use the tip to heat up the area where you want the solder to go, and then push solder-wire onto the hot soldering iron tip
 - What you want in this act is to flow the solder off the iron straight onto the hot tinned surfaces. It'll flow beautifully once you get the knack of it, and you'll get perfect solder joints.

Practise makes Permanent!!

This older mantra from my violin teacher at university⁴ is the much improved and true-to-life version of ‘practice makes perfect’. ‘Practice makes perfect’ is wrong, if you think about it - if you keep practicing something wrong, then you’ll always get it wrong!

The trick is to do what anyone who gets good at something does - learn a technique, and then practise it slowly, and concentrate on getting the technique down. You’ll soon find that your brain takes over in a sub-conscious way, and you’ll get the perfect solder joint every time!

To do this with soldering, you can buy ‘solder Training Kits’ that give you various footprint and some old/obsolete IC’s and components to solder to a PCB, and two PCB’s to train on. It’s worthwhile, and once you have the technique nailed down, it’s a fairly straightforward exercise.

Preparing Solder Contact Areas

Here are some tips on preparing wires and contact points for soldering:

- Use an appropriate wire - if the thing you want to solder to is thin/delicate, use a thinner wire
- Strip the wires with plenty of space to get some contact when you join it to a board

⁴Yes, I did my UG in Music, and Tony was my violin teacher who drilled this idea into me!

- When tinning plastic coated wire:
 - Strip the plastic and make sure the inner wire is clean
 - Twist the wire between finger and thumb to turn the loose copper wires inside hold together in a ‘spiral’
 - Place the wire in a holder, with the metal pointing up, and with a tinned soldering iron tip and fluxed solder, cover the end of the wire from the bottom of the stripped section to the top in a smooth quick motion.
 - Allow the wire to cool before touching it - it will only take a few seconds.
- When tinning enamel-coated copper wire:
 - If the wire is 0.2mm thick or less, put a blob of solder on a wedge iron tip, and holding the wire in pliers, dab the end in the blob of solder to remove the enamel and tin the wire, cleaning the iron tip on a sponge thoroughly afterwards.
 - if the wire is 0.2mm-1mm thick the enamel can be removed using a lighter or rubbing the end in extra fine sandpaper. Tin as normal before using.
 - If the wire is >1.0mm thick, it is best to use fine to extra fine sand paper to remove the enamel before tinning the end of the wire.
- Apply solder to contact points you want to solder to:
 - Often a plastic layer is on top of the pad, so use an emory board or scratch the surface with a knife before applying an iron
 - Using a blob of solder on a wedge tip, gently apply the iron to the contact point for around a second, and pull back to see if the solder has adhere.
 - careful repetition of this will prepare the contact pad to be soldered to with a jumper wire
- if soldering to a leg of an IC, line the wire up carefully, and gently rub a tinned iron tip over the pin and wire to bond them. **Be very careful** not to short the pins, and check with a voltmeter to ensure no erroneous continuity (short circuits) have been introduced.

It is also possible to use the ground-plane (the sheet of copper that covers most of a PCB and acts as ground shielding) to anchor some test pins in the following manner:

1. Cut an appropriate number of pins to have an anchor at either end of the jumper (so, if you need 5 pins, cut 7 from the strip of pins)
2. Find a free area of the PCB edge that can accommodate the new pinout, and mark the pins at either end.
3. Scratch off the solder/silk masks from the PCB in order to expose nice shiny orange copper.
4. Solder the pins onto the ground plane:
 - First, apply solder blobs to the scratched-off areas
 - Tin the first and last pins as you would a wire
 - Solder the pins onto the solder contact points we carved out, and check it's secure.
5. Now gently put small blobs of solder onto PCB side of the pins
6. Connect the tinned ends of the jumper wires you want to secure to the pins.

[Pictures to follow]

Doing this after you've identified and soldered jumper wires onto the relevant contact points (usually for a JTAG that doesn't have a development header on the PCB) will produce a 'makeshift' pinout for the interface you want to access.

[More to follow regarding how to solder, soldering tips, and what's going on when you are soldering.]

Preparing BGA for removal and replacement

BGA's are all the rage! It's because of these being used for everything from eMMC to SoC's that we need to have effective means for dealing with these. Thanks to Deral Heiland⁵ from Rapid7 for his work on this!

BGA Removal

To achieve this we will need a heat-gun soldering station, some IPA⁶ if necessary, and a little patience.

1. First off, clean the area around the BGA chip.
2. Set the heat-gun to an appropriate temperature for the thickness of PCB

⁵https://twitter.com/percent_x

⁶Isopropyl Alcohol

- **Note** - fibreglass, the substrate of a PCB, and copper, the conductive overlay on fibreglass, are both pretty decent conductors of heat.
 - As such, it is worth remembering that you need a higher temperature on the heat-gun in order to compensate for this differential.
 - A reasonable guideline temperature is to start at 350 °C and go up from there if nothing happens after a minute or so.
3. carefully saturate the area around the BGA with heat from the heat-gun, making sure to spread the heat *evenly* around the target BGA.
 - This is important - it's possible to damage the PCB by causing the copper pads to delaminate (*i.e.* pull off from the PCB) or damage the BGA component if you focus all the heat in one area.
 4. using tweezers, test the component by nudging it from the side to see if it will move.
 5. If it moves, keep the heat present and remove the component.

You can now load this into an eMMC reader or a custom board you made earlier for the purposes of reading the chip off (using, say, pogopins or a read/write clip for the BGA package layout).

BGA replacement

Now suppose that, after all this, you want to put the BGA chip back.⁷ This is where Deral's work really comes into play - the following is derived from his talk at Nuit du Hack in Paris, 2018.

For this, you will need:

- A BGA to solder onto a PCB
- An infrared reflow oven (with temperature curves, ideally)
- A vial of solder balls the width of your BGA solder pads
- A microscope
- Some solder removal wick
- A flux pen and IPA for removal of excess flux
- A soldering station with a heat-gun and temperature controlled soldering iron

⁷It's always a good idea to try and disassemble a target device in such a way that you can reassemble it in a working state later on. Many hardware hackers subscribe to this ideal.

- Nerves of steel/soothing music⁸

Here's the rough workflow:

1. First of all, use the solder braid and a broad soldering iron tip with a touch of flux to remove excess solder off the BGA footprint on the PCB
2. Apply a thin layer of flux from the flux pen to the surface of the BGA
3. Get the BGA component under the microscope, and using tweezers, carefully place a solder ball on each pad.
 - Although this sounds like a task dreamt up as a special for of water-torture, it's actually a reasonably straightforward process. With a well-anchored steady hand, a reasonably sized BGA can be 'reballled' in this sense within ~20mins. No, really.
 - Also worth noting, the flux is going to act as a kind-of adhesive to stop the balls from moving around. Too much, however, and in the next step (putting the chip in the oven) the solder will join together and form a clump. If this happens, you have to clean up, and start again.
4. Carefully move the BGA to the reflow oven, and set the temperature to have a maximum just above the melting point of your solder for ~1min. On removal from the oven, it should be completely reballed.⁹
5. Apply some flux to the PCB footprint, and carefully place the reballed (and cooled) BGA on the PCB footprint in the proper alignment - there are usually silkscreen symbols to help here.
6. with the heatgun or in the reflow oven - reflow the BGA onto the footprint.
 - for a reflow oven, set the temp to a max of 250 °C, and if using a heat gun, pre-warm the area and then set to 350 °C
 - if doing this with a heat gun - keep an eye out for the BGA 'dropping' onto the PCB. As soon as this occurs, you're done, so STOP APPLYING HEAT!
7. Clean up as usual - with IPA and a sponge/cotton pad for the flux, especially - and do whatever continuity testing you can.

⁸If it's your thing, now is the time to break out the Enya.

⁹This is usually done in factories by a machine that drops a ball of solder onto a pad and zaps it with a laser to melt it on... so this is the long-winded by-hand low-fi DIY version of that.

With these methods, you should be able to remove and replace BGA chips with relative ease. There are other techniques, but these ones are reliable and proven to work - though the replacing technique does require some more initial outlay.

Chapter 3

Protocol Overviews

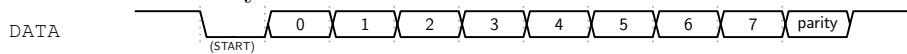
Lots of protocols are in play in hardware, so this chapter will give an overview of the main ones used by electronics IC's/systems to communicate with each other. Ultimately, the aim will be to use these to let us interface with the devices we want to, at a low level. This is the window that hardware hacking opens up to pentesting.

A ‘serial’ protocol is simply one that acts sequentially, usually in an ‘asynchronous’ manner (that is, like a turn-based game of chess, rather than a synchronous game of Command & Conquer). We will cover the main ones here, and have some special mentions at the end.

3.1 UART

Universal Asynchronous Receiver/Transmitter, or UART for short, is *so ubiquitous* that of all the many and varied serial protocols out there, we just call it ‘serial’. It’s a very easy protocol to understand - there’s no clock signal (it relies on a preset baud rate), and the only thing to remember is that $Tx \leftrightarrow Rx$ and $Rx \leftrightarrow Tx$ when you wire it up. But first, here’s how the protocol is laid out.

Like many protocols, this is a pull-up voltage protocol. The lines sit at 3.3V until they are pulled down to 0V to indicate a signal. Let’s suppose we want to send a full byte of data - 8 bits - over the wire from one end to the other. Here’s the layout:



Here is a breakdown of this presentation of the protocol:

- All clock speeds are set independently of the signal - the ‘baud rate’ is the rate at which bits will be transmitted. On this kind of setup, it is just the number of bits per second.
- There is one ‘start’ bit that gets the ball rolling.

- There is one stop bit (a low voltage cycle) that signals the end of the chunk of data
- 8 data bits in the middle
- 1 parity bit at the end of the data

Normally, the settings for (number of data bits/number of parity bits/number of stop bits) is (8/N/1) (**8** data bits, **No** parity bit, and **1** stop bit). This is generally shortened in screen and minicom to just ‘8N1’.

There are four connections to make this protocol work:

1. VCC - usually 3.3V positive
2. GND - ground
3. Tx - transmit, which is wired to Rx on the other device
4. Rx - receive, which is wired to Tx on the other device

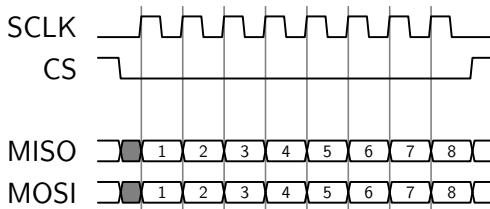
So you can see that there are 2 common pitfalls. The first is wiring the devices up wrong. At worst, this will fry your USB or the device itself. At best, the devices just won’t communicate. **Take care when connecting and powering these pins!** Second main issue is finding the baud rate. Thankfully, you can try many baud rates using the recorded traces from a logic analyser. The most common baud rates are as follows (in order of what I’ve seen the most):

- | | |
|---|--|
| <ul style="list-style-type: none"> • 115200 • 9600 • 57600 | <ul style="list-style-type: none"> • 19200 • 38400 |
|---|--|

It’s straightforward to script testing these, but if you try them in this order by hand (when you run screen in BASH for example), you’ll probably find the right one more quickly.

3.2 SPI

Serial Peripheral Interface (SPI) was designed by Motorola to permit the conversation between different devices based off a master/slave or controller/source relationship. As opposed to UART, SPI is synchronous, and specialised for high-speed, short distance communication. The timing diagram (due to Martin Scharrer)



Here's what the various parts mean:

- **MISO** - Master In Slave Out - the bus that lets data flow from the slave units to the master unit
- **MOSI** - Master Out Slave In - the bus that lets data flow from the master unit to the slave units
- **SCLK** - The clock signal pin - the rising edge of the clock triggers the level of MISO/MOSI to be read as the current bit by the target device.
- **SS/CS** - in order to select a device, this pin is grounded from 3.3V to 0V, telling the particular device to listen up.
- **VCC/GND** - don't forget that the chip will still need power and ground to operate!

So, as you can see, it's hard to wire SPI connections wrong! MOSI on the target maps to MOSI on the connecting device, and the same for MISO, as well as all the other pins.

NB - for the real geeks who already know this stuff I have chosen the clock phase and polarity to both be CPOL=0 and CPHA=0 in the above diagram.

3.3 I^2C

Where one manufacturer goes, another follows. Motorola may have had their SPI protocol, and Philips designed their I^2C !

Very unlike SPI, I^2C operates as a packet-switched protocol, permitting multiple master units, multiple slave units, and is engineered towards slower, short distance communication.

[Put a diagram here - give a proper description]

Initially, a license fee was associated with using this protocol, but this was waived in 2006. That said, the allocation of ID's to devices is still chargeable.

3.4 JTAG

The Joint Test Action Group (JTAG) had a problem to decide how to test PCB's, especially with the advent of Ball Grid Array and related packages,

where the pins are hidden away underneath the package? How do you check if two of the pin connections are bonded with a solder short circuit? What is the best way to solve this problem in a general way?

In 1990, after 5 years of work, the JTAG published IEEE Standard 1149.1-1990, entitled *Standard Test Access Port and Boundary-Scan Architecture*.

This solution is roughly as follows - each leg that comes from the main silicon die in the middle of a chip has on its path a small module. These modules are joined up into a ring, and collectively form the testing apparatus for the IC. Given this particular level of access, three main functions were defined for what a 'JTAG' should perform:

1. Debugging - allowing a developer to debug the chip, much like you would debug a local machine with GDB or related
2. Managing Firmware - erasing, uploading, and downloading
3. Boundary Scan Testing - this is testing the Boundary Scan Register (BSR) over the Test Access Port (TAP), or in normal terms, testing the IC pins for problems or issues.

[Put a diagram of how the JTAG looks vis-a-vis chaining chips together, etc]

JTAG doesn't just permit access to an individual chip, but rather, to all the chips on a JTAG chain. This makes debug tasks work more slowly, but gives wider access to the IC's on the target PCB. Most instances, however, are singular JTAG's for a particular IC.

JTAG is nominally a 5-pin protocol (not counting VCC or GND):

1. **TDI** - Test Data In, this is the input pin.
2. **TDO** - Test Data Out, this is the output pin, and when you daisy-chain IC's with JTAG, the TDO of one goes to the TDI of the next, until it loops back to the debug header.
3. **TCK** - Test Clock, continuing the 'test' naming convention, this is just the JTAG clock signal, the rising edge triggering a read operation. TCK is not chained, but rather forms a 'test clock bus' along with TMS - each IC can see the clock and TMS signals.
4. **TMS** - Test Mode Select, This is read as the clock signal rises, and determines the next state of the internal JTAG controller
5. **TRST** - Test Reset, an optional pin that can reset the internal test controller, but this isn't required.

The internal mechanics are as follows; There is an internal Test Access Port (TAP) Controller, just a mini-chip inside the main chip that manages the JTAG signals the IC receives. There minimally 3 registers - an Instruction Register, and two or more data registers. There is a defined state-machine that uses the TMS level to decide what to do after each clock cycle. This TAP controller connects to the boundary cells (the little bits of silicon I alluded to earlier, that are on every leg coming from the main silicon die) and drives their behaviour. Boundary Cells can raise/lower a leg's voltage to influence the behaviour of the chip.

As such, you can see that JTAG basically controls the chip, and a JTAG unit within an IC, if accessible to a skilled assessor (that's you), can be leveraged to get full control of the CPU stack, access to memory, as well as the ability to dump/upload firmware without pesky restrictions.

[Signal diagram and in depth descriptions to follow]

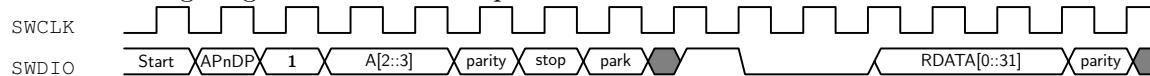
There is also a cut-down version of JTAG called 'reduced pin count' and this just has TMSC (Test Serial Data) and TCK (Test Clock) pins in play. This, however, was phased out by more accomplished two-wire protocols such as SWD.

3.5 SWD

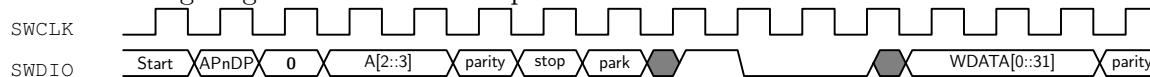
Serial Wire Debug bus was designed by ARM in 2004 to replace JTAG on IC's that use their processors (so chips that have inside a Cortex MCU, for example). Thankfully, most JTAG modules (like the J-Link by Segger) support this as a two-wire protocol from the data and clock lines (YMMV, so check their documentation).

[PTO for some timing diagrams]

Here's a timing diagram for a READ operation:



Here's a timing diagram for a WRITE operation:



The differences should be apparent - the 'grey blobs' are clock cycles where the host shuts up and lets the target talk. The three bits '100' that occur after every host opening byte sent is an ACK from the target - if that isn't received the whole thing loops back to the beginning.

The first byte has the third bit deciding if the READ or WRITE is required, and the other bits decide on the mode (AP ('1') or DP ('0') modes) and the ADDR bits decide which register is being talked to/requested from.

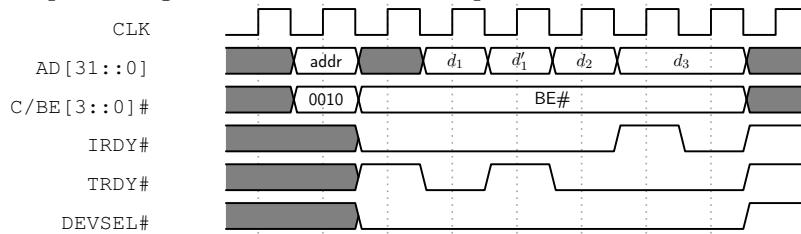
SWD has its own little state machine, much like JTAG, and you use the single wire to talk to it, and get access to the debug data (stack traces, etc.) for ARM chipsets. Neat-o.

It is worth noting that for both JTAG and SWD we will explore more examples in ??.

3.6 Special Mentions

Context Area Network (CAN) used in cars. Texas Instruments CC21xx series custom Debug Protocol.

template diagram - for use with ther protocols...:



Chapter 4

Things you need to know

Before we get stuck in, here are the things you need to know. This is the only chapter you **must** read.

4.1 Safety First

So, some safety pointers from a physical and legal point of view:

4.1.1 Hardware Safety

I remember a handyman at a job I once had in a factory telling me the following maxim for getting stuff done with hardware:

Right tool for the right job.

That's a very important idea. Although he meant don't use a MIG welder for something a dab of glue would fix, it's a very pervasive notion. Using tools properly and safely is hugely enjoyable! But improper use can lead to some real problems for you.

Here, we'll go through some safety tips - both for your and your person, as well as for the electronics you'll be using.

Soldering

First off, soldering. That thing is hot. Very hot. It's roughly the temperature inside your kitchen oven on its highest setting, except that it's metal and you're holding it in your hand. Soldering irons are very efficient at heating things to this temperature - it's what they're designed for! So, here are some tips when soldering:

1. Use an appropriate rest for the iron when it's not in use.
2. If you're done with the iron, turn it off.

3. Keep children and animals well away.
4. Don't be daft - if you think you shouldn't use an iron on something, then you probably shouldn't.
5. Safety eyewear is advised - if solder flicks into your eye, it's painful and difficult to get out.

PCB Handling

You remember at school - you were either shown a video of, or even a real Van de Graaff generator. Someone, usually a girl with long hair, would get on an insulated plinth and put their hand tentatively onto the large metal dome of the generator, and as it whirred her hair would slowly become charged and the strands would float away from each other. Remember that?

Good. The science you saw there was that the human body can hold enough charge as a voltage to literally make your hair stand on end, and not kill you. Sources vary, but it is safe to say that the body can hold hundreds to thousands of volts of charge, and you won't die so long as the current, the number of electrons involved, isn't high. So, as a rule of thumb, volts don't kill you, amps do.

Now think to our poor little IC's. Volts love them, and they love volts. But too many, and they die, never to work again. This is why proper PCB handling is important - it's trivially easy to do the wrong thing, gain a charge, and brush up against a PCB and annihilate all the lovely chips with their secrets.

1. 'Earth' yourself before you touch any electronics - it's as simple as just touching the exposed copper on a radiator.
2. Use an earth bracelet if you have one - attaching it to the exposed copper of a radiator is a good idea.
3. Socks on carpet is a bad idea.
4. Keep pets at bay, lest they do the job of zapping your target PCB on your behalf.
5. Get an anti-static mat and use it properly - see manufacturer's setup guide.
6. Buy anti-static bags on eBay/amazon and use these to store your previous target PCB's in. **Don't** leave them lying around unprotected in a box if you can avoid it.

4.2 Electronics Intelligence Gathering

Pro tip - being proficient in reading/basic knowledge of Russian/Chinese/German is very useful. Here are a few reasons why -

- Chinese has lots of open manufacturer forums,
- Russia has lots of forums where people have already broken into stuff or gathered information.
- freifunk is a really good resource, along with DeviWiki and OpenWRT wiki/forums.

[Expand on this list]

4.2.1 Handy Phrases

Here are some handy phrases to search for when doing intelligence on IC's on the internet:

German Phrases

Russian Phrases

Chinese Phrases

filename 下载 "Download filename"

数据表下载 "Download the datasheet"

Part II

Basic Assessment Techniques

Murphy's Law of IoT

"The router/camera/internet toaster in your home, with access to all your data, was made by the lowest bidder..."

Kirckhoff's Principle for IoT

"In any IoT system, some level of security should be assured, even assuming the complete compromise of at least one device."

Due to Mark C., inspired by Andrew Tierney.

3 Step Plan to IoT Nirvana

We present the following '3 Step Guide to IoT Nirvana'¹, originally by Mark C.:

1. Ensure you have per-device secrets.
2. Make firmware updates possible.
3. Make the security problem a named-person's problem.

¹No guarantee of actual Enlightenment is given nor implied by the author.

To include:

- UART locating and connecting walkthrough
- SPI Flash ROM dump walkthrough
- pentesting tips/techniques using UART and dumped firmware.

Chapter 5

The Hardware Pentesting Methodology

To get the most out of security assessing hardware, you should ideally be proficient in numerous areas of security testing - infrastructure, web, mobile apps - as these are all presented attack surfaces on modern devices. Hardware hacking can greatly enhance the testing of these surfaces, and offers up some new attacks and a level of detail that is completely missed if the device's debug/terminal functionality is not explored.

Here, we'll explore the motivations, realistic threats, and give an overview of the testing possibilities we will build on practically in the next chapter.

5.1 The Hardware Threat Model

5.1.1 Motivation

A matter of longevity

The ‘half-life’ of a vulnerability - the length of time for half the users to have their service/software patched against an identified threat - is nowadays very short; minutes for a hosted web service, even hours for a maintained executables and services, to weeks with clients who have strict roll-out testing requirements. However, a vulnerability in hardware will be around for **YEARS**. Once a device is installed, it’s there for the long run, and if there is one guy who procured it, went on the training course, but then left under a cloud and never wrote down how to update it - you’re really stuck until you upgrade. These are the true facts of hardware based security models.

Looking past just getting a shell

There are 1,001 guides online for doing ‘Hardware Hacking’ online, and they mostly stop at ‘...and now you’re connected’, leaving the reader to come up

with the relevant attacks, tests, and experience for themselves. Few places deal with what the risks really are, and as has been noted below, the risks are only accumulating. They are not going away.

Yes, I understand (more than most) the value of having a r00t shell on a device, delivered to you through UART, but let's be honest for a moment - **this is not a viable attack vector!** In the majority of cases, there is no way that you'll actually execute access to a UART on a victim's device without either their noticing you, or travelling around fully equipped to perform the installation and assessment of the development header.¹

It is important to fully appreciate what the threat model is from an attacker's perspective, looking at realistic scenarios for attack.

5.1.2 What hardware hacking means

First, a quick digression into what **I**, the author, mean by 'hardware hacking'. This is a very broad topic - it covers circuit bending, repurposing tech, breaking DRM on games consoles, accessing debug functionality on IC's built into product PCB's, designing and getting an arduino/RasPi/(generic microcontroller) to achieve a certain task for a project you're making, accessing/writing to protected ROM areas, reverse engineering IC's that don't have public datasheets, *etc. etc. etc.*

Hardware hacking is a *vast* topic, and I don't want to get bogged down with definitions - they're all valid uses of the term 'hacking'.

Specifically, for me, the phrase hardware hacking is linked to my work as a pentester and security researcher where I would analyse IoT/hardware from the standpoint of an attacker looking to get some malicious leverage on a system - so, my point of view for 'what to do when you have access to a debug/terminal interface with a chip?' is to use this deep level of access to test for vulnerabilities and security flaws. This is what I'm interested in, and so I'm making this explicit from the start.

As such, here are the kinds of questions I ask of a device I come into contact with:

- What else can you do besides (generic product description)? (you could argue this is a definition of 'hacking' in a generic way...)
 - Are you really a microcontroller being used to flash a light?
 - Do you have a firmware-controlled radio that could be repurposed?
 - Do you have some advanced functionality that goes above and beyond what you need for the task at hand, and that I can use for other stuff?

¹The only examples I have ever found where this is not the case is when such a debug port is publicly available, such as Cisco switches/ASA's in car parks without adequate physical security being provided.

- Can I turn you into a malicious device for blackmail?
 - e.g. bugging, video recording, disclosing sensitive user data or Personally Identifiable Info (PII) - that kind of thing
- Can I get data out of you that relates to your owner?
 - e.g. PII again, usernames, passwords, secret keys for 2FA, private keys for encryption, etc.
- Can I add you to a botnet/BTC miner?
 - easy ways to monetise pwned devices - get them to mine BTC to a wallet you control, or add them to a botnet and sell DDoS services a la Mirai clone botnets
- Can I think of any other malicious activity that would benefit me in some real way, that I can leverage to this device and others?

5.2 Pentesting Hardware - an Overview

We want to get you to understand the attack scenarios and ways in which you can pentest so as to cover off these attack vectors effectively in your work.

So, let us suppose you're a pentester and you've been given a device to hack as part of a test - what now?!? In fact, many categories of attacks relate to Hardware hacking from this security standpoint, in particular:

Outdated Software

There is a major supply chain issue in IoT. A recent mass hack² in the UK showed that there are some serious flaws in the generic underlying software presented on the devices. *But how the hell are these vulns here?* Well, in part it's a supply chain problem.

Often (incorrectly) referred to as *ZyXEL*, Allegrosoft³ manufacture software for embedded systems, including the very popular RomPager server software. It was clear that there was some issue with this software in the wake of this hack. However, according to sources⁴, Allegro had *fixed this vulnerability already in 2005*, which was revealed when checkpoint disclosed the ‘misfortune cookie’ issue, which was similar in nature.

²https://www.theregister.co.uk/2016/12/08/talktalk_routers_may_be_botnet_imperova_says/

³<https://www.allegrosoft.com/>

⁴<http://www.ispreview.co.uk/index.php/2014/12/masses-broadband-routers-hit-misfortune-cookie-security-scare.html>

So what went wrong? Well, Allegro make money, like any other company, so they sell their software to IC manufacturers directly. In turn, these manufacturers develop a chip, and then sell/give away an SDK with the chips so that companies can quickly get products to market. As such, the problem was specifically that the manufacturer had never upgraded, so the product designers never upgraded, and these devices are in place for years and years - so a vuln from 2005 can and will be a major pain in the arse in 2017 because of this major problem in the supply chain **delay-line effect**.

You will find web servers you've never heard of (because they were discontinued in 2003, like D-Link's Boa webserver), DNS servers older than most of your underwear, kernel versions that belong in museums, and best not to mention what versions of busybox they are running.

Exposing to the WAN what shouldn't be on the internet

The above hack was caused by a UPnP payload working on a TR069 endpoint. *But how did this happen?* Well, the binary that ran the TR069 service also ran the local (unauthenticated) UPnP service on the LAN. As such, you could exploit the one through the other, as they were, in fact, the same.⁵

Other D-Link routers expose the 'file sharing' feature; any USB stick plugged into the router has its contents displayed for retrieval from the WLAN/LAN over a web service on port 8181. This service, however, is also displayed on the WAN device. It has a default login of user 'admin' and blank password.⁶ It also has an auth bypass by forced browsing to [/folder_view.php](#) endpoint.

As is hopefully obvious, there were no need for the design and implementation 4AM decisions that lead to these issues.

Wi-Fi woes

Many attacks have been found against commercial Wi-Fi setups. But from this to this, it's clear that all your Wi-Fi experience as a teenager has not gone to waste.

Remote Command Exec (RCE)

Many instances where you can send a command over the internet and run commands on the local device as r00t. This is quite obviously bad - and there are many sources for 'what to do' if you find this kind of vulnerability. Many devices also run everything (even web pages) through BASH scripts because then they don't have the overhead of including a (probably as vulnerable) PHP service. As such, injecting RCE strings all over is likely to find something, somewhere, on the web-app used to manage the devices.

⁵Echoes of the Hesperus vs. Phosphorus argument from antiquity rings loudly, here...

⁶To see this, a shodan search of '[WebServer port:8181](#)' will display online devices.

Remote Code Exec

Likewise, connecting to some exposed service on the external device presentation and send crafted payloads to exploit BOF, format string vulns, etc. These services, if vulnerable, can easily be found using [shodan](<https://shodan.com>) and the like.

Cross Site Request Forgery (CSRF)

Using CSRF to fire in a user's browser, connect to the ubiquitous web based apps/API's used to manage the devices, and either get RCE or directly manipulate settings is a common vulnerability I have found on firmware-based devices. The required CSRF-Token checks needed to mitigate this can produce lines and lines of code, which is a hefty price to pay on storage-light firmware chips (remember - the max for most firmware chips is 16Mb). This vuln is going to be a common theme.

Cross Site Scripting (XSS)

XSS is always bad, and if there is a way to compromise the device with, say, a stored or reflected XSS, the user's browser as well as the device can be compromised. I'm not going to write here on the evils of XSS as many others have done this before me! But your honed bug-bounty/web-CTF skillz are not lost here!

File Upload vulns

Linking into my nightmare CSRF scenarios, file parsing is almost always comically bad on devices, and I will put money on that if there's a file upload on a new-to-market device, it'll have some vulnerability related to either Local File Inclusion, Directory Traversal, or straight up pwning (as seen with a device that just unzipped a file as root into /, including the malicious /bin/evil-ELF I included).

XML/SOAP based attacks

More and more devices are using mobile based android/iOS apps to manage these devices. As cool as this sounds, these devices open up a SOAP/XML/JSON based API to talk to the app (rarely implementing any authentication), and UPnP attacks, enumeration, TR069 attacks and enumeration (and others) fall into this category, too - but see 'cloud' for more.

Authentication / Session management issues

Default credentials are everywhere - literally! admin:admin or similar über-weak credentials used to be a joke in the office, until it was everywhere on

every device you could talk to that used some embedded platform. It is also a myth that most users even *can* change their passwords to these devices (many don't support this anymore - see 'cloud' below) so this is also worth reporting. Authentication information enumeration (usernames/emails) is also rampant, as well as authentication bypass - I remember doing a test on one device where just having a valid cookie (which you could guess or brute-force to get) was enough to bypass username/password requirements.

Information disclosure

If the device passes through much of the user's data (as in CPE's) then this data is at risk if the device is - and so these devices can act as a source for further malicious activity (dumping CC data, seeing if someone is in the house for IPCams, etc.). Don't forget, setting malicious DNS servers is very, very easy on CPE's once you've compromised one, and most users (even tech savvy IT people) won't notice, not ever.

Additionally, if a device offers up access to other services that they then manage (SIP/VOIP, FemtoCell or IPTV access) then disclosure of the internal mechanisms on these devices can reveal private keys for encryption, common shared secrets or credentials (often unencrypted/unhashed), certificates for authentication and encrypted access, DRM information and access methods, etc.

The 'Cloud'

...is just someone else's computer. But aside from all these issues (which are dealt with much better elsewhere) there is another concern. Amid this melange of written-at-4am-with-no-sleep and updated-when-West-Virginia-last-voted-Democrat software (see, US-relevant jokes... how international of me), introducing new and modern services is as bad an idea as you might think.

All I can say is that there is little reason for a user to be able to change their firewall settings on their router from the other side of the planet (read; over the internet) but that hasn't stopped manufacturers' efforts to 'solve problems people didn't know they had', thereby creating more problems that nobody whatsoever needs.

5.2.1 Hardware Specific Issues

There are also some specific issues relating to hardware devices themselves:

More Supply chain woes

Example: U-Boot doesn't support secure booting. At all. Secure boot mode has been in ARM processors from 2005/6 time, but the most common boot-

5.3. HARDWARE PENTEST TIMELINE - FROM SCOPING TO REPORTING51

loader in existence doesn't support it. This means that your bootloader probably doesn't support checking if the firmware is malicious. Again, tie this into a CSRF vuln (see above) and you can upload a malicious firmware file (I'll probably do a blog on how to make these).

But how is this a supply chain issue? Well, ARM is sold under license, and as such, many of the cheapest chips use licenses of older chip designs. So we see the precise same issue as above - Jazelle mode is an old ARM exec mode that supports native Java bytecode execution (*YES YOU READ THAT RIGHT*). It was 'replaced' by a new mode in ARMv7 called 'ThumbEE', but you'll still see Jazelle in millions of chips in the wild worldwide.

Development Headers

Kinda obvious from the outset, but the reason these are left on the boards is that the cost of tooling a 'consumer safe' board is very high. Prototyping can take weeks, and when you double that time to create a 'secure' version of a PCB layout, most project managers would shelve such an idea on the grounds of cost alone.

Also, 'omitting' a connecting resistor when going from dev board to production board *isn't fooling anybody...*

Bootlog Oracles

Bootlogs are great. They're on the list a little down from post-it notes and using your MAC address as the default Wi-Fi password⁷. They tell you processor versions, chip ID's (they are often covered up with heat-sinks you need light munitions to get off), kernel versions, precise software versions, and tell you all the modes that things are running in/at/on. We will do some detailed commentaries on some bootlogs to show how this works later on.

5.3 Hardware Pentest Timeline - from Scoping to Reporting

This is a proposed 'ideal' timeline that you should work towards in order to do the most effective penetration test and security assessment of an embedded device that you can. This timeline is the product of professional experience and talking to other experts/testers in the field. Although you may not be able to execute all these steps exactly, it should serve as an overview for the 'ideal' hardware/embedded security assessment.

⁷<https://twitter.com/LargeCardinal/status/682591420969029632>

5.3.1 Scoping

Scoping Time

Have an allotted amount of time within which you can scope the device in question. It can take up to a week to fully analyse and research a board, and at least one example device must be provided for scoping. As such, this is the first step and only takes a few minutes to determine how long you will need to scope out the testing of a device.

PCB Analysis and research

Once you have this time, there three tasks, the first of which is PCB Analysis. The aim here is to gather information about the devices and IC's in play, and how they work at a hardware level. Actions should include:

- Identifying all IC's on a PCB
- finding and storing the datasheets for all identified IC's (no matter how obscure!)
- making black-box diagrams for anything that isn't obvious
- Identify developer headers and assess what they could be
- Using logic to solve identity problems

An example of the last step is as follows; An SoC had a digital signal out trace sent to a 'black blob' epoxy IC (these are the ones where the IC is smothered in black epoxy, and no clue of what the IC underneath does is available), and a line came straight out of here to an 8-pin SOIC speaker driver. It was therefore reasonable to assume, given this black blob spoke to the 'signal in' pin that drives the speaker, that this IC was a Digital-to-Analogue Converter (DAC). This made sense later when it was determined that the SoC had a built in ADC (Analogue-to-Digital converter) for recording sound, but no inverse audio-output channel.

Recording all this information - IC identitites, datasheets, diagrams, etc. will be very useful later down the line.

5.3.2 Testing

Testing should ideally have the following provided by the client:

- At least 2 devices for testing
- Firmware files
- Source Code for custom services (if code-assisted pentesting is in scope)

- Relevant technical documentation as available
- Disclosure of any known issues, should the client be willing

Ideally, the client should also determine any trophies - things that they are concerned about specifically. This is especially important if testing time is limited.

PCB Analysis and Research

It is highly likely that the device you eventually test is completely different from the one that you scoped. This is due to the fact that satisfying the requirement that a pentest needs to be provisioned is often put ahead, in time, of the device being finally ready - the idea is that you will test the device as soon as it completes its basic testing.

As such, it is possible that the entire PCB analysis may need to be redone once the test devices have been acquired. Allow time for this to be the case.

Dumping Firmware

Dumping the firmware, or analysing prepared firmware, is an essential step in assessing embedded systems.

Debug/Terminal access

This is going to be especially useful, and time should be allotted to specifically trying to access this functionality. If there are any in place restrictions, these should be documented, and their effectiveness should be determined - if you got past them, the manufacturer should know how it was done.

5.3.3 Reporting

Found Secrets/Trophies

Common secrets between different devices should be highlighted as a matter of priority, as they are probably in the wild. If it was possible to determine a unique secret from artefacts on the device (such as the MAC address of an interface giving rise to the default Wi-fi password), this should be made clear. Usual sensitivity should be exercised when dealing with these secrets.

'Testing Narrative' of Hardware Hacks

With regard to giving visibility, if required, the client should be able to get ahold of a rough testing narrative of all the things that were attempted against the hardware. This can be particularly useful for their own risk assessment exercises, so should be encouraged.

Chapter 6

PCB Teardown and R.E.

PCBs are fundamental to our work - every piece of electronics utilises a PCB, and so understanding how they work, what they are made of, and more importantly, how to Reverse Engineer them is a very important skill to cultivate.

6.1 The Basics - PCB's in the Wild

We'll cover the basic parts of a PCB, how they are manufactured, and what their makeup is.

6.1.1 Anatomy of a PCB

PCBs come in all shapes, sizes, and colours, but what are they really made from? Well, life starts with a sheet of bonded fibreglass with a layer of copper coating on the top and bottom. The thickness of these fibreglass sheets can range from 0.2mm through to 2.0mm and beyond - though the extra thin and extra thick layers tend to have special purposes. Most PCB manufacturers will make you boards in the 0.6mm to 2.0mm thickness range, with the nominal thickness being pitched at 1.6mm.

Once you have your board, the PCB is ready to be produced from the schematic files you have created. If you are using a tool like Eagle or KiCad (the latter being FOSS and actually very useful, when you get used to the controls), you will follow a process along these lines:

1. Design the schematic of your PCB - define how the parts will connect together, etc.
2. Take this schematic and generate the netlist - this lists the parts you want to place, and you can select the appropriate footprints to be placed (most components have more than one option)

3. This data is then loaded into a board design where the devices are placed and joined up with traces - where traces need to pass through to different layers of the PCB, drilled conductive-filled holes called *vias* are placed to facilitate this.
4. Once all the components are located, an edge-cut is defined, and the silkscreen can be checked/added to the design. This is usually intended to have part numbers, indicators, fiduciary marks, etc. that are useful during board construction.

Once the PCB has been designed, then the software will generate files called ‘Gerber’ files, along with some drill files that essentially contain all the information needed to construct the board to your specifications.

In order to add the traces, resist is applied and UV exposure processes are used to cure and then etch away the waste copper (this is recovered for recycling). Once the traces have been drawn, large automatic drilling machines drill all of the holes for vias and sockets.

The board then passes through various plating processes which fill the holes, join things up, and make the board ready for use. At this point, a soldermask - the trademark green colour of a PCB is just this soldermask - is applied in order to prevent the solder from running amok around the board during assembly. Finally, the skilscreen is added, the boards are tested in large automatic machines called Flying Probe Testers¹ or FPTs.

PCBs can be manufactured as two-layer boards - this is the most common, with one layer on top and the other layer on the bottom of the board. However, you can have multi-layer boards - the iPhone 4 had a 10-layer PCB, and the Tesla CID unit has a 14-layer PCB! This is one way in which you can reduce the footprint of your PCB, by moving things from layer to layer. This is, however, considerably more expensive.

This, in a nutshell, is a whirlwind tour of the anatomy of a PCB.

6.2 Reverse Engineering a PCB

There are two fundamental ways of reverse engineering a PCB - destructive analysis, and non-destructive analysis. We’ll cover both here, though the best resource is KT Ng’s “The Art of PCB Reverse Engineering” - a *spectacular* book on the subject, and well worth tracking down a copy.

In addition, Joe Grand’s paper on PCB RE is also worth reading to inform this material:

http://www.grandideastudio.com/wp-content/uploads/pcb_deconstruction_techniques_paper.pdf

¹These are also useful for board reverse engineering, too!

6.2.1 Non-Destructive Analysis

We will cover here some techniques for locating what goes where without removing components from the board (or removing minimally-many components), and without damaging the board (or causing minimal/repairable damage).

IC Number Reading

Sometimes, owing to grease, dirt, or damage, the IC numbers are hard-to-impossible to make out and read. Here are a few quick fixes to get around this:

- **Chalk** - apply a small amount of chalk to the IC by rubbing it completely over the top. Then, with a fine brush (like a makeup brush) remove the excess and the part numbers will be revealed. The chalk fills in the indentations on the top of the board which are the identifying marks.
- **Scotch Tape** - this tape's finish lends itself to revealing the indentations very clearly when placed atop an IC. The matt finish make brighter the edges caused by the markings.
- **IPA/Acetone Clean** - sometimes all you need is some isopropyl alcohol or acetone on the tip of a cotton bud to remove the crap that is preventing you reading off the information you seek!

Steel Wool Binary Search

Working out which trace goes where can be very time consuming. However, there is a simple trick that will let you do this very quickly - HT to Robert Leale Jnr.² who taught me this trick:

1. Start with the pin or test point or trace you are interested in, and set your multimeter to ‘continuity’ (or ‘diode’ mode) - ideally with an audio alert turned on. Place the positive lead at this point (the board should not be powered for this exercise)
2. On the other test probe, hold some steel wool against the pin, and ‘dab around’ on the PCB until you hear a «BEEP» - the rough area where you just dabbed is where the trace goes to.
3. halve the amount of steel wool you have, and repeat step 2 until you have it down to a small enough area to do a manual search.

²<https://twitter.com/carfucar/>

This technique is simple, effective, and reduces the time-to-find-a-trace down from potentially tens of minutes (with by-eye examination and probing) to minutes, even on large PCB targets, as you reduce the problem to a spacial binary search, once you have the general area.

Flatbed Scanner Imaging

Taking images of a PCB is a great way to work out what is going on by tracing over lines in photoshop or GIMP.

However, getting precise and reliable placement and angles using a standard phone camera or even a DSLR on a tripod mount can be challenging - you find yourself with planar distortions of the PCB under analysis.

In order to work around this, Andrew Tierney³ has made great use of a flatbed scanner. This gives you a perspective-neutral image of the PCB, with which you can easily load the images into photoshop/GIMP, set the transparency to 50% and you can see easily where traces go when they swap to different sides of the board.

It is advised to remove any large or high components, like power decoupling capacitors, first in order to maintain the focus and details on the resulting scan. However, this method is incredibly reliable and very powerful for getting high-quality, useable images of PCBs for reverse engineering.

6.2.2 Destructive Analysis

This is the fun stuff!

³<https://twitter.com/cybergibbons/>

Chapter 7

Terminal Access through UART

Now that we've covered the motivation for *why* we should include low-level hardware hacking in our pentesting, and we've covered off the specific details elsewhere, it is now time to do a full run-through of what to do when you want to connect to a piece of hardware.

We'll start with the most ubiquitous and easy example - UART.

7.1 What we will need

Here's a list of the stuff you'll require to do everything in this chapter:

- Multimeter
- Soldering Iron (with solder and flux as required)
- Header Pins - usually set in plastic, 4mm apart
- Plastic Coated wire - esp. if you don't have the pins, cut to lengths with tinned ends
- Logic Analyser - connecting over USB, compatible with Saleae Logic
- USB to UART adapter - a CP2102-like adapter is what I use, or a bus-pirate
- A target board.
- Datasheets for the SoC/IC you want to connect to (may not be available)

7.2 Identifying UART

7.2.1 Identify and Connect to Development Headers

So, when they manufacture a PCB, they don't necessarily remove the development headers from the production models. What are these? Well, they are the parts of the PCB that would carry interfaces designed for use with test equipment the designers use in development of the firmware.

Remember, the PCB design comes in two stages - the first decides the layout of the parts, and correlates this to the layout requirements of a PCB that will fit inside the proposed product. The second stage is where the firmware is written and developed so that the board works properly, and performs all required tasks.

Obviously, the first stage is costly - it costs a lot of money to prepare a PCB design, then have the prototypes made and populated, to then have access to a board you can develop on. Often, removing something like a debug header is not an option, owing a requirement to re-test the board. So to keep costs low, many prototype boards make it into production, minus the costly plastic headers that actually go on the board.

As such, this is what we want to find, so that we can do a little debugging of our own.

Identifying a Development Header

[Insert pictures]

Given we are looking for UART, we would expect there to be four pins, with two fine traces coming from two of them. As such, the ones with thick traces/connected to ground are likely VCC/GND, and the two with fine traces are signal lines, indicating they are likely Rx/Tx.

You can also look to the datasheets - usually available online - for the chip you want to connect to. Search the document for 'serial' or 'debug' or 'uart' and you'll usually find some listing for the UART connection. Get the name of the pins, e.g. P_2_6 for Tx, say, and then search for the 'pinout' and look for that pin. There is usually a large dot, notched corners, or some other fiduciary mark on the chip to tell you which way round the diagram matches to the pins. You should then be able to follow the PCB traces and see if the development/debug access leads to a development head or not.

If you don't see what looks like a gap for a development head, then it is very possible that your SoC connects to UART through a 'bed of nails', and so you will need to examine the test pads. A bed of nails test rig is one where little springy brass collared pins called 'pogopins' are arranged in a fixed way to allow the developers to connect to a fixed set of test pads. As such, if your UART is sent there, then you need to solder onto a test pad, as described in the soldering skillz chapter.

Once you've identified the development header you want, solder a wire into each hole in the PCB, or onto each test pad, and get ready to analyse the signals with the logic analyser.

Analysing the Header Signals

We will need to use a multimeter to get a sense of which pin is what. For the sake of argument, we will assume that we have four pins, which we will write $[a | b | c | d]$. Find some ground point (either labelled, or use the one connected to the main ground off the power supply unit (PSU)).

1. Use the multimeter continuity setting (or the diode setting) to check which of our candidate pins is GND. You should see full conductivity to the ground point, as if there is nothing in the way (should be a close reading to touching the probes together). Label this on some diagram, for example if it's the last pin, label as $[a | b | c | -]$
2. Next identify the VCC line - this should be the pin with the highest voltage whilst the device is running. Chips run at 3.3V, but generally you'll see 3.2V or thereabouts. As such, two have a lower voltage than the third, that is a stable 3.3V from the PSU assembly. Label this on your diagram, e.g. if it is the port far left, we update it to $[+ | b | c | -]$.
3. We can now make the safe assumption (which we will test) that the remaining two pins are Tx and Rx. To work out which is which, you may have noticed some variance on the voltage of one of the pins when finding VCC. If you did, then that's probably a Tx line.¹

We are now in a position to use the logic analyser with the Logic software to The general methodology for analysing the signals from a device is as follows.

Connect up the analyser - Connect GND to the relevant pin on the analyser. You can connect VCC to a channel if you aren't certain it's the power line. This will now let you use the logic analyser to work out what is what.

Connect the pins to analyse - We have two, so we'll use channel 1 and 2 in our example. Now configure the sampler as in Figure ???. This will work for most UART implementations in the wild.

¹Multimeters use Root-Mean Square (RMS) chips to 'even out' the peaks and troughs of a signal to give a true level reading. As such, it would show lower voltages, such as 2.7V or briefly 0V, as it tries to interpret an digital signal.

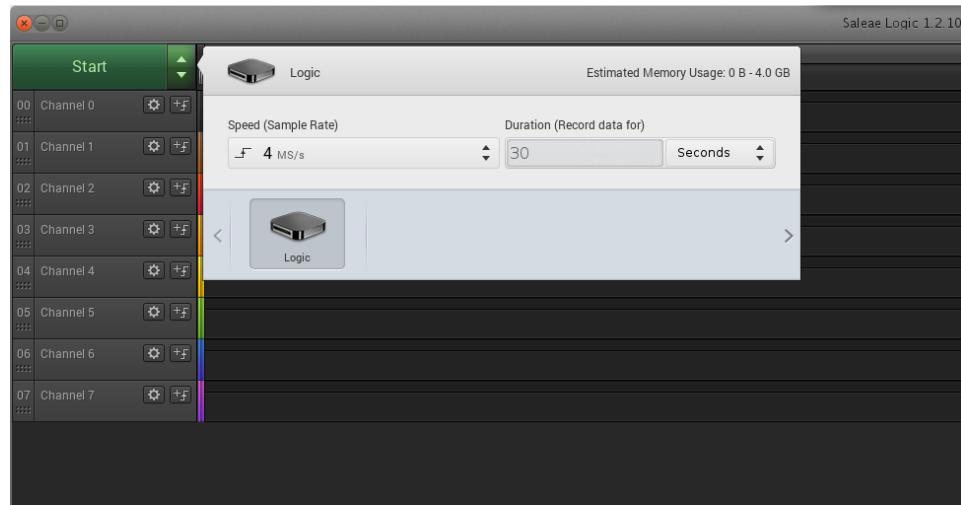


Figure 7.1: Setting the Sample Rate in Saleae Logic

Once this is setup, start the capture, and then immediately start the target device. Once it is finished you will see a captured digital signal, as in ???. This will then need to be run through the analyser to get out the baud rate.

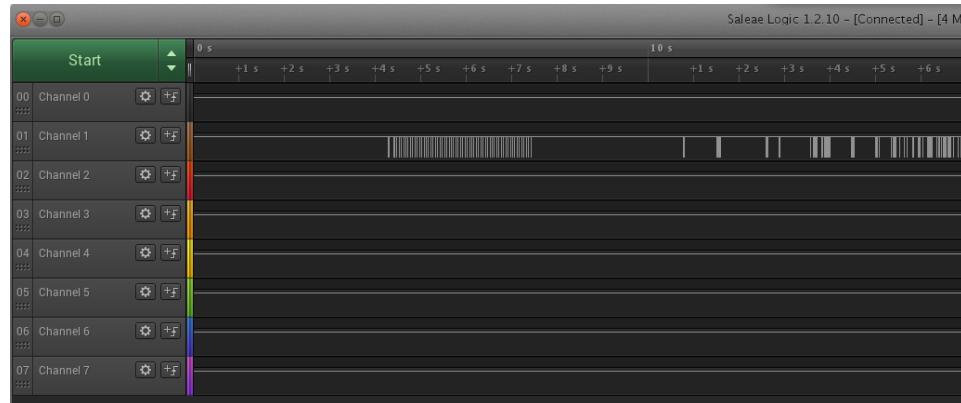


Figure 7.2: Captured Digital Signal

ProTip: If you get no signal from either, or signal from one but you can't connect, double check the traces. They may have omitted a resistor from the production model to stop people connecting to the debug interface (or just to save costs).

Double check the traces leading to/from the SoC to your, and look for gaps. If you find one, use a dab of solder to carefully short the gap, allowing you to connect.

Run the Analyzer Now we are ready to analyse the signal we have. On the right of the main screen in Logic you will see a small box for 'analyzers'. Click the cog icon, and you will see the analyser setting screen, as per Figure ???. This will then allow you to load an analysis of the signal for different baud rates and settings. Nominally, all UART systems in the wild that you will encounter will tend to use 8N1 as the data settings (see the UART protocol description for more details about that). As such, our only task is to work out what the baud rate is.

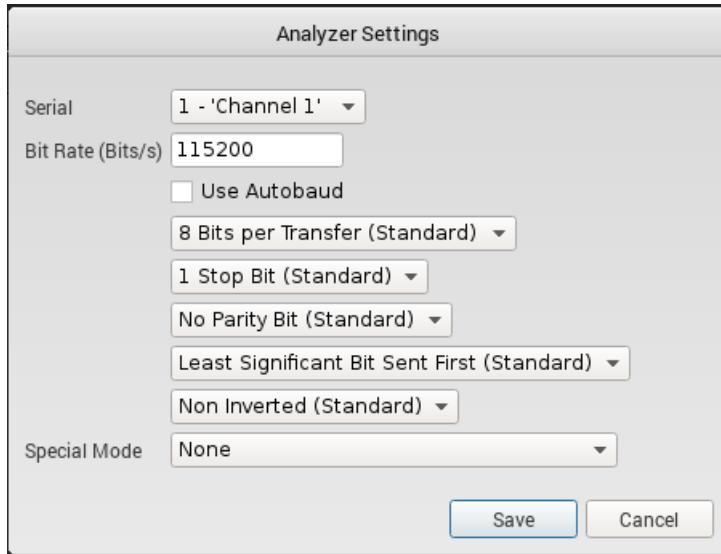


Figure 7.3: Analyser Settings - note that the baud rates are free-form text.

The most common baud rates, as mentioned in a previous chapter, are as follows:

- 115200
- 9600
- 57600
- 19200
- 38400

Once the successful one is found, you should see some meaningful output mapped above the signal when you zoom in on Logic. A successful analysis with the correct baud rate will yield an output that looks sensible, as found in Figure ???

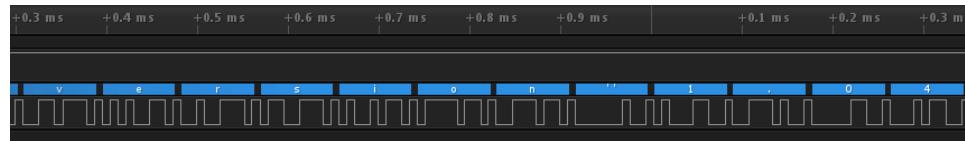


Figure 7.4: Successfully analyzed signal showing version 1.04 being read from the data.

We have now successfully found our pinout for the development header, as well as our connection settings.

7.3 Connecting to UART

NB - For the sake of this section, we will assume that the pinout was $+[\text{Tx}|\text{Rx}-]$, and that the baud rate was 115200 (the most common).

7.3.1 Testing Tx

We are now ready to connect for the first time to our device. We now want to wire up our USB to UART adaptor - for ease, I will assume it's a CP2102 (the most common IC used on these boards) and will refer to it by the chip name. **NB** - ensure that the device is off!

Connect GND to the GND pin on the CP2102, and take the Tx from the target pinout, and connect that to the Rx on the CP2102. Now run the following command (as root):

```
|| screen /dev/ttyUSB0 115200
```

Here, `/dev/ttyUSB0` is the CP2102 device² and 115200 is the baud rate we identified in the previous section.

Now power on the device and you should see text output in the screen terminal. It should look like a boot sequence log - that's because it is! If you're seeing this, then congratulations, you've successfully connected over UART!

7.3.2 Connecting fully with Tx/Rx

We are now ready to try and get a shell. Power off the target device, and disconnect the CP2102. Connect the remaining Rx pin on the pinout to the Tx on the CP2102. Reconnect the CP2102, start screen (the command will be the same) and power on the device. If you see some generic Please press a key to interrupt boot process... try it. You should

²If you're not sure, you can check by reading the `dmesg` log after connecting the CP2102.

either get dumped into a bootloader shell, and can access memory, or if you allow it to continue (power cycle the device if needed) you may get dumped into a root shell.

This is the debug shell that the developers used when building the firmware, so now you're ready to start testing.

7.4 **UART Situational Awareness**

There is a wealth of information you can glean from the boot sequence dump (called a ‘bootlog’). Items worthy of note include:

- Bootloader name and version (Common ones are U-Boot and HiZ)
- CPU/MCU architecture (ARM/MIPS usually)
- Linux Kernel version
- busybox version (discussed below)
- Other software names and versions - check against known CVE’s
- Look for locations of important scripts (like `/etc/rc.d/rc0.sh` or the like)
- Device mappings should also be checked (like whether the UART console is mapped to `/dev/ttys0`, `/dev/ttys1`, `/dev/console`, etc.)

If you are a little lucky (not all devices do this, but most do) then at the end of this you will see a root shell. You may have to login, but that login is usually the same as the default for the management web application, as a rule of thumb. We will discuss how to use this shell and what to do with it once we have given a methodology for dumping flash ROM data using SPI over a bus pirate.

Chapter 8

Fimware Dumping through SPI

SPI (Serial Programming Interface) is as ubiquitous as UART when it comes to prevalence in the IoT/embedded hardware world. When used as a communications protocol, its primary use has been to interface an SoC with a Flash ROM IC. Although it can be used in other areas - many LCD screen breakout boards use SPI as their communications protocol with Raspberry Pis, for example - this is the most common form you will find SPI in use on production PCB's.

8.1 What we will need

Here is a list of the stuff you will need to do everything in this chapter:

- Multimeter
- SOIC Clip or IC Leg Taps
- Bus Pirate - or other SPI compatible protocol board¹
- `flashrom` installed - or the software that works with another chosen SPI protocol board
- `binwalk` installed
- Optional items:
 - Logic Analyser with Saleae Logic software

¹A ‘protocol board’ is just a board that speaks to your computer and lets you communicate with a target over the defined protocol.

8.2 Identifying SPI

8.2.1 Standard Analysis

Analysis of the PCB should identify one chip that is an 8- or 16-pin IC that has a part number corresponding to a datasheet that tells you this chip is a Flash ROM Storage device.

It must be noted that not all 8-pin SOIC IC's are ROM chips. Many power converters, speaker driver/op amp IC's, small MCU's (Microcontrollers) all come in these packages. You **must** confirm that the chip is indeed a ROM chip and supports SPI before continuing to connect it to a Bus Pirate. If you don't, bad things will happen.

ProTip: Flash ROM SOIC's tend to be a little 'thicker' and more square than the other types of chip using an SOIC package. This is due to the silicon die needing a large square NAND flash storage area. As such, the IC's themselves tend to look a little different from the other 8-pin SOIC's in particular.

Once you have identified the pinout from the datasheet, you're ready to connect to the device and dump the contents!

8.2.2 What to do when Standard Analysis fails

First, with the device powered on, identify the VCC pin with a multimeter - this can be done by putting the negative (black) probe on GND (say the collar of the power supply connector that goes to ground) and trying each of the pins in turn. As a rule of thumb, it is generally the pin opposite the pin next to the little circle (the fiducial marker) - if the circle is 'top left' then this pin is 'top right'.

Now disconnect the unit and use a multimeter in continuity mode (or diode test mode) to identify the GND pin. This has to be continuous from all GND terminals to main ground, so using a known GND point (like we did earlier) is best. As a rule of thumb, this is generally the pin at the opposite end of the chip, on the side with the circle (fiducial marker). If the marker is 'top left' the this pin is 'bottom left'.

You now have a decent indicator for the layout of the IC. This should help you work out if the pinout for this chip is similar to others - it is best to see if you can get the datasheets for chips in the same series, if not the exact one. Searching online for these datasheets may be as fruitful. Atmel ROM IC's, by and large, do not follow the diagram below.

Following from my 'rule of thumb', here is a standard pinout that works on most flash ROM IC's I have encountered personally:

Bus Pirate		Flash chip				Bus Pirate	
CS	1	CS		VCC	8	3.3V	
MISO	2	DO (data out)		HOLD	7	3.3V**	
3.3V**	3	WP (write protect)		CLK	6	CLK	
GND	4	GND		DI (Data In)	5	MOSI	

Figure 8.1: Generic ROM IC Pinout - the numbers are just IC pin numbers

**See the IC datasheet for whether these should be high or low in order to read the contents.

8.3 Dumping firmware with SPI

Once we have our proposed pinout, we are ready to wire up the chip!

Ensure that the device is powered off! Else, you will have the SoC to compete with and all manner of mayhem will be set loose.

Following the wiring that you have identified, as the SPI pinout descriptions line up with the markings on the bus pirate.

- MOSI goes to MOSI
- MISO goes to MISO
- SS/CS goes to CS (chip select) on the bus pirate
- CLK goes to CLK.
- Now wire up VCC for power and GND.

ProTip: Flash ROM's can operate at 3.3V or 5V. Check the datasheet or use a multimeter to determine which is the correct voltage, else you either won't get a signal, or worse, could damage the chip by over-volting.

Now we can attempt to dump our firmware. Using a bash prompt, run the following as root:

```
||>$ flashrom -p buspirate_spi:dev=/dev/ttyUSB0,spispeed=1M -o ./firmware-dump.bin
```

NB - if at first it doesn't work, don't worry. Try it again. If it still doesn't work, adjust the spispeed parameter. The supported speeds for bus pirate in the flashrom source are:

- | | |
|--------|--------|
| • 30k | • 2M |
| • 125k | • 2.6M |
| • 250k | • 4M |
| • 1M | • 8M |

So, if you are having problems with a particular ROM chip, start by lowering the `spispeed`. Upping it to 2M is supported by most ROMs (and may be a hard limit on the bus pirate itself), and will dump the storage faster.

If it is still not possible to dump the ROM file out, then there are a few possibilities:

1. You may have made an error in your wiring
2. You may have followed the wrong datasheet
3. You may have a ROM chip that has specific requirements - you may need to adjust the buspirate firmware settings
4. Worst case, you may have killed the chip by accident during handling
- use a logic analyser when the device is running (if it runs)

If all has gone well, you now have a file `firmware-dump.bin` that you are ready to analyse.

8.3.1 Firmware Dumping through UART and `tftp`

If all else fails, then it is sometimes possible to dump firmware through the UART connection. This is documented here only for the U-Boot bootloader, although similar things are possible on others.

The basic methodology is as follows:

1. Start the interface for UART (screen or minicom) in logging mode - this is important.
2. Load the bootloader into debug mode (usually by pressing any key during boot).
3. Determine (through the `help` menus) whether there is a way to display memory.
4. Ensure that the firmware image is loaded in memory.
5. Use the memory display functionality to dump the contents of memory to screen.
6. As we're logging, this will now be logged to your local drive.

Worked Example

Here's a walkthrough for U-Boot, using screen, with the decoding at the end.

1. As root, run `screen -L screenlog.log /dev/ttyUSB0 115200`
2. press any key to interrupt booting
3. either using the bootlog or `sf probe`, get the addresses and offsets for your flash image
4. Load flash partitions into memory using `fload`, or if this fails, try `sf`
 - for example `sf probe 0; sf read 0x82000000 0x40000 0x370000`
5. Dump memory to screen using `md` (memory display)
 - e.g. `md 0x82000000 0x40000`
6. Copy the escaped hex into a separate file.
7. Hex that is in a hexdump format can be decoded into a binary file using `xxd -r`
 - e.g. `xxd -r text-dump.txt > firmware-dump.bin`

Further options

If this fails, see if your u-boot supports tftp, and follow these steps:

1. As above, use `sf` and/or `fload` to ensure that the ROM images are in memory.
 - To be clear, I mean; `sf probe 0; sf read 0x82000000 0x40000 0x370000`
2. using the `tftp` command, load the firmware image you want over the wire
 - Start a local tftp server
 - start a local DHCP server if you can, else, set static IP's that are compatible with the subnet used by the device
 - Attempt to load the file over tftp; `tftp 0x82000000 romfs.cramfs 0x370000`

Chapter 9

Firmware Analysis

We will assume that the dumped firmware file is called `firmware-dump.bin`. This can either be an in-memory dump off a device, but can also be from a firmware file provided by the client or one that you found online - generally the ones online are just replacements for the ones that are held locally; the binary file gets written as-is to the ROM chip.

9.1 Extracting Firmware

First, run binwalk on the file to see what the likely interesting offsets are. This can be done via

```
|| binwalk -e firmware-dump.bin
```

This will create a subdirectory `./_firmware-dump.bin` in which you can find any files that binwalk recognized and dumped or the contents of compressed directories and filesystems that binwalk was able to recognise and extract.

If this produces nothing, then run `binwalk ./firmware-dump.bin` to get an overview of the technical details and offsets. You can then use dd like so:

```
|| dd if=./firmware-dump.bin skip=<offset address> count=<number  
      of bytes to dump>
```

Remember, the count parameter takes decimal input, not hex, so make sure to use these. Thankfully, binwalk outputs this as well, which is handy.

ProTip: If you are not sure what to make of a file and binwalk is not making sense of it, look at the bootlog. This will usually contain the filesystem format, start addresses and block sizes, as well as compression schemes involved.

9.2 Trophies and Loot

Here are the things you should keep an eye out when parsing through extracted files from a firmware image. These artefacts are particularly useful when testing, as they may give you some clue as to how to overcome a particular login/parameter, and are usually universal for all devices of the same model number or series.

General tips when hunting for such treasure:

- Keep a note of *what* was found, so you can check if it changed later
- Keep a note of *where* it was found for use in later testing

Configuration Files

Many firmwares, especially for RomPager based firmware files, there is a file `/etc/romfile.cfg`, that can also be found at `/var/romfile.cfg` in some instances. This file contains all the details required to setup the system on the fly. When the system is booted this file is parsed out, and its contents are used to populate various BASH script `.sh` files, as well as populate configuration files for `wpa_supplicant` settings, DHCP ranges, etc.

As such, this is a good example for understanding the kind of thinking that takes place when firmware files are put together. There are near-infinite ways of achieving the same thing with the same setup, so it would be remiss to say “This is how it’s done!”, when it clearly may not be the case. Instead, here are some tips for

Passwords/Secrets

The following command, run in the directory where you store the extracted files, will find all sorts:

```
|| grep -RiE --color 'pass|usern|config|private|key'
```

Pulling `/etc/passwd` is a good idea - cracking the passwords¹ therein may reveal default user accounts that can access SSH/Telnet/management ports/management web application/etc. It also may give you a clue as to how authentication is done on the web applications used to manage the app - very often, the `/etc/passwd` file, or some other localised `passwd` file, is used for all authentication systemwide, or for multiple subsystems.

Likewise, parsing for common developer variable names such as `secret_key` or `api_key` may reveal uniform secrets that are used across all these devices.

The same can be done for DECT phone secrets on CPE’s that can talk to wireless phones (which would enable call snooping), and for media access/Digital copyright keys on devices that access media, like IPTV’s. Also,

¹Generally passwords are hashed with one of the earlier `crypt` implementations - usually `3` or `3`. No joke.

keep track of what private keys for certificates are kept, and what they look like they are used for. For example, some systems require authenticated access to some private subnet for hosted services, with access being restricted to any device that has a particular certificate - if you now have said certificate, you have open access to these ‘restricted’ systems.

NB - Recording the locations for these will be useful in later testing, too, so keep note!

BASH Scripts

When an embedded device is loaded, chances are everything is managed/updated/configured by BASH scripts. As such, knowing where scripts are, and where they get their parameters from is really important. This is how many RCE’s are found in embedded systems.

Binary files

Knowing the location of important binaries is very useful in later testing. Keep track of what binary files are where. Also, if a binary file does something specific but is not named properly, note stuff like that down. You *will* forget, and *will* get frustrated trying to remember later!

We will cover the disassembly of ARM/MIPS/XTensa/etc. binaries later on.²

9.3 Correlating with UART

Having access to a UART interface to some SoC will greatly increase the chances of you making full sense of the way a system works, and increase thereby the likelihood you will make significant advances in determining the vulnerabilities on the attack surface of the device.

When checking the actual locations of files, it can be handy to use the UART interface to make sure that files are indeed where they are in the initial firmware. Often, files are copied and then manipulated - in our example above, the `/var/romfile.cfg` file is soft symlinked from `/etc/romfile.cfg`. This can take hours to decipher looking at the firmware dump itself, but is abundantly clear when you have access to the UART terminal.

The UART will also tell you what devices are where - mappings such as `/dev/ttys0` is a common mapping for the serial port. Using `echo` to send some data to that location should reflect back to you. Once confirmed this can be used when finding RCE vulnerabilities (covered in the next chapter).

²Fun Fact: if you see lots of hex opcodes in a file that begin `0xDx` or `0Ex` that are close to each other (or line up in a hex editor into columns), then chances are it’s ARM in thumb-mode, as these are the opcodes for the branching instructions.

9.4 Reverse Engineering Firmware and Binaries

Firmware generally splits into three categories; those running just straight firmware (single purpose code, without any OS or filesystem to speak of), those running as some sort of RTOS (e.g. FreeRTOS), or those running some known kernel and filesystem (sometimes running in an RTOS-y manner - as found with Linux on many home routers/CPEs).

In this section we'll cover basic reverse engineering of both - in IDA and Radare2/Cutter where appropriate.

9.4.1 Reverse Engineering ARM/MIPS Binaries

These are actually quite straightforward to reverse engineer. Both ARM and MIPS instruction sets are well supported by both IDA, R2, and other disassemblers. As such, we will leave this to the chapter (to be included) which will deal with binary reverse engineering for embedded systems more specifically.

9.4.2 Reverse Engineering a Firmware ‘Blob’

Sometimes, you just get a load of firmware - no filesystem, just code that runs and that is it. This is by far the most common thing to find in IoT/embedded devices. But how to reverse engineer it? Here are some pointers for the curious:

- Find the datasheet for the microcontroller you want to target and get the memory address values - these tell you where, for example, the first UART port or the third I²C controllers are placed in memory for the device.³
- Once you have attained this information, seek out functions and name them appropriately - here are a few possible examples to get you thinking along the right lines:
 - If a function reads from a UART Rx buffer and puts this data into another memory location, then this is a `UART_Read()` kind of function.
 - If a function writes a byte from a register to a UART Tx buffer, then this is a `UART_Send()` type of function.
 - If a function reads many blocks from an SPI bus, then this is possibly a `firmware_copy()` kind of function

³For an example IDA memory mapping config file for AVR microcontrollers such as ATMega328, see: <https://gist.github.com/extremecoders-re/8d3e9b846a6ec883e5ae3b2bccf5cc88>

- If this is called at the start and then has multiple calls to the ARM `CRC32W` opcode, then this is likely where a firmware check is being done...

A very nice bug that was found in the NXP LPC1348 bootloader (and will be discussed later for how it can be used to bypass CRP) by Chris Gerlinsky⁴ - here, it would only retrieve the security bits one time, and had fixed values for firmware protection modes. This lead to his using this to enable a glitch at the moment this read is undertaken - thereby disabling CRP protections. But more on this later...

Once you have this mapping done, and have worked out a few of the main functions, you are ready to commence reverse engineering in the ‘normal way’ - as you would if you had a file with an ELF header that you were more familiar with. Again, these specifics will be dealt with more fully elsewhere.

⁴<https://twitter.com/akacastor>

Chapter 10

Radio Hacking

10.1 Wireless Communications - the Backbone of IoT

One of the things that makes IoT devices (and indeed, any embedded system) so useful is the capability of acquiring some sort of data and then relaying that to some central assessment destination for analysis, or taking some sort of command and then actuating that into a real action across multiple devices in multiple locations/spaces.

Although we can achieve much with wires, the advent of wireless communications has facilitated a multitude of easy-to-include, ready-made wireless solutions for sharing data between devices. As such, any consideration of IoT security must necessarily include a treatment of radio hacking, from basic WiFi vulnerabilities we all know and love, through to BLE ubiquity and hacking, and rounding off on the more IoT/embedded specific protocols; LoRaWAN, SigFox, ZigBee, etc.

10.1.1 Radio Basics

What is a ‘Radio’

Wired communication protocols (such as the serial ones like UART, SPI, etc. we saw earlier) have their range determined by properties of the wire, the power used to send the signal down said wire, and external parameters such as noise, interference, temperature, etc. Taking all these together defines how well your ethernet will work, whether you can run that serial cable from the basement to the top floor terminals (back in the day...), how fast you can communicate down that wire, etc.

Similarly, wireless protocols are governed by similar parameters, except we don’t have actual wires. Instead, we have a suite of frequencies - from ~10MHz through to 5GHz and beyond. As society, we have split this up into bands, and some of these are free to use for anyone - such as 433MHz,

868MHz, 2.4GHz bands - and some are very strictly controlled - such as 800, 1800, and 2600MHz bands for use by 4G/LTE communications, for example.

The propagation of a wireless signal is governed by a few easily identifiable factors - these are factored into equations for what is called ‘Path Loss’:

- Antenna length/height
- ‘line of sight’ between transmitter and receiver - how clear is the path between Tx and Rx
- The distance between a transmitter and some receiver

Signal strength is measured in dBi, or Decibel Isotropic - and yes, this is a cousin of the decibel system used for measuring sound signals in audio electronics, usually dBu or dBv in that case. dBi is the gain measurement of an antenna compared to the ‘isotropic antenna’ - although this sounds like something you got threatened with in pre-calculus, the isotropic antenna is the ‘idealised most perfect antenna’ - given such a beast is physically impossible, we measure antennas on how close they come to imitating the gain of this perfect antenna.

Worth mentioning common terms you’ve seen before, e.g. in UART, come from radio - such as using TX for ‘transmit’ and RX for ‘receive’.

Software Defined/Configurable Radios

In order to access this seemingly magical world of EM waves traversing the air, you’re going to need a few tools. These tools take the form of Software Defined Radios (SDRs), or sometimes the related (but distinct) Software Configurable Radios (SCRs). In short, there are a few classes of tools available, ordered from cheapest to most expensive:

- **VGA to USB SDR** - costing as little as €5, these little VGA chips don’t have any cache, meaning that the chip can be pulsed raw from the USB socket by sending carefully crafted packets to the IC. This rudimentary transmit-only radio can achieve TX of multiple frequencies, though its reliability is clearly better at the lower end - sub 1GHz.
- **RTL-SDR radio** - usually packaged as ‘TV Tuner SDRs’, these RX-only devices are very cheap, usually around €10, and have a reasonable range and reliability. Can be used to read from around 100MHz through to 2GHz reasonably reliably (meaning they can be used as passive IMSI catchers, or passive 1090 transceiver aviation radar, for example).
- **Sub-1GHz SCR** - based around the TI CC1111 chipset (found in the [in]famous IM-Me devices, and many others) the Yard Stick ONE

comes in at around €110 and are very effective. Being an SCR means you cannot configure it to use any modulation, but rather the ones it has pre-programmed - but it has enough for serious usage in testing the radios of common household devices.

- **Desktop USB SDRs** - A ‘proper’ SDR can be picked up for as little as €80 in sales, up to a few hundred Euros for a HackRF ONE or BladeRF device, all the way up to €2000+ for an Ettus Research SDR. These classes of SDR have full configurability, and ranges from low MHz all the way up to 5+GHz - serious amounts of usage for serious amounts of money.

It is worth noting that you should acquire the SDR that is most appropriate for what you want to do. A good example of this is the HackRF vs. BladeRF - on the face of it, the HackRF has a better frequency range and features, but it is half-duplex - not so useful if you want to MITM a protocol such as GSM/LTE - to do so, you would need 2 HackRFs, which is a bit pricey. However, the BladeRF, though not as wide a frequency range, is full duplex - it can simultaneously transmit and receive. As such, It is always best to choose the right SDR tooling for the job at hand, so survey the market and assess your exact requirements early on in the assessment preparation phase.

ProTip: When it comes to any non-ISM^a band, there are legal restrictions. For 3G/4G/LTE, these frequency bands are heavily controlled in every country on the planet. Interfering in these bands is particularly illegal - mainly due to the fact that you incapacitate people from making emergency calls easily. In the UK, this is governed by the Wireless Telegraphy Act 2006, for reference. **Proceed with caution.**

^aThe ISM bands are the ones that we mentioned earlier - where anyone can transmit or receive without a license.

10.1.2 Radio Terminology

General Terms

Just so we’re clear, here are some common terms used in radio communications, and some brief explanations about what they mean:

- **Frequency** - literally, the number of cycles (specifically, the full up-and-down of a waveform) per second, measured in Hz - 1Hz = 1 cycle (one up-and-down) per second.
- **Modulation** - the way in which we communicate data over a given frequency/frequencies

- **Channels** - For any given frequency *band*¹ there are usually defined several frequency ‘channels’ which are distinct, if sometimes overlapping, sub-bands identified to prevent interference between devices
 - The ISM 2.4GHz band - from 2400MHz through to 2490MHz - can be divided up into bands in many ways.
 - WiFi divides the 2.4GHz band into 14 channels; channels overlap with each other, but are designed that channels far apart will never interfere - as such, channels 1, 6, and 11 are most commonly used as they have no overlap with each other, so devices can talk freely within the ISM band.
 - Bluetooth and BLE split the 2.4GHz ISM band into 40 distinct channels, and uses lower-capacity modulations to communicate over such narrow bands. They also channel hop - Bluetooth negotiates this, but BLE has it predefined - between the channels in order to minimise interference (both from and to the bluetooth/BLE communications).

Modulation Methods - an Overview

There are broadly 4 types of ‘keying’ - ways of representing ‘1’ and ‘0’ in radio communications. We’ll summarise them here: [TODO - add diagrams for each...]

On-Off Keying (OOK)

This is probably the simplest - you represent a ‘0’ with silence, and a pulse at a given frequency indicates a ‘1’. Clearly, this has to be clocked/sampled independently of the signal, but it is the simplest digital radio transmission to enact - it’s the digital version of analogue morse-code.

Very simple, very low power, but also, not incredibly reliable, and prone to being messed up by reflections (off a building or other generic Large Thing near the receiver). This is the kind of modulation used by garage doors, doorbells, etc.

Amplitude Shift Keying (ASK)

Amplitude Shift Keying, or ASK, is the digital version of Amplitude Modulation in analogue transmission - yes, the same AM that you get transmitting on radios (normally received on a ‘crystal’-based (read: diode based) radio set ‘back in the day’...). Technically, OOK (above) is a slightly extreme version of ASK - where the gain is pulled down to zero to represent ‘0’.

As the name would suggest, the amplitude shifts up for a ‘1’ and down for a ‘0’ - this means, quite literally, that the 1’s are louder than the 0’s, and that’s how you can decode the signal. These are coded by adding gain amplifiers to the carrier wave when a 1 is present on the data bus.

¹Not as in rock bands - more as in ribbon bands...

Frequency Shift Keying (FSK)

This is the digital version of Frequency Modulation in analogue radio transmission - and yes, this is the FM that you listen to radio stations over, normally.

This time, there are two carrier frequencies - the lower frequency usually representing '0', and the higher frequency representing '1'. These are effectively switched by the radio given an input of 1 or 0, and as such, it is very efficient and easy to implement as well - however care has to be made to ensure the phase of the two oscillations is maintained to prevent 'broken'/discontinuous waveforms.

There are 2 variants you will find - 2FSK and 4FSK. The one I just described is 2FSK. 4FSK is the same, but with 4 frequencies, one each for '00', '01', '10', and '11'. There are also several ways of generating the waveform - one is called Gaussian FSK or just GFSK, where the waveform is passed through a Gaussian filter to 'clean it up a bit' and smooth it out before transmission. GFSK is the primary modulation used in DECT phones, Bluetooth, BLE, LPRF, Z-Wave, and more.

Phase Shift Keying

Phase shift keying, or PSK, may sound like how door locks work on a sci-fi teenager's bedroom, but in actuality it's another kind of data keying for digital transmission.

First, a note about phase - this is the measurable difference between the peaks of two waveforms of the same frequency. If you have two waveforms 180° out of phase, then you get no sound at all from mixing them - waveforms that are only slightly out of phase give all sorts of weird and wonderful tonalities and timbres to instruments; this is a problem (or a feature) when recording electric guitars sometimes.²

So, with BPSK - Binary PSK - we have that two phases of waveform represent '1' and '0', and these waveforms are shifted around based on the data input values. As such, data can be transmitted at one frequency and at one amplitude, yet data can still be recovered from the (discontinuous) waveform.

The big-brother of BPSK is QPSK - Quadrature PSK. This has a 4-node complex-plane constallation mapping that defines data transmission, not too dissimmilar from the 4FSK above.³

PSK-based modulations are very reliable, but harder to implement in hardware. As such, they tend to be used in applications where higher data rates are required. That said, it is a strict 'subset' modulation to QAM.

Quadrature Amplitude Modulation

²That said, the iconic sound of the Greeny/Moore Gibson Les Paul was entirely due to this property - they'd put the pickups on the guitar out of phase to get the sound effect they popularised.

³If that sentence made no sense, then don't worry - either google more, or forget you read it :P

QAM is the hybrid between QPSK and ASK. But how does this work? Well, first we need to do understand a little about complex numbers.

Complex numbers are just numbers of the form $a + ib$ where $i = \sqrt{-1}$ - this last bit is the ‘imaginary’ numbers bit, but really it is just a number space that doesn’t ‘make sense’ in the normal number line (as no number times itself gives -1) - but the solutions to equations often find themselves branching out into complex number spaces. As such, you can make an x, y graph, just like in school, only with the ‘real’ numbers going left-to-right, and the ‘imaginary’ numbers going up-and-down.

To spare you a *lot* of maths, I’m just going to tell you that ‘spinning stuff’ on a complex plane is the same thing as shifting waveforms forwards and backwards - that is, rotating around the unit circle on the complex plane is akin to manipulating the phase of a wave.

Last thing to realise is that the amplitude of our wave corresponds to a ‘bigger/wider’ circle on the complex plane.

That’s it - maths over. That’s all you need to know.⁴

Now that you know this, you can see how the phase of a wave (rotation around the central origin) and the amplitude of a wave can be coded to retrieve bits from the diagram in figure ?? - called a ‘constellation diagram’.

NB - For reference, the QPSK we discussed earlier would be equivalent to the 4 inner nodes in this graph (labelled 0101, 1101, 0111, 1111) alone - without the surrounding extra 12 nodes.

⁴Trust me, I’m a mathematician...

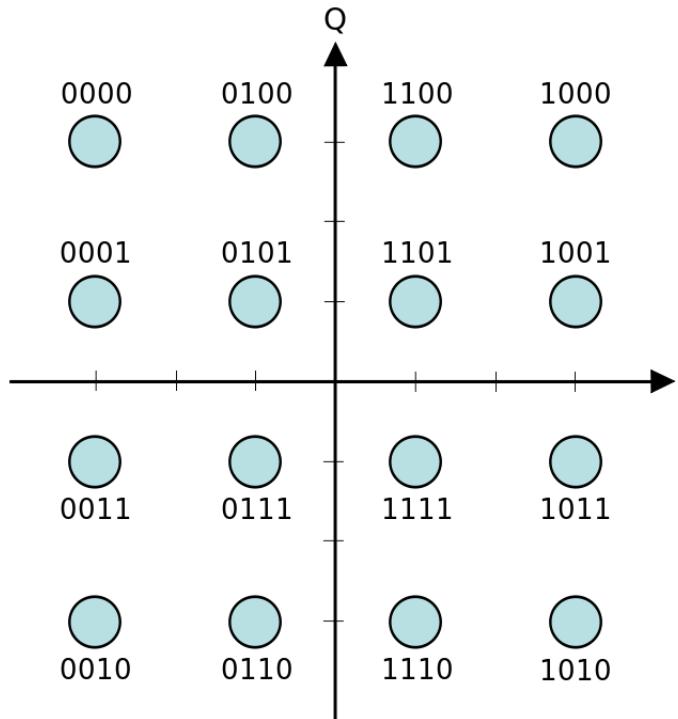


Figure 10.1: QAM16 diagram - c.o. Wikimedia; https://commons.wikimedia.org/wiki/File:16QAM_Gray_Coded.svg

As such, we can combine the ‘rotational’ aspects of QPSK with measurements of gain level from ASK in order to have a setup corresponding to every node on this diagram. Thus, each snippet of our radio transmission corresponds to 4-bits in this new setup - hence the name, QAM16.

Fun Fact: QAM, in its various guises, is the base modulation used by WiFi, 3G, 4G/LTE, etc.

Additionally, you can have a 4×4 grid in every quadrant, giving you what is known as QAM64 - this is used in LTE, in particular.

10.1.3 Other Nomenclature

Here are some other phrases you will come across in radio-hacking:

- **RSSI** - Received Signal Strength Indicator, this is a power-level output for a signal that is mandated by so many standards, inc. IEEE 802.1x, that it is standard kit for pretty much any radio IC.

10.2 The Radio Attack Surface

Now that we have given a solid coverage of (quite a bit of) radio theory, we should get down to business - *How can we use this knowledge to hack radios used in IoT?*

10.2.1 An overview of attacks

With IoT/embedded systems, you will find that you have to get closer for some radio implementations (such as BLE) or can attack the presented radio interface from some distance away (as with SigFox or LoRa). Here will will outline some of the attacks that are possible against radios of any type:

- **Replay Attacks** - from opening garage doors to opening cars to deactivating alarm systems, radio replay attacks are often found on systems that have a simple switch logic that reads a specific code in order to switch state - that is, the same radio pattern will lock as unlock a door, say.
 - It is worth emphasizing that, even with encrypted communications, unless your key changes or you use a stream cipher, packets will always look the same. An attacker might not know that some packet indicates a temperature of 15°C , but they can deduce it means ‘less than 20°C .
- **No Encryption** - this peculiar folly is where the designers, assuming nobody would think to use an SDR (or not being aware how easy they are to acquire), neglect to make use of any encryption during transmission, rendering the information readable to anyone who cares to listen in. This is very common on older systems, such as early pager networks, for example.
- **MITM Attacks** - as with other hacking, radio hacking has ways of man-in-the-middle’ing radio traffic in order to defeat rolling codes⁵, usually for the express purpose of fooling the receiver into thinking a command is legitimate by capturing one code, whilst jamming the receiver - for more, see Samy Kamkar’s ‘RollJam’ project.⁶
- **Information Disclosure** - some radio protocols are very helpful, and sometimes data can be disclosed or access can be granted unwittingly. This is common with naive BLE stacks, for instance.

⁵Codes that are sequential, usually based on a shared key known only to the receiver and transmitter.

⁶<https://samy.pl/defcon2015/>

10.3 Attacking Known Radios

There are two very common radios used in IoT, and a few other radios that are specific to the IoT space. We'll cover what these are, how to hack them, and link to more detailed resources where we can.

10.3.1 WiFi Hacking IoT

WiFi hacking is one of the first things most hackers learn these days.⁷ We are well used to using `airodump-ng` for logging data and then optimising `aircrack-ng` to crack WPA2 handshakes to get the WiFi key out of it. And yes, all of these apply to IoT/embedded systems and devices.

But, assuming you're already familiar with this stuff, what is specific to IoT? Here are some considerations:

- **DeAuth/DoS Conditions** - Devices do not always maintain connectivity, with or without an attacker nearby causing trouble. So, the question becomes this; *How does the device and/or system deal with outage?*
 - This might not be a security concern on face value, but recent work on the Amazon Key⁸ found that if you disable the camera, the stream would display the last known image... as opposed to some indicator that the service was down. This lead to the Wired article linked below - someone could be there, in view of the camera, but a simple WiFi deauth attack would mean you didn't know.
- **Insecure Default WiFi Credentials** - it is not unknown for generating weak passwords.
 - Perhaps one of the best examples is in figure ???. Yes, the default password is just the MAC address. Yes, that is transmitted in the wireless frame, so you can just sniff the password out of the air.⁹
- **Admin Functions to WLAN** - as we will see with WAN, often it is possible to perform administrative functionality over the WLAN. Whilst this is necessary in most cases, there are some (such as firmware update) that might be considered to have an option to be restricted to LAN-only connections (though this is very rare requirement and found implementation)

⁷To show my age, I remember building a cantenna in high school out of a Pringles can, and soldering it to my very-nice-for-the-time PCMCIA WiFi card to see stuff for *miles* from the top floor of the school. Oh, and warchalking... that was a thing... Good times.

⁸<https://www.wired.com/story/amazon-key-flaw-lets-deliverymen-disable-your-camera/>

⁹This example has been the author's pinned tweet for a few years, and people double-take it every time.



Figure 10.2: See if you can spot the weak WiFi password?

10.3.2 BLE Hacking IoT

Bluetooth Low Energy has overtaken Bluetooth as the default PAN (Personal Area Network) technology of choice. It has pervaded through the tech industry and found itself in everything from mobile phones, speakers, and headphones to medical instruments, children's toys, and even cars.¹⁰

Overview of the Structure of BLE

Fig ?? gives the overall structure of the BLE software stack¹¹:

¹⁰More on this later.

¹¹Image Credit: <https://www.digikey.com/en/articles/techzone/2011/mar/bluetooth-low-energy-for-wireless-sensors-and-actuators>

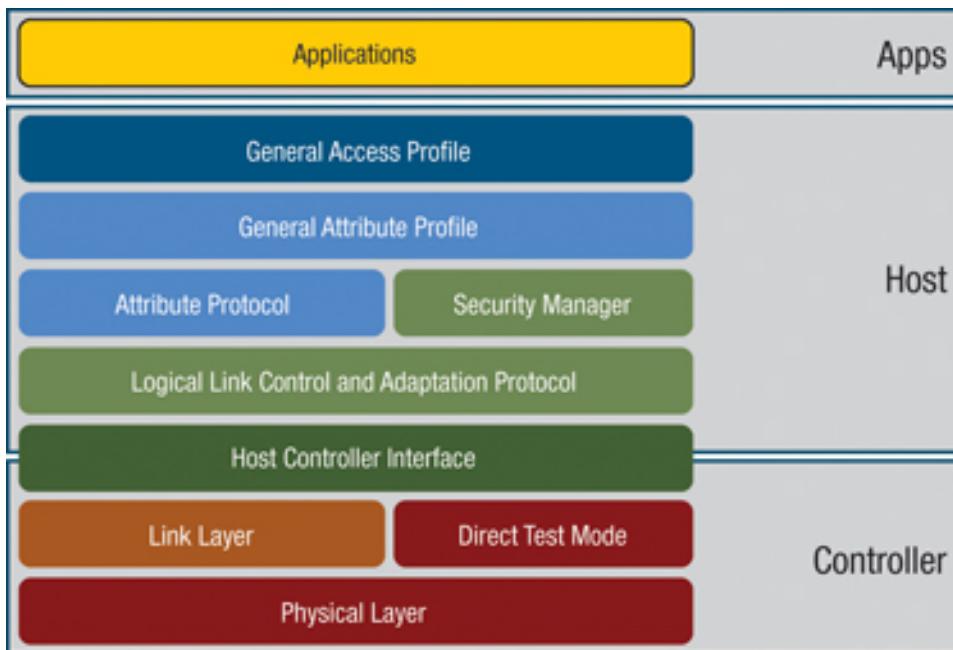


Figure 10.3: BLE Structure Overview

Nomenclature - EVERYTHING in BLE Is little-endian. Aggressively so. Also, we don't have 'bytes' - we have 'octets' of 8 bits. However, I will probably mess up and say byte, and you're allowed to, too.

Additionally, UUIDs we will see later can be either 128-bit (e.g. `dcc828ef-5799-4072-b6c3-ad7abf5b72eb`), or 32-bit (the 'xxxxxxxx' in `xxxxxxxx-0000-1000-8000-00805F9B34FB`) or 16-bit (the xxxx in `0000xxxx-0000-1000-8000-00805F9B34FB`).

Let's walk ourselves through the overall structure of this.

First off, we see that the application layer is exactly as any OSI model would give - it is just encapsulated data that is interpreted by some firmware. As such, much of our hacking will focus here, with understanding BLE as a means of delivery.

Next, we have the Generic Access Profile (GAP) - this is something we will never see, but this marries up the software interpretation of the GAP to the Generic Attribute Profile (GATT) which then connects to layers beneath this in the software stack.

GAP layer also handles addressing; for data, the hardware addresses are random, and for advertising, they are specific. Advertising packets have a specific structure:

- 1 octet of preamble
- 4 octets of address

- PDU 2-to-57 octets in length - PDU stands for ‘Protocol Data Unit’, which the standard describes as ‘A message’
 - 2 octets of header - this defines whether it is an advertising or data packet, amongst other things
 - up to 55 octets of data
- 3 octets of CRC check

PDU’s have 4 bits of type, of which the most important are:

- SCAN_REQ - scanning request
- SCAN_RSP - scanning response
- ADV_ID - connectable undirected advertising
- ADV_DIRECT_IND - connectable directed advertising
- ADV_NONCONN_IND - non-connectable undirected advertising
- ADV_SCAN_IND - scannable undirected advertising
- CONNECT_REQ - Connection request

Once we have a connection request, and then negotiate our connection, we are then dealing with the Generic ATTribute Profile, or GATT. The GATT may be familiar to anyone who knows Bluetooth prior to 4.0 - it was one of the (many) options for organising data. Now, in BLE at least, it is the fundamental way in which we sort data into a hierarchy:

- **Profile** - This is the entire profile for the device
- **Service** -
 - Primary or Secondary
 - This layer is a generic description of what characteristic data we want to be dealing with.
 - Referenced Services can point to other services - unlimited nesting
 - Examples would be ‘battery’ or ‘heart monitor’.
- **Characteristic** - where the magic happens:
 - Maps to an attribute in ATT
 - These can be referenced by UUID (128-bit, 32-bit or 16-bit as above)
 - Can be referenced by a handle in the ATT table

- Has set security options (rarely used)
- Has defined ‘properties’ which are really permissions
- Data!
- **Descriptor** - these simply describe the characteristic.

The full list of Characteristic properties are:

- Broadcast
- Read
- Write without response
- Write
- Notify
- Indicate
- Authenticated signed writes
- Extended properties

HCI

Tools for Identifying BLE characteristics

So, now we have seen what the overall structure is - the natural question is *Well, how can we assess it?* Here’s a rough guide to getting information from a BLE target, and the various ways you can interact with them.

First off, here’s a list of hardware:

- Any BLE dongle - most BLE dongles are fully featured for doing BLE assessment
- An Ubertooth - versatile, long range, and very flexible piece of kit, they are certainly worth the investment for coverage alone
- BLE Module - intended for use with a uC and then soldered on as a PCB add-on, these are useful for hacking if you are designing a tool for BLE hacking.

Next, here is some software packages for gathering information:

- On Android:
 - NRFConnect - an app by Nordic Semiconductor, is very useful and can afford low level control of a BLE interface (though not the easiest to use when testing)

- BLE Scanner - and app that is another option from NRFConnect. Free, easy to use, and just as powerful, with some nifty features (like proximity scanner that shows you devices by RSSI)
- hcisnoop - any android phone can log the data sent to the physical HCI layer of BLE. The output is sent to `/data/sdcard` as `hci_snoop_XXX.log` where XXX is variable. These files can be pulled off with `adb pull` and then opened up in wireshark for inspection!
- `bleah` - A superb tool, with a very nice interface, ideal for screenshots and putting in reports.
- `gatttool` - now a standard tool in the BLE hacking arsenal, this tool has command line switches for getting data out of the full BLE stack.

And finally, tools that can be used for assessing/PoC BLE:

- BluePy - a neat python library that makes it really easy to script BLE interactions and design tests/fuzz.
- `bleah` - As much as it is a tool for intel gathering, it is also a decent platform for injecting data
- `gatttool -I` - the Interactive mode for gatttool is very handy and easy interface for exploring a device's BLE interface

Worked BLE Hacking Example - the iTAG

[Put the full worked example here!]

Should lead up to the creation of the Annoyatron 2000! <https://gitlab.com/snippets/1720857>

10.3.3 ZigBee

10.3.4 LoRa, LoRaWAN, and SigFox

10.4 Attacking an Unknown Radio

To include:

- Using and SDR to record traces in gnuradio
- Using gr-inspector¹² to analyse radio signals
- Using this information to deduce what data is being transmitted
- How to use this to hack the radio interface

[Much more to come as I write it!!]

¹²<https://github.com/gnuradio/gr-inspector>

Chapter 11

Putting It All Together

11.1 Background

Now that we have explored the two main basic techniques of accessing UART and dumping firmware over SPI, as well as covering how to analyse firmware and what to look for, we should make plain how this low-level access can inform your pentesting.

I will assume general knowledge, familiarity, and experience with the following types of testing:

- Infrastructure testing - nmap/nessus scanning, and running basic exploits through MSF as a baseline
- Wi-Fi testing - DEAUTH attacks to get handshakes, understanding of WEP/WPA/WPA2, and more recent attacks like Pixie-Dust
- Web Application testing - familiarity with XSS, CSRF, LFI, RCE, SQLi, and the OWASP Top 10 in general
- Mobile Application testing - Testing web API's (XML/SOAP/JSON) and assessing mobile devices

The reasons for all of these testing methodologies being relevant comes from the vast attack surface presented by embedded devices in a modern world. Many of these were alluded to in the first chapter of Part II, so we will assume that this has been understood for this more practical chapter.

11.2 Hardware Specific Vulnerabilities

Here follows a list of common vulnerabilities that affect only the hardware:

1. UART Access - yes, this should actually be restricted. Many SoC's support the limiting of this service, and it should not be possible to get a root debug shell from accessing the UART on boot.

2. Unsigned firmware images or SecureBoot disabled on ARM SoC's

11.3 Pentesting Embedded Devices - Technical Overview and Attack Models

11.3.1 Web Application

This is perhaps the most fruitful of all the testing areas, so I will mention it first. The first two vulnerability classes are so numerate, that they have their own sections.

Remote Command Exec (RCE)

RCE is a high risk, high probability vulnerability on embedded systems, owing to the fact that the main way that performance overheads are mitigated is by the implementation of BASH scripting to perform almost every function on the device. It is not unknown for entire '.asp' websites on embedded systems to be nothing more than renamed BASH scripts that get run on access by the webserver, passing parameters in as you would any other script.

Given the prevalence of this *modus operandi*, it should be tested for exhaustively. Here is a good procedure to gain familiarity with RCE's in embedded systems; Determine whether any user input ends up in any BASH scripts we identified earlier. Common ones include DNS settings, DynDNS settings, NTP timeserver settings, etc. Use the UART connection to test whether a submitted payload has reached the script unencoded. Once this has been confirmed, different payloads can be created.

It is handy to bear in mind the following RCE payloads that do not use spaces (lines 2,3): (We assume that the UART terminal is mapped to /dev/ttys0)

```
||"; echo 'test' > /dev/ttys0;
||";echo$IFS'test'>/dev/ttys0
||";{echo,'test'}>/dev/ttys0;
```

NB - when testing with \$IFS, beware when trying certain commands or using pipes. Sometimes you need to add an escape at the end - e.g.

```
||nslookup$IFS.google.com
```

would need to be

```
||nslookup$IFS\google.com
```

on some systems.

Cross Site Request Forgery (CSRF)

CSRF tokens are generated by many embedded platforms, but are almost never checked. The general appearance is that developers of embedded man-

agement web-apps appreciate the need for a fix, but fail to understand fully the nature of the vulnerability. Given that most devices will sit on the same local IP address (usually in the 192.168.0.0/16 range) and users almost never change this base address, CSRF has a wide ranging capability. The ability to CSRF a setting/configuration by embedding a javascript payload that talks to the management application has been found to permit the following actions to be performed by payloads loaded over the internet:

- RCE through the local browser
- Setting malicious DNS servers
- Deactivating protection mechanisms (built-in AV and firewalls)
- Uploaded malicious firmware files
- Loading a known stored XSS that then fires to disclose information back to the attacker

(**NB** - this is by no means an exhaustive list, and is included for illustrative purposes.)

Example code to use for testing a payload's effectiveness can be found in the [Appendix I am yet to write].

OWASP Top 10

Taking excerpts of the most significant issues, the following should be kept in mind when pentesting:

Authentication Control and Management

Perhaps the third most common issue on embedded systems, Authentication is often handled in a very lax manner. Things to watch out for include:

- Weak Password Policies or client-side only password strength enforcement
- Passwords not required to be changed after first login to management panel
- Access to management screens without authentication (by forced browsing, for example)
- Platform Authentication Bypass - if the system has NTLM or Basic Auth implemented, could this be bypassed, e.g. by having a valid cookie present

- Predicatable Cookies/Tokens - often, to save on processing and I assume large crypto libraries, hashing is always weak, random numbers generally aren't that random, and so cookies are predominantly time-based and very predictable.

SQLi

Sadly, not much of a thing on embedded devices. Although some do use an sqlite implementation for data storage (very, very few), it is generally not much of an issue owing to a Database server causing significant performance and space overheads that are unnecessary.

It is thereby not something to be overly concerned about if it is not found - a static analysis of the firmware will show you that there probably is no database service to be talking to in any case. **NB** - this does not preclude it being present on externally hosted services!!

Insecure Direct Object Reference (IDOR)

This type of issue is also generally not found on embedded systems, however, their cloud-based counterparts for externally hosted services are often vulnerable, opening up access to other customers' devices in the worst cases.

11.3.2 Infrastructure Vulnerabilities

When pentesting any network service presented on the external WAN or internal LAN/Wi-Fi interfaces there are a number of things you are looking for. Here is how low-level access and firmware analysis can help when assessing this attack surface:

1. Analysis of the bootlog (from UART) will yield precise service versions that nmap/nessus can miss.
2. Analysis of configuration files will yield default passwords or required keys to access and test running services.
3. Analysis of the firmware can determine what SSL Certificates are in play and how they operate (are the private keys universal?).
4. If testing a VOIP service, the VOIP secrets should be analysed to see if they are generic/universal. It is also worth doing SIP testing with tools like [sipvicious](#) in order to ascertain if there are any issues there.
5. If testing a UPnP service, this should be assessed to see if unauthorized access can be granted through some other (web based) vulnerability (they have been known to be CSRF-able).

6. When testing any service, if you successfully generate an exception, then the stack trace has been known to be dumped to serial (in the spirit of debugging). As such, this can be useful for trying/fuzzing payloads to binary or other services on open ports.

11.3.3 Wi-Fi and Wireless testing

Here, we focus on primarily CPE's (home routers or switches with Wi-fi capability), but these can apply more generally to any Wi-Fi interface.

1. Uploading Firmware over Wi-Fi - this should not be possible. Else the attacker who cracks the password and comes back can set a malicious firmware without ever getting physical access to the device.
2. It should not be possible to determine the default Wi-Fi password or WPS pin by sniffing for publically available information (such as the default SSID or MAC address)
3. It should not be possible to bruteforce WPS pins nor trigger a pairing without physical access to the device.
4. Attacks that are more recent, such as the Pixie-Dust attack I mentioned earlier, should be guarded against, as well. This is unlikely in many cases, owing to the lack of updating base software in the firmware.

See the next chapter for more information about WiFi, BLE, and other radios in use on embedded/IoT systems.

11.3.4 Mobile App/Cloud Integration

We assume for this testing that there is a tie-in mobile app, or management mobile app that may also have some 'cloud'-based connectivity.

1. Many apps permit management of the device itself - and the credentials to allow this access should not be pushed to `adb logcat`. This should be checked during standard operation.
2. The cloud-based service should be checked for all the usual problems, with special attention paid to username enumeration and bruteforcing - such a compromise can lead to actual compromise of real devices.
3. If the cloud-based platform or mobile application permits the 'adding' of devices, it should be confirmed that no Insecure Direct Object References (IDOR) vulnerabilities could be used to add devices arbitrarily.
4. Tying into RCE vulnerabilities, if any are identified locally, then these should be tested to see if they can be chained with a CSRF or over the app API itself to trigger such a vulnerability over the internet.

We'll cover more details regarding mobile app testing for IoT and cloud hosting tests relative to IoT.

11.3.5 Denial of Service (DoS)

This can be brought about in a myriad of ways, and it is very hard to give general guidance. However, the following examples have been found in the wild:

- Bad Error Handling
- Unexpectedly large parameters on web headers (such as User-Agent or Referer)
- Incorrect Basic Auth headers - where the header is invalid Base64 has crashed a device completely
- TR069 Port on external interface - DoS through fuzzing.

In more general terms, it should be assessed wherever some level of instability of processing is found in the system.

11.3.6 Additional Testing Points

Bluetooth/BLE

Many BLE chips have custom debug stacks that are now well known, and tend to run their own firmware. As such, it should be seen if the devices can be connected to without any keys, or with default/universal keys. They should be segregated from the rest of the device's functionality, and this should be confirmed (it shouldn't be possible to cause a system-level exception from the BLE data being sent to the SoC).

See the chapter on radio hacking for more information about BLE hacks.

Part III

Advanced Testing Topics

Ideas include:

- running firmware files in virtual machines (qemu)
- JTAG/SWD in practice
- Glitching
- IC Decapping - with and without acids - and silicon die reverse engineering
- Reverse Engineering bootloaders
- worked example of bit-banging a protocol - CC2541 custom dbg??
- CRP Bypasses

Chapter 12

Virtualizing Firmware

12.1 The Problem

Devices are fickle and delicate and prone to, oh, blowing up or dying or not working anymore, or even not being provided by a client in time for an assessment, sometimes. This is annoying, and hard to deal with, but part of the many facts of life.

However, we have virtual machines, and indeed, virtual appliances (such as virtual firewalls appliances, etc.) - these are just some virtual machines that run the firmware for the desired security appliance in a sandbox/VM/-container. Additionally, it is well known that QEMU, for example, supports emulation for ARM/MIPS/x86/etc., and makes them available for access on an x86 machine.

Wait, so... why can't we virtualise firmwares then?

Well, the fact of the matter is, *you can...* -ish... kinda... permit me to explain.

12.2 The Solutions

As we have seen earlier, there are fundamentally two main types of firmware - those running a Linux-like environment, or those running straight code. We'll run through how you can virtualise both - with examples linked, but more importantly, we will aim to include enough theory that you can take something that hasn't been virtualised before and actually have a reasonable chance of achieving this.

12.2.1 Emulating Linux-based Router Firmware with Firmadyne

Firmadyne can be cloned from the firmadyne repo: <https://github.com/firmadyne/firmadyne> ← You'll need this to understand this part

of the chapter.

Firmadyne was developed to perform the following actions automatically:

1. Search vendor websites for firmware, and download them
2. Unpack the firmware and extract the filesystem.
3. Virtualise this in qemu for either MIPS or ARM as needed.
4. Run exploits/nmap scripts/metasploit modules against the running firmware.

However, for our purposes, it will serve to be able to run some various firmwares virtually, meaning we can do auxiliary/automated testing of our own, without needing to sacrifice the device or restrict our access to it.

The workflow for firmadyne internally is as follows:

1. Use `scraper` to find files on the vendor site and download them.
2. Let `extractor` extract the relevant binary data from the firmware file.
3. Setup `qemu` using scripts and two custom binaries:
 - `libnvram` - an NVRAM peripheral emulator compatible with qemu
 - `console` - a serial console device emulator compatible with qemu
4. Choose the appropriate architecture and run the firmware virtually
5. Invoke custom and metasploit based exploit scripts, and run nmap scans on the running target

We wish to essentially use the automatic part of this process for steps 3 and after in our work, having done the first two steps manually before now.

Roughly, we would do this as follows:

1. Extract the firmware and filesystem as we have shown previously - judicious use of `binwalk`, `unsquashfs`, `jefferson`, and other tools will make this possible and reasonably straightforward.
2. Prepare the firmware extracts to run in qemu; initialise the filesystem, and setup the NVRAM emulator and console device.
3. Run the relevant QEMU script in the firmadyne repo.

With a lot of luck, it should now be running, however, usually you can expect to have to fix the following things:

- Memory mappings for UART devices
- Memory mappings for network devices

- Memory mappings for GPIOs
- Locating the entrypoint for the firmware's kernel

NB - most of these can be fixed with some patience, some thought, and some steady immersion in the relevant data sheet for the SoC this is supposed to be running on (if it is available).

12.2.2 Emulating STM32 Firmware with [qemu-stm32](#)

Anyone who has ever fought with GRUB2 or a Linux ditribution that doesn't have a binary in its package manager (or even not having a package manager¹) - Linux internals are scary for beginner, but there are enough people around who have solved the problem (and posted it to StackOverflow) for us to not be too worried.

But what about those pesky firmwares that *just run code??!!*

Well, for the STM32 there is the [qemu-stm32](#) fork of QEMU that permits the ARM-SoftMMU to run an STM32 firmware, and even emulate parts of a development board for it!

You can download qemu-STM32 from the repo here: https://github.com/beckus/qemu_stm32

For a walkthrough and tutorial for how to use these see here: https://github.com/beckus/stm32_p103_demos

[More details to come... would be nice to show how to take a memory map and implement that for virtualisation]

¹I leave it for history to decide whether my love of Slackware was the result of Stockholm Syndrome or not...

Chapter 13

Practical JTAG/SWD Debugging

13.1 What is JTAG and SWD?

13.1.1 Motivation

We saw JTAG and SWD before in this volume - in ??, for reference. However, in order to really get to grips with hardware hacking, we need to dive in-depth and explore the protocols in full. So, let's go!

JTAG, or ‘Joint Test Action Group’¹ is a consortium of manufacturers who defined the standard that we now refer to by the same name. *But what is it??*

Quite simply, it’s a debug and testing protocol designed to add various troubleshooting and programming features to IC’s of any package. It was envisioned that such a thing was required with the advent of BGA packages in particular. Why BGA? Well, principally it comes down to soldering errors. With packages with external pins - like LQFP for example - it’s easy to either visually or physically test it with probes. With a BGA you can’t see under the component, and you can’t access it with probes - so what do you do? You use something like JTAG. Hopefully we’ll explain how in the next section.

This said, the functionality of JTAG expanded significantly, so we need to talk about JTAG functions, the state machine, and how it ties together. We will also discuss Serial Wire Debug (SWD) which is an implementation of JTAG used by ARM-based processors. For the purposes of practical action, they are the same, so we will treat them as such (and highlight any differences as and when they arise)². As such, when I say JTAG, you can substitute SWD, and vice versa, without any loss of generality.

¹Resist the urge for ‘legal state quality assurance’ jokes, please...

²The ARM Debug Interface Architecture v5 is the definitional point for SWD - it can be found here: <http://infocenter.arm.com/help/index.jsp?topic=/com.arm.doc.ihi0031c/index.html>

13.1.2 JTAG/SWD Functions

There are three principle functions performed by JTAG/SWD:

1. Loading Firmware
 2. Debugging the system
 3. Boundary Scan Testing

We'll cover these in reverse order, especially as people are probably not very familiar with the term 'boundary scan'.

Boundary Scan:

JTAG has the boundary scan silicon literally all around the boundary of the silicon die of an IC. As such, it can talk to every pin directly, and set its level to high or low, as well as read the level of every pin. Kinda neat! Very useful for hacking, as I'm sure your brain is thinking right now.

The boundary scan is when the JTAG state machine sends a binary string around all the pins, sets the boundary nodes to read, and sees how the string differs on return. What use is that? Well, suppose we send out

- to set pin 26 to high - and we get back

- with both pins 26 and 27 high. This would indicate a possible solder short or some other physical issue with the board.³

Of course, we can also use this to manipulate the operation of the chip or read pin values as we need during regular operation - it has even been shown that you can emulate UART over JTAG, but this is non-trivial and very unreliable compared to the easier methods you could employ.

³Yes, this is a slightly over-simplified and maybe contrived example, but it illustrates the point and one of the original intended uses of JTAG.

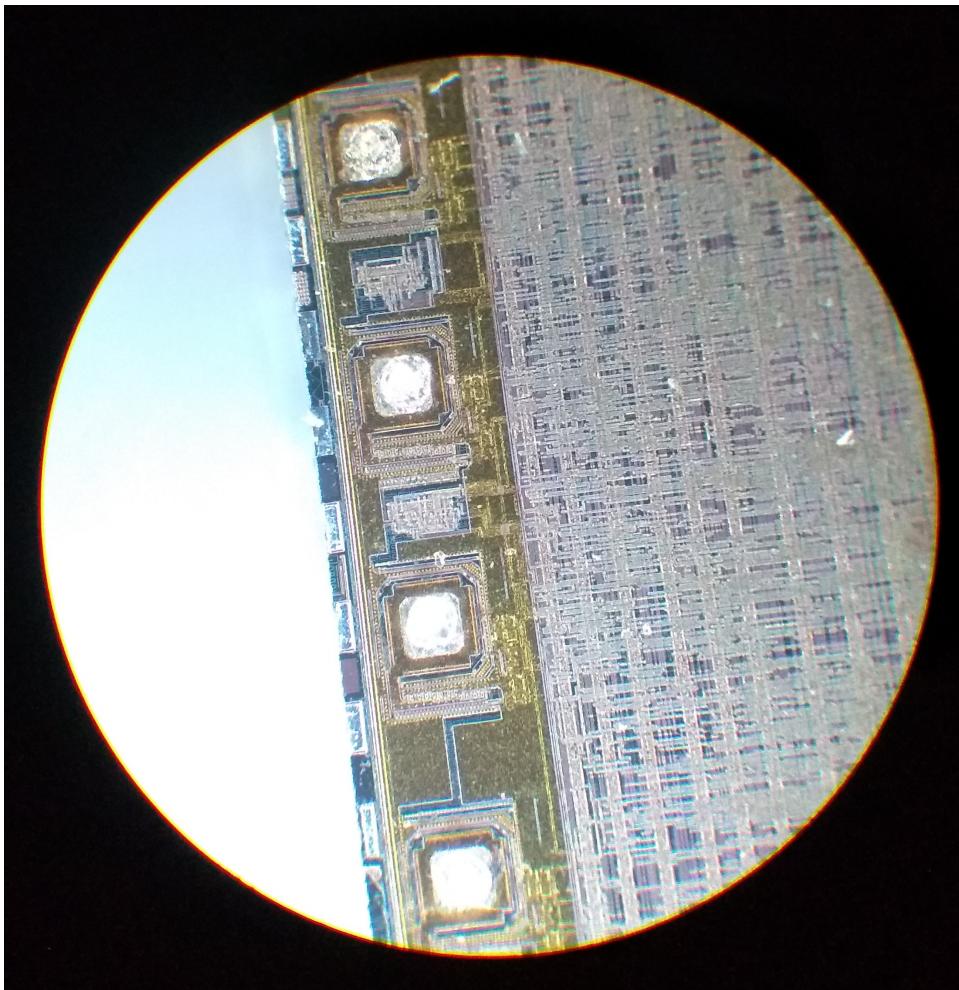


Figure 13.1: Showing the boundary modules next to pins on a silicon die under a microscope

Debugging:

JTAG was also designed to allow a developer to debug the CPU and peripherals⁴ onboard the microcontroller. As you would expect, you can perform all the usual debugging operations:

- Setting breakpoints
- Halting the CPU
- Setting register values/manipulating the \$PC

⁴Recall that SPI or CAN or UART interfaces are all ‘peripherals’ included in the microcontroller silicon - they really are just placed and wired up in a similar fashion to peripherals in macro computers.

- Dumping memory contents (for disassembly, say)

OpenOCD as part of its normal operation will spawn a GDB server that will happily interface GDB to our target with ease (see a worked example later). Thus, you can reduce the problem of debugging a chip quite easily to debugging any binary in the familiar GDB context.

Loading Firmware:

JTAG can load firmware on, and off, the microcontroller. Yes, you read that right. You can often retrieve firware using JTAG/SWD on a target IC. This will, of course, depend on the settings of the security bits on the flash (See the chapter on CRP bypass for how these can sometimes be circumvented).

OpenOCD has, for many microcontroller architectures, implemented flash download and flash upload functionality, and this is the same for other software packages for specific debuggers, such as J-Link.

13.2 Fully worked JTAG/SWD example

The road to JTAG/SWD is paved with good intended features.

Here we will describe the process of getting SWD on a STM32F103 chip by means of a STLink-V2 - specifically, a cheap chinese knock off, and a ‘blue pill’ development board as our target.

First we wire up the devices. Here are the pin mappings:

- SWDIO goes to SWDIO
- SWDCLK goes to SWDCLK
- GND goes to GND
- VCC goes to VCC - assuming we want to power the chip through the debugger
 - If the IC gets power elsewhere, don’t connect this.
 - If connecting to a target currently *in situ*, double check for any ground loops before powering anything.

We will use OpenOCD - we cover the installation and configuration of this in ??, so we will assume that this is installed.

First off, we run this command:

```
|| sudo openocd -f interface/MyConfig.cfg -f target/stm32f1x.cfg
```

This will run the OpenOCD servers in the background. In order to pull the firmware, we will need to connect to the telnet server using this command:

```
|| telnet localhost 4444
```

[Maybe put some stuff in here about ELF-ing a firmware for STM32?
It's a bit niche, though...]

Once we want to start debugging, we can open the appropriate GDB (such as `arm-none-eabi-gdb`) and connect to the local GDB server that will interface us to the IC as follows:

```
|| target remote localhost:3333
```

Once we are all connected we are ready to begin debugging our target as if it were any other binary process we are familiar with debugging. Success!!

13.3 JTAG Disabled Bypass

(Material in this section is based on <https://medium.com/@LargeCardinal/how-to-bypass-debug-disabling-and-crp-on-stm32f103-7116e7abb546>

13.3.1 Overview

There is a good case to state that modern devices are moving towards a very slightly more homogenised ecosystem. This is signalled by the number of cheap, easily available ARM-based microcontrollers and microprocessors available on the market.

ARM offers many features ‘as standard’ across various applications, and with low-price IC’s on the market, there is a good chance that the uptake of ARM architecture will continue. The STM32 series of microcontrollers are quite popular with developers owing to good physical/performance properties, as well as a low price point.

We will consider the STM32F103 - specifically the STM32F103C8T6 IC - and how you can disable debug protections and CRP. These are nifty little Cortex-M3 based chips with 64kb flash, 20kb RAM, and clocked up to 72MHz. You can see, perhaps, why they might be popular!

13.3.2 The Problem

Implementations of ARM can vary wildly. There are some standard features that are either missed out or reconfigured completely - each IC manufacturer is left to their own devices (→pun) for how they produce ARM chips.

There are two protection mechanisms that are usually recommended by hardware/embedded/IoT security people:

1. Disable Debug ports - JTAG and SWD disable
2. Enable firmware protections e.g. Code Readout Protection or CRP

We want to look to see if we can disable these protections on the target IC - the STM32F103.

13.3.3 Disabling Debug

The first is quite straightforward - for STM32F103 you simply set `__HAL_AFIO_REMAP_SWJ_DISABLE()`; in the standard file `stm32f1xx_hal_msp.c` (if you are using STCubeMX generated output - you can do something similar with libopencm3, too).

This is the default configuration of STCubeMX-generated firmware source projects. Meanwhile, the IC's themselves ship with SWD Enabled but the flash being locked - but this is trivial to unlock over SWD.

But once such a firmware with SWD/JTAG debugging disabled is applied, then this device no longer can be debugged. At least, that is the idea.

So, our initial scenario is an STM32F103 device with Debug Disabled.

13.3.4 Bypassing Disabled Debug

All this said, there are two pins that should be brought to attention - the BOOT0 and BOOT1 pins. These are designed to be used to trigger fall-back/default modes on the IC - and as usual, this can spell disaster if you're trying to hide some IP.

Specifically, we first want to use BOOT0, which is pin 44 on the LQFP48/UFQFP48 packages. Here is the attack for this:

13.3.5 Attack Outline

- We need to raise this BOOT0 pin to 3.3V and then power cycle the IC
- Acquire a copy of the firmware
- Patch the firmware to enable SWD debugging
- Re-flash the new firmware to have the device with debugging enabled

So, first, we raise the voltage of pin 44 to 3.3V. If you're playing along at home with a Blue Pill development board, you'll note that this is one of the two jumper pin options on the top of the board.

Here is the board before and after we apply the change to BOOT0:

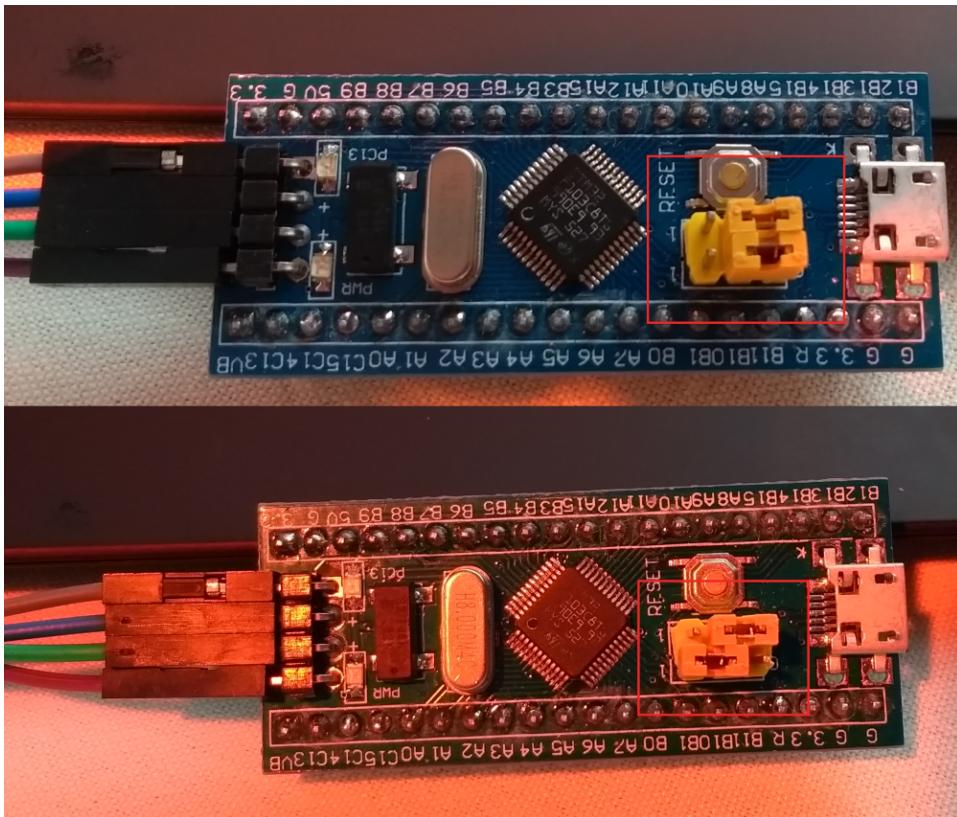


Figure 13.2: Showing the change of pin.

We will now see that we can debug the IC, and dump firmware using SWD. That's it. That's the first part of the attack. Having issues? You can just use mbr to cycle through all the addresses on the flash using a script, or dump the image, both using the openOCD client, an STLink-V2 USB-to-SWD interface, and a little bit of time/SSD space.

Once you have acquired the firmware, how can we enable debugging on the code? Well, first, we reassemble the firmware into some `firmware.bin` file.

With this done, we need to do some binary editing. Here is the diff between two otherwise identical firmware files - the first has debug disabled, and the second (lower) has SWD Enabled.

```

./blinky1-SWDDisabled.bin
0000 0E70: E3 84 20 46 FF F7 7E FD E4 E7 00 2A EB D0 00 23 .. F..~. ....*....#
0000 0E80: E3 85 20 46 FF F7 7F FD E5 E7 90 F8 39 30 90 F8 .. F.... ....90..
0000 0E90: 3A 00 18 43 70 47 00 00 00 B5 83 B0 21 4B 9A 69 :..CpG.. ....!K.i
0000 0EA0: 42 F0 01 02 9A 61 9B 69 03 F0 01 03 01 93 01 9B B....a.i .....
0000 0EB0: 03 20 FF F7 45 FA 00 22 11 46 6F F0 0B 00 FF F7 ...E..".Fo.....
0000 0EC0: 51 FA 00 22 11 46 6F F0 0A 00 FF F7 4B FA 00 22 Q..".Fo. ....K..".
0000 0ED0: 11 46 6F F0 09 00 FF F7 45 FA 00 22 11 46 6F F0 .Fo..... E..".Fo.
0000 0EE0: 04 00 FF F7 3F FA 00 22 11 46 6F F0 03 00 FF F7 ....?..".Fo.....
0000 0EF0: 39 FA 00 22 11 46 6F F0 01 00 FF F7 33 FA 00 22 9..".Fo. ....3...".
0000 0F00: 11 46 4F F0 FF 30 FF F7 2D FA 07 4B 5A 68 22 F0 .F0..0.. ...KZh".
0000 0F10: E0 62 5A 60 5A 68 42 F0 80 62 5A 60 03 B0 5D F8 .bZ`ZhB. .bZ`...].
0000 0F20: 04 FB 00 BF 00 10 02 40 00 00 01 40 02 68 18 4B .....@ ...@.h.K
0000 0F30: 9A 42 00 D0 70 47 30 B5 87 B0 03 F5 58 43 9A 69 .B..pG0. ....XC.i
0000 0F40: 42 F4 80 42 9A 61 9B 69 03 F4 80 43 01 93 01 9B B..B.a.i ...C....
0000 0F50: 4F F4 00 73 02 93 02 23 03 93 03 23 05 93 0D 4D 0..s....# ...#.M
0000 0F60: 02 A9 28 46 FF F7 8E FB 4F F4 80 63 02 93 00 24 ..(F.... 0..c....$.
0000 0F70: 03 94 04 94 02 A9 28 46 FF F7 84 FB 22 46 21 46 .....(F ...."F!F

./blinky1-SWDEnabled.bin
0000 0E70: E3 84 20 46 FF F7 7E FD E4 E7 00 2A EB D0 00 23 .. F..~. ....*....#
0000 0E80: E3 85 20 46 FF F7 7F FD E5 E7 90 F8 39 30 90 F8 .. F.... ....90..
0000 0E90: 3A 00 18 43 70 47 00 00 00 B5 83 B0 21 4B 9A 69 :..CpG.. ....!K.i
0000 0EA0: 42 F0 01 02 9A 61 9B 69 03 F0 01 03 01 93 01 9B B....a.i .....
0000 0EB0: 03 20 FF F7 45 FA 00 22 11 46 6F F0 0B 00 FF F7 ...E..".Fo.....
0000 0EC0: 51 FA 00 22 11 46 6F F0 0A 00 FF F7 4B FA 00 22 Q..".Fo. ....K..".
0000 0ED0: 11 46 6F F0 09 00 FF F7 45 FA 00 22 11 46 6F F0 .Fo..... E..".Fo.
0000 0EE0: 04 00 FF F7 3F FA 00 22 11 46 6F F0 03 00 FF F7 ....?..".Fo.....
0000 0EF0: 39 FA 00 22 11 46 6F F0 01 00 FF F7 33 FA 00 22 9..".Fo. ....3...".
0000 0F00: 11 46 4F F0 FF 30 FF F7 2D FA 07 4B 5A 68 22 F0 .F0..0.. ...KZh".
0000 0F10: E0 62 5A 60 5A 68 42 F0 00 72 5A 60 03 B0 5D F8 .bZ`ZhB. .rZ`...].
0000 0F20: 04 FB 00 BF 00 10 02 40 00 00 01 40 02 68 18 4B .....@ ...@.h.K
0000 0F30: 9A 42 00 D0 70 47 30 B5 87 B0 03 F5 58 43 9A 69 .B..pG0. ....XC.i
0000 0F40: 42 F4 80 42 9A 61 9B 69 03 F4 80 43 01 93 01 9B B..B.a.i ...C....
0000 0F50: 4F F4 00 73 02 93 02 23 03 93 03 23 05 93 0D 4D 0..s....# ...#.M
0000 0F60: 02 A9 28 46 FF F7 8E FB 4F F4 80 63 02 93 00 24 ..(F.... 0..c....$.
0000 0F70: 03 94 04 94 02 A9 28 46 FF F7 84 FB 22 46 21 46 .....(F ...."F!F

```

Figure 13.3: vbindiff view - look at the bits in RED

As you can see, at 0x00000F18, the 2-byte word there is the only thing that is different. So, search for 0x8062 at the same address in our target firmware and patch it to be 0x0072, and then re-flash the firmware onto the STM32F103 chip.

We can then lower the BOOT0 pin voltage, power cycle the IC, and then we should have (possibly with a little more hacking) SWD enabled on our target firmware.

13.3.6 Remarks

Although quite simple, this is a valid attack when you consider that

‘disabling debug’ == ‘firmware protection’

in the minds of some developers and frameworks. As such, attacks like this both help a researcher/assessor potentially access a chip, and demonstrate by Proof of Concept⁵ that this alone is not an effective countermeasure.

⁵Remember: PoC||GTFO.

See Chapter ?? for details on other hardware control bypasses on firmware!

Chapter 14

Glitch Attacks - Real World Fault Injection

Breaking firmware controls and hardware security with expertly timed fault injections.

14.1 What is Glitching?

We will introduce glitching - what it is, and how to do it, and what you will need. The majority of this material is greatly expanded on Colin O'Flynn's New AE Technology Wiki pages: <https://wiki.newae.com/> We strongly suggest that the interested reader head over there to gain much more background and specifics about this material.

14.1.1 Glitching - Theory and Praxis

A fault injection, put simply, is the intentional introduction of a system fault. Whilst that may sounds like a cynical definition, think about it for a moment - if we can control when a system goes into 'undefined behaviour', whatever that may be, then we can control a 'glitch' or fault on that system.

But why would you want to do this? Well, when a system is operating normally, it may be difficult to influence its behaviour (if, say, it has a limited input capability or user interface), or have some controls placed on it, e.g. activated CRP¹ on the firmware. Though not a guaranteed technique, and a technique that is as much an art as it is a science, glitching is a very useful tool to have in the arsenal of a security assessor/researcher.

¹Code Readout Protection

The general overview of a fault injection, or ‘glitch attack’², is that we want to influence the execution of code on some particular target by pushing the system into some unstable state. We can then use this unpredictability to cause the system to do any of the following:

- Break out of a control loop and into other areas of code
- Bypass some defined control by causing a mis-read of a security setting
- Bypass a silicon-level control by directly influencing the operation of a circuit
- Utilise a misread to read data we should not have access to

We’ll discuss how we achieve these outcomes in the next section.

14.1.2 Types of Fault Injection

There are ostensibly only a few major types of Fault Injection. We’ll list the most common here:

- Voltage Fault Injection
- Clock Fault Injection
- Electro-Magnetic Fault Injection
- Photo-Fault Injection

As is probably apparent, these are named based around the method by which the fault is introduced into a system.

But how exactly does this work?

This is a tricky question to answer in a general way, but the overall method of operation relies on the general operating principles of a system.

Voltage Fault Injection

A voltage fault injection generally relies on the capacitance of the operational circuit - by cutting out the voltage for microseconds, either completely or by cutting it out completely or to a lower level (say from 3.3V to 1.8V), we force the target CPU to begin using up its stored charge.

If we time our injection just right, there will not be enough power to write values to registers or memory locations, or for the CPU to properly track its program counter. If we properly time the glitch - both in the ‘when’ to glitch and the ‘for how long’ we glitch - we can give the CPU full power just as we have introduced the fault.³ It is because of this exact timing

²These names are pretty much synonymous, so will be treated as such in this text.

³This is why I say that glitching is ‘half art, half science’ - the science is quite straightforward, but working out where and when and for how long we perform a glitch to get the desired result varies tremendously - even between two identical devices.

requirements across microsecond intervals that we usually make use of an FPGA or high-speed uC, with high-speed switches for activating a glitch.

Clock Fault Injection

A clock fault injection is a little different. These introduce spurious rising and falling edges to a clock signal in an attempt to cause misalignment of the internal timings for the CPU operation.

Most CPU's operate on a 'load/execute' simultaneous pipeline structure. So, instruction one is loaded on the first clock cycle, then the second instruction is loaded whilst the first is executed on the second cycle, and so forth from there.⁴

As such, we can break out of the designed timing for execution by injecting clock faults at just the right moment. The effects of this kind of fault is generally to get the CPU to mis-read a function output incorrectly - say, a password function does not manage to return its 'incorrect password' error before the comparison is performed.

EM Fault Injection (EMFI)

EMFI is where we use coils placed above various parts of the IC in order to induce a glitch current directly into the IC by means of a burst generator.

Methods vary for this, from Colin O'Flynn's ChipSHOUTER to more manual techniques involving hand-wound coils and burst generators run over precise triggers. We won't cover this particular method in depth here, but instead refer the reader to these papers:

- <https://pdfs.semanticscholar.org/522f/14843a03372588397b9de88b088e50e6638.pdf> - This paper, from Riscure, is well worth starting with.
- <https://eprint.iacr.org/2012/123.pdf>
- <https://arxiv.org/pdf/1402.6421.pdf>

Photo-Fault Injection

Similar to the above, and inspired by similar techniques found in CRP Bypass techniques with UV light, it has been found that glitches can be introduced by introducing a large amount of light (say, from a camera flash) into an IC.

The most notable example of this was found on the Raspberry Pi 2 some years ago. Here is the 'free physics lesson' from the Raspberry Pi Foundation in relation to this 'bug': <https://www.raspberrypi.org/blog/xenon-death-flash-a-free-physics-lesson/>

14.1.3 Equipment Needed

Let's discuss the equipment we're going to need in order to perform glitch attacks on a variety of targets.

⁴This is a greatly simplified view of what is really happening, but this is 'good enough' for our purposes.

Minimum Setup

Minimally you will need the following:

- Fast-enough processor to drive the glitching
 - Normally, an FPGA is used to generate the glitch signals - usually owing to even cheap FPGA's having PLL clocking capability up to 255MHz
 - A fast XMega can also be used, however, it very much depends on the clockspeed of your target.
- Some fast-response switch - a suitable MOSFET will generally have very fast switch time, and to switch between two voltages a MAX4619 switch IC (that also comes in a handy DIP16 package)
- Misc. electronics - connector wires, breadboard, bench PSU, multimeter, etc.
- A decent Oscilloscope is a must for this work

Ideal Setup

The ideal setup for both this and the Side-Channel analysis work that is discussed elsewhere is a ChipWhisperer from NewAE Technology⁵. Either a Level 1 kit - which is ideal for this work, but also a very safe and well structured environment in which to learn - or a ChipWhisperer Lite kit make glitching targets very straightforward.

The software, written in-house in python, is very easy to use, and very easy to adapt to what you need for specific requirements. We will see Micah Scott (ScanLime) make excellent use of this later.

14.2 Theory into Practice - Worked Example with iCEStick FPGA

14.2.1 Overview

The idea of this chapter is to identify the two most important parameters of a glitch:

- *Width of a glitch* - or, how long the glitch should last. Too long, and the IC will reset. Too short and you will have no effect on its operation.

⁵Run by Colin O'Flynn and Hilary Taylor - both are wonderful people!

14.2. THEORY INTO PRACTICE - WORKED EXAMPLE WITH ICESTICK FPGA119

- *Placement of a glitch* - too early, and you'll affect the wrong part of execution. Too late, and you'll miss your window to cause the effect you want.

The process of finding these two values is what is most important in performing a successful glitch - and yet, this process is so often overlooked in the literature. We hope to offer two examples, below, that clearly elucidate some methodologies for locating these parameters, and the immediate example giving a sense of the minimal tooling required to perform a glitch.

We will follow the construction (due to the author) of the materials in this GitHub repo: <https://github.com/unprovable/glitch-stm32>⁶

There is also a tentative writeup I did on twitter as a thread: <https://twitter.com/LargeCardinal/status/1006286104939593730>

14.2.2 Setting Up

First off, you'll need the following materials for this particular hack:

- An iCEStick⁷ - this is a relatively inexpensive FPGA development board that permits you to learn verilog programming through some FOSS software
- STM32F103 development board - I usually just use the ‘blue pill’ or derivative uC development boards that feature this IC.⁸
- An STLink-V2 or similar SWD interface for the STM32F103.
- Breadboard and Cables
- A reasonably high-speed MOSFET - I’m using a 2N7000 MOSFET, that cost around €0.10
- A Soldering iron, and ideally ESD protected tweezers - these are for removing the decoupling capacitors from the board.
- An oscilloscope - yes, for this stuff, it really is quite necessary that you have access to one that has enough resolution for doing power analysis. I would recommend a picoscope⁹ if you don’t have access to a bench scope.

⁶This project by the author is itself based on another project: <https://github.com/deanjerkovich/avr-glitch-101>

⁷ Available from many retailers - here’s Farnell’s page: <https://uk.farnell.com/lattice-semiconductor/ice40hx1k-stick-evn/ice40-hx1k-icestick-eval-kit/dp/2355207>

⁸Rather than link to an eBay listing, you can find these on; eBay, AliExpress, numerous online retailers possibly local to you. You can use any board that features the STM32F103, and you don’t mind desoldering things from! :-P

⁹The 2206B is an excellent introductory model - it can be acquired here: <https://www.picotech.com>

Once you have this, you will need the following software - I usually recommend the latest version from the github, but these are not always easy to compile, so with this in mind, the code should be minimally compatible; You might have to make a few tweaks, but it should ‘just work’.

Here is the software list - note, I run Linux¹⁰:

- Ice-storm package - This *incredible* package has been spearheaded by the excellent Clifford Wolf and Mathias Lasser¹¹. You can download the software here: <http://www.clifford.at/icestorm/>
- The `arm-none-eabi` toolchain for compiling the ARM firmware if you want to make changes.
- An installation of OpenOCD - see Section ?? for details on installation and testing.

¹⁰For windows support, please install a linux VM on your PC.

¹¹It is very hard to put into a small paragraph how important their work is - but in short, thanks to them, you don’t need to register nor pay for the ability to generate bitstreams for the ICE40 FPGAs, and they are working hard on doing the same for Xilinx and other FPGAs. This is true Open Sourcing, and hails the advent of openly available FPGA platforms and software. It really is a game-changing feat of reverse engineering.

14.2. THEORY INTO PRACTICE - WORKED EXAMPLE WITH ICESTICK FPGA121

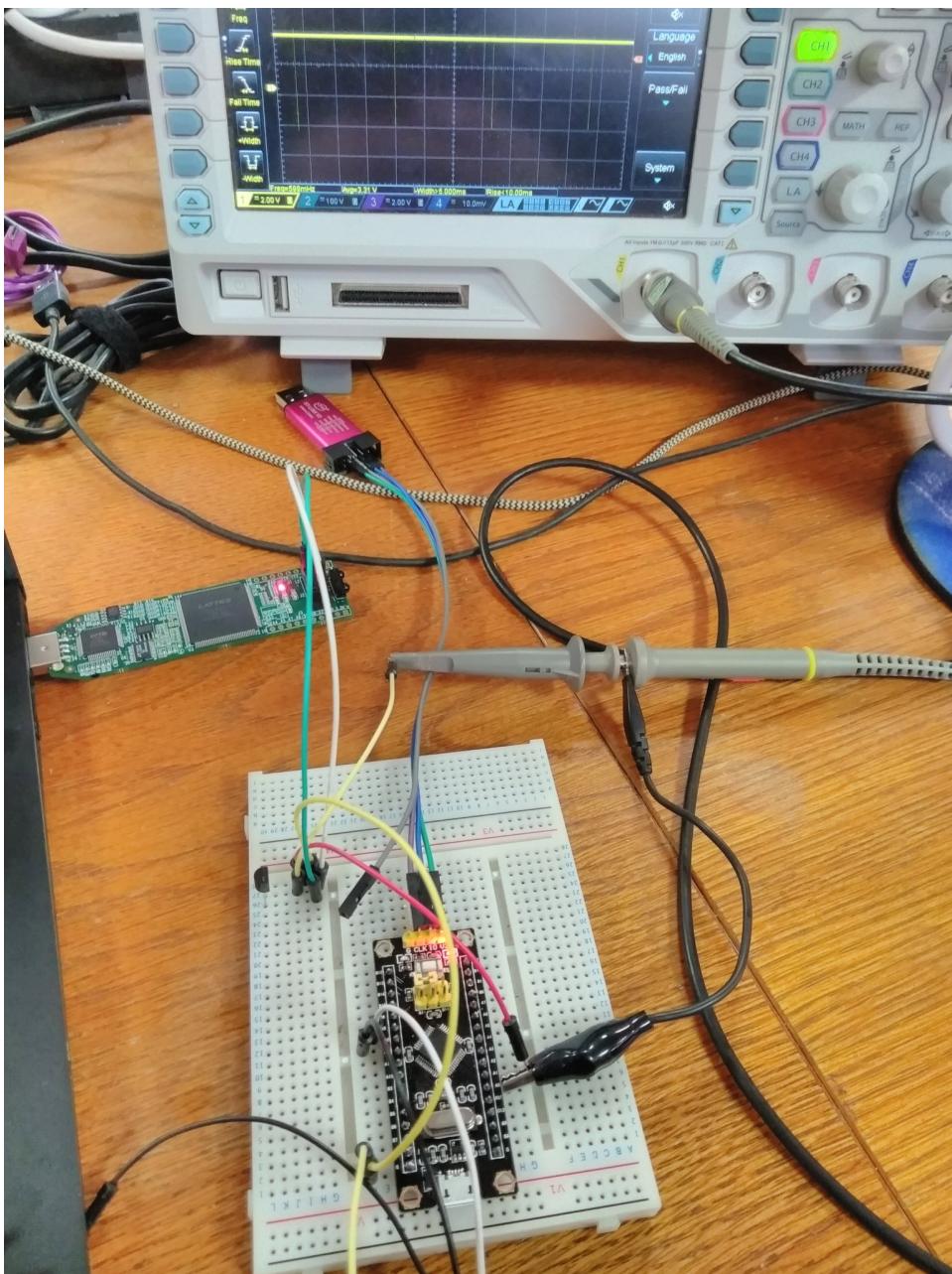


Figure 14.1: Overview of the setup

We will assume you have these materials to hand - here's the rough setup procedure:

1. Setup and load the firmware onto the STM32F103. The rough setup should be something along the lines of:

(a) using OpenOCD, the command should be something like this:

```
|| sudo openocd -f target/stm32f1x.cfg -f interface/stlink-
  v2.cfg -c "init" -c "reset halt" -c "stm32f1x unlock
  0" -c "flash write_image erase unlock /path/to/
  firmware.bin 0x08000000" -c "exit"
```

2. Edit the hardware and wire up:

- Remove the power decoupling capacitors - not all the caps need to come off, but the decoupling ones should be removed. See fig ??.

3. Wire up the uC board and the FPGA s.t. the following wiring is followed:

- The board should be powered from the UART to USB interface (CP2102 boards usually have a low power output) or bench PSU - the SWD connections can now be removed after programming.
- The ground should now be wired to the drain of the MOSFET¹²
- the source of the MOSFET should go to reference ground¹³
- The gate should be wired to the FPGA pin we will be using - this is pin 78 on the board, or pin 1_02 on this diagram: <http://www.pighixxx.net/2016/02/icestick-pinout/>

¹²Always check datasheets: <https://www.onsemi.com/pub/Collateral/2N7000-D.PDF>

¹³Note that this means the FPGA has to be programmed and powered for the uC to work. We will cover this shortly.

14.2. THEORY INTO PRACTICE - WORKED EXAMPLE WITH ICESTICK FPGA123

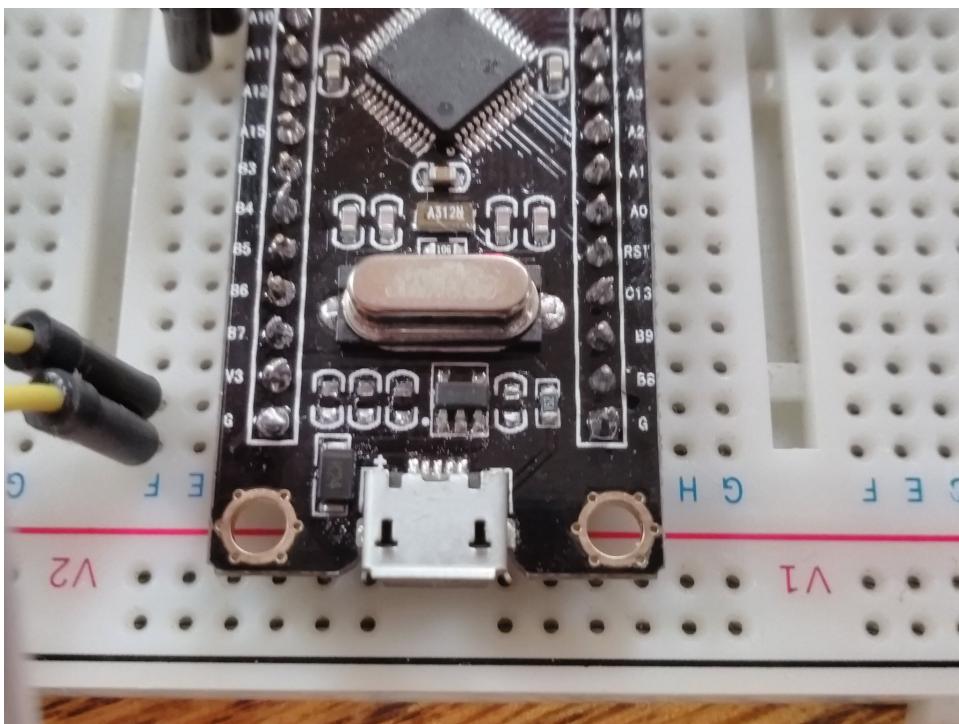


Figure 14.2: Powerline decoupling capacitors removed off a ‘blue-pill’ like dev board

Assuming that we have set things up correctly when the device is powered we should be able to connect to the UART pins A10 and A11 and see UART traffic. We are now ready to program the FPGA. The rest of this section will be setting up the FPGA so that we can control the glitch through UART.

14.2.3 Programming the FPGA for glitching the STM32F103

We will start off with the basic example code from <https://github.com/deanjerkovich/avr-glitch-101> - here is the verilog in full:

```
'define GLITCH_DURATION 36

module top (
    input  clk,
    output LED1,
    output LED2,
    output LED3,
    output LED4,
    output LED5,
    output J1_1
);
```

```

assign LED1 = outreg;
assign LED2 = ~outreg;
assign J1_1 = outreg;

reg outreg;
reg [31:0] counter;

initial begin
    counter <= 0;
    outreg <= 0;
end

always @ (posedge clk) begin
if (counter == 20000000 && outreg) begin
    counter <= 0;
    outreg <= 0;
end
else if (counter == `GLITCH_DURATION && !outreg) begin
    counter <= 0;
    outreg <= 1;
end
else begin
    counter <= counter +1;
end
end

endmodule

```

If we load this onto the FPGA and ensure the GND is controlled through the FPGA, then this code should generate small glitches visible on the oscilloscope. However, in order to change the glitch duration, we need to recompile and reload the bitstream. This isn't optimal for testing and working out the best glitch length timing.

So, let's implement UART. Rather than do this by hand in verilog, we'll use the following library that implements it on the icestick: <https://github.com/cyrozap/icestick-UART-Demo> - we've included the library in the repo. This allows us to communicate with the FPGA over UART and give it bytes of data. Each byte will do two things:

- Indicate that a glitch is wanted, and set the bit to true
- The byte value will be multiplied by 256 and then this value will be used as our glitch length timing

So, we'll take this data and set the PLL to be the same freq as the target processor¹⁴. So, here is the verilog for this, with comments:

¹⁴Using the PLL on an FPGA, we can use the actual freq or some subdivision of it - doubling it or more is useful for clock fault injection, but we don't need that granularity here.

14.2. THEORY INTO PRACTICE - WORKED EXAMPLE WITH ICESTICK FPGA125

```

always @(posedge pll_clk) begin
    // increment clock...
    counter <= counter +1;
    // check if glitch_signal is on, and if yes,
    // reset counter and set glitch flag
    if (glitch_signal && !glitch) begin
        // turns out we don't know the interface register
        // uart_var1... take rx_byte and multiply it by 256
        - this
        // gives us roughly +3.5 microsecond increase in
        // glitch time
        // for every +1 increase in UART byte value... ish
        ...
        var1 <= rx_byte << 8;
        glitch <= 1;
        counter <= 0;
    end else if (glitch) begin
        // if we are in glitch, set outreg to 0
        outreg <= 0;
    end
    // if we are in a glitch (outreg is 0) and we reach the
    // desired
    // duration, then unset glitch and put outreg to 1
    if (counter == var1 && !outreg) begin
        outreg <= 1;
        glitch <= 0;
    end
// end always
end

```

Figure 14.3: Example glitch output on the left, with input from python script on the right

With this now done you should see the kind of output visible in fig ?? - with some experimentation, it should be clear that for an input of `0x7A` (ASCII ‘z’) the device will reset, whilst an input of `0x61` (ASCII ‘a’) will cause intermittent glitches, whilst an input of `0x41` (ASCII ‘A’) the device will not reset nor glitch.

14.2.4 Overview

Hopefully this section has been instructive on how to generate glitches with a reasonably cheap and easy to acquire FPGA - this glitch example is fairly straightforward, and this example shows how to break out of control loops in uC software.

In our above example, although our control of the width of a glitch is reasonably good, it isn't too accurate for controlling the 'when' of our glitch - that is, we are dependant on the speed of our UART to define the resolution of slots within which we can have a glitch.

For much finer control - and built-in support for doing powerline analysis - a ChipWhisperer is worth the investment if you are serious about pursuing doing these attacks on a wide range of targets.

14.3 An Introduction to the ChipWhisperer

The ChipWhisperer is the brainchild of Colin O'Flynn of NewAE Technology - and his work has made glitching and SCA attacks as straightforward as plugging in a logic analyser. The tooling is designed to work perfectly hand-in-hand with the python software used to both perform measurements for SCA as well as glitching for attacks. The user interface also includes such useful facets as a glitch explorer (which we will see later in scanlime's work).

For those who prefer a video introduction - an overview of the ChipWhisperer's capabilities, features, and workings can be found in this video from Scanlime, starting around 5mins in: <https://www.youtube.com/watch?v=TeCQatNcF20>

14.3.1 Settting up the ChipWhisperer

[The ChipWhisperer Lite](#)

[The ChipWhisperer Level 1 kit](#)

[The ChipWhisperer Nano](#)

14.3.2 ChipWhisperer's First Glitch

14.4 Real-World Examples

As with the chapter on SCA¹⁵ it is all to easy to consider glitch attacks as 'too academic' and 'not real world' attacks. However, as we hope to demonstrate, these are very much real-world and applicable attacks, that are feasible against targets found in products.

To show this, we will utilise two excellent examples¹⁶ of glitch attacks:

¹⁵Side-Channel Analysis

¹⁶Also, two of my personal favourite examples...

1. Chris Gerlinsky's glitch attack to bypass CRP on LPC1343
2. Scanlime's glitch attack to extract firmware over USB off a Wacom graphics tablet

14.4.1 Gerlinsky's CRP Bypass on LPC1343

This section is part is based on Chris Gerlinsky's hack of the LPC1343 CRP. Resources can be found in his talk <https://www.youtube.com/watch?v=98eqp4WmHoQ> and his slides.¹⁷ There is also a github for the Xmega code - found at <https://github.com/akacastor/xplain-glitcher>

Scoping the Target

We want to be able to pull a firmware off an LPC1343 IC that has CRP enabled. This should not be a possible task, but Gerlinsky found a bug in the implementation of the bootloader's CRP control.

Looking at the datasheet, we see that there are 4 settings for CRP:

1. `NO_ISP` - Using hex `0x4E697379`, this disables sampling on ISP pins
2. `CRP1` - Using hex `0x12345678`, this mode limits flash updates and reads
3. `CRP2` - Using hex `0x87654321`, this mode prevents flash and RAM reads and bootloader
4. `CRP3` - Using hex `0x43218765`, this mode disables debug and fully locks down the flash, and does all that CRP2 does.

However, note that these settings only declare for 4 specific values. Ergo, for any of the other 4,294,967,292 possible 32-bit words, the system will be open. So, if we can get *literally any other value* to be read as the CRP mode value, then we will have no CRP. Nada. Nowt. Nicht. We will have bypassed the flash protections.

ProTip: This is a good rule of thumb indicator that a chip will be vulnerable to a glitch - where there is only a strict subset of a larger search space that defines a restriction or defined behaviour, we can cause the chip to fall back on its default mode (usually fail-open systems) in order to gain the control we want.

Now we have our target, we need to do some analysis - for this, Gerlinsky saw that the CRP value is actually *written into memory only once on boot* and referred to from there. Quite why, we may never know. But the upside is that we don't need to glitch at the point of requesting flash contents, we only need to glitch for when the bootloader is reading those values and then passes execution on to the firmware entry point.

¹⁷https://recon.cx/2017/brussels/resources/slides/RECON-BRX-2017-Breaking_CRP_on_NXP_LPC_Microcontrollers_slides.pdf

The last thing that we need is an indicator to know when we have successfully glitched the chip. We know that the chip has two modes - bootloader mode and regular execution. We also know that, under CRP, bootloader is restricted.

To do this, Gerlinsky put a 10Ω resistor in series with GND¹⁸, and read the voltage across it. As we will see in the chap ??, we can use the voltage across this and its known resistance to work out the current that is going across it. This current corresponds to the power consumption of the microcontroller.

So, if we see 3V across a 10Ω resistor, then we know that $3/10 = 0.3$ Amps are flowing through the circuit. Gerlinsky was able to find identifiable traces using this power analysis to know whether the chip has gone into normal execution or into bootloader mode.

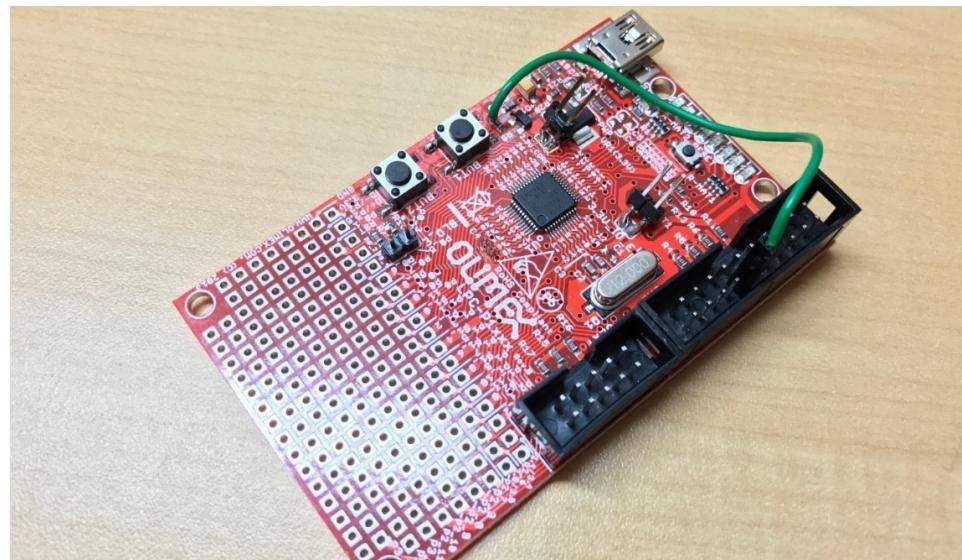


Figure 14.4: Gerlinsky's Setup for glitching the LPC1343

Setting up the Glitch

We are now ready to tool the actual glitch.

Gerlinsky opted for a level shift glitch to cause a brownout on the LPC1343. To do this, you switch between the operating voltage of, say, 3.3V, and a lower voltage, here 1.2V. This causes a *brownout*, as opposed to a full blackout of the chip. This is a good glitching option to keep in mind when setting

¹⁸We use the GND measurements for power analysis, as opposed to VCC measurements, in order to avoid ground loops. This is another reason why it is helpful to remove decoupling capacitors, esp from ground. For more information, see the section in Chap ??.

up glitches - if the target gets too unstable under a blackout, consider a brownout - they are much more forgiving.

To inject our glitch, Gerlinsky made the modifications seen in fig ??, similar to the modifications we saw previously in ?? except that his are on the positive side of the board¹⁹. This gives us control of the power for the target.

In order to switch between the two voltages (supplied by a bench PSU) Gerlinsky made use of the MAX4619 chip - a high-speed switching IC with some good characteristics. By hooking up the two power supply levels to the MAX4619 chip, and then using a glitch trigger on the switch the on/off pin.

Glitching the Target

To actually setup the glitch, Gerlinsky used an XMega-A1 on an Xplained development board. This uC runs at 32MHz and he wrote some custom ASM to allow for single-cycle delays. This gave him enough control over the chip to be able to break out of the CRP.

In order to work out the width of a glitch, Gerlinsky used the following test code on an LPC1343:

```

||| while(1)
||| {
|||     b = NUM_GLITCH_LOOPS
|||     // Wait for synchronise with glitcher
|||     while( GPIO_PORT0 & (1<<13))
|||         ;
|||     // Wait for PIO_13 to go high
|||     while( !(GPIO_PORT0 & (1<<13)))
|||         ;
|||     // TIME FOR A GLITCH! :D
|||     // Pull pin low for time in for loop
|||     GPIO_CLRPO |= (1<<17)
|||     for( int a=0; a<NUM_GLITCH_LOOPS; a++) {b--;}
|||     // Pull pin high again
|||     GPIO_SETP0 |= (1<<17)
|||     // test a and b... should be the same...
|||     if (a == NUM_GLITCH_LOOPS && b == 0 )
|||         pp_printf(",");
|||     // if not, we had a glitch!! :D Tell the user...
|||     else
|||         pp_printf("\n! a: %08X b: %08X\n", a, b);
||| }
|||

```

This code allowed Gerlinsky to scope out the appropriate width of a glitch on the target by looking for glitch values that would give the output values of a and b, when they didn't meet. It would also allow Gerlinsky to see the severity of different values, and to see if the target reset.

¹⁹One of the many reasons I included this example.

```
chris@archer:~$ cortex-m3/tools/dump 0 32768 > LPC1343_FLASH.bin
Opening serial port /dev/ttyACM0
Forcing boot loader mode
Synchronizing... done
Identifying... done
part number: 3d00002b
LPC1343, 32kB Flash, 8kB RAM
.....
```

Figure 14.5: Successful Glitch of LPC1343 to bypass CRP

As can be seen in Fig ??, Gerlinsky was able to successfully glitch the LPC1343 by using these methods.

Some important takeaways are as follows:

- You don't need the clock speed of the glitcher to exceed the clock speed of your target - you just need enough granularity in order to successfully cause a glitch.
- You can glitch a brownout between voltages, not just blackouts
- Glitches work well on targets where they can be made to fail-open or 'miss' a specific condition that is restrictive.
- Glitching is as much about the analysis and setup as it is about the precise timings and outcome - Gerlinsky's presentation shows how you can solve the problems of identifying the width and placement of a glitch on a real target.

The Aftermath

It would seem that NXP have noticed this, possibly before Gerlinsky did this work. The new USB bootloader was found by Gerlinsky to contain an explicit `NOCRP` value, in hex it is `0x11223344`, added to the list of CRP modes. This would mitigate this particular glitch, but it is an open question whether there are other opportunities for CRP bypass.

14.4.2 Scanlime's USB Glitch on Wacom Graphics Tablet

Micah Scott, a.k.a. scanlime²⁰, documented quite a whirlwind journey towards dumping the firmware off a Wacom Graphics tablet she had access

²⁰<https://twitter.com/scanlime>

to. As such, scanlime was able to play with the hardware and create some really cool tooling that permitted the chip to be glitched very precisely by a chipwhisperer by creating a control board add on to the ChipWhisperer, the FaceWhisperer²¹, which would permit the ChipWhisperer to reset the USB device, and then time glitches precisely during the USB loading sequence.²²

Having done a full teardown and reverse engineering of the tablet²³, the aim was to see if the tablet could be used as an RFID reader, given that the operating frequency of the tablet's sensors was within the range of some RFID cards.

The FaceWhisperer designs and code can be found here: <https://github.com/scanlime/facewhisperer> and a video explaining this attack can be found on scanlime's YouTube channel here: <https://www.youtube.com/watch?v=TeCQatNcF20>.

Scoping the Target

First analysis of the IC on the Wacom target was that it was a LC871A00 series 8-bit microcontroller. however, it was not possible to read the firmware off the IC (see the video for details on this).

Scanlime found that the USB descriptors referenced in the initial USB communication give data such as version numbers, which are really just strings in the firmware image, pointers for which are dereferenced at the time the USB handshake occurs.

The idea for the glitch is to 'trick' the uC into dumping details or chunks of the firmware from the flash by glitching this dereference.²⁴

After making various hardware modifications, such as:

- Removing the 12MHz clock crystal and replacing it with our own clock signal
 - By using a signal generator and an oscilloscope, we can determine the input and output pins for the crystal by observing the impedance of the circuit under this injection - the one that doesn't 'fight back' is the input pin.
 - Once we have identified this, scanlime added a small 50Ω resistor in series with the input signal line in order to reduce reflections²⁵

²¹Name inspired by Travis Goodspeed's excellent FaceDancer USB hacking device - which is a subject we will cover elsewhere.

²²Although the analysis of this USB traffic is beyond the scope of this section, it will be covered elsewhere. [TODO: Appendix B???

²³An excellent account of which can be found also on scanlime's YouTube channel: <https://www.youtube.com/watch?v=nPab7pbOhBY>

²⁴As we shall see, however, the full firmware ended up being dumped in this process, straight out of the USB port!

²⁵Both of these tricks are really cool, and worth experimenting with yourself for auto-didactic completeness.

- Removing the decoupling capacitors immediately around the uC - much as we did in our example above
- Adding a stack of resistors to act as a current shunt resistance.
 - Some theory that is important for this - the resistance of resistors in parallel is given by $\frac{1}{R_t} = \frac{1}{R_1} + \frac{1}{R_2} + \dots + \frac{1}{R_n}$ - so, if you have, say, a pile of 50Ω resistors, then for each one you add in parallel (in the stack, in scanlime's video) you decrease the total resistance exponentially. As such, 4x 50Ω resistors in a stack gives an overall resistance of $\frac{1}{R_{total}} = \frac{1}{0.02+0.02+0.02+0.02} = 12.5\Omega$ (using the solved version of the equation). These very small resistances are most useful for current sensing, and this is an easy way to achieve them!
 - Some more capacitors were added on the non-measurement side of this resistor stack in order to make the power going into the target as clean as possible - much like is done on the ChipWhisperer target boards. The same was done by the power regulators for the board.
- A header was added in order to breakout the reset and debug pins such that we can both see when the target is at different stages of booting, and trigger a reset as and when we need to.

The mechanics of these modifications may look a bit gnarly, but they are actually very straightforward and useful for the attacks we can consider against the target.

Adding an Interface - the FaceWhisperer

The FaceWhisperer, although exclusively designed for this project, is actually a nice, extensible board that permits integration of ChipWhisperer control signals into the USB stack on a device. The device, in all its glory, can be seen in figure ??.

It features clock output pins for synchronisation, and a USB plug that connects to an on-board MAX3241E. This nifty little chip lets you control USB via SPI from a microcontroller. Given pretty much any uC supports SPI, you can add USB support to more-or-less any IC you desire. In this case, we can control the USB stack for the target via an XMEGA128D.

Ostensibly, this device works the same as the facedancer but replacing the MSP430 with the aforementioned XMEGA128D. This is then programmed with the USB `GET_DESCRIPTOR` command and some code for ensuring synchronisation of the clocks.

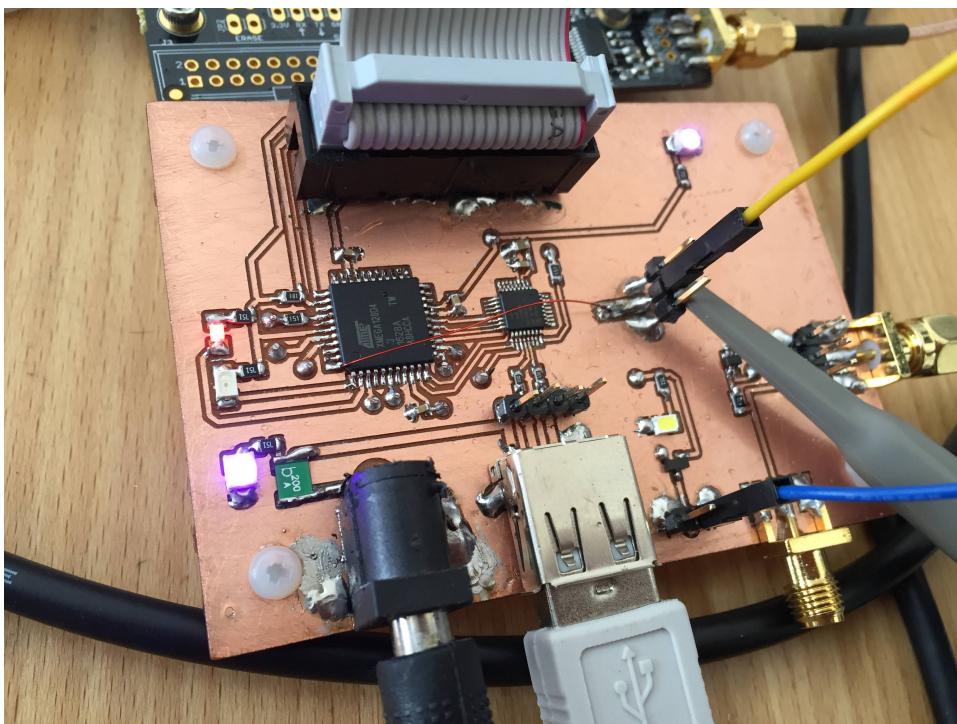


Figure 14.6: The facewhisperer add on for ChipWhisperer, from scanlime’s github

Glitch Exploring

Now that the target is setup and the facewhisperer is programmed and plugged in, we are ready to explore the glitches! Normally for this, you would use the Glitch Explorer window (see fig. ??) - in this, you can select what ‘normal’ responses are, and define any specific ‘error’ responses you are looking for.²⁶

However, scanlime opts to script the glitchexplorer’s functionality, but follows the same logic - vary the glitch timings and look for any USB traffic that is longer than 34 bytes - the length of the standard response from the `GET_DESCRIPTOR` USB request to the target.

The video shows the process in fantastic details, but a summary of the findings are as follows:

- By analysing the measurements from the target, we can identify rough timings for when the USB replies to requests

²⁶This is very similar in functionality to Burp Suite’s Intruder window, for those who are familiar with that.

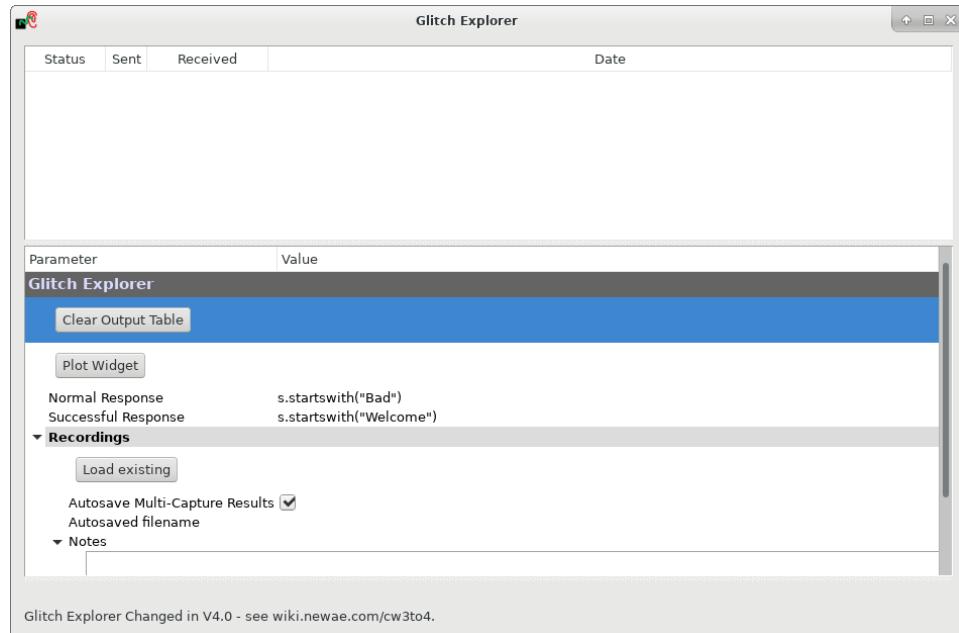


Figure 14.7: The normal glitchexplorer window

- The USB request is set so that the control request allows for up to 64kB of response from the device (even though this is usually limited by the on-board USB stack, it doesn't hurt)
 - **NB** - this is part of the core of the concept of this glitch; we hope that the reading of this value might overwrite the actual read length used by the chip when it fetches values.
 - This is a core part of the glitch - to setup the USB traffic as much as possible to help us with the glitch.
- The glitch explorer is then parameterized with the timings from analysis of the USB traffic and asked to scan the ranges of timings from a defined offset (the start fo the USB call/response)

In the end, it was found that the device will gladly spit out the firmware and/or memory dump in glitched responses to the USB `GET_DESCRIPTOR` requests. By analysis on these, it was found that the full firmware had been dumped via this method.

Target acquired - objective achieved!

The rest of the task is piecing together the firmware file from the dump, but this is quite straightforward. The dump usually begins in the data segment, where strings and other artefacts are stored. From this assumption,

it is straightforward to then locate the end/beginning of memory, and from here, we can rearrange the file into a proper firmware image for the target.

Overview of the Attack

This attack is an excellent example of how you can create custom tooling for a target, and then explore the glitches you can generate, and the effect this process can have on a target. Hopefully this has served as a good example for how careful thinking and planning can lead some very nice tooling and some incredible results!

14.5 Further Reading

This is a huge subject, and we hope that this has given some help in getting going! Here are some more resources for fault injection attacks:

- Anything by Skorobogatov at Cambridge Computer Lab: https://www.cl.cam.ac.uk/~sps32/ECRYPT2011_1.pdf

Chapter 15

Side-Channel Attacks

Learn how to break
cryptography on hardware
targets by just listening really,
really closely...

15.1 Side-Channel Attack Overview

A pre-requisite for this chapter is the sections on AES and RSA in ??, and/or a reasonably cogent understanding of their basic mechanics.

15.1.1 What is a Side Channel Analysis Attack?

Put bluntly, it's an attack that makes use of analysing a side channel.

FFS, Mark...

What? :P But more seriously, they are attacks concerning small amounts of information leakage from a system in order to 'fill in some blanks' - enough that we can make a logical inference with some degree of certainty about the internal mechanics of an otherwise black box device. A 'side channel' is just a 'channel' of information, like any other, but one which is adjacent to, and not explicitly defined in the operating design/software outputs. What do I mean? Well, here's an example from folklore:

A windowless room contains three identical incandescent lightbulbs in fittings, also identical. Each light is connected to one of three switches outside of the room. Each bulb is switched off at present. You are outside the room, and the door is closed. Before opening the door you may play around with the light switches as many times as you like. But once you've opened the door, you may no longer touch a switch. After this, you go into the room and examine the lights. ***How can you tell which switch goes to which light?***

Have a think about it before moving on. Ok, ready?

So, we get a one-shot chance to work out which bulb is controlled by which switch. How can we do it? Easy... Let's assume the switches are numbered 1, 2, and 3.

1. Set switches 1 and 2 to 'on' and wait a few minutes.
2. Turn off switch 2 and then open the door. We now assess the outcomes:
 - (a) The bulb that is on is controlled by switch 1
 - (b) The bulb that is off but warm is controlled by switch 2
 - (c) The bulb that is cold is controlled by switch 3

It's obvious when you see it, and assumes a world pre-energy saving light-bulbs, but the side channel here is the bulb getting hot. It isn't its intended purpose to get hot - it just does. But the setup permits us to enumerate this behaviour into our logic and bypass the 'control' that is "you can't use the switches once the door is opened". We use some thermodynamics in conjunction with Ohm's Law to work out which bulb is which in one sitting.¹

15.1.2 UART Pin Impedance as semi-side channel

Recall the quick'n'dirty, logic analyser-free method we described earlier for identifying UART pins with only a standard multimeter - assuming 3.3V operating voltage:

- The GND pin is the one that has continuity to ground.
- The VCC is the pin at full 3.3V.
- The Tx is the pin that is at ~3.2V - owing to the pull up resistor, the voltage is usually lower than VCC - and that, when in operation (say, during boot) fluctuates as the multimeter calculates RMS.
- The Rx is the pin that is left.

In a similar way, we shall see unintended outcomes regarding information disclosure that arise naturally from the operation of circuits.

15.1.3 SCA attacks in Web apps

To further show that you may already be familiar with SCA's - consider this attack. Suppose we have a login form with two outputs:

- Error of "Incorrect Username"

¹In fact, using Ohm's Law and some other property is a pretty common method for disclosing data from within a blackbox.

- Error of "Incorrect Password"
- Logging the user in

Now, this is a standard ‘user enumeration’ vulnerability we see all the time. But suppose the error messages are unified to be “Incorrect username or password” - it’s fixed right? Well, maybe not.

Suppose the same site now upgrades to using `bcrypt` or some other strong password hashing algorithm, with salting and all the good stuff. They’re secure, right? Well, in some (admittedly quite rare) circumstances, this can cause a time-based username enumeration. It would look as follows:

- Error of “Incorrect Username or Password” with a page return delay of 5ms for a known-invalid username
- Error of “Incorrect Username or Password” with a page return delay of 15ms for a known-valid username
- Logging the user in

Clearly, with some further testing, we can confirm whether the site has a username enumeration vulnerability using page timings.

But how can we fix this?

The fix for this is simple - always compute the password hash, and *then* check if the username and password combination is valid. This example will be useful in understanding the attack below in sec ??.

15.1.4 Types of Side Channel

As such, we can generalise that there are fundamentally 4 main types of side-channel:

1. **Time** - variance in timings can reveal information about the background processes
2. **Heat** - Heat changes and thermal variance can determine the amount of work something is doing
3. **Sound** - Similar to heat, sounds can usually indicate or fingerprint ongoing hidden activity²
4. **Power** - the amount of power a target is consuming directly correlates to its activity, and the values of its parameters (as Hamming weights - we will see this later).

²The RSA attack below was performed against a target from across a room with a microphone listening to the *coil whine* of the heatsink expanding and contracting whilst doing a computation, by a team in Tel Aviv. Amazing work! See their paper for more: <https://www.cs.tau.ac.il/~tromer/papers/acoustic-20131218.pdf>

The most useful for what we want to do - break cryptography - is power analysis. So, from here on in, we will focus almost entirely on this side channel, and demonstrate the various methods for analysis and assessment that are most pertinent in the literature.

15.2 Types of Power Analysis

So, let's take a look at the three types of Side Channel Analysis attacks:

- Simple Power Analysis
- Differential Power Analysis
- Correlation Power Analysis

We'll cover how these work and what part of the encryption schemes they target. But first, let's look at the physics of what is going on.

15.2.1 Power Analysis - an Introduction

For these examples we will add a current shunt³ resistor to the GND pins of our targets. The idea is that the introduced resistance is very low - around 10Ω or less - so that it doesn't affect the operating impedance of the circuit too much, but enough that, by taking measurement of the voltage across the resistor, we can know how much current is flowing across the resistor - knowing that $V = I \times R$ we get that $I = \frac{V}{R}$.

This gives us an idea of the power consumption of the target device. But how/why is this useful? Well, there are two facts we will use most later:

- Operations and calculations all use varying amounts of power.
- The Hamming Weight of a bitstring correlates to the amount of power required to represent that bitstring - a string of bits '11101010' requires more power than '10000010'⁴

As in our examples above, we will use this information to intuit what is going on when we know that the chip is doing a specified operation.

Note that in general - we are interested in using SCA to break *cryptography*, so that is what our examples will focus on. You can use this technique in other ways, but this is the most useful and most powerful.

³Also called a 'current sense' resistor.

⁴Note that this string '10000010' has the same weight as '11000000', so power consumption to Hamming weight is not indicative of value, but of the amount of power being consumed. As such, this isn't the whole story, so we have to do enough work otherwise to work the rest out.

15.3 Simple Power Analysis Attack on RSA

Here we will present how SCA can be used against the RSA cryptosystem.

15.3.1 RSA Decryption - a reminder

Recall how RSA decryption works - we take our private key d , the modulus value n , and the ciphertext c generated by the person sending us a message (they having made use of our public key, of course). We then know that

$$c^d \equiv (p^e)^d \equiv p \pmod{n}$$

where p is our plaintext message.

So what good does a side channel analysis of this process do us? Well, that comes down to the fact that ‘exponents are hard to compute’...

15.3.2 Computing Exponents - the Square-Multiply Method

Exponents - or, taking some number n to the power k , usually written n^k - generate huge numbers, and are generally considered a hard computational task. However, there is an efficient workaround called the ‘square-multiply method’ (SMM) that will let us perform exponents in a breeze! Here’s some pseudocode of the process:

```

|| power(x, e, n):
||   int z=x;
||   for(int i=1; i< bitlen(e); i++) {
||     int b = (e & ( 1 << i )) >> i;
||     z = (z*z) % n;
||     if (b == 1) {
||       z = (z*x) % n;
||     }
||   }
||   return z;

```

Just to make it explicit, here’s a verbal description of the process:

1. Take an input x , and exponent e and modulus n .
2. For each bit of e , square the current result, starting with x , take the modulus.
3. If the current bit of e is 1, also do a multiplication with the original x .

As you can see, each bit of the exponent, which is of course our private key, guides the execution process. We will use the fact that the traces for a square-only operation are very different from a square-then-multiply operation. As such, monitoring the power consumption we can derive a sequence:

sqr	sqr-mult	sqr-mult	sqr	sqr	sqr-mult	sqr	sqr-mult
0	1	1	0	0	1	0	1

...and on the second row we can easily derive the bits of the private key, and we can do so exactly - we don't need to guess when there is a one or zero, as both start with the same squaring operation. If we see 'square square square' we know the first two bits are 00, and only need to see if the fourth trace shows a square (third bit is now a '0') or a multiply (third bit is now a '1').

Noooo - it can't be that simple?!?

Well, usually it isn't *quite* so simple, but in its most basic form, this is a 'simple power analysis' attack against RSA, and this is a practical attack you can code up on an arduino.⁵ For more information see this paper:

<https://onlinelibrary.wiley.com/doi/pdf/10.4218/etrij.10.0109.0249>

15.3.3 Can we protect against this?

Protections against this attack can be quite trivial so long as the attack is this basic; always do a square multiply - just use the value of the bits of e to decide whether you assign that value or not. Other attacks don't quite scale up in their resilience to this attack.

15.4 Differential Power Analysis Attack on AES

This is all very well and good against RSA - the attack requires both knowledge of the specific exponentiation algorithm and no obfuscation. The truth of the Real World™ is that many vendors will supply '**MILITARY GRADE ENCRYPTION**' to their customers withing their devices.

That's right, the kind of encryption good enough to secure military applications. However, there are sometimes a few snags with this - as has been documented by Travis Goodspeed⁶, AES-128 in CTR mode⁷ running on an 8-bit hardware counter is effectively just XOR with a 512-byte key⁸

However, what if we are met with a full-blooded AES-128 in CBC mode? Or AES-256 in CBC mode? Is there anything we can do against block ciphers?

As it happens; 'yes'

⁵[TODO - totes write up this example, Mark.]

⁶<https://www.youtube.com/watch?v=GX8-K4TssjY>

⁷For an explanation of modes, please see the associated section in chapter ??

⁸We will cover this and other hardware-specific crypto snafu's later on in chap ??.

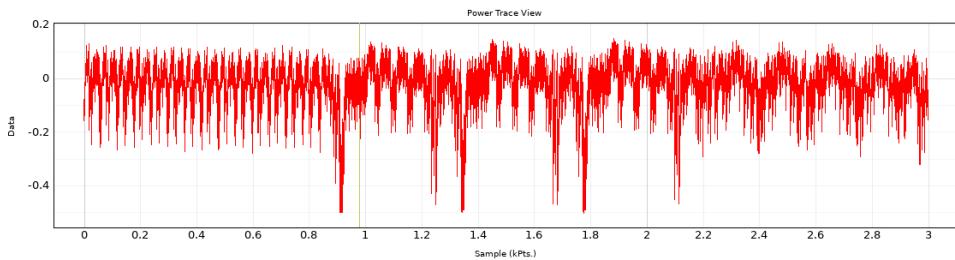


Figure 15.1: Example Power Trace of an AES128 encryption - can you tell why it is AES128?

15.4.1 SCA against AES - Background Theory

The theoretical background for this section comes from a few papers, but the most useful is probably Owen Lo, Buchanan, Carson (2016) <https://www.tandfonline.com/doi/full/10.1080/23742917.2016.1231523>

So, here is the basic idea:

- We know that the power consumption will vary between a sequence of random inputs and encryption.
- We monitor the power consumption of the target whilst recording the cryptographic input and output ciphertext.
- If we don't know it already, use Simple Power Analysis to identify key lengths - see fig ???
 - If we see 10 rounds, we know that the key length is 128
 - If we see 12 rounds, we know that the key length is 196
 - If we see 14 rounds, we know that the key length is 256
- We collect as many samples as we can - 50 is a minimum, but up more will increase accuracy. Somewhere between 50 and 250 is usually sufficient.

Once we have the power traces we are ready to move onto the analysis phase of the attack.

15.4.2 Differential Power Analysis - in Practice

We will now discuss how to perform a differential power analysis attack. It relies only on knowing the traces and the output those traces correlated to.

Background Theory

The idea of this attack is essentially as follows:

1. Guess both the plaintext and the key.
2. For each guess of both, compute at least one of the round keys. and compare LSBs of these guesses with the known LSBs of the traces
 - We construct a selection function that will do this comparison for us.
 - We compare the LSBs recorded, and then sort the traces into buckets according to whether the LSB would be 0 or 1 for that guess
3. We take the averages over these buckets and look for spikes - these indicate that we have found a likely candidate for the key!

But.... WHY?!? I hear you ask? Well, this is because of a fact about the sorting. We're making guesses, and then calculating the outcomes, and sorting our traces appropriately.

Remember, cryptography makes our data reasonably indistinguishable from noise. As such, if we made an incorrect assumption, we are sorting 'noise' by random so the output will be... well... more noise. Taking the differences of a set of traces will simply yield no discernible difference in power consumption, as we've essentially sorted the traces by random.

However, if we guess the key correctly, then we will have caused the traces to have been sorted according to some correlation of power consumption. *Why?* Well, this is because, if you recall, the power traces correspond exactly to the amount of power required to perform that computation. As such, we have made the correct guess, we will have aligned some part of each of the power traces, meaning that when we take the average *some part of the traces will line up and give us a non-zero result!*

And this, my friends, is how the magic happens...

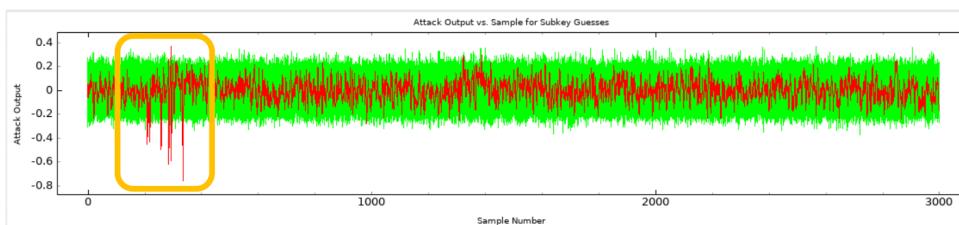


Figure 15.2: Example Differential Trace Analysis - you see the spike indicating we make a correct guess, highlighted on the left.

Putting it into Code

So, this is all very well and good - but what does it look like in code?! Well, something like an extension of this pseudocode:

```

for KeyGuess 0 to 255:
    for plainText 0 to 255:
        sBox_lookup = s[plainText XOR keyGuess]
        if (LastSignificantBit(sBox_lookup) == 1):
            Add trace relative to plainText input to subset1
        else if (LastSignificantBit(sBox_lookup) == 0):
            Add trace relative to plainText input to subset0
        getAverage(subset0)
        getAverage(subset1)
        calculatePointToPointDifference(subset0, subset1)
    
```

As you can see, we use the sBox lookup function from AES in order to work out what the SBox lookup of our key would be - given the key schedule is completely deterministic, it doesn't really matter which key we get, we can always recover the original key.

However, there is a better way of achieving this - by means of a Correlation Power Analysis...

15.4.3 Correlation Power Analysis Attack on AES

This is an attack that is very much the same as the Differential Power Analysis, except with a few key differences - the main ones being the use of Hamming Weights for the SBox lookups, followed by an assessment of the sorting that takes into account the overall power profile of the target; the Pearson Correlation Coefficient.⁹

⁹The Pearson Correlation Coefficient is given by

$$\rho_{X,Y} = \frac{cov(X, Y)}{\sigma_X \sigma_Y}$$

where $cov(X, Y)$ is the covariance, and σ_X is the standard deviation of X.

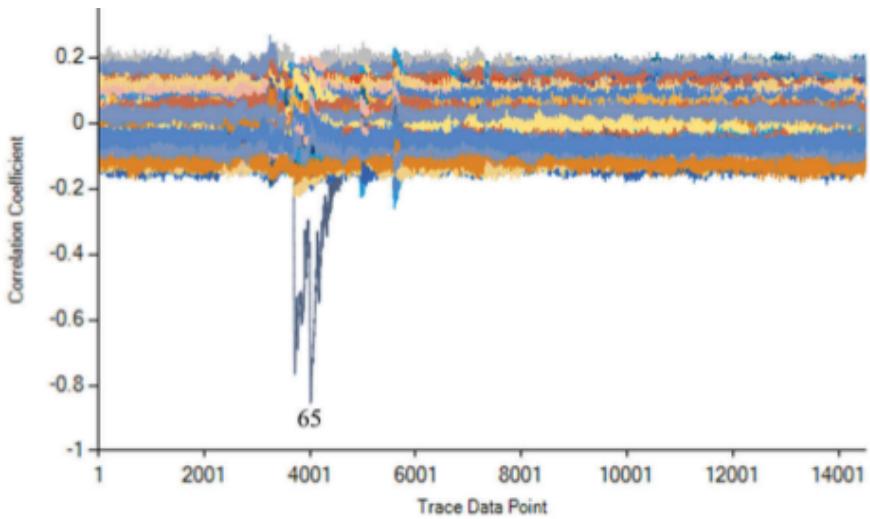


Figure 15.3: Example Differential Trace Analysis - you see the spike indicating we make a correct guess, highlighted on the left.

Lets take a look at some of the code¹⁰:

```

for KeyGuess 0 to 255:
    empty hammingPower array
    reset pos1 to 1
    for plainText 0 to 255:
        sBox_lookup = s[plainText XOR keyGuess]
        hammingPower[pos1] = getHammingWeight(sbox_lookup)
        increment pos1 by 1
    empty actualPower Array
    for p 0 to points:
        actualPower = getEachTracesDataPointAtPosition(p)
        calculatePearsonCoefficient(actualPower, hammingPower)

```

Here is a rough outline of the process:

1. We make guesses, just like before.
2. We compute the hamming weight of our SBox lookup of the guess and commit this as the ‘Hamming power’ of the guess at that position.
3. We then take the actual power of each trace at that given point, and calculate the Pearson Correlation at that point.

What this allows us to do is see how well our ‘guess’ fits with the actual power profile derived from the traces recorded off the device. The positions will give us an actual power rating, and our guess will give us a power

¹⁰Both snippets come from the reference paper, Owen Jo, et al, link above.

consumption guess based on the fact that the Hamming Weight of a bitstring correlates exactly to the power consumption.

Performing this across the entire power trace will yield various data points and spikes, and we use the Pearson Correlation Coefficient as our guide - this is fundamentally the mean average of the product of the mean-adjusted ‘random’ variables. It is designed to bring out variances within a ‘random’ population, which is essentially what we have. By computing this for every power point against our guesses, most of the time there will be no correlation - it will be comparing random noise to random noise as above in DPA. However, when our guess is correct, we will have a distinct spike in correlation, meaning we can identify our guesses!

Although they look quite similar, the trace analysis output can be found in fig ??.

It should be said that in practical examples, there is a lot of lining up to do of the positions, but this is a good outline for how the attacks can be implemented. For more, see the NewAE tutorial on this attack here:

[https://wiki.newae.com/Tutorial_B6_Breaking_AES_\(Manual_CPA_Attack\)](https://wiki.newae.com/Tutorial_B6_Breaking_AES_(Manual_CPA_Attack))

and the original paper here:

<https://www.tandfonline.com/doi/full/10.1080/23742917.2016.1231523>

15.5 Real World and Worked Examples

Our examples will make use of the ChipWhisperer setup, as mentioned in section ??.

[Do a writeup of the AES attack... in all its scripted beauty!]

Chapter 16

IC Decapping

16.1

16.1.1

Chapter 17

Bootloader Reverse Engineering

17.1

17.1.1

Chapter 18

Custom Protocol Implementations

18.1 How do you solve a problem like a custom protocol?

By and large, in the modern technological world, we find that most devices will use common protocols - the ones you have seen covered in previous chapters.

However, it is not uncommon to see some custom protocol implemented - whether it be some debug or custom method transferring serial data, we will present some ideas around how these can be reverse engineered, and then bit-banged on some uC we have lying around.

18.1.1 Bit-Banging - a beginner's guide

'Bit-banging' may sound like something from a *Gibson-esque cybersexual noir* book, but in actual fact, it's a way of using GPIO pins in order to emulate any protocol we want.

Similar to when we defined a radio in software and then used an SDR to then transmit and receive, GPIO pins are flexible ways to read digital levels and output digital levels. As is hopefully reasonably apparent by now, many serial protocols work on measuring voltages and comparing against a threshold, and being able to set a line to be some particular voltage.

By creative use of internal clocking and setting GPIO pins to appropriate levels, and then reading these levels, we can effectively emulate a protocol without having a specific protocol silicon peripheral on our microcontroller.¹

¹After all, these are just doing the same as our software would do, but hard-coded into silicon hardware.

18.1.2 Basic Example - UART Rx

So, for a first example we're going to big-bang a UART receiver on an Arduino Nano (though this will work on any arduino). We're going to keep it simple, and here's the rough design; we'll listen for a UART byte being transmitted, receive it, and then send it out the regular UART so we can see it.

Let's recall what happens during a regular (non-inverted) UART transmission²:

1. We will assume a baud rate of 9600, with a setup of 8N1.
2. The line is held high until a stop bit is sent by holding the line low.
3. Once we see this, we then wait 104 microseconds (the time it takes for one bit on 9600 baud)
4. Next we enter a for loop where we see if the line is high, and if so, put a '1' in the right position in the byte
 - If the line is held low, we will do nothing, effectively recording a zero.
5. We then print what we saw out the regular UART so we can check.

We will work that we connect the grounds on the transmitter and receiving arduino, and will connect the Tx of the UART source to pin D2 on our arduino. The code that accomplishes this task is as follows: (**Note** - this will compile in Arduino IDE for pretty much any arduino)

```
// get delay time for 9600 baud UART
// take 1/9600 and this will give you the time in microseconds
int DELAY = 104;

// setup char for UART
char a;

void setup() {
  pinMode(2, INPUT); //setup pin as input
  Serial.begin(9600);
}

void loop() {
  if (!digitalRead(2)) {
    // we have a start bit!
    // so, wait for the first bit...
    delayMicroseconds(DELAY);
    for (int i = 0; i < 8; i++) {
      // check if the line is high...
    }
  }
}
```

²We covered this earlier in a timing diagram, but we'll describe it here in words.

```

if (digitalRead(2)) {
    // if so, put a one in the correct position
    // remembering the transmission order of the bits
    a += (1 << i);
}
// if not high, then leave it as 0 and wait...
delayMicroseconds(DELAY);
}
// print the byte we got...
Serial.print("We got: ");Serial.println(a);
// reset a to null...
a = NULL;
}
// let's go round again!
}

```

We can follow a similar procedure for UART transmit, but we leave this as an exercise for the reader.

This is a very straightforward example of how to bit-bang a protocol, and as such, hope that this example goes some way to demystify what we are doing.

Protocols that are relatively easy to bitbang include SPI, I²C, and debugging protocols such as JTAG/SWD not being particularly more complicated once you understand the workings of the internal state machines.

18.2 Implementing a Custom Protocol w/ Datasheet

Here we will show how we can build a protocol library for an aruidno by making use of the excellent CCLib by Ioannis Charalampidis <https://github.com/wavesoft/CCLib> - an arduino implementation of the Texas Instruments CC Debugger protocol. Handy to use with BLE chips such as the CC2541 or CC11xx series.³

To do this, we will reference one of TI's own documents about implementing this protocol, with examples drawn from the CCLib code as appropriate: <http://www.ti.com/lit/an/swra410/swra410.pdf>. We also cite the debug protocol as another useful reference: <http://www.ti.com/lit/ug/swra124/swra124.pdf>

18.2.1 Analysing the Protocol Wiring

Many protocols have reasonably straight forward wiring, but understanding how these things can interact can be helpful in deciphering what we are seeing

³Historical note: CC stands for ‘ChipCon’, a company that made lots of low-power IC’s for short-range radio. They were acquired by TI in 2006 for around \$200 million. Although they are owned entirely by TI, they still seem to operate as a separate unit, releasing chips with their trademark “CC” prefix, and making use of their debug protocol.

on an oscilloscope. The reference document gives us three pins required by CC debug protocol:

- `RESET_N` - Reset line, that is required by the protocol.
- `DC` - Debug Clock
- `DD` - Debug Data

Note, things to keep in mind when analysing the wiring for any given protocol include:

- Whether the protocol uses pull-up or pull-down resistors.
- Whether the protocol uses inverted logic (or some hybrid, like CAN bus).
- Whether any lines are redundant - such as the reset lines on JTAG/SWD.
- Whether any usually-redundant lines are actually useful - such as DTR and DSR for full RS232.

Note that CCLib splits the `DD` line in two - it has one pin acting as input, `CC_DD_I`, and one acting as output, `CC_DD_O` - this saves time on the arduino switching modes on a single pin. The wiring for this is a fairly common hack; a $100\text{k}\Omega$ resistor is placed between `CC_DD_O` and `CC_DD_I` (with another $200\text{k}\Omega$ resistor going from this connection to ground), with the actual `DD` line connected to `CC_DD_I`. Similar $100\text{k}\Omega$ resistors are placed in front of the other output lines.

This allows us to work more easily with the signals and processing on the slightly limited arduino.

18.2.2 Analysing The Protocol Timings

As we hope is clear from the UART Rx bit-bang example, much of this requires an understanding of the timings - what lines need to be raised or lowered, and when. Important aspects of timings when working out how to implement a protocol for use include:

- Making sure that the line is held high in time for the target to read a rising clock edge (or falling edge, if this is what is used)
- Dealing with accumulating timing errors for continuous streams
 - In the above example, we should have used 104.1666... as our microsecond delay.

- The 0.1666... error will accumulate, so we will lose a microsecond every 6 bits or so.
- This is fine for the 10 bits we are concerned with for each UART transmission, but other protocols are much more sensitive to timing errors
- Protocols such as CAN bus that have sometimes confusing artefacts, like bit-stuffing, as part of the protocol.

However, let's consider the timing requirements for the CC Debugger protocol we are concerned with here.

- The clock signal drives all of the communications - including the `DEBUG_INIT` start.
- Data is driven on `DD` by the rising/positive edge of the `DC` line.

18.2.3 Using GPIOs to Implement a Protocol

This is the bit we have been building up to - how to use the information we have gathered and then create an implementation that works for what we need.

It is worth keeping in mind, CCLib is quite a mature library - it has a python interface to communicate to the Arduino, and has implemented error logging/checking and a custom protocol to signal that back to the python interface and alert the user. This is quite a lot of work - and sometimes it will be required (for which CCLib is a good example).

However, some situations will only require that you implement one function in order to prove a PoC, hook in a fuzzer, or create a manual testing interface you can use in security assessment of a device.

We will not give a full teardown of the CCLib implementation - we will, instead, use three illustrative examples; `DEBUG_INIT()`, `GET_CHIP_ID()`, and `CHIP_ERASE()`.

All examples are taken from `CCDebugger.cpp` from the CCLib library, in the Arduino/CCLib folder.

`DEBUG_INIT()` Implementation

Kicking a CCxxxx target into debug is fairly straightforward:⁴

1. Pull `RESET_N` low
2. Pump `DC` twice⁵

⁴See page 6 of the TI reference; <http://www.ti.com/lit/an/swra410/swra410.pdf>

⁵In Travis Goodspeed's debug interface for CC430 in the GoodWatch project, this line is pumped a few times more, just to be sure.

3. Pull `RESET_N` high

This is fairly straight forward in code:

```
/**  
 * Enter debug mode  
 */  
byte CCDebugger::enter()  
{  
    // ----- SNIP -----//  
  
    // Reset error flag  
    errorFlag = CC_ERROR_NONE;  
  
    // Enter debug mode  
    digitalWrite(pinRST, LOW);  
    cc_delay(200);  
    digitalWrite(pinDC, HIGH);  
    cc_delay(3);  
    digitalWrite(pinDC, LOW);  
    cc_delay(3);  
    digitalWrite(pinDC, HIGH);  
    cc_delay(3);  
    digitalWrite(pinDC, LOW);  
    cc_delay(200);  
    digitalWrite(pinRST, HIGH);  
    cc_delay(200);  
  
    // We are now in debug mode  
    inDebugMode = 1;
```

`GET_CHIP_ID()` Implementation

Getting the Chip ID is a useful way of seeing how data transfer works. First, we will summarise some of the utility functions used in this part of the CCLib code:

- `read()` - this function sends 8 clock pulses, and reads `DD` line for incoming data. Here is a code snippet

```
// Send 8 clock pulses if we are HIGH  
for (cnt = 8; cnt; cnt--) {  
    digitalWrite(pinDC, HIGH);  
    cc_delay(2);  
  
    // Shift and read  
    data <= 1;  
    if (digitalRead(pinDD_I) == HIGH)  
        data |= 0x01;  
  
    digitalWrite(pinDC, LOW);  
    cc_delay(2);
```

- `write()` - This does the inverse action, sending 8 clock pulses and then raising or lowering the `DD` line according to the byte we want to send. Here's a snippet:

```
// Sent bytes
for (cnt = 8; cnt; cnt--) {

    // First put data bit on bus
    if (data & 0x80)
        digitalWrite(pinDD_O, HIGH);
    else
        digitalWrite(pinDD_O, LOW);

    // Place clock on high (other end reads data)
    digitalWrite(pinDC, HIGH);

    // Shift & Delay
    data <= 1;
    cc_delay(2);

    // Place clock down
    digitalWrite(pinDC, LOW);
    cc_delay(2);
}
```

- `switchRead()` and `switchWrite()` - these helper functions let the code switch modes between reading and writing data to the `DD` line.
- The code has the ‘read chip ID’ command, `0x68`, in a variable we can call.

Here is a snippet of the relevant code:

```
write( instr[I_GET_CHIP_ID] ); // GET_CHIP_ID
switchRead();
bRes = read(); // High order
bAns = bRes << 8;
bRes = read(); // Low order
bAns |= bRes;
switchWrite();
```

We can read the response, and then get an idea of which chip is in play. For example, a response of `0xA5` would indicate a CC2530 target, whilst `0xB5` would indicate a CC2531. See page 7 of the technical reference.⁶

This example shows very well how we can bit-bang a custom protocol in software to both send data, and read the responses back. It is worth noting that the CCLib code has other commands available, not just `GET_CHIP_ID`, such as:

```
// Default CCDebug instruction set for CC254x
```

⁶<http://www.ti.com/lit/an/swra410/swra410.pdf>

```

||| instr[INSTR_VERSION]      = 1;
||| instr[I_HALT]            = 0x40;
||| instr[I_RESUME]          = 0x48;
||| instr[I_RD_CONFIG]        = 0x20;
||| instr[I_WR_CONFIG]        = 0x18;
||| instr[I_DEBUG_INSTR_1]    = 0x51;
||| instr[I_DEBUG_INSTR_2]    = 0x52;
||| instr[I_DEBUG_INSTR_3]    = 0x53;
||| instr[I_GET_CHIP_ID]      = 0x68;
||| instr[I_GET_PC]           = 0x28;
||| instr[I_READ_STATUS]      = 0x30;
||| instr[I_STEP_INSTR]       = 0x58;
||| instr[I_CHIP_ERASE]       = 0x10;

```

CHIP_ERASE() Implementation

This is the last in our basic examples of implementing the CC Debug protocol in software on an arduino.

The `CHIP_ERASE()` function has a slightly different operation than the read/write we have seen so far:

1. Send the command `CHIP_ERASE()` by writing the value `0x10` to the target
2. Wait for the `DD` line to go low.
3. Send 8 clock pulses and check if the 7th bit is high - if so, this means `CHIP_ERASE_BUSY`
4. Once this is unset, the flash is erased

It should be noted that the third and fourth steps in this description are not observed in the CCLib code strictly. The command `CHIP_ERASE()` is sent to the device and then the status is read back to the user immediately:

```

||| byte bAns;

||| write( instr[I_CHIP_ERASE] ); // CHIP_ERASE
||| switchRead();
||| bAns = read(); // Debug status
||| switchWrite();

||| return bAns;

```

As such, an understanding of how this code works may help in troubleshooting such code, should a mass erase operation fail.

*18.3. REVERSE ENGINEERING AND IMPLEMENTING A CUSTOM PROTOCOL*161

18.3 Reverse Engineering and Implementing a Custom Protocol

[Find an example IRL of a custom protocol. Write up should be - RE it, and then bit bang on an ARM or arduino uC...]

Chapter 19

CRP Bypass Techniques

19.1 Overview of CRP

19.1.1 What is CRP?

CRP, or Code Readout Protection, is a collection of various methods employed by IC manufacturers to ensure that access to the firmware loaded onto a chip is limited once the chip has been flashed.

The intention is to prevent the acquisition of firmware by unauthorized parties - this is more to prevent copying of devices and IP¹/copyright/patent infringement. However, as we shall see, there are a variety of ways in which CRP has been found to be bypassed.

19.1.2 CRP Bypasses - the Basics

CRP is generally only applied to microcontrollers - and is not found on microprocessors, where the RAM and ROM are in separate packages. As such, you can't really control the firmware as the ROM IC is open for an interested party to access and pull off the firmware (as has been covered earlier).

CRP is a control specific to microcontrollers - and as such, there are as many ways of implementing it as there are manufacturers of uC's! There are broadly several 'levels' of CRP:

- No control applied
- Protect from read, but allow erase/rewrite
- Protect from read, disallow erasure

As such, we want to achieve one of the following situations:

¹Intellectual Property, not the other one.

- Disable the security control altogether, giving free access to the ROM
- Find some other way of accessing the information through another channel, say from RAM
- Leverage fail-safe or fail-open modes on IC's

We shall cover Glitch-based firmware recovery or CRP bypass in Chapter ??, so please see there for details. This chapter will assume that you have access to the target device, and suitable equipment for gaining debug access - UART interfaces, JTAG/SWD adapters, etc.

This chapter is broadly influenced and inspired by Andrew Tierney's² talk given at BSides Leeds entitled 'Armadillos'³

19.2 CRP Bypasses

19.2.1 UV Light to Reset Security Fuses

This attack was formulated by Bunnie Huang for opening up firmware on a locked PIC 18F1320 - Source: http://www.bunniestudios.com/blog/?page_id=40

How UV affects storage

In order to understand this attack we need to go down to the very small and understand what a transistor looks like under a microscope.

Figure ?? shows the structure of a floating gate MOSFET (FGMOS) - these use a neat effect called Fowler-Nordheim tunnelling⁴ to 'store' a charge in the floating gate.

There are two ways in which this information can be reset - by either applying a large inverse voltage (which at these scales is in the 3-5V range) or by applying UV light (typically c. 205nm wavelength is used).

²<https://twitter.com/cybergibbons>

³A recording can be found here: https://www.youtube.com/watch?v=DTuzuaiQL_Q

⁴I tell ya, this stuff is magical physics...

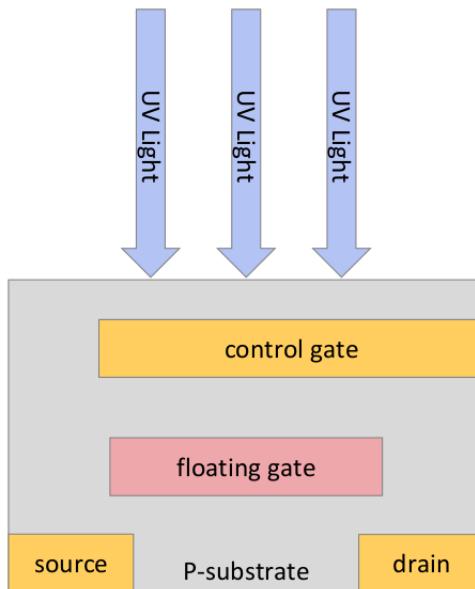


Figure 19.1: Structure of a FGMOS - the UV light resets the state of the floating gate

To make this a viable way of resetting, say, security bits that prevent us from reading firmware, we can follow the following method:

- Decap the chip or expose the silicon
 - This subject is dealt with in depth in Chap ??
- Put a mask over the silicon to prevent blanking the flash with the UV
- Carefully expose the relevant areas we want to reset
- Place the chip under UV light
- Enjoy reverse engineering your firmware!

So, eventually companies have grown wise to techniques such as this, and have begun to put some controls in place. One such example is found on the PIC 18F1320 - the security bits that control the read/write access to the firmware are covered by small metal squares. These are designed to prevent anyone resetting them and extracting the firmware. See figure ?? for a die photo and diagram about how this is supposed to work.

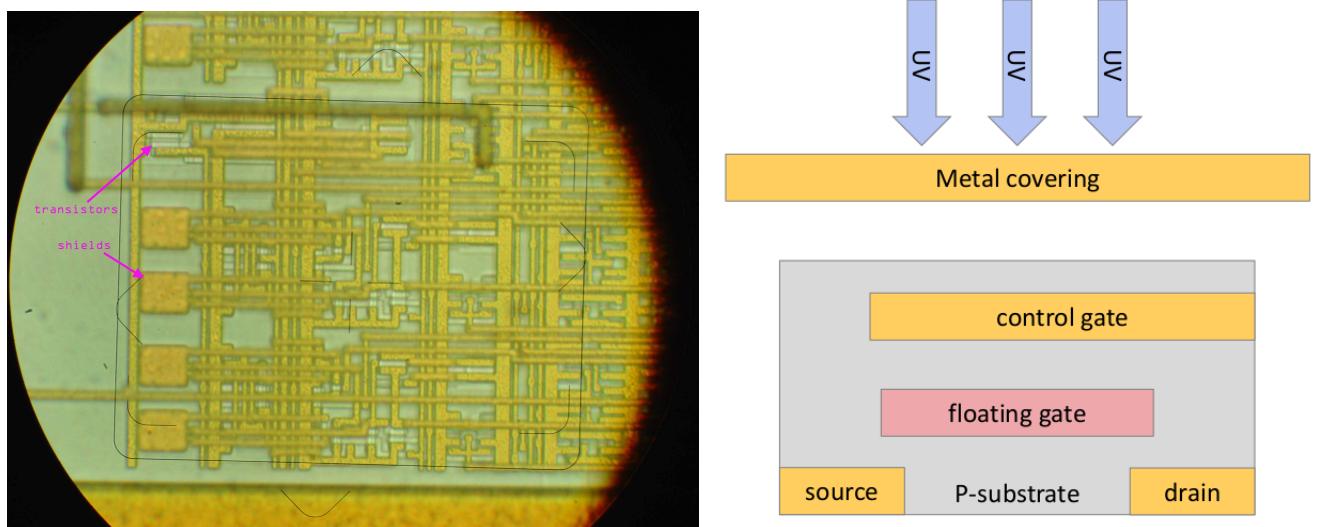


Figure 19.2:

Resetting via A Tilt

These controls, however, overlooked one key thing - light can be reflected by metal. It seems very obvious when you say it, but this is a very subtle attack - we proceed as the steps outlined above, however we now tilt the IC by 45° in order to cause the light to reflect off the bottom of the IC, then up to the metal covering, and reflect off this into the substrate in order to reset the FGMOS. Figure ?? has an illustration about how this is envisaged to work. Photographs of the attack can be found on Bunnie's blog - the link above.

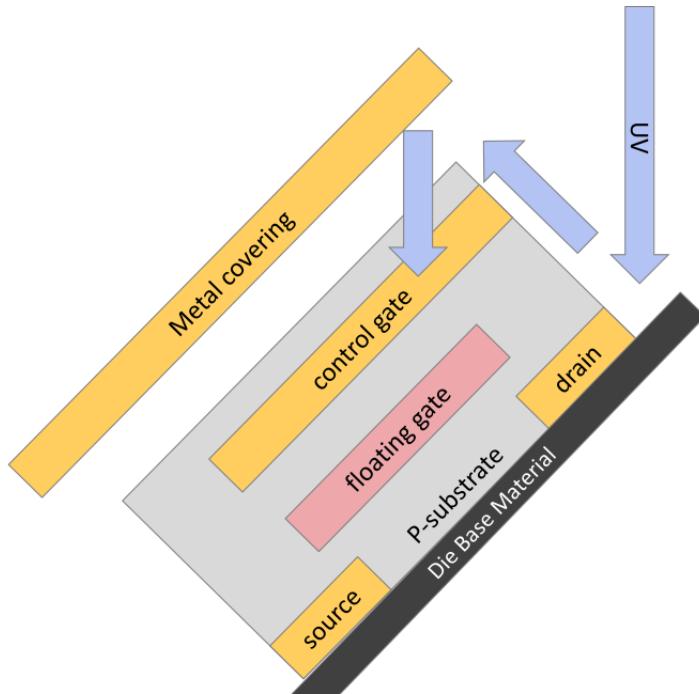


Figure 19.3:

This is a very elegant and simple solution to this particular problem of bypassing these controls - and although it requires some other work (such as decapping) in order to be effective, it is a very elucidating example of the kinds of approaches that have been employed to bypass CRP and other hardware controls.

19.2.2 LDR Register Readout on NRF51822

Background to the Attack

This attack comes from the Include Secuirty blog - Source: <http://blog.includedsecurity.com/2015/11/NordicSemi-ARM-SoC-Firmware-dumping-technique.html>

ARM-based microcontrollers are all the rage. After 8051, ARM is the most in-use RISC architecture in embedded systems. In any given computing unit (a laptop, mobile phone, desktop PC, etc.) there are several embedded cores to any major CISC-based CPU.

Being a RISC architecture, there are some things that are done very differently - one of these things is how data is loaded. In x86, you can reference memory locations in opcodes without any issue. But on ARM, if you want to store or load a string you cannot simply access memory - you

have to use two specialist opcodes that will store and fetch data for you. These are `STR` and `LDR` - the first one stores data, and the second one loads it.

`LDR` is a very straightforward opcode - here's the rough structure for those less familiar with ARM:

```
|| LDR r4, [r3, #4]
```

Here, the value in `r3` is increased by 4, and then the contents of memory at the address corresponding to that value is loaded into `r4`.

This is a ubiquitous opcode, so being able to leverage it to bypass something like CRP is a real boon. Especially given the popularity of the target in this case, the NRF51822.⁵

Attack Method

To exploit this attack, we look to the following method:

1. Access debug on the target NRF51822
2. Locate a `LDR` opcode, and record its location
3. Work out which register it is reading from
4. Use the debug control to set the `PC` register and input register (`r3` in the above example)
5. Step the processor and retrieve the data from the output register (`r4` in the above example)
6. Repeat until the whole firmware is acquired

Assuming we have found the relevant opcode location and parameters, we can script this attack over OpenOCD's telnet server as follows in this python script:

```
import sys
import telnetlib
import re

HOST = "localhost"
PORT = "4444"
timeout = 500

tn = telnetlib.Telnet(HOST, PORT)
f = open("firmware.txt", "w")

tn.read_until("> ")
```

⁵This IC is used in devices like the BBC Micro:bit, and devices that need low power BLE, suchas fitbit clones.

```

tn.write("reset halt\n")
for i in range(0,0x200):
    # replace this '3' with the location
    # of the LDR opcode in the firmware
    tn.write("reg pc {0}\n".format(3))
    tn.write("reg r3 {0}\n".format(i))
    tn.write("step\n")
    tn.write("reg r4\n")
    num,_,lines = tn.expect([r'r4.*\w'], timeout)
    f.write(lines)
    #print(lines)
f.close()

```

Clearly this script will only get the first 512 bytes of firmware, but can be used to acquire whole images from microcontrollers such as the NRF51822.

Note - This attack works on ICs in the NRF51xxx series, but is known to not work on the NRF52xxx series ICs, so it may well be that this issue has been mitigated on newer chips.

19.2.3 ‘Heart of Darkness’ Attack

This attack was originally due to Milosh Meriac <https://www.openpcd.org/images/HID-iCLASS-security.pdf> against the HID iCLASS.

This particular attack concerns the PIC18F452 IC - the CRP controls are quite strict. You can bulk-erase the bootloader or all of the storage blocks - but no further granularity of control. However, this is enough for us to proceed.

First we create a bootloader that prints out the entire firmware:

```

for(int i=0; i<0x8000; i++) {
    printf("%x: %x", i, &i);
}

```

We now compile this and proceed as follows:

1. Bulk erase the bootloader, leaving the firmware intact
2. Load our own custom bootloader
 - Note - this will now dump the entire firmware through a UART connection, so we have the bulk of the firmware
3. On a second device, bulk erase the bootloader, and write the dumper code to the end of the firmware
 - Bulk erase replaces the firmware with NOPs, so the bootloader will send execution to the firmware entrypoint, but this will be a NOP sled that will be ridden all the way down to the dumper code
 - We have now acquired the bootloader firmware from the device

Figure ?? has two diagrams that illustrate the two main phases of this attack.

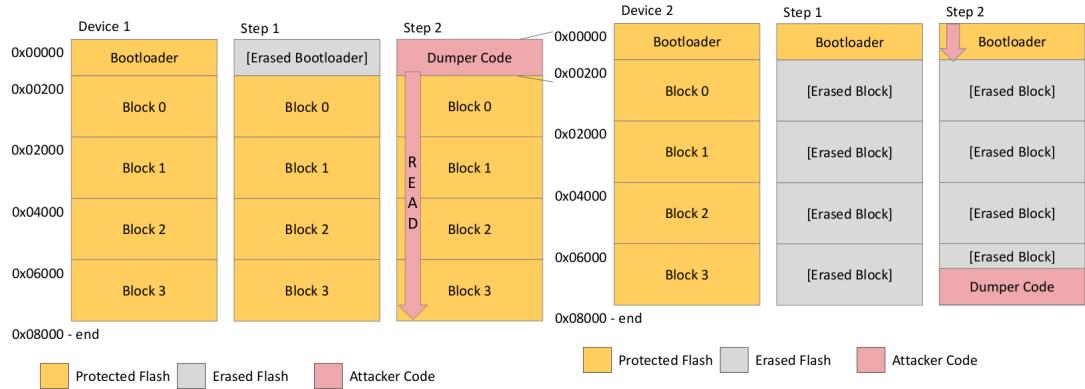


Figure 19.4: Phase 1 is on the left, phase 2 is on the right

This attack is a very nice example of how, with some careful management, a CRP can be bypassed even with a bulk erase being the only real tool we can use. The paper this attack first featured in is well worth a read!

19.2.4 Cold Boot Stepping

This attack was first put forward by Obermaier: <https://www.usenix.org/system/files/conference/woot17/woot17-paper-obermaier.pdf>

Background to the Attack

This attack instead chooses to focus on the memory of a microcontroller, as this can sometimes be left unaffected by CRP restrictions. As such, if a firmware does some sort of integrity check on itself, or some other process that ranges over the memory, this can possibly be leveraged to extact the firmware.

In the example presented, a CRC32 check is done on the firmware on the chip, and we abuse the level of access granted by a debug connection to control the flow of this function and pull the firmware, one double word at a time, out of memory.

This is the case on the STM32F051 - the memory is left unaffected from being read in full, even when CRP is fully enabled on the chip to prevent reading the firmware. As such, this is the example used in Obermeier's paper.

Attack Method

Here is an outline of the attack - firmware examples and code are all available in the paper and accompanying materials linked to above. Here is the general methodology of the attack:

1. Get debug control of the CPU
2. Identify when an integrity check is in progress
 - This is tricky, but fingerprinting the number of steps needed and the register behaviour for a CRC32, for example, and then comparing this to a target would be relatively straightforward.
3. Once this has been identified, we then setup a similar process to the LDR attack:
 - (a) Let the check run for the first address `0x0001` and retrieve the double word of firmware from that address from memory.
 - (b) Reset the chip, and let it run to the second address, `0x0002`, and then retrieve the memory contents again
 - (c) Repeat this incrementally for the entire firmware.

Figure ?? has a diagram demonstrating the rough mechanics of this attack. It is quite slow - Obermeier's example has the attacking STM32 chip emulating SWD in order to control the clocking very accurately, as well as the actual debugging. However, there is nothing that would explicitly prevent this from being effective over OpenOCD in some circumstances.

This attack highlights very well the kind of creative thinking that is very useful to aim to achieve - indeed, the kind of hack-y thought process, tenacity, and resourcefulness that is successful on limited systems such as those found in IoT/embedded devices hones many skills that are very transferrable within the tech industry.

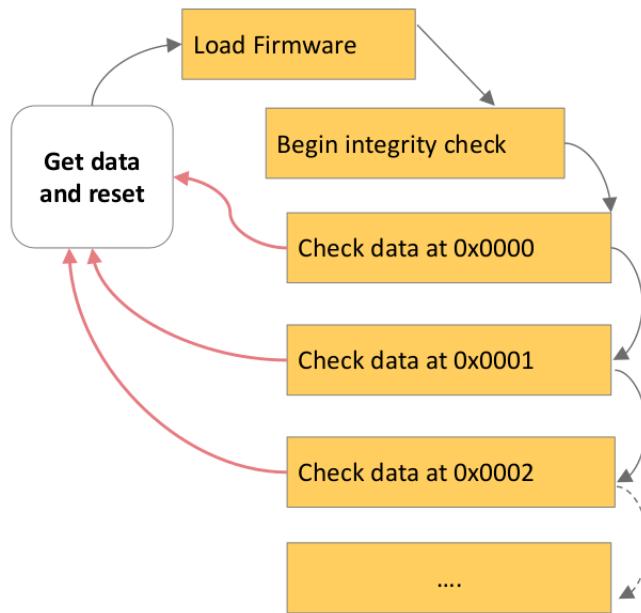


Figure 19.5:

Chapter 20

Cryptography on IoT/Embedded Systems

20.1 Cryptographic Overview

This section will give some practical overviews of ciphers important to and used within this book, and embedded systems in general. We will focus on two ciphers: AES (128, 196, and 256 variants) and RSA.

If an in-depth understanding of a cryptographic protocol, scheme, or attack is required, we recommend J. P. Aumasson's excellent *Serious Cryptography* - Available from No Starch Press.¹ This is a truly excellent and accessible volume for those who need to do serious cryptography.

However, for our purposes, we need only a general overview of the specific cryptoschemes found in modern embedded and IoT devices.

20.1.1 Overview of AES

[Overview to follow!]

20.1.2 Overview of RSA

[To follow!]

20.1.3 Overview of SSL/TLS

[This section should cover SSL/TLS issues found in IoT/embedded systems - bad cert checking, processing, randomness, etc.]

¹<https://nostarch.com/seriouscrypto>

20.2 Randomness in IoT

[I'll put my research here - what went wrong, and how it can be assessed.
]

Appendix A

GDB Cheatsheet

Here is a GDB Cheatsheet of useful command/ways to do things in GDB:

Command	Description
info <var> or i <var>	print some information
i b	Get info about breakpoints
i r	dipslay all registers and their values
i proc map	display the process map information for the code being analysed
break label	break on function label, e.g. break main
break *0x1000	break on 0x1000
delete 2	delete breakpoint 2
set \$r0=5	set the value of register \$r0 to 5
stepi or si	step forward one instr.
nexti or ni	step forward one instr., and STEP OVER any calls
stepi 5	Step forward 5 instr.
nexti 5	step forward 5 instr., stepping over any calls
continue or c or cont	continue on with the process
fg	as in BASH, foregrounds the process (push to background with c&)
bt	'backtrace' - show the call stack. NB -useful when analysing exceptions
run arg1 arg2	execute the debug target (from cli) with arg1 and arg2 as arguments.
print system	print the system symbol table data (stored as the system variable in GDB) NB - very handy for ROP exploring.
p system	Shorthand for print system
p *(array-var)@20	print 20 bytes of array-var

<code>x 0x1000</code>	short for ‘eXamine’, examine 0x1000 in this case
<code>x/nfu 0xADDR</code>	<p>Starting from 0xADDR, examine the;</p> <ol style="list-style-type: none"> 1. (N)umber of bytes, in a specified 2. (F)ormat (any of <code>x/d/u/o/t/a/c/f/s</code> from C’s print formatting, with additional <code>i</code> for instructions - defaults to <code>he(x)</code>adecimal), 3. Specified (U)nit size (any of (b)yte, (h)alfword, (w)ord, or (g)iant,double words) <p>e.g. <code>x/4dw \$ps</code> is ‘examine address in \$ps, and print 4 words as signed integer values’ [Yeah, that one’s a mouthful...]</p>
<code>x/10i \$pc</code>	disassemble the next 10 instructions from the addr value in PC
<code>x/s 0x1000</code>	display the string at memory addr. 0x1000
<code>x/10xb 0x1000</code>	display 10 bytes, as hex, starting at addr 0x1000
<code>x/10xw \$sp</code>	Show the 10 32bit words, in hex, starting at \$sp
<code>set * (int*) 0x6000 = 42</code>	Set an int (32bits) at addr 0x6000 to 42
<code>disp <expression></code>	do <expression> each time a breakpoint is hit, e.g. <code>disp x/10i \$pc</code>
<code>disassemble (label)</code>	disassemble the function with the label specified, e.g. <code>disassemble main</code> will disassemble main.
<code>disassemble 0x1000 0x2000</code>	disassemble from 0x1000 to 0x2000 (replace with addresses in the program)
<code>set print pretty on</code>	set print formatting of C structures on (use <code>off</code> to turn off)
<code>set logging file /path/to/file</code>	set logging to go to <code>/path/to/file</code> from your session
<code>set arm force-mode arm/thumb/auto</code>	Set the arm forced instruction mode (arm/thumb/auto) - overrides the symbol table
<code>show arm force-mode</code>	Show the current force instruction mode

Useful Commands

I also find these commands useful when debugging (on linux posix like environments):

Command	Descr.
<code>gdb file</code>	Debug file in gdb (You never know, I might forget...)
<code>gdb a.out 789</code>	attach GDB to a.out running as process 789
<code>gdb a.out core</code>	open core dump for a.out in GDB
<code>cat /proc/789/maps</code>	get address space layout for 789
<code>echo 0 >/proc/sys/kernel/randomize_va_space</code>	turn off the pesky ASLR...
<code>execstack -s foo</code>	Disable XN for foo
<code>ulimit -c unlimited / ulimit -c 0</code>	turn core dumps on / turn core dumps off

Glossary

This will be a basic glossary of commonly looked-up terms.

BGA - Ball Grid Array. The IC package where blobs (balls) of solder underneath the IC are used to bond it to the PCB contacts.

CPE - Consumer Premises Equipment. Often, routers, but can also be IPTV boxes, DECT phones with base units, etc. Anything supplied by some service provider and left in a consumer's home.

PCB - a Printed Circuit Board. The thing with all the components on. Can be multilayered, or simply double-sided. Power and signals are communicated via traces on the PCB that link up components.

SMC - Surface Mount Components - the little tiny components that sit atop/underneath a PCB, and you daren't put down in case they get blown away by a flea's fart.

SOIC - Small Outline Integrated Circuit. A very common IC Package you will see on many consumer electronics.

Solder Mask - The part of a PCB that's coloured (green/black/red/whatever). This actively resists solder being put on it, driving it through surface tension to stick to the copper contacts and components instead.