# Data Mining
# Hardware Descriptions

**from Vendor Code, Configuration Tools, and Documentation**

**Niklas Hauser, emBO++25**
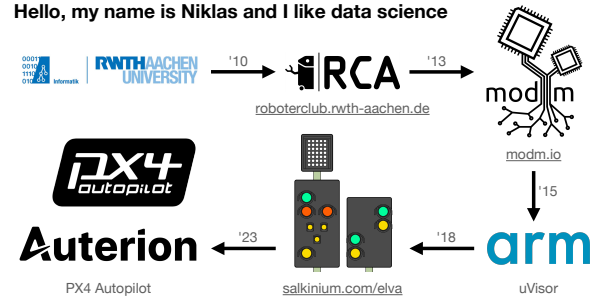
emBO++

---

Thank you for the introduction.
Thanks to emBO for the opportunity to talk here.

---

## Who?
**Hello, my name is Niklas and I like data science**



RCA — roboterclub.rwth-aachen.de

modm.io

PX4 Autopilot

Auterion
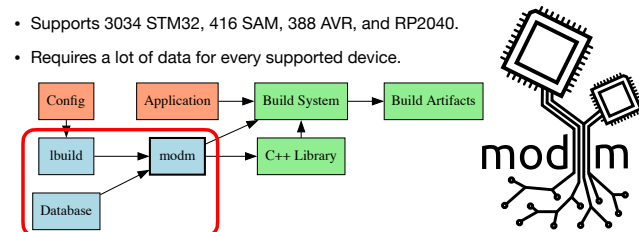
salkinium.com/elva

uVisor

---

My name is Niklas.
- I started studying Computer Science some time ago.
- I began building autonomous robots in 2010.
- We created a C++ library which is known as modm.io, a C++23 library generator that supports several thousand Cortex-M devices.
- I then started at ARM working on Cortex-M sandboxing, before returning to the university to study for my masters degree.
- There, I worked on a digital modular signalling system for railways.
- I'm currently working at Auterion on the open-source PX4 Autopilot.

---

## modm.io C++23 barebone embedded library
**Modular, data-driven HAL and build system generator**

- Generates startup code, linker script, peripheral drivers for microcontrollers.
- Supports 3034 STM32, 416 SAM, 388 AVR, and RP2040.
- Requires a lot of data for every supported device.



modm
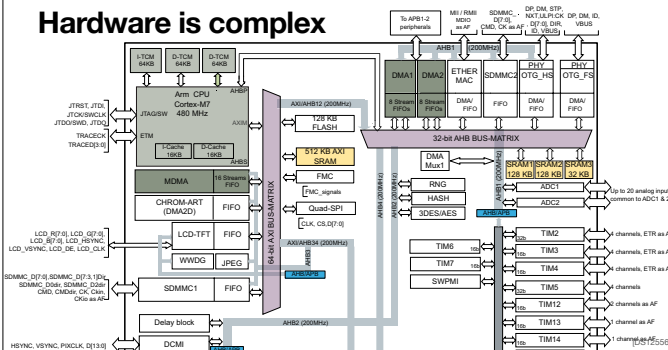
---

modm is a C++23 embedded library generator.
- The core of modm is a code generator written in Python called lbuild.
- It queries a database of device data and formats the results into C++23 code.
- The HAL is highly modular and configurable and it allows a very small maintainer team to support thousands of microcontrollers.
###
- Today we'll talk about the database part of this construct.

## Hardware is complex
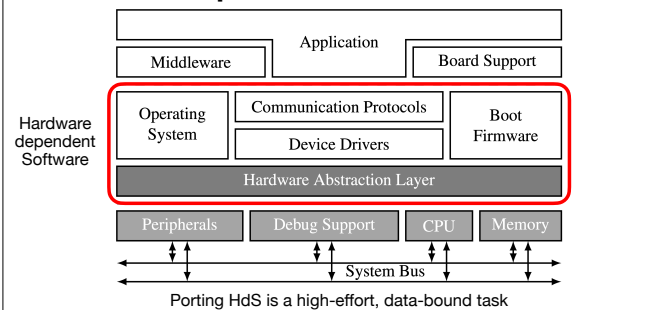


Microcontroller hardware is quite complex nowadays.
Here is a STM32H7 with its many internal busses.
You can see distributed memories in yellow, you can see lots of peripherals, and many DMA engines.
Everything also needs to be externally connected via the pins.
That's a lot of hardware to abstract.

## Hardware-dependent Software



The HAL is actually part of the hardware-dependent software and there's a lot of it.
###
It's also operating systems, external sensors, communication protocols, and bootloaders.
So it's a fairly large topic, not just about microcontrollers itself.

## Combine and Share all Data Sources



So the idea is to parse every data source I can find and merge it into a single database.
Then I can share this among all my embedded friends: Zephyr, modm and embassy.
And then I would benefit from any of their improvements to the database as well.

## data.modm.io Conversion Pipelines

And I decided to make this an open-source project on GitHub.
It's split into individual pipelines, where each data source is converted eventually into Python.
###
Let's first focus on the trivially machine-readable data sources.

---



## Configuration Tools: CubeMX
**STM32_open_pin_data contains all packages, pinouts, memories**

```
<Pin Name="PB13" Position="26" Type="I/O">
    <Signal Name="ADC3_IN5"/>
    <Signal Name="COMP5_INP"/>
    <Signal Name="I2S2_CK"/>
    <Signal Name="OPAMP3_VINP"/>
    <Signal Name="OPAMP3_VINP_SEC"/>
    <Signal Name="OPAMP4_VINP"/>
    <Signal Name="OPAMP4_VINP_SEC"/>
    <Signal Name="SPI2_SCK"/>
    <Signal Name="TIM1_CH1N"/>
    <Signal Name="TSC_G6_IO3"/>
    <Signal Name="USART3_CTS"/>
    <Signal IOModes="Analog,EVENTOUT,EXTI"
</Pin>
```
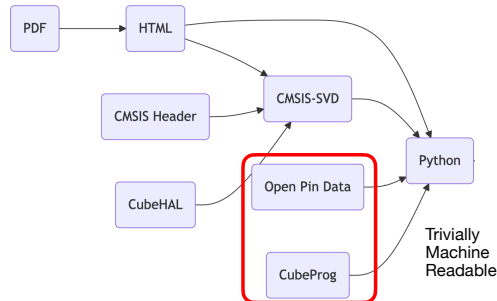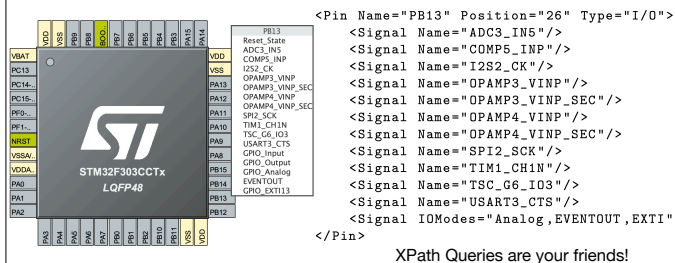XPath Queries are your friends!

The most well known is the CubeMX GUI application, which allows you to configure the pin functions of the STM32.
This is actually backed by a XML database that STMicro actually publishes on GitHub with a BSD licence.
It contains the entire catalog of STM32 ever made, their package, their pinout, and all alternate functions.
It's undocumented but you can get very far with simple XPath queries.
Many people already use this, including Zephyr, embassy and KiCad to generate HALs and footprints!

---



### CubeMX Clock Tree

However, the CubeMX database also contains a fully annotated graph of the entire STM32 clock tree.
###
A typical configuration is to have an external clock source fed into the PLL, which then increases the clock frequency and feeds it into the system clock, from which most peripherals are powered.

**Clock Graph**

salkinium.com/stm32/clock

We can also render this clock graph as graphviz graph, and you can see that it contains all frequency limitations that are used to solve the problems of the clock tree in CubeMX.
Here we can follow the same configuration: external clock source gets fed into the PLL and comes out into the system clock.
But now there is a lot more detail visible.
You can also see that this is not really a tree, it's really a graph.
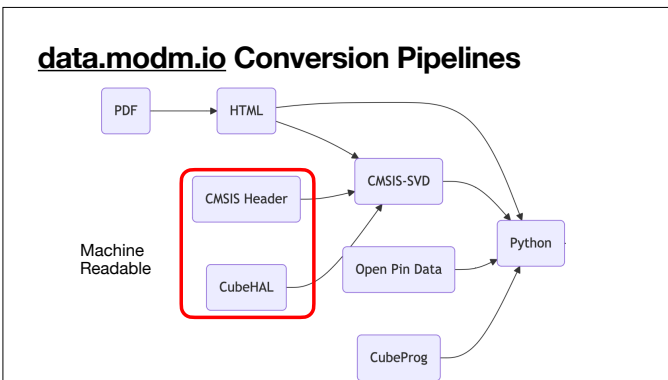You can find more of these rendered clock graphs on my homepage.



**data.modm.io Conversion Pipelines**

Machine Readable

So that was the easy part, let's now focus on more difficult data sources: source code.

**Parsing CMSIS Header Files**
**and converting them back to CMSIS-SVD**

1. Parse peripheral structs to reconstruct register order and offset.

2. Resolve numeric values of macros to reconstruct bit field order and offset.

3. Connect TypeDef instantiation (=peripheral) with macros via name matching.

```
#define PERIPH_BASE        0x40000000UL
#define AHB1PERIPH_BASE    (PERIPH_BASE + 0x00020000UL)
#define CRC_BASE           (AHB1PERIPH_BASE + 0x3000UL)
#define CRC                ((CRC_TypeDef *) CRC_BASE)

/****************** Bit definition for CRC_CR register
#define CRC_CR_RESET_Pos   (0U)
#define CRC_CR_RESET_Msk   (0x1UL << CRC_CR_RESET_Pos)
#define CRC_CR_RESET       CRC_CR_RESET_Msk
```

```
typedef struct
{
    __IO uint32_t DR;
    __IO uint8_t  IDR;
    uint8_t       RESERVED0;
    uint16_t      RESERVED1;
    __IO uint32_t CR;
} CRC_TypeDef;
```

We can convert the CMSIS header files back into a register map:
We know the order and width of the registers from the typedef struct.
We know the order and width of the bit fields from the macros.
And we know the peripheral instance and address from the typedef cast.

This does not give us enumerations of any bit fields unfortunately, since they are simply not in the header files.

## Parsing CubeHAL Header Files
### CMSIS files are missing Bit Field Enumerations

- Neither the STM32 CMSIS-SVD nor CMSIS Header define Bit Field Enumerations.

- We need to parse the Low-Level CubeHAL header files to reconstruct.

```
/******  Bit definition for RCC_CFGR register  ******/
#define RCC_CFGR_SW_Pos      (0U)
#define RCC_CFGR_SW_Msk      (0x3UL << RCC_CFGR_SW_Pos)
#define RCC_CFGR_SW          RCC_CFGR_SW_Msk
#define RCC_CFGR_SW_0        (0x1UL << RCC_CFGR_SW_Pos)
#define RCC_CFGR_SW_1        (0x2UL << RCC_CFGR_SW_Pos)

#define LL_RCC_SYS_CLKSOURCE_HSI    0x00000000U
#define LL_RCC_SYS_CLKSOURCE_HSE    RCC_CFGR_SW_0
#define LL_RCC_SYS_CLKSOURCE_PLL    RCC_CFGR_SW_1
#if defined(RCC_PLLR_SYSCLK_SUPPORT)
#define LL_RCC_SYS_CLKSOURCE_PLLR   (RCC_CFGR_SW_1|RCC_CFGR_SW_0)
#endif
```

For some of the bit field enumerations we need to parse the CubeHAL low-level header files.
Same procedure, we interpret the macros.
###
We can do a reverse lookup to see which macros use the register bit field definitions and then work backwards from that.
Annoying, but doable.
But does not give every bit field enumeration possible.

---

## data.modm.io Conversion Pipelines



Machine Renderable

Human Readable

Now for the really hard stuff: parsing PDF datasheets.
PDFs are machine-renderable, but not machine-readable.
There's a lot of research out there on information extraction from PDFs, mostly relating to financial statements.

---

## STMicro PDF Documentation
### You can look, but you cannot parse

- STMicro publishes >2600 PDFs for documentation: ~15GB on disk.

- You must consult multiple PDFs with thousands of pages: STM32H7A3/B0/B3.

- How hard could it possibly be to make all these PDFs machine-readable?



STMicro publishes a lot of PDFs: We are only looking at active components, microcontrollers, sensors, memories.
And there are over 2600 PDFs available: ~15GB.

For one microcontroller, a lot of PDFs apply, here the STM32H7 family has 7 PDFs involved.
Nobody reads them all.

## PDF Datasheets: Text

**17**     Cyclic redundancy check calculation unit (CRC)

**17.1**     Introduction

The CRC (cyclic redundancy check) calculation unit is used to get a CRC code from a 8-, 16- or 32-bit data word and a generator polynomial.

Among other applications, CRC-based techniques are used to verify data transmission of storage integrity. In the scope of the functional safety standards, they offer a means of verifying the Flash memory integrity. The CRC calculation unit helps compute a signature of the software during runtime, to be compared with a reference signature generated at link time and stored at a given memory location.

**17.2**     CRC main features

- Uses CRC-32 (Ethernet) polynomial: 0x4C11DB7
  $$X^{32} + X^{26} + X^{23} + X^{22} + X^{16} + X^{12} + X^{11} + X^{10} + X^8 + X^7 + X^5 + X^4 + X^2 + X + 1$$

How can we access PDF data?

For text its relatively simple: each glyph is individually positioned on the page.

There's no semantics for headings or lists or superscript. It's all just individually positions characters.

THERE IS NO NEED TO OCR PDFs!

---

## PDF Datasheets: Figures



Figures are a mix of vector graphics and text. There's no special indication that this is a figure, it must be detected.

---

## PDF Datasheets: Tables

Table 13. STM32F303xB/STM32F303xC pin definitions (continued)

| Pin number | | | | Pin name (function after reset) | Pin type | I/O structure | Notes | Pin functions | |
|---|---|---|---|---|---|---|---|---|---|
| WLCSP100 | LQFP100 | LQFP64 | LQFP48 | | | | | Alternate functions | Additional functions |
| J3 | 52 | 34 | 26 | PB13 | I/O | TTa | (4) | SPI2_SCK/I2S2_CK/USART3_CTS/TIM1_CH1N/TSC_G6_IO3/EVENTOUT | ADC3_IN5/COMP5_INP/OPAMP4_VINP/OPAMP3_VINP |
| J2 | 53 | 35 | 27 | PB14 | I/O | TTa | (4) | SPI2_MISO/I2S2ext_SD/USART3_RTS_DE/TIM1_CH2N/TIM15_CH1/TSC_G6_IO4/EVENTOUT | COMP3_INP/ADC4_IN4/OPAMP2_VINP |

A special case of a figure is a table, where the table cells are drawn in vector graphics and the text is placed inside that.

That's why if you just attempt to copy the text of the table into an editor, you usually get garbage.

Note the rotated text in the header, in which order is that copied? It's up to the PDF reader how to copy this text.

**Table 10. STM32F413xG/H pin definition**

| UFQFPN48 | LQFP64 | WLCSP81 | LQFP100 | UFBGA100 | UFBGA144 | LQFP144 | Pin name (function after reset)(1) | Pin type | I/O structure | Notes | Alternate functions | Additional functions |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | - | NC | 1 | B2 | A3 | 1 | PE2 | I/O | FT | (2) | TRACECLK, SPI4_SCK/I2S4_CK, SPI5_SCK/I2S5_CK, SAI1_MCLK_A, QUADSPI_BK1_IO2, UART10_RX, FSMC_A23, EVENTOUT | - |
| - | - | NC | 2 | A1 | A2 | 2 | PE3 | I/O | FT | (2) | TRACED0, SAI1_SD_B, UART10_TX, FSMC_A19, EVENTOUT | - |
| - | - | NC | 3 | B1 | B2 | 3 | PE4 | I/O | FT | (2)(3) | TRACED1, SPI4_NSS/I2S4_WS, SPI5_NSS/I2S5_WS, SAI1_SD_A, DFSDM1_DATIN3, FSMC_A20, EVENTOUT | - |
| | | | | | | | | | | | TRACED2, TIM9_CH1, SPI4_MISO, SPI5_MISO, | |

**Table 12. STM32F413xG/H alternate functions**

| Port | AF0 SYS_AF | AF1 TIM1/2/LPTIM1 | AF2 TIM3/4/5 | AF3 DFSDM2/TIM8/9/10/11 | AF4 I2C1/2/3/I2CFMP1 | AF5 SPI1/I2S1/SPI2/I2S2/SPI3/I2S3/SPI4/I2S4 | AF6 SPI2/I2S2/SPI3/I2S3/SPI4/I2S4/SPI5/I2S5/DFSDM1/2 | AF7 SPI3/I2S3/SAI1/DFSDM2/USART1/USART2/USART3 | AF8 DFSDM1/USART3/4/5/6/7/8/CAN1 | AF9 I2C2/I2C3/I2CFMP1/CAN1/2/TIM12/13/14/QUADSPI | AF10 SAI1/DFSDM1/DFSDM2/QUADSPI/FSMC/OTG1_FS | AF11 UART4/UART5/UART9/UART10/CAN3 | AF12 FSMC/SDIO | AF13 | AF14 RNG | AF15 SYS_AF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PA0 | - | TIM2_CH1/TIM2_ETR | TIM5_CH1 | TIM8_ETR | - | - | - | USART2_CTS | UART4_TX | - | - | - | - | - | - | EVENTOUT |
| PA1 | - | TIM2_CH2 | TIM5_CH2 | - | - | SPI4_MOSI/I2S4_SD | - | USART2_RTS | UART4_RX | QUADSPI_BK1_IO3 | - | - | - | - | - | EVENTOUT |
| PA2 | - | TIM2_CH3 | TIM5_CH3 | TIM9_CH1 | - | I2S2_CKIN | - | USART2_TX | - | - | - | - | FSMC_D4/FSMC_DA4 | - | - | EVENTOUT |
| PA3 | - | TIM2_CH4 | TIM5_CH4 | TIM9_CH2 | - | I2S2_MCK | - | USART2_RX | - | SAI1_SD_B | - | - | FSMC_D5/FSMC_DA5 | - | - | EVENTOUT |
| PA4 | - | - | - | - | - | SPI1_NSS/I2S1_WS | SPI3_NSS/I2S3_WS | USART2_CK | DFSDM1_DATIN1 | - | - | - | FSMC_D6/FSMC_DA6 | - | - | EVENTOUT |
| PA5 | - | TIM2_CH1/TIM2_ETR | - | TIM8_CH1N | - | SPI1_SCK/I2S1_CK | - | - | DFSDM1_CKIN1 | - | - | - | FSMC_D7/FSMC_DA7 | - | - | EVENTOUT |
| PA6 | - | TIM1_BKIN | TIM3_CH1 | TIM8_BKIN | - | SPI1_MISO | I2S2_MCK | DFSDM2_CKIN1 | - | TIM13_CH1 | QUADSPI_B K2_IO0 | - | SDIO_CMD | - | - | EVENTOUT |
| PA7 | - | TIM1_CH1N | TIM3_CH2 | TIM8_CH1N | - | SPI1_MOSI/I2S1_SD | DFSDM2_DATIN1 | - | - | TIM14_CH1 | QUADSPI_B K2_IO1 | - | - | - | - | EVENTOUT |
| PA8 | MCO_1 | TIM1_CH1 | - | - | I2C3_SCL | - | DFSDM1_CKOUT | USART1_CK | UART7_RX | - | USB_FS_SOF | CAN3_RX | SDIO_D1 | - | - | EVENTOUT |
| PA9 | - | TIM1_CH2 | - | DFSDM2_CKIN3 | I2C3_SMBA | SPI2_SCK/I2S2_CK | - | USART1_TX | - | - | USB_FS_VBUS | - | SDIO_D2 | - | - | EVENTOUT |
| PA10 | - | TIM1_CH3 | - | DFSDM2_DATIN3 | - | SPI2_MOSI/I2S2_SD | SPI5_MOSI/I2S5_SD | USART1_RX | - | - | USB_FS_ID | - | - | - | - | EVENTOUT |
| PA11 | - | TIM1_CH4 | - | DFSDM2_CKIN5 | - | SPI2_NSS/I2S2_WS | SPI4_MISO | USART1_CTS | USART6_TX | CAN1_RX | USB_FS_DM | UART4_RX | - | - | - | EVENTOUT |

Here you can see the first page of a datasheet. We detect the double column layout manually, then convert each side. We need to simplify the problems, so first we

- Convert all 2D information into an abstract syntax tree.
- Then modify that AST to detect the hierarchy of the document and then normalize page breaks.
- Then format it as HTML.

If this sounds like a compiler, it's basically a PDF frontend, then a number of AST passes, then a HTML backend. And this actually works really well.

This is the result, for example the pin definition table in the datasheet.
This is a pure HTML table with minimal CSS to look similiar to the PDF.
All of the data is converted as is including line breaks.

Here is the alternate function table.
This is normally broken up across many pages, in the HTML its just one long table.

## Table 24. RCC register map and reset values for STM32F413/423

| Addr. offset | Register name | 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0x00 | RCC_CR | Res. | Res. | Res. | Res. | PLLI2SRDY | PLLI2SON | PLLRDY | PLLON | Res. | Res. | Res. | Res. | CSSON | HSEBYP | HSERDY | HSEON | HSICAL[7:0] | | | | | | | | HSITRIM[4:0] | | | | | Res. | HSIRDY | HSION |
| 0x04 | RCC_PLLCFGR | Res. | PLLR[2:0] | | | PLLQ[3:0] | | | | Res. | PLLSRC | Res. | Res. | Res. | Res. | PLLP[1:0] | | Res. | PLLN[8:0] | | | | | | | | | PLLM[5:0] | | | | | |
| 0x08 | RCC_CFGR | MCO2[1:0] | | MCO2PRE[2:0] | | | MCO1PRE[2:0] | | | Res. | MCO1[1:0] | | RTCPRE[4:0] | | | | | PPRE2[2:0] | | | PPRE1[2:0] | | | Res. | Res. | HPRE[3:0] | | | | SWS[1:0] | | SW[1:0] | |
| 0x0C | RCC_CIR | Res. | Res. | Res. | Res. | Res. | Res. | Res. | CSSC | PLLI2SRDYC | PLLRDYC | HSERDYC | LSERDYC | LSIRDYC | Res. | Res. | Res. | Res. | PLLI2SRDYIE | PLLRDYIE | HSERDYIE | LSERDYIE | LSIRDYIE | CSSF | Res. | PLLI2SRDYF | PLLRDYF | HSERDYF | LSERDYF | LSIRDYF | | | |
| 0x10 | RCC_AHB1RSTR | Res. | Res. | Res. | Res. | Res. | Res. | Res. | Res. | DMA2RST | DMA1RST | Res. | Res. | Res. | Res. | Res. | Res. | Res. | Res. | CRCRST | Res. | Res. | Res. | Res. | GPIOHRST | GPIOGRST | GPIOFRST | GPIOERST | GPIODRST | GPIOCRST | GPIOBRST | GPIOARST | |

---

We also find the register layout information again for each peripheral. Note that the text is rotated only by CSS, so the table data is still easily accessible in HTML.

---

### 6.3.3 RCC clock configuration register (RCC_CFGR)

Address offset: 0x08

Reset value: 0x0000 0000

Access: 0 ≤ wait state ≤ 2, word, half-word and byte access

1 or 2 wait states inserted only if the access occurs during a clock source switch.

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MCO2[1:0] | | MCO2 PRE[2:0] | | | MCO1 PRE[2:0] | | | Res. | MCO1[1:0] | | RTCPRE[4:0] | | | | | PPRE2[2:0] | | | PPRE1[2:0] | | | Res. | Res. | HPRE[3:0] | | | | SWS[1:0] | | SW[1:0] | |
| rw | rw | rw | rw | rw | rw | rw | rw | | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | | | rw | rw | rw | rw | r | r | rw | rw |

Bits 31:30 **MCO2[1:0]:** Microcontroller clock output 2

Set and cleared by software. Clock source selection may generate glitches on MCO2. It is highly recommended to configure these bits only after reset before enabling the external oscillators and the PLLs.
00: System clock (SYSCLK) selected
01: PLLI2S clock selected
10: HSE oscillator clock selected
11: PLL clock selected

```
#define LL_RCC_SYS_CLKSOURCE HSI
#define LL_RCC_SYS_CLKSOURCE HSE
#define LL_RCC_SYS_CLKSOURCE PLL
#if defined(RCC_PLLR_SYSCLK_SUPPORT)
#define LL_RCC_SYS_CLKSOURCE_PLLR
#endif
```

Bits 1:0 **SW[1:0]:** System clock switch

Set and cleared by software to select the system clock source.
Set by hardware to force the HSI selection when leaving the Stop or Standby mode or in case of failure of the HSE oscillator used directly or indirectly as the system clock.
00: HSI oscillator selected as system clock
01: HSE oscillator selected as system clock
10: PLL selected as system clock
11: not allowed

But what is the enumeration name?

---

And I can even convert the invisible table of the bit field and their enumerations description as a HTML table.
###
And indeed this is accurate, the PLLR does not exist for this device, so the guard in the CubeHAL header is actually correct.
###
Unfortunately we have the enumeration value and description, but not a name. That would need to be generated from the description and that not always easy to do automatically.

---

## PDF to HTML conversion
**Open-sourced at data.modm.io**

- Manually written Python3 code based on pypdfium2.

- ~157k PDF pages in 65mins on a MacBook Air M2 => ~25ms per page!

- Works on all PDFs from STMicro: also sensors, not just STM32!

- Most valuable data is inside tables, but table processing is hard and fuzzy.

- Not easily portable to other vendor data sheets due to content segmentation!

- Figures and images are ignored, math formulas are not recognized.

---

I'm very happy with this pipeline. It's written in Python3 using native bindings for pdfium (PDF renderer in Chrome). It's entirely deterministic, so the translated HTML is byte reproducible. It's also very fast with 25ms per page. All STMicro PDFs are supported, including sensors.
Some compromises: it's not easily portable to other vendors, since the format recognition is hardcoded. I'm only interested in tables and text, so figures are completely ignore (should be converted to SVG) and math formulas are turned into garbage.

## PDF Formatting Mistakes

- The PDF sometimes have formatting mistakes: tables with missing cell borders.
- Apply git patch to HTML result: works, but fragile.

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Res. | Res. | Res. | Res. | Res. | PEC BYTE | AUTOE ND | RE LOAD | | | | NBYTES[7:0] | | | | |
| | | | | | rs | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw |

| 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| NACK | STOP | START | HEAD1 0R | ADD10 | RD_ WRN | | | | | | SADD[9:0] | | | | |
| rs | rs | rs | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw | rw |

Even though it's deterministic and reproducible, some formatting mistakes are not easy to fix.
A classic is missing border cells in tables.
We could try to infer cells from whitespace analysis between text, but it's fairly unreliable for such issues.
I just apply a git patch to the HTML, which works because the HTML is so reproducible.

## Interpreting Datasheet Tables
### Substitution hell to fix typos in PDFs

```python
package = package.replace("UFBGA/TFBGA64", "UFBGA64/TFBGA64")
package = package.replace("LPQF", "LQFP").replace("TSSPOP", "TSSOP")
package = package.replace("UFBG100", "UFBGA100").replace("UBGA", "UFBGA")
package = package.replace("UQFN", "UFQFPN").replace("WLCSP20L", "WLCSP20")
package = package.replace("UFQFPN48E", "UFQFPN48+E").replace("UFQFN", "UFQFPN")
package = package.replace("LQFP48 SMPS<br>UFQFPN48 SMPS", "LQFP48/UFQFPN48+SMPS")
package = package.replace("LQFP<br>64", "LQFP64").replace("LQFP<br>48", "LQFP48")
package = package.replace("UFQFPN<br>32", "UFQFPN32").replace("UFQFP<br>N48", "UFQFPN48")
package = package.replace("WLCSP<br>25", "WLCSP25")
```

But, still interpreting the HTML tables was actually way more annoying that converting the PDF.
You can to clean the data, because of many typos and random line breaks.
Here I'm using text substitution.

## Interpreting Datasheet Tables
### Substitution hell to fix typos in PDFs, now with more RegEx

```python
patterns = {
    r" +": "", r".*?\(([A-Z]+|DMA2D)\).*?": r"\1",
    r"Reserved|Port|Power|Registers|Reset|\(.*?\)|_REG": "",
    r"(\d)/I2S\d": r"\1", r"/I2S|CANMessageRAM|Cortex-M4|I2S\dext|^GPV$": "",
    r"Ethernet": "ETH", r"Flash": "FLASH", r"(?i).*ETHERNET.*": "ETH",
    r"(?i)Firewall": "FW", r"HDMI-|": "", "SPDIF-RX": "SPDIFRX",
    r"SPI2S2": "SPI2", r"Tamper": "TAMP", "TT-FDCAN": "FDCAN",
    r"USBOTG([FH])S": r"USB_OTG_\1S", "LCD-TFT": "LTDC", "DSIHOST": "DSI",
    "TIMER": "TIM", r"^VREF$": "VREFBUF", "DelayBlock": "DLYB",
    "I/O": "M", "I/O": "", "DAC1/2": "DAC12",
    r"[a-z]": ""
}
```

And then I decided to use Regex to fix many patterns in the register definitions.

## Interpreting Datasheet Tables
### RegEx hell to fix typos for bit field reconstruction

```
off_replace = {r" +": "", "0x000x00": "0x00", "to": "-", "x": "*", r"\(\d+\)": ""}
dom_replace = {r"Register +size": "Bit position"}
reg_replace = {
    r" +|\.+": "", r"\(COM(\d)\)": r"_COM\1",
    r"^[Rr]es$||0x[\da-fA-FXx]+|\(.*?\)|-": "",
    r"(?i)reserved|resetvalue.*": "", "enabled": "_EN", "disabled": "_DIS",
    r"(?i)Outputcomparemode": "_Output", "(?i)Inputcapturemode": "_Input", "mode": "",
    r"^TG_FS_": "OTG_FS_", "toRTC": "RTC", "SPI2S_": "SPI_",
    r"andTIM\d+_.*": "", r"x=[\d,]+": ""}
fld_replace = {
    r"\] +\d+(th|rd|nd|st)": "]", r" +|\.+|\[.*?\]|\[?\d+:\d+\]?|\(.*?\)|-|^[\dXx]+$|%|__|:0\]": "",
    r"Dataregister|Independentdataregister": "DATA",
    r"Framefilterreg0.*": "FRAME_FILTER_REG",
    r"[Rr]es(erved)?|[Rr]egular|x_x(bits)?|NotAvailable|RefertoSection\d+:Comparator": "",
    r"Sampletimebits|Injectedchannelsequence|channelsequence|conversioninregularsequencebits": "",
    r"conversioninsequencebits|conversionininjectedsequencebits|or|first|second|third|fourth": ""}
bit_replace = {r".*:*": ""}
glo_replace = {r"[Rr]eserved": ""}
```

And this then got a little out of hand for the bit field enumerations.
I do not recommend using regex for this, there needs to be a better way.

---

## Evaluation of Data Sources 🧑‍🔬
### Actual Science! OMG

Extracted 4 datasets with increasing complexity for ~2700 STM32 devices:

1. Interrupt vector table: PDF vs CMSIS Header
2. Package and pinout: PDF vs CubeMX Database
3. Pin functions: PDF vs CubeMX Database
4. MMIO register map and descriptions: PDF vs SVD vs Header

Compare PDF against machine-readable sources: Headers, SVD, CubeMX

Ok, but enough regexing around. Let's do some actual science!
We want to find out how accurate our data import pipelines actually are.
So we're going to compare the machine-readable data against the PDF data.
We evaluated in detail four data sets for this.
We fixed obvious spelling mistakes, but only as long as the fix is unambigious.

---

## PDF Interrupt Table vs CMSIS Header 🧑‍🔬
### Device → Reference Manual → Table → Position + Name

| Position | Priority | Type of priority | Acronym | Description | Address |
|---|---|---|---|---|---|
| - | - | - | - | Reserved | 0x0000 0000 |
| - | -3 | Fixed | Reset | Reset | 0x0000 0004 |
| - | -2 | Fixed | NMI | Non maskable interrupt. The RCC clock security system (CSS) and the RAM parity check are linked to the NMI vector. | 0x0000 0008 |
| - | -1 | Fixed | HardFault | All classes of fault | 0x0000 000C |
| - | 3 | Settable | SVCall | System service call via SWI instruction | 0x0000 002C |
| - | 5 | Settable | PendSV | Pendable request for system service | 0x0000 0038 |
| - | 6 | Settable | SysTick | System tick timer | 0x0000 003C |
| 0 | 7 | Settable | WWDG | Window watchdog interrupt | 0x0000 0040 |
| 1 | 8 | Settable | PVD_VDDIO2 | PVD and V_DDIO2 supply comparator interrupt (combined EXTI lines 16 and 31) | |
| 2 | 9 | Settable | RTC | RTC interrupts (combined EXTI lines 17, 19 and 20) | |

**98.8% match (N=190 109)**

This is fairly easy: it's the interrupt vector table for STM32 microcontrollers.
Quite good.

## PDF Pinout vs CubeMX Database
**Device → Datasheet → Table → Pin Position + Name**

| WLCSP100 | LQFP100 | LQFP64 | LQFP48 | Pin name (function after reset) | Pin type | I/O structure | Notes | Alternate functions | Additional functions |
|---|---|---|---|---|---|---|---|---|---|
| J3 | 52 | 34 | 26 | PB13 | I/O | TTa | (4) | SPI2_SCK,I2S2_CK,USART3 _CTS, TIM1_CH1N, TSC_G6_IO3, EVENTOUT | ADC3_IN5, COMP5_INP, OPAMP4_VINP, OPAMP3_VINP |
| J2 | 53 | 35 | 27 | PB14 | I/O | TTa | (4) | SPI2_MISO,I2S2ext_SD, USART3_RTS_DE, TIM1_CH2N, TIM15_CH1, TSC_G6_IO4, EVENTOUT | COMP3_INP, ADC4_IN4, OPAMP2_VINP |
| H4 | 54 | 36 | 28 | PB15 | I/O | TTa | (4) | SPI2_MOSI, I2S2_SD, TIM1_CH3N, RTC_REFIN, TIM15_CH1N, TIM15_CH2, EVENTOUT | ADC... |

**99.88% match (N=247 756)**

The package pinout was extremely accurate. This is just the pin position and name on the package.

---



## PDF Pin Functions vs CubeMX Database
**Device → Datasheet → Table → Pin Name + Function**

| Port & Pin Name | AF0 | AF1 | AF2 | AF3 | AF4 | AF5 | AF6 | AF7 | AF8 | AF9 | AF10 | AF11 | AF12 | AF14 | AF15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PA12 | - | TIM16_CH1 | - | - | - | - | TIM1_CH2N | USART1_RTS_DE | COMP2_OUT | CAN_TX | TIM4_CH2 | TIM1_ETR | - | USB_DP | EVENT OUT |
| PA13 | SWDIO -JTMS | TIM16_CH1N | - | TSC_G4_IO3 | - | IR_OUT | - | USART3_CTS | - | - | TIM4_CH3 | - | - | - | EVENT OUT |
| PA14 | SWCLK -JTCK | - | - | TSC_G4_IO4 | I2C1_SDA | TIM8_CH2 | TIM1_BKIN | USART2_TX | - | - | - | - | - | - | EVENT OUT |
| PA15 | JTDI | TIM2_CH1_ETR | TIM8_CH1 | - | I2C1_SCL | SPI1_NSS | SPI3_NSS, I2S3_WS | USART2_RX | - | TIM1_BKIN | - | - | - | - | EVENT OUT |

Data cleanup for cells requires interpretation of newlines, commas, hyphenation

**96.2% match (N=1 107 035)**

The signals are a bit more interesting, this is our first 2D data structure.
We looked at over a million signals in our dataset, didn't find any issues with our PDF-to-HTML pipeline, but many issues in the CubeMX database, as well as formatting issues in the PDF.
Still very accurate.

---



## Register Map: Header vs SVD vs PDF
**183 three-way comparisons of occupied linear address space**

Legend: CMSIS Header, CMSIS-SVD, Reference Manual

Memory Map Size [Byte] vs Unique Three-way Comparison Sorted by Device Identifier and Grouped by Family (F0, F1, F2, F3, F4, F7, G0, G4, H7, L0, L1, L4, L4+)

| Hierarchy Level | Conflict Size | Total Map Size | Conflict-Free Locations | Overlap Map Size | Matching Locations |
|---|---|---|---|---|---|
| Peripherals | 2 748 B | 55 376 B | 95.0 % | 42 752 B | 93.6 % |
| Registers | 35 406 B | 1 188 994 B | 97.3 % | 891 044 B | 96.0 % |
| Bit Fields | 379 180 bit | 5 711 619 bit | 93.3 % | 3 425 903 bit | 88.9 % |

And finally we compared the register maps reconstructed form the CMSIS Header, vs the CMSIS-SVD vs the PDF.
And this was the most interesting part, because it shows that STMicro has three slightly different datasets for their hardware.

As a proxy for completeness we can look at the size of the register map. How many bytes are occupied by the registers.
You can see that the register map reconstructed from the reference manual is very accurate.

BUT the device resolution is not great:
- the CMSIS headers create 183 register maps,
- The CMSIS-SVD files only 100 register maps, and
- The PDFs only 53 register maps.



**Register Map: PDF vs SVD vs Header**
Mismatches and fixing them with 2-1 majority voting

| Hierarchy Level | Resolvable by Majority Vote | Reference Manual + CMSIS Header | CMSIS Header + CMSIS-SVD | CMSIS-SVD + Reference Manual | Matching and Resolved Locations |
|---|---|---|---|---|---|
| Peripherals | 44.1 % | 78.9 % | 19.8 % | 1.3 % | **96.4 %** |
| Registers | 41.7 % | 40.1 % | 15.5 % | 44.4 % | **97.7 %** |
| Bit Fields | 63.0 % | 38.8 % | 22.9 % | 38.3 % | **95.9 %** |

Here is the conflict rate in more detail. We can see that the complex families like F7, H7 and L4 have the most conflicts overall.
Since we have three differing data sources, we can do majority voting and see how many differing registers we can fix.

It works well for simple families, and improves the matching data quite a bit, but we can also see that the combination of CMSIS header and CMSIS-SVD is the least successful in majority voting.

This is very weird since the CMSIS header files are supposed to be generated from the CMSIS SVD files.

## Results Overview
### It's almost great!

- We didn't find any systemic issues in our PDF-to-HTML pipeline!
- STMicro maintains three slightly different datasets for register maps???

| Dataset | Sources | Method of Comparison | Result | N |
|---|---|---|---|---|
| Device Identifier | | Datasheet ⊇ CubeMX | 93.2 % | 3024 |
| Package | ③ Datasheet vs. ⑦ CubeMX | Datasheet = CubeMX | 99.68 % | 2819 |
| Pinout | | Matching pin name at package position | 99.88 % | 247756 |
| Pin Function | | Matching index for function name at pin | 96.2 % | 1107035 |
| Interrupt Vector Table | ③ Reference Manual vs. ⑤ Header | Matching vector name at table position | 98.8 % | 190109 |
| Peripheral | | | 96.4 % | 42752 |
| Register | ④ Reference Manual vs. | Matching peripheral, register, or bit field name | 97.7 % | 891044 |
| Bit Field | ⑤ Header vs. ⑥ SVD | at byte or bit address after majority voting | 95.9 % | 3425903 |
| All Datasets | All Sources | Weighted average over all data points | 96.5 % | 5910442 |

Overall, the machine-readable data is very accurate with 96.5% match at 5.9 million data points.
As a result, I would not use the PDF or the CMSIS-SVD files as primary data sources unless necessary.
Extract as much as possible from the CubeMX database and CMSIS headers instead.

---

## data.modm.io : Data Interface



Knowledge Graph is the Interface!

So the question is how to we make this data accessible?
We have a highly heterogeneous dataset, which includes clock graphs, so why not use a graph database?
The graph database then also acts as the interface to the external world.

---

## Knowledge Graph as Interface
### Kuzu Embedded Database

- Perfect for heterogeneous dataset like hardware description.
- Implements Cypher Query Language for many languages.
- Serialization into sorted plaintext allows trivial archiving.
- Simple to install and use: `pip install kuzu`

```
kuzu> MATCH (:Package)-[po:hasPin]->(pi:Pin)-[af:hasAlternateFunction]->(s:Signal)<--(pe:Peripheral)
      RETURN pi.name, po.position, pe.name, s.name, af.index;
```

| pi.name STRING | po.position STRING | pe.name STRING | s.name STRING | af.index UINT8 |
|---|---|---|---|---|
| PA0/WKUP | N3 | ETH | CRS | 11 |
| PH2 | K4 | ETH | CRS | 11 |
| PA0/WKUP | N3 | TIM1 | CH1 | 1 |

I chose Kuzu, which is a small and fast embedded graph database that implements the Cypher language.
It's easy to install and use and comes with many language bindings: C/C++, Rust, Python, Web Assembly.
You can see a part of the schema on the right and a cypher query at the bottom showing alternate functions.
The query returns a table which you can then use to generate code.

There a browser based explorer tool including graph visualization.
Here you can see the package node, surrounded by all the pins and their corresponding signals.
This is a bit chaotic

**Knowledge Graph as Interface**
`MATCH p=(:Peripheral)-->(:Signal) RETURN p`

We can also query only a part of the graph, like specific relations.
Here we query all relations between peripherals and signals.

**Package and Pinout Shenanigans**

```
┌STM32G071GBU6─────────────┐
        PB8 PB7 PB6 PB5 PB4 PB3 PA15
PC14                            PA14
PC15                            PA13
VDD                             PA12
VSS                             PA11
PF2                             PC6
                                PA8
PA9 PA2 PA3 PA4 PA5 PA0 PA7 PB0
```

`pip install stm_layout`

```
┌Regex─────────────────────┐
ADC\d_IN
```

```
┌Pin Info─────────┐  ┌Alternate Functions──┐  ┌Additional Functions─┐
      Name: PA2      [ ] 0: I2S1_SD/SPI1_MOSI   ADC1_IN2
       Pos: 8        [ ] 1: USART2_TX           COMP2_INM
    Config: Default  [ ] 2: TIM2_CH3            RCC_LSCO
      Mode: [ ] GPI  [ ] 3: -                   SYS_WKUP4
            [ ] GPO  [ ] 4: -                   UCPD1_FRSTX1
            [ ] Alternate [ ] 5: TIM15_CH1      UCPD1_FRSTX2
            [x] Analog    [ ] 6: LPUART1_TX
     Speed: Low      [ ] 7: COMP2_OUT
            Med       [ ] 8: -
            High      [ ] 9: -
            Very High [ ] 10: -
      Type: Push-Pull [ ] 11: -
            Open-Drain [ ] 12: -
  Resistor: None     [ ] 13: -
            Pull-Up  [ ] 14: -
            Pull-Down [ ] 15: -
```

But there are also many other things you can do with this code.
For example REGEX your alternate functions!
This is a smol TUI tool based on the CubeMX database, here showing all the pins that have an ADC input signal.

## Package and Pinout Shenanigans



This also works for BGA pins, here searching for all pins with I2C data and clock signals.

This is very useful to quickly find alternate functions, since the CubeMX gui is not great for searching like this.

---

## Memory Map Shenanigans



Why not simply™ regex your register map?

I apologize to your eyeballs. This is code from modm.

modm uses Python Jinja templates to generate startup code, where we need to enable the clock to the system config and power peripherals.

And it's very annoying since the bit is in different registers depending on the family.

###

So instead, why not just regex the register?

###

This abomination actually works really well…

---

## More Use Cases
### Everything is a Query when all you have is Data

- Modularize and generator your own HAL much easier.
- HTML version of pinout and clock configurator. No more CubeMX.
- Optimizing constraint solver for pinout and clock limitations.
- Diffing HTML versions of Datasheets and Reference Manuals.
- Much more accurate CMSIS-SVD files from the CMSIS Header + PDF.
- Testing AI models against PDF-to-HTML-to-Knowledge-Graph pipeline.

There are many more use cases that I didn't go into.

A nice one would be to create a simpler CubeMX application as a HTML page, something that can configure and generate code for other HALs.

You can apply a SAT solver to the database of course to solve design constraints and help with parts selection.

You can now diff PDFs via the HTML version.

And you can generate much more accurate SVD files than the official ones.

If you think your AI model can do better, I've basically built you a benchmark.

BEWARE.

## Conclusion
### and Future Work

- STMicro publishes several machine-readable data sources on GitHub!
- Parsing machine-readable data sources is easy and very accurate.
- Parsing PDF/HTML is difficult due to typos and formatting mistakes.
- modm-data: PDF2HTML works well, rest is "academic" code quality.

- Knowledge Graphs are a good database for heterogenerous data sets.
- Documentation and discoverability of Knowledge Graph Ontology is difficult.

There's a lot of machine-readable data on GitHub, it needs to be put in a good database. Knowledge graphs are still pretty niche.
Parsing PDFs is hard because of humans, rather than accessing the PDF. Some fuzzy matching required.
The PDF2HTML pipeline works really well in modm-data, the rest needs to be rewritten.

## Data Mining Hardware Descriptions
### Questions?

Niklas Hauser likes data science.

| | |
|---|---|
| Homepage: | salkinium.com |
| Fediverse: | @salkinium@chaos.social |
| Code: | github.com/salkinium |
| Thesis: | salkinium.com/master.pdf |
| Paper: | salkinium.com/hp23.pdf (peer-reviewed!) data.modm.io |
| GitHub: | github.com/modm-io/modm-data |

If you want to know more details, including citations, check out master thesis and my peer-reviewed paper.
All the code is public and somewhat documented.
Do you have Questions?